



HELLENIC REPUBLIC
National and Kapodistrian
University of Athens

Re-defining Quality in Journalism and Audience Engagement in the Digital Era: A Computational Approach Using Big Data and AI

Aikaterini - Alexandra Sotirakou

Department of Communication and Media Studies

School of Economics and Political Sciences

National and Kapodistrian University of Athens

December 2022

Re-defining Quality in Journalism and Audience Engagement in the Digital Era: A Computational Approach Using Big Data and AI

Supervising Professor: Constantinos Mourlas

Supervisor 2: Hajo Boomgaarden

Supervisor 3: Antonis Armenakis

Intellectual Property

The candidate confirms that the work submitted is her own and that appropriate credit has been given where reference has been made to the work of others.

This copy has been supplied on the understanding that it is copyright material and that no quotation from the dissertation may be published without proper acknowledgement.

© 2022 National and Kapodistrian University of Athens, Aikaterini - Alexandra Sotirakou

Signed



Acknowledgements

First and foremost, I am thankful to Professor Constantinos Mourlas, who inspired me to write this dissertation and gave me confidence that I could get through this difficult topic, as well as for continuing to believe in me throughout the years. His mentorship and guidance were crucial for my scientific work and it was my privilege working alongside him. I express my gratitude to numerous individuals, as the research presented in this dissertation stems from extensive collaboration. I wish to thank my committee, Professors Antonis Armenakis and Hajo Boomgaarden who offered illuminating feedback on my methodological approach, analysis and the theoretical framework of the research, enriching my scientific knowledge. I am indebted to Dr. Panagiotis Germanakos who helped me find my footing as a young researcher and to the Ph.D candidate Katerina Mandenaki with whom we had countless conversations about the theoretical part of my dissertation. Furthermore, I gratefully recognize the help of the NKUA students A. Karampela, D. Sinnis and E. Koutromanou for their assistance on the research papers we wrote together. Also, I wish to thank Professor Spiros Moschonas who read and commented on draft versions of this work.

I have greatly benefited from the scholarships at Columbia Journalism School, University of Vienna and London School of Economics and I want to thank Stavros Niarchos Foundation, the British Embassy of Athens and the Greek State Scholarships Foundation for providing these opportunities to advance my education. Lastly, I gratefully acknowledge the financial support by the General Secretariat for Research and Technology (GSRT) and the Hellenic Foundation for Research and Innovation (HFRI) - (Scholarship Code: GA. 14540). Finally, I owe everything to my mother and my sister to whom I dedicate this dissertation.

Abstract

The way news is consumed and circulated has undergone significant changes, with the proliferation of the internet and social media platforms providing new avenues for news organizations to explore. Publishers are under increasing pressure to demonstrate their reach and to engage with their audiences online, allowing news consumers' preferences to influence their daily agenda. Additionally, technology giants and algorithmic curation challenge the survival of news organizations that must compete with various actors on social media, all while remaining committed to reporting the truth and producing high-quality journalism.

In today's hybrid media environment where traditional values like objectivity have shifted allowing space for more intimacy in journalism, social media engagement has become an important metric for the success of a news piece. In the pursuit of the "magic recipe" for writing quality and compelling news stories, this dissertation will conceptualize quality and engagement based on various communication theories, and utilize data journalism and artificial intelligence techniques to identify certain attributes in news articles that could influence their perceived quality and engagement.

Throughout the studies and discussions, the research explores various quality dimensions in different settings such as blogging platforms, online news outlets, and social media. Furthermore, textual, visual and contextual information are taken into consideration to provide a deeper understanding of how the elements of a news story influence its quality and its potential to become engaging.

This dissertation explores the use of quantitative and computational methods, specifically explainable artificial intelligence, in the social study of journalism. The focus is on iden-

tifying hidden relationships within data and using those insights to inform the creation of high-quality news stories that are engaging for readers. The findings show that characteristics such as the depth of a story, its diversity, the conveyed emotions, readability, and the choice of images can predict both the quality and social media engagement of a news article.

The objective of this research is to provide guidance for journalists on how to create captivating content using effective techniques and images while maintaining journalistic standards. This work has the potential to be applied in academia and the media industry, both advancing the theoretical foundations of journalism studies and providing practical applications for news production.

Keywords News, Data Journalism, Digital Journalism, News Quality, Computational Social Science, Audience Engagement, Emotions, Artificial Intelligence, Machine Learning, Explainable AI, Image Recognition.

Contents

1	Introduction	1
1.1	Quality and Engagement in the Digital News Era	1
1.2	Objective and Research Questions	4
1.3	Research Design	5
1.4	Structure of the dissertation	6
2	(Re)defining Quality	11
2.1	Navigating the Changing Landscape of Digital Journalism	11
2.2	Quality in Journalism	13
2.3	Previous Efforts to Measure News Quality	15
2.4	Data Journalism	19
2.4.1	Computational Journalism	23
2.4.2	Data journalism as a Profession	25
2.5	Exploring the Possibilities of AI in Journalism	28
2.5.1	Existing Initiatives	35
2.5.2	Algorithmic Accountability	40
3	Exploring the Impact of Social Media	45
3.1	Adapting to the Digital Age	45
3.1.1	From Imagined Audience to Micro-targeting	47
3.1.2	Definition of Engagement	51
3.1.3	Popularity Metrics	54

3.1.4	Shareability	57
3.1.5	News Values and Engagement	59
3.2	The Power of Emotions	61
3.2.1	Definition of Emotions	62
3.2.2	Emotions and Journalism	64
3.2.3	Emotions for Disinformation Detection	82
3.2.4	Disinformation	83
4	The Role of Images in Journalism	88
4.1	The Importance of Visual Content	88
4.1.1	Pictures as the Main Way of Expression	89
4.2	Social Media and the Power of the Image	91
4.2.1	What Makes an Image Newsworthy and Engaging?	93
4.2.2	Brand Related Images	99
4.2.3	Machine Learning Approaches	101
5	Methodology	103
5.1	Computational Social Science in Journalism Studies	103
5.1.1	Knowledge Discovery in Databases	106
5.1.2	Cross Industry Standard Process for Data Mining	108
5.1.3	Data Mining Methods in CSS	111
5.1.4	Machine Learning Algorithms	114
5.2	Research Questions	118
5.3	Research Design	118
5.3.1	Business Understanding	119
5.3.2	Data Understanding	122
5.3.3	Data Preparation	122
5.3.4	Modeling	124
5.3.5	Supervised Machine Learning Algorithms	124
5.3.6	Unsupervised Deep Learning for Text	130

5.3.7	Evaluation	132
5.3.8	Deployment	134
6	Studies	136
6.1	Study:1 Audience Engagement Metrics and Perceived Quality	136
6.1.1	Study Overview	137
6.1.2	Audience Engagement	141
6.1.3	Model & Feature Extraction	144
6.1.4	Method & Dataset	149
6.1.5	Phase A - Evaluating the Importance of the Proposed Model Engage- ment Metrics in Relation to the Perceived Quality	151
6.1.6	Phase B - Proposing a Set of Rules with Respect to the Engagement Met- rics of each Category	156
6.1.7	Conclusion & Future Work	165
6.2	Study:2 Predicting the Quality of News Articles	168
6.2.1	Toward a Model of Quality in Journalistic Texts	168
6.2.2	Methods and Data	173
6.2.3	Analysis	177
6.2.4	Discussion & Future work	184
6.3	Study: 3 An Analysis of News Engagement using AI	186
6.3.1	Methods and Dataset	187
6.3.2	Analysis and Results	189
6.3.3	Discussion and Conclusion	193
6.4	Study:4 The Impact of Images on News Quality and Engagement	195
6.4.1	Methods & Data	199
6.4.2	Analysis & Results	201
6.4.3	Discussion	203
6.5	Study:5 News Quality and Engagement for Fake News Detection	204
6.5.1	Method & Dataset	208
6.5.2	Data Analysis & Findings – in two Distinctive Phases	209

6.5.3	Discussion of the Results	212
6.5.4	Conclusion & Future Work	213
7	Conclusion and Future Work	216
7.1	Conclusion	216
7.1.1	Discussion and Limitations	225
7.2	Knowledge Transfer and Future Work	229
	References	232
A	Biographical Sketch	292
B	Code	294
B.1	Study: 1	294
B.2	Study: 2	298
B.3	Study: 3	305
C	Lexicons for Feature Engineering	309
D	Dataset Cleaning	312

List of Figures

5.1	Interdisciplinary nature of CSS. Graphic by J. Zhang et al. (2020)	105
5.2	Overview of KDD Process, from the Fayyad et al. (1996)	107
5.3	The life cycle of a data mining project, from the Chapman et al. (2000)	109
5.4	Example of unstructured data from the <i>Sun</i> .	121
5.5	A classifier is trained by using examples that are labeled, and then used to predict labels for a test set that has not been seen before. The predicted labels are then compared with the true labels and the accuracy of the classifier can be determined by finding out the fraction of examples that the prediction was correct. Graphic by Pereira et al. (2009)	125
6.1	Features that affect the perceived quality of the readers	145
6.2	Feature importance score for high claps bucket	153
6.3	The importance measures of Top-5 features in six different article categories	154
6.4	Feature Importance from the Random Forest Classifier	180
6.5	Confusion Matrix of the XGBoost Classifier for the small dataset	181
6.6	Visualization of one decision tree of the XGBoost classifier	182
6.7	Visualization Example of Quality Dimensions using the Theoretical Framework	184
6.8	Visualization of one decision tree of the XGBoost classifier for “Total Interactions” prediction.	191
6.9	Contributions to the low- and high-engagement bucket for “Total Interactions” prediction	192

6.10 Contributions to the low- and high-engagement bucket for “Total Interactions” prediction	193
6.11 Visualization of the feature importance of the Random Forest Classifier for “Likes” prediction.	202
6.12 Visualization of the feature importance of the Random Forest Classifier for Quality prediction.	203
6.13 Feature importance score for the content-based features	211
B.1 Feature importance for the high-claps class.	297
B.2 The contribution of each feature to the prediction made by the model.	298
B.3 The two buckets of the Quality variable.	299
B.4 Dendrogram	302
B.5 Feature Importance from the Random Forest Classifier for the whole dataset . .	304
B.6 Visualization of the leaf samples using the dtreeviz library.	306
B.7 Visualization of the correlation matrix.	307
B.8 Explanation of a prediction using LIME.	307
B.9 Permutation importance using ELI5 library.	308
C.1 Formula for “Flesch-Kincaid Grade Level”	310
C.2 Formula for “Flesch Reading Ease Score”	310
C.3 Flesch Reading Ease Scores	310

List of Tables

6.1	Descriptive statistics	150
6.2	Permutation Importance for the top 10 features	152
6.3	Precision, Recall, and F1 scores	153
6.4	Comparison of Top-10 and Top-5 Articles Categories' Features Similarity	155
6.5	Creation of the features.	175
6.6	Baseline Models	178
6.7	Accuracy scores of the different classification algorithms	179
6.8	The rules extracted from the large final leaves of the XGBoost classifier for the High Quality class	183
6.9	The rules extracted from the large final leaves of the XGBoost classifier for the Low Quality class	183
6.10	Creation of the features	189
6.11	Permutation Importance for the top ten features	190
6.12	The creation of the image-based features	200
6.13	The accuracy of the models	201
6.14	Permutation Importance for the top 10 combined features	212
6.15	Accuracy of Machine Learning Classifiers	212
B.1	Variable Correlation	300

Chapter 1

Introduction

1.1 Quality and Engagement in the Digital News Era

Quality in journalism is defined by a combination of norms including accuracy, objectivity, information quality, and diversity, which have long been considered essential to good journalism but now they have become difficult to maintain in the fast-paced, highly competitive environment of online news. In addition to traditional criteria social media networks and their increasing influence on how news is reported by journalists and consumed by online audiences have introduced new dynamics, with engagement becoming one of the key factors of news organizations' success. News outlets are continuously under pressure to establish a sizeable online audience and bring more visitors to their websites that rely on platforms such as Facebook, Twitter, Instagram, and Google, for better visibility and profit generation. That in turn, has led to an overall prioritization of low-quality content like, for example, sensational stories with witty headlines over news that is of high journalistic quality. On the brighter side, crowd-sourced journalism is now prevalent because of social media, which gave journalists access to a broader audience and a huge pool of user-generated content. Additionally, as citizen journalism has grown, more viewpoints and diverse voices have been heard in reporting.

As technology advances, so, too, does journalism as a profession, with data journalism and

robot journalism to be progressively integrated and accepted by traditional news organizations. More importantly, the rise of data journalism and its resulting popularity has had a profound impact on the practice of journalism around the world. The use of digital technologies, big and open data, and the computerization of many aspects of life has changed the way journalism is conducted (W. Weber et al. 2018). Additionally, elaborate visualizations allowed for more in-depth analysis and opened up complex subjects to a wider audience. Similarly, the increasing use of artificial intelligence in the newsroom is a trend that has been gaining traction in recent years. Machine learning algorithms are increasingly being used to create stories from data such as financial reports or weather forecasts, to identify hidden patterns in big data, automate trivial tasks, and offer a personalized experience to the audience with the help of machine learning algorithms.

The practical application of data science and artificial intelligence in journalism is not only an issue of having access to new, intelligent software, but also of revolutionizing the media ecosystem. In order to incorporate data in their reporting, journalists have become more accustomed to using sophisticated tools, working with massive volumes of unstructured data, and learning how to code. Such a shift in the everyday routines of many newsrooms was seen by communication scholars as having the potential to revolutionize the field of journalism by making it more creative, democratic and transparent. However, the use of data, predictive algorithms, and robot systems in journalism presents certain challenges, such as the requirement for journalists to possess a broad spectrum of advanced technical skills, an understanding of how algorithms work and how to interpret the results. Journalists should also have the ability to recognize the ethical issues that may arise from working with those tools and address them sufficiently.

As online audiences become more accustomed to consuming media in different ways, the field of journalism is also forced to adapt. New methods of reporting and storytelling are being developed to keep up with the changing habits of online audiences. This requires journalists to be both adaptive and aware of their audience's needs, as well as to embrace change and experiment with new tools and methods. Journalists are striving to make their reporting more engaging and interactive, leveraging the advantages of social media to con-

nect with their audience and promote their work. In doing so, traditional norms may be re-defined, while new principles such as transparency emerge. For instance, often on social media the journalist's voice is heard while journalists are trying to provide an honest and objective view of the story they are reporting. Additionally, in data journalism projects, the sources of information and the analysis is provided openly to the audience in contrast with traditional practices. As a result, it is obvious that in order to stay up-to-date journalists must keep up with these changes.

Research shows that while mainstream media are still playing an important role in society, new forms of media and journalism are emerging across the world. According to Deuze and Witschge (2020), news startups do not necessarily mean the dissolution of a shared journalistic ideology, since the vast majority of those startups share the same commitments to the principles of journalism. Rather, it suggests that quality journalism needs to be continually re-evaluated in order to ensure the highest standards of reporting in the digital era, taking into account the new opportunities provided by technology.

Prior attempts to define what constitutes "quality" journalism draw on direct observations, audience preferences, expert judgments, and indirect indicators, considering factors such as the use of language, the accuracy and reliability of news content, the ethical principles guiding journalism, and the commitment of journalists to their profession (Arapakis et al. 2016; Urban and Schweiger 2014; Rosenstiel et al. 2007; Lacy 2000). Even though there is a rich body of literature, that provides definitions of quality journalism and how it can be measured, there is a lack of research from a computational perspective.

In the current era of "Big Data", the proliferation of data and the technological developments have caused traditional social science to implement data-driven approaches to investigate how people behave online. Kennedy et al. (2015) define the act of quantifying aspects of the world that had not previously been digitally represented as "datafication". They claim that this method improves data gathering, analysis, and interpretation. Social science can present a thorough and holistic analysis of data using modern computational tools, allowing for a more accurate understanding of underlying social phenomena, more trustworthy

decision-making, and more effective use of resources. This emerging field is called Computational Social Science (CSS), and has the potential to change social science research and provide a better and more effective approach of dealing with complicated societal challenges. Similarly, given the digitization and datafication of the media environment, computational methods can arguably support and enrich the interpretation of news quality.

1.2 Objective and Research Questions

This dissertation draws from Journalism Studies, and Computational Social Science and employs quantitative and computational methods to examine the notion of journalistic quality in digital journalism. The first overarching research question is: How can AI provide a nuanced understanding of what constitutes quality journalism in the digital age? To answer this question, two further research questions need to be posed. Can subjectivity, emotionality, entertainment, and quality of language predict journalistic quality in online news? And what is the exact contribution of the quality criteria to quality prediction? For instance, one hypothesis is that lower subjectivity predicts high-quality news stories while emotional coverage contributes to stories being of low quality. Additionally, research questions are posed concerning audience engagement, such as: Can a machine learning model accurately predict social media engagement of news? And what images are significant for predicting the engagement or the quality of a news story? Finally, a subsequent question is: Is the quality predictive model able to detect online disinformation? To answer these research questions and hypotheses five studies took place.

Therefore, this dissertation makes a significant original contribution to knowledge by answering those questions using computational social science. The findings of the five studies along with the resulting guidelines on how to write a news piece that is of high journalistic quality and resonates with readers can be used on various industry applications. Tools that utilize predictive models can be useful for publishers by helping them to better understand how to reach and retain their audience and to appreciate the factors that contribute to audience engagement. The outcomes of this dissertation could be of interest to journalists,

communication scholars, and news audiences alike.

1.3 Research Design

Computational Social Science involves a wide range of methodologies that can be used to answer specific research questions. For this dissertation, the Python programming language was used and for the data analysis pipeline, the Cross Industry Standard Process for Data Mining (CRISP-DM) model (Chapman et al. 2000) was followed as a framework for all the different studies. Specifically, this reference model provides a comprehensive set of guidelines for the development of data mining projects and consists of six distinct phases: Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment. Each phase includes specific activities and tasks that should be completed in order to ensure the success of the data life circle.

The first phase of the study involved understanding the problem at hand and forming hypotheses. To do this, the relevant literature was reviewed, key gaps were identified, and research questions were formulated. The theoretical framework was then developed and the research approach was discussed. For all the studies, news stories from blogs and news organizations were scraped from the internet using computer software to create the datasets. After that, for Data Understanding, the unstructured textual data was examined and the right strategies were designed and deployed to convert it into a structured dataset. In the Data Preparation stage, multiple techniques were implemented to clean and prepare the text for further analysis. Then the raw data had to be transformed into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy on unseen data. This is a key step in machine learning, as it can have a significant impact on the performance of the models.

The modeling phase involves creating and refining a model to solve a problem. This involves choosing an algorithm and running experiments to determine the best model. Different classification algorithms have been used such as naive-Bayes, nearest neighbors, decision trees, logistic regression and support vector machines. Additionally, evaluating the results

of a machine learning model is necessary to assess its performance and determine if further improvements are needed. Therefore interpretable models were preferred since they are designed to be more transparent and explainable in their decision-making process, as compared to deep learning models which are more difficult to understand. There are various techniques that can be used to extract insights on the decisions made by explainable models, and in the evaluation stage all of the available Python libraries were utilized to extract insights on the models decisions. Finally, in the deployment phase, the results were analyzed, the AI models were interpreted, and the value of each study was discussed.

1.4 Structure of the dissertation

In the introduction, an overview of the research is provided. In chapter 2, the basic terms are defined, which will serve as a reference point for the research that follows. Specifically, the definition of journalism as a discursive institution is adopted according to Vos and Thomas (2018). In addition, in section 2.3 the concept of journalistic quality is examined through the lens of the supply side, which focuses on the media's performance of content (McQuail 1992), as opposed to the demand side (Ruggiero 2000), which is concerned with satisfying the needs of individual members of the audience. The discussion of quality is based on a normative understanding of the term, and particular attention is given to the research of Kovach and Rosenstiel (2014), who identified the key principles of quality journalism, and McQuail (2005), who proposed basic values of media quality that define the performance of a media organization. A number of other studies on journalistic quality are also reviewed and this information is later used to construct a theoretical framework for measuring quality.

In section 2.4, the rise of data journalism is examined. Data journalism is not just “doing journalism with data” (Rogers 2014), but can be further divided into two main categories: investigative data journalism and general data journalism (Uskali and Kuutti 2015). Investigative data journalism requires a significant time investment, advanced data and coding skills, and a focus on uncovering hidden information within data through investigative reporting methods. In contrast, general data journalism tends to have shorter production

times, requires less advanced skills, and has a less ambitious professional approach.

Furthermore, besides analyzing numbers, research tools for text analysis and coding are growing in popularity among journalists. Journalists have traditionally used qualitative methods such as interviews and observation to gather information. These methods, nevertheless, have been enhanced in the digital age with the use of text analysis and coding, opening up new ways for the discovery of important stories. Additionally, journalists can quickly and reliably examine enormous amounts of unstructured data using natural language processing (NLP). Moreover, NLP can be used to write news articles, follow stories as they develop online, and keep track of shifts in public opinion, giving reporters access to the most recent information, so they can always stay updated with the latest developments.

Additionally, the use of artificial intelligence in news organizations is discussed in section 2.5, with a focus on the perspective of Beckett (2022), who emphasizes the importance of AI in modern newsrooms. Beckett suggests that AI can improve the work of human reporters by producing a more comprehensive journalistic output and help ensure the survival of the organization. Examples of “robot journalism” include the automatic production of content for social media, the creation of personalized newsletters, and the automated publication of sports results. Automation can also aid in the transcription of interviews and enhance advertisement planning at a business level. While AI may have certain advantages, it also has some drawbacks, thus it is crucial to use it responsibly. As a result, algorithmic accountability is vital for ensuring that algorithms remain transparent, impartial, and fair employed in a manner that complies with regulations and ethical principles.

Chapter 3 is about social media, which had a significant impact on the media industry and the way journalism is practiced. In section 3.1 The architecture of those platforms promotes participation, creation and dissemination of content, enabling users to express themselves, share ideas, report wrongdoing, and act as watchdogs, similar to the role of journalists. The media ecosystem has evolved from a one-way communication model to a participatory process between news producers and the public, transforming the process of a centralized news monologue into a conversation. What is more, social media has decentralized journalistic

authority and changed the rules of the media game, resulting in the emergence of social media journalism (Kuyucu 2020), which is more open and dynamic, and relies on two-way communication. Naturally, publishers have adopted social media in their daily journalistic practice as a means of marketing, and profit-generation to avoid being left behind in the digital disruption of the media industry. At the same time, social media has also provided journalists with access to a large pool of free, textual and visual information to select from, as well as the ability to discover breaking news and eyewitness footage, particularly in crisis situations.

Section 3.1.1 focuses on the changing role of journalists in selecting and publishing news, with the intervention of social media and the public in the news-making process. The rise of social media and changes in news consumption patterns have disrupted the traditional journalistic model, known as “gatekeeping” in which journalists acted as gatekeepers who controlled the flow of information and influenced social reality. In this model, journalists had an asymmetric relationship with the public and shaped the news based on their own assumptions about what the public would find interesting. Some journalists now act as content aggregators by linking and reposting news, and rely on social media users for feedback. This has caused a new form of reporting to emerge, called “gatewatching” (Bruns 2018) in which journalists must navigate and distinguish newsworthy events from the internet.

Section 3.1.2 is a thorough discussion of audience participation. Specifically, engagement with the news refers to the attention and interaction paid by the audience to a news story. It can be seen as a means to encourage audience consumption and participation in the news, and can be evaluated in terms of the time spent, shares, likes, and other indicators of involvement. From a journalistic perspective, engagement can be seen as confirmation that the work is relevant and important to the audience.

Furthermore, the incorporation of emotions in news articles is an important topic in section 3.2 of this chapter. Previous research has shown that the inclusion of emotions in reporting can make readers more immersed in the narrative, as it elicits an emotional response. While this may conflict with the traditional principles of impartiality in journalism,

some scholars have called for a shift towards an “emotional turn” in journalism, particularly given the current hybrid media landscape and the presence of networked publics (Wahl-Jorgensen and Pantti 2021; Lecheler 2020; Chadwick 2017).

The chapter concludes in section 3.2.4 by referring to disinformation, since the use of emotional language has also been identified as one of the major strategies employed by fake news creators. While journalists and professional fact-checkers can verify the accuracy of potential threats through their expertise, the field of disinformation has attracted significant attention from computer scientists, who use machine learning and other automated methods to identify fake news. The chapter explores previous studies in this area with a focus on linguistic approaches.

Chapter 4 begins with the “iconic turn” of journalism and how the role of images has evolved through the years. An image that is emotionally powerful, unique, well-composed, timely, and focused on a human-interest story is more likely to be memorable and perform well on social media than an image that lacks these qualities. The value of the image is also examined through a marketing and advertising perspective, with attributes such as colors, image aesthetics, human faces and emotions to be very significant.

In chapter 5 the research design and the methods used are thoroughly discussed, and in chapter 6 each of the five studies is presented. Every study investigates a specific aspect of news stories, although all of them follow the same data mining framework and machine learning methods. The first study 6.1 investigates the use of artificial intelligence to evaluate the perceived quality of stories on the blogging website Medium.com. The results showed that machine learning models, including decision tree, random forest, and XGBoost, can accurately predict the success of an article on the blogging platform based on features related to the author, style, content, and context. The second study 6.2 examines the use of AI to determine the quality of news stories based on a theoretical framework derived from the literature.

The third study 6.3 looks at the use of machine learning algorithms to predict engagement on social media. The findings demonstrated that tree-based approaches, support vector

machines, and logistic regression can accurately predict audience engagement on Facebook. The fourth study 6.4 focused on the use of the visual elements of journalism to predict both the engagement and the quality of a news article. The featured image of each news story was first used for image recognition and then as an input to the classifiers. The findings showed that images alone can reach an accuracy of 0.70 (F1-score), but when combined with textual features they do not improve the performance of the models. The final study 6.5 focused on the application of AI for disinformation detection. The results showed that textual attributes can predict fake news with high accuracy. Furthermore, the accuracy of the models was improved when combined with additional engagement features from Facebook such as likes and comments.

The dissertation's final chapter 7 discusses in detail the findings and limitations of the studies and provides a quick review of future research. One important component of this work is its ability to transfer knowledge from the research findings to ongoing projects in academia and the media sector, with the goal of advancing both the theoretical foundations of the research and its practical implementations. This chapter also analyzes the potential influence of the research on the Greek media industry, discusses possible future uses of the findings, and gives a road map for further research and applications.

Chapter 2

(Re)defining Quality

2.1 Navigating the Changing Landscape of Digital Journalism

When discussing journalism in today's diverse media landscape one has to take into account the massive transformation in the field due to shifts in social, political, cultural, and economic planes to which journalism struggles to adapt. It is also paramount to consider, the important manifestations of technological change in the field of journalism such as data science (data journalism) and artificial intelligence (robot journalism) (e.g. Carlson (2015)). Regardless of whether we discuss watchdog journalism or lifestyle stories, we understand journalism as a social institution (Vos and Thomas 2018), as a "set of values, principles, and practices enacted in different ways and settings with a "sense of wholeness and seamlessness" (Hallin 1992, p. 14), and as such the media organizations act as "an ordered aggregate of shared norms and informal rules that guide news collection" (Sparrow 2006, p. 155). Following this approach, Hanitzsch and Vos (2017, p. 121) described the journalistic roles as "the discursive articulation and enactment of journalism's identity as a social institution. Hence, journalistic roles set the parameters of what is "appropriate" or "acceptable" action in a given context". In this sense, journalistic roles are not only seen as a set of rules and norms but also as a set of values and principles that enable journalists to engage in mean-

ingful activities and practices, such as investigative journalism, that can potentially lead to positive social change.

The concept of occupational ideology of journalism has also transformed its discursive institutionalism encompassing a worldview, values, practices, ethics and principles to which professional journalists in elective democracies adhere to. Communication scholars have frequently studied in great detail the role of journalists and editors in news making, their beliefs, and strategies in a universal context to spot similarities for creating a consensus shared among professional journalists. In fact, multinational studies have shown that even though social and political differences exist and the practice of journalism differs from country to country, professional principles such as detached reporting, immediacy, reliability, factuality of information, impartiality, and neutrality are universal denominators around the globe (Papathanassopoulos and Miconi 2023; Hanitzsch et al. 2011; Weaver and Wu 1998). Although it seems as if journalists are following a universal ideology, the way they interpret those principles and the way and the degree in which they apply them is not the same, depending on the country, the individual and the media outlet (Deuze and Witschge 2020; Hanitzsch et al. 2011; Shoemaker and Reese 1996).

Deuze (2005) after rigorous investigation of the strong global coherence ruling the journalistic roles and views conceptualized the ideology of journalists that can be defined by values such as the public service ideal, credibility, autonomy, immediacy, and legitimacy, and provide journalists with an “exclusive role and status in society” (Deuze 2004). Ten years later along with his colleague Tamara Witschge, she started the project “Beyond Journalism” which focuses on journalistic startups around the world. The results showed that even though journalists and editors criticize legacy news organizations, the “commitment to the ideology of journalism remains firmly in place” (Deuze 2019). Although journalistic ideology remains significant for the study of journalism, they emphasize the need to go “beyond” the industry and include entrepreneurs in the “what journalism is” discussion (Deuze and Witschge 2020).

In this work, the definition of journalism as a discursive institution (Vos and Thomas 2018)

will be used. In this thesis, journalism will not be defined based on different models of democracy, genres, types of outlets, or different regions of the world, but rather the research will focus on some dominant theories of the study of journalism such as occupational ideology (Deuze and Witschge 2020; Deuze 2005), news values (Harcup and O’neill 2017; Harcup and O’neill 2001), journalistic roles (Hanitzsch and Vos 2017), and best practices (Carlson and Lewis 2015) that serve as a benchmark, and have provided a general understanding of the way journalists think and produce journalism (Deuze and Witschge 2020; Weaver and Wu 1998).

2.2 Quality in Journalism

Current technological developments are shifting the professional identity of journalists, who are now described as the “robotic reporter” (Carlson 2015), “exo journalist” (Tejedor and Vila 2021), “networked” (Van der Haak et al. 2012), “reinvented” or “social” journalist (Olausson 2017) and an audience-oriented editor (Ferrer-Conill and Tandoc Jr 2018), while “real” journalism is depicted as “ambient”, “beyond” or “niche” journalism (Ruotsalainen 2018). In the past, news editors used to be autonomous and neglected their audience, thinking they already knew what people expected of them. However, audiences nowadays are more critical towards all media, and demand in-depth analysis (Ruotsalainen 2018) and as J. B. Singer (2005, p. 179) puts it, they want to know “how the sausage is made”. The fact that the internet encourages users to help sort out the most engaging journalistic products, provides a unique opportunity for more relevant stories (Hanusch 2017), online participation, transparency, and overall better journalism (Ruotsalainen 2018, p. 19). This reciprocity between citizens that made up the audience and the press is evident in recent endeavors from the *New York Times* (T. A. Evans 2019) to identify exactly which elements of their news stories are highly engaging in order to “shape (their) output” in the digital world.

Recent studies have demonstrated that online journalism is taking a data-driven approach, congruent with a shift in the press to provide more quality news in an effort to regain public trust. Quality news has the ability to attract more readers and establish new, sustainable

subscription models, like paywalls, to help compensate for the financial instability in the industry which has been caused by changes in the media landscape, and the decline in advertising revenue (Hanusch 2017; Papathanassopoulos 2002).

Google and Facebook have been transforming the media landscape by hosting news content without paying for the rights and monetizing it as the new gatekeepers of news. They have prioritized cheap and engaging content over more costly, but high-quality material, resulting in clickbait and platforms for disinformation campaigns (Bell 2017). In an effort to counter this development, Google recently announced its intention to pay certain news outlets for paywalled high-quality stories to provide access to its users through its Google News Showcase product¹ that launched in Brazil and Germany (Cellan-Jones 2020; Drozdiak 2020). This move is aimed at supporting quality journalism and appeasing its long-standing conflict with publishers; however, it raises the concern that it could make Google even more powerful in controlling news on a global scale.

Technology platforms and social media networks disrupted journalism by changing the behavior of media professionals as well (Tandoc Jr and Ferrucci 2017). Even though the cornerstones of institutional journalism remained the same throughout history (Deuze and Witschge 2020), when one delves deeper into the newsroom's everyday routines one realizes that journalism as an activity depends heavily on the working environment (Hanitzsch and Vos 2017). This environment might curtail some of the principles of journalism, such as the freedom to decide how and what to write about, over being more relevant to the audiences (Nygren 2012). The incorporation of web performance metrics in journalists' working routines and more specifically the reliance on audience metrics can lead to a lack of nuance in reporting, and lower quality of journalism (Agarwal and Barthel 2015). Such practices in the day-to-day work include: reorganizing the website (Tandoc Jr 2019), choosing a story to write about as a means of increasing social media engagement (Tandoc Jr and Ferrucci 2017), deciding which story to follow up (Tandoc Jr 2014) or even altering an article that does not perform very well online (Hanusch 2017; Vu 2014), such practices are nowadays commonplace. This digital-first journalism, except for driving many editorial decisions based

¹<https://blog.google/products/news/google-news-showcase/>

on internet-oriented and data-driven practices, also requires a constant flow of information that makes revisions and reappraisals otiose.

Within the seeming friction between audience-driven online journalism and a move towards more quality, journalism, normatively speaking, it is still crucial for democracy, given its “primary purpose to provide citizens with the information they need to be free and self-governing” (Kovach and Rosenstiel 2014, p. 61). Values of media quality like freedom, equality, diversity, the obligation to report the truth, information quality, and so on (Kovach and Rosenstiel 2014; McQuail 2005) define the performance of a media organization. It is important to differentiate three levels of how media organizations operate (McQuail 2005), namely their structure (media system), conduct (organizational level), and performance (content). The first level is media structure which includes aspects like ownership, finance, and infrastructure, the second is conduct which relates to policies, editorial decisions, and public relations, and the third is media performance which is the product that ends up to the individual members of the audience.

2.3 Previous Efforts to Measure News Quality

The quality of journalism has long been a topic of interest among scholars and media professionals alike. Many different perspectives, formats, and levels have been used to measure journalistic quality, including asking both journalists and audiences for their understanding of good journalism, as well as analyzing the content of news outlets (Arapakis et al. 2016; Urban and Schweiger 2014; Rosenstiel et al. 2007; Lacy 2000). In journalism studies, a differentiation is often made between the demand and supply sides of the field. The demand side refers to the behavior and feedback of the audience or intended audience (e.g. individuals, households, and communities), and the content provided to them to satisfy their needs. In contrast, the supply side focuses on the services provided by journalists and can include characteristics of the content and human resources.

Much of the research on the *demand side* of journalism focuses on the consumer as an active participant in mass media, and is based on the theoretical framework of uses and gratifica-

tions (Ruggiero 2000; McQuail 1994). This approach examines the motivations of audiences for consuming news and has identified different types of media use, such as diversion, personal relationships, personal identity, and surveillance (Blumler 2019; Katz et al. 1973), as well as decision-making (Lacy 2000), social interaction, civic responsibility, and improving their lives and professional skills (Purcell et al. 2010). In this framework, content that satisfies these gratifications is considered to be of high quality.

In a study by Gladney (1996), the author sought to move away from measuring news quality based solely on user utility and instead focus on editorial quality. To do this, Gladney used a set of 18 quality criteria divided into two dimensions: organizational standards (e.g. integrity, staff enterprise, community leadership, editorial independence, staff professionalism, editorial courage, decency, influence, impartiality) and content standards (e.g. news interpretation, lack of sensationalism, strong local coverage, visual appeal, accuracy, strong editorial page, community press standard, comprehensive coverage, good writing). The findings showed that, in general, newspaper readers and editors agreed on how to evaluate the journalistic quality and considered important characteristics such as integrity, impartiality, editorial independence, strong local news coverage, accuracy, and good writing to be key indicators of quality. However, the study also found that editors placed greater value on staff enterprise and visual appeal, while readers appreciated decency and a lack of sensationalism more. To better understand readers' preferences for news content, Urban and Schweiger (2014) tested a set of normative dimensions with the audience to evaluate the degree to which readers value diversity, relevance, ethics, impartiality, objectivity, and comprehensibility. The results revealed that readers rated relevance, impartiality, and diversity highly; however, the audience was unable to differentiate quality based on ethics, objectivity, and comprehensibility. This could be due to the audience's focus on the (manipulated) media brand images.

With the digitization of the media landscape, news organizations welcomed into the practice of journalism web analytics companies (Q. Wang 2018) that offered sophisticated tools and structured datasets for further experimentation into the issue of audience engagement with news. Some recent studies harness the power of data science and use engagement met-

rics as a proxy to measure quality from the audience's point of view. News values (Harcup and O'neill 2017; Harcup and O'neill 2001) have been used to shed light on what drives users to prefer some news stories over others, with characteristics such as geographical proximity, referring to Western countries, conflict, human interest, sentiment, and exclusiveness (Trilling et al. 2017) proven to well predict preferences. The platform "Metrics for News"² by the *American Press Institute*, offers publishers a way to identify which news is perceived as more valuable for the readers by correlating the qualities of a single news story against engagement data such as page views.

With the present contribution's focus on textual features to explain journalistic quality, a more thorough discussion of the *supply-side* perspective is in place. Adopting a supply point of view, the scholarship is focused on a product-based approach and is mostly grounded on the notion of media performance (McQuail 1992). The basis for this approach lies in a conception of the journalistic product as containing a set of unique attributes indicative of quality rather than content with the only intent to satisfy the individual audience members' social and psychological needs.

Some of the first and most supported attempts to measure quality was the study of Bogart (1989), which surveyed 746 newspaper editors and created a scheme of 23 high-quality identifiers for newspapers, such as high "readability" score, the ratio of illustrations to text, and length of an average front-page story. Another scheme of quality scores was created by Rosethiel and colleagues who conducted a five years examination of local TV broadcasts in the United States that provided them with a list of quality identifiers that they correlated with television ratings to see if they resonate with the audience (Rosenstiel et al. 2007). The project included extensive content analysis, along with annual surveys of news directors, focus groups with viewers, and panels with veteran journalists. The measurements used were the significance of the topic, breadth of issues covered, the narrow or broad focus of the story, level of effort by the journalists, accuracy, and fairness, number of sources and viewpoint diversity, authoritativeness (meaning the level of expertise of the sources cited), clarity of information and sensationalism. The authors used these indicators to predict the audi-

²<https://www.metricsfornews.com>

ence ratings based on this quality criteria and demonstrated “that it’s not what you cover, it’s how you cover it that matters” (Rosenstiel et al. 2007, p. 117). They created a “magic formula”, a set of strategies for attracting viewers by focusing on high-quality elements, such as divergent and balanced viewpoints, investment in investigative journalism, quoting multiple authoritative sources, local relevance, and providing longer stories on important matters.

In their book “The Elements of Journalism”, Kovach and Rosenstiel (2014) identified 10 key principles of quality journalism. Amongst these are the obligation to report the truth, loyalty to citizens, verification, independence, monitoring those in power, and making the important news interesting and relevant. Similar to their approach McQuail (2005) suggested basic values of media quality to define the performance of a media organization. These include freedom, equality, diversity, truth, and information quality, and social order and solidarity. Undeniably, information quality and diversity are two of the most frequently used quality criteria in the literature.

The current contribution is methodologically closer to the work by Tang et al. (2003) who created automatic textual features based on a set of information quality criteria of news articles, such as objectivity, depth, readability, and grammatical correctness. They found that the automatic quality measurements they used captured successfully four of the information quality attributes, namely depth, objectivity, multi-view, and readability. Arapakis et al. (2016) consulted media professionals and computational linguists to identify content-based linguistic characteristics of news stories that can define an article’s editorial quality. They proposed a taxonomy of 14 attributes of the editor’s perceived quality designed around five main dimensions namely, readability, informativeness, style, topic, and sentiment. Experts annotated a sample of 561 news articles from a single source, creating a corpus that was used for correlation analysis and quality prediction. The findings showed that characteristics such as fluency, completeness, and richness were the most strongly correlated to quality, while subjectivity and polarity showed the least strong correlations. Furthermore, the researchers used a generalized linear regression model to predict the quality of a news story. Even though they admit that automatic prediction of news quality proved to be a challenging task, it showed that their news content decomposition was pointing in the right

direction. More recently Choi et al. (2021) combined news values manually annotated by humans, automatically constructed linguistic features, and audience attributes to predict quality via a neural network. They arrived at the conclusion that the stronger predictor was the group of journalistic values because when running the deep learning model without it the results worsened significantly. The news values used were believability, depth, diversity, readability, objectivity, factuality, and sensationalism. Following the literature reviewed in sections 2.2, 2.3, this study proposes a normative theoretical framework of journalistic content quality and subsequently applies this framework to online news stories using machine learning algorithms.

2.4 Data Journalism

Digital transformation can describe the zeitgeist of the last decade and refers to any radical transformation of the traditional ways of interaction, creation of value, and communication in both social and professional contexts caused by digital technology advancements. Undoubtedly, the Internet, social media platforms, data availability, Technology Giants such as Meta (former Facebook) and Google, and various new technologies disrupted the mass media ecosystem and changed both the journalistic routines inside the newsrooms and the business models of the news outlets. The introduction of computer-assisted reporting (CAR) in the 1960s has revolutionized the field of journalism. Initially developed by Philip Meyer of the Detroit Free Press, CAR has allowed journalists to efficiently and accurately analyze data related to their reporting. Meyer utilized a computer to automatically analyze the demographics of Detroit's 1967 riots (Royal 2010). Since then, CAR has been used to integrate social science-based methods of data collection and statistical analysis into the journalistic process (Hannaford 2015). As technology has advanced, journalists have been able to add additional capabilities such as archival and online research as well as conducting interviews via email (Coddington 2015). The integration of CAR into the journalistic process has presented a great opportunity for journalists to gain a deeper understanding of their topics and provide more accurate and comprehensive reporting. Furthermore, CAR has enabled jour-

nalists to access a larger and more diverse range of data sources, leading to greater diversity in news content. The development of data journalism has been a major boon for the field of journalism, allowing for more accurate, comprehensive, and diverse reporting. With further advancements in artificial intelligence, data journalism will continue to evolve and grow.

In the late 1990s, the Internet and its associated services had a marked impact on the news industry, resulting in the emergence of new models of journalism. BBC Online was one of the first to take advantage of this new technology, forming a team of professional and inexperienced journalists who embraced digital practices and adopted the identity of the digital journalist (Hermida and Young 2019). This digital transformation of news delivery allowed for human-computer interaction which resulted in the creation of new forms of journalism, such as multimedia, visual, participatory, and data journalism. One example of how this new model has been embraced is in the case of data journalism, where datasets are analyzed to reveal hidden information and stories are accompanied by visualizations to make complex subjects accessible to everyone (Flew et al. 2012). Through this way of data representation and in-depth explanation, important points can be highlighted and communicated to readers. As Boyles and Meyer (2016) suggest, this type of journalism has the potential to democratize journalistic practices since data journalists need to communicate and explain their data analyses to the audience and be transparent about the whole process (Boyles and Meyer 2016; Lewis and Usher 2013). Also, the new ways of presenting data to interested citizens could potentially lead to a more inclusive and participatory form of journalism (Palomo et al. 2019).

Data journalism is seen by many scholars as a result of digital technologies, big and open data, and the computerization of many aspects of life. As W. Weber et al. (2018) note, this has given rise to novel ways of conducting journalism. Big Data lies at the root of data journalism, and in this environment, the software and data journalism tools were built to allow journalists to process and analyze large amounts of data. Although, a wide range of skills is required to conduct data-driven investigative journalism such as programming, statistics, visualization techniques, and database management (Ausserhofer et al. 2020; Parasie and Dagiral 2013). Subsequently, the use of these tools and methods has changed the way that

journalists can gather information, generate news stories, and convey them to their audiences.

Data journalism is a complex field of work that requires a range of skills and knowledge. While previous definitions of data journalism simply note that it is done with data (Rogers 2014; Gray et al. 2012), Uskali and Kuutti (2015) provide more insight into the two main streams of data journalism, namely Investigative data journalism and General data journalism. Investigative data journalism is characterized by a large time commitment, advanced data, and coding skills, and an investigative reporting ethos focused on discovering hidden “gems” within data. By contrast, General data journalism is characterized by shorter production times, basic skills, and a less ambitious professional attitude. Both streams can operate within newsrooms, providing a range of data journalism services. However, while the distinctions between the two streams are useful, they should not be seen as mutually exclusive. For example, there may be cases where Investigative data journalism is produced in short timescales. Furthermore, while General data journalism may be seen as less ambitious, it can still provide valuable insights and stories, and may even uncover important stories that would otherwise go unnoticed.

In recent years, data journalism has become increasingly popular, with major news outlets such as the New York Times, Reuters, the Guardian, and the BBC leading the way in 2012 (Pellegrini 2012). These pioneering news organizations dedicated resources and personnel to the development of data stories, a trend that has continued to grow in the years since. Now, data journalism is commonplace, with many news outlets using data to tell stories that would otherwise go untold. Data journalists have also been using data to investigate and report on a variety of topics, from public health and environmental issues to criminal justice and political developments. Furthermore, data journalism has become an invaluable tool for uncovering patterns and trends, providing an objective and quantifiable look at the world. Undeniably, the success of data journalism has been largely attributed to the emergence of open data initiatives and the availability of data sets from governmental, non-governmental, and corporate sources (Brugger et al. 2016; Stoneman 2015; Anderson 2013). This has given journalists access to a wealth of information and the ability to create sto-

ries from that data. However, there are still many challenges associated with the practice that need to be addressed in order for it to be used efficiently. First, smaller newsrooms do not have the financial capacity to support data-driven investigations and rely on their employees' own will to learn, with many data journalists to be self-taught using free online courses and open source materials (Porlezza and Splendore 2019). Another challenge is the availability and accessibility of public datasets (Porlezza and Splendore 2019, p. 1230) which is not always a given (Borges-Rey 2016), driving many newsrooms to create datasets by themselves through crowdsourcing, conducting online or offline surveys, scraping and so on (Zamith 2019). Scraping is the practice of utilizing software to extract and collect data from websites or other online sources, and it could serve as an effective method for gathering and organizing massive volumes of data. It is important to note that a recent ruling by a U.S. appeals court, (Whittaker 2022), found that scraping publicly accessible data is acceptable in the United States. This is advantageous for researchers, scholars, and archivists, who frequently use publicly available data in their work.

When datasets are not open or do not exist, journalists take on the task of creating their own datasets in order to produce a more in-depth analysis of a particular subject. The Guardian's "Behind the Blade"³ investigative multimedia series is a prime example of this type of data journalism. The purpose of this series was to look into the issue of knife crime involving children and adolescents in the United Kingdom. To accomplish this, the journalists meticulously gathered all the relevant information and data about these crimes, including reported locations, victim ages and genders, and the kinds of knives used. Readers were able to explore the specifics of each crime by using an interactive map created from the thoroughly analyzed data that has been acquired. Interviews with specialists and the relatives of the victims were included in this map to provide further awareness of the knife crime problem. By combining data analysis with personal stories, the "Behind the Blade" series was able to present an in-depth look at an important social issue. The success of "Behind the Blade" demonstrates the potential of data journalism and manual data collection.

³<https://www.theguardian.com/membership/ng-interactive/2017/mar/28/beyond-the-blade-marking-the-death-of-every-child-and-teen-by-a-knife-in-2017>

The use of online activist crowdsourcing platforms by newsrooms is also an effective way of collecting data. ProPublica's use of the gamified collaborative news platform, Vozdata, for their "Free the Files"⁴ investigation is a notable example. Through Vozdata, ProPublica was able to collect election spending data from thousands of receipts and having readers sorting it. Similarly, The Guardian used the platform to collect data from 700,000 documents and 5,500 pdfs concerning the House of Commons and its 646 members for the story "How to crowdsource MPs' expenses"⁵. La Nacion⁶, meanwhile, was able to investigate 95,000 telegrams published by the Argentina government authorities from each polling station. The data collected through these crowdsourcing platforms is a valuable resource for investigative journalism. It provides an original source of information and an essential foundation for further investigations. This data is crucial in uncovering the truth and providing the public with an accurate overview of the electoral process. Moreover, by collecting, analyzing, and presenting large datasets in an accessible and interactive format, journalists are able to create compelling stories that provide valuable information for democracy (Boyles and Meyer 2016; Stoneman 2015; Baack 2015).

2.4.1 Computational Journalism

The use of technologies such as machine learning, natural language processing, artificial intelligence, and big data analysis in journalism has the potential to revolutionize the way stories are created and is referred to as computational journalism (Coddington 2015). These technologies provide journalists with unprecedented capabilities to uncover hidden trends, detect patterns, and provide predictions. Hamilton and F. Turner (2009) have highlighted the dramatic shift in the landscape of journalism due to the proliferation of "ubiquitous computation". They note that this shift has weakened traditional business models, altered the power dynamics between journalists and their audiences, and enabled the rapid dissemination of news. While computational journalism is unable to completely resolve the economic issues present in contemporary journalism, it has the potential to create tools that

⁴<https://www.propublica.org/series/free-the-files>

⁵<https://www.theguardian.com/news/datablog/2009/jun/18/mps-expenses-houseofcommons>

⁶<https://blogs.lanacion.com.ar/projects/data/vozdata-2015-unveiling-argentinas-elections-system-failures-with-impact/>

reduce the cost of investigative reporting, optimize the new information environment, and facilitate the continuation of watchdog work during the current technological revolution.

Big data analysis can be used to uncover correlations between different variables, predict future trends, identify patterns and outliers, and explore relationships between different data points. By understanding these correlations, journalists can better inform their audience and assist them in decision-making (Flew et al. 2012). For example, a journalist could use big data analysis to uncover how different factors, such as age, gender, and income, impact voter turnout in an election. This could inform their reporting of the election by providing more context and revealing deeper insights to the readers. Big data analysis can also be used to uncover insights about a particular population or used for audience analysis. For example, a journalist could use big data analysis to understand the impact of a particular policy on a certain demographic or to uncover correlations between different types of media consumption and people's political beliefs. In addition, big data analysis can also be used to improve a newsroom's workflow and operations. For instance, a newspaper might employ big data analysis to examine story readership patterns or to check on the effectiveness of its advertising initiatives. This might make it easier for the media company to comprehend its audience and modify its material to better fit its readers.

The introduction of computational solutions and coding into the newsroom has allowed for greater levels of openness and collaboration (Zamith 2019; Tandoc Jr and Oh 2017). This has ensured that journalistic efforts are no longer completed by a single individual, and instead, the process is now completed through collective efforts. This move to collaboration has been seen in some of the most notorious news leaks, such as the Edward Snowden revelations, Julian Assange's WikiLeaks platform, the Paradise Papers (ICIJ 2017), and the Work crimes in Afghanistan and Iraq (Stray 2010b). These stories all resulted from investigations into anonymous news leaks and inaccessible material, such as emails, secret web addresses, videos, bank accounts and transactions, military or secret service documents, and more. Also, the collaborative nature of data journalism has invited experts from other fields such as data scientists, developers, and other professionals to become involved in the news production process (Hermida and Young 2017). This has allowed for the sharing of resources

and the completion of time-consuming tasks that would have otherwise been completed by journalists alone. Additionally, the presence of computer scientists provided a greater depth of understanding and analysis when it comes to automated classification tasks like in the case of Panama Papers⁷. Additionally, the open-source nature of data journalism allows for greater levels of transparency and collaboration (Tandoc Jr and Oh 2017; Stoneman 2015) that led to the formation of interdisciplinary communities such as Hacks and Hackers⁸ (Lewis and Usher 2014).

Automated web applications and tools for journalism appeared in the last decade aiming at assisting reporters in various stages of their work. Many tools focus on visualization providing ways to produce appealing graphics and maps with only a few steps. Specifically, visualization and mapping software such as QGIS⁹, Leaflet¹⁰, Datawrapper¹¹ and Geolocation services like Google Earth¹² are used from investigating journalists to collect, analyze, verify and visualize data and geographical information and discover stories. Furthermore, by analyzing location-based data, local news outlets can publish personalized news related to the preferences of the residents (Flew et al. 2012). In addition, some tools focus on data gathering and scraping. For instance, ScraperWiki¹³ and Import.io¹⁴ are two popular web crawlers that allow journalists to automate the extraction and collection of data from webpages and other sources. Moreover, automated data-driven journalism tools such as Tableau¹⁵, and Processing¹⁶ allow users to quickly build interactive, visual stories using large datasets.

2.4.2 Data journalism as a Profession

The boundaries between the three forms of journalism, namely, computer-assisted reporting (CAR), data journalism, and computational journalism are not so discrete, and as Usher

⁷<https://www.icij.org/investigations/panama-papers/>

⁸<https://www.hackshackers.com/>

⁹<https://qgis.org/en/site>

¹⁰<https://leafletjs.com>

¹¹<https://www.datawrapper.de>

¹²<https://earth.google.com/web>

¹³<https://scraperwiki.com>

¹⁴<https://import.io>

¹⁵<https://www.tableau.com>

¹⁶<https://processing.org/>

(2016) assumes, they depend on the education of the data journalist or the country they live in. More specifically, after discussing with newsroom employees from different outlets, such as Al Jazeera, Washington Post, FiveThirtyEight, Guardian and others about the labels associated with data-driven journalism, she found that the type of training received by practitioners shapes how they perceive their work. Specifically, Usher notes that those from a journalistic background usually classify the work as data journalism, while those with a programmer or computer science background typically refer to it as computational journalism. Additionally, she asserts that European practitioners refer to “data journalism” as what Americans define as “computational” or “programmer journalism”.

Many studies examined the role of data journalists, in countries like the United Kingdom (Borges-Rey 2016; Hannaford 2015), Norway (Karlsen and Stavelin 2014), Italy (Porlezza and Splendore 2019), the United States (Pavlik 2013; Royal 2010), China (Shuling Zhang and Feng 2019) and elsewhere. They demonstrate the importance of continuing to monitor the development of data journalism in order to understand how it is changing journalism and how it is impacting the public’s understanding of the news. Furthermore, the studies show the potential for data journalists to become increasingly important in newsrooms, as well as the need for journalists and programmers to collaborate in order to achieve their ultimate goals (Trappel and Tomaz 2021).

To ensure that journalists are able to keep up with the pace of technological advancements, it is important for newsrooms to provide adequate training and education in both programming and how to use data for journalism. The comparison of two different studies on data journalism reveals an interesting perspective on the changing nature of the profession. Royal (2010) found that journalists in the New York Times interactive news technology department, emanate from diverse backgrounds but share the same goal to learn how to code, use data for journalism, and master Django and Ruby to present news in an interactive way; Hannaford (2015) however found that journalists in the BBC and Financial Times are not trained in coding, but solely in the way automated data journalism tools work and how to make sense of the outcomes. This reveals that the training and skill sets required for data journalists are not the same and each newsroom has its own work practices in response to

data journalism (Splendore et al. 2016).

The study by Borges-Rey (2016) offers valuable insight into the adoption and development of data journalism in the United Kingdom. This survey of media professionals in fourteen British news organizations showed that high-quality news providers were quick to recognize the potential of data journalism to provide interactive, non-linear content to their audiences. The results of the survey provide an overview of the three main forms of data journalism present in the United Kingdom: “daily brief” data journalism, with a focus on the visual aspect; extensive and more “investigative” data journalism; and “soft”, entertaining, and gamified data journalism. The survey also found that almost all of the newsrooms had at least one data journalist on their staff. Furthermore, the survey highlights the value of data journalists and the role they play in newsrooms.

On a similar note, Karlsen and Stavelin (2014) investigated the role of the computational journalist as a rhetorical craftsperson by utilizing the Aristotelian view of *techne*. The authors interviewed data journalists in six of the largest Norwegian newsrooms and found that there is a “journalistic-programmer” (role additional to editor etc.), who collects, analyzes, and communicates the results of the investigation to the public. This notion was further explored by Pavlik (2013) in Chicago newsrooms, who distinguished between the roles of journalist and programmer, positing that the journalist’s ultimate goal is to uncover the truth and advance democracy, a task that cannot be completed by programmers working with data alone (Boyles and Meyer 2016). In Porlezza and Splendore (2019), a study of 15 data journalists in Italy was conducted, revealing that the field is highly interconnected within and beyond the country’s borders, producing quality journalism. The study also found that most of the financial support for data journalism in Italy comes from European grants and awards and that the field has a strong commitment to openness and transparency.

Undoubtedly, data journalism’s potential to fulfill (Karlsen and Stavelin 2014) or revolutionize journalism (Borges-Rey 2016; Flew et al. 2012; S. Cohen et al. 2011) and the role of data journalists as digital watchdogs (Felle 2016) is evident in many studies as is the integration of innovative technologies in news production. However, drawing from the nexus between

journalism and computer science, the scholarship on computational journalism has been primarily concerned with building digital tools to supplement, routinize, and algorithmically expand traditional news practices, and not so much with the impact of computation on larger social, political, organizational, and cultural currents in journalism (Anderson 2013). The current approach of many computational journalists has been to adopt a “just build it” attitude, that overshadows social institutions and aspects such as public policies, open datasets, and unequal access to technological resources and only focuses on how useful these computational tools are (Anderson 2013). It is therefore essential for data journalists to adopt a more socially conscious approach to their craft, including the wider implications to institutions outside journalism (Anderson 2013). This is especially true in countries such as China, where open datasets are scarce, and the ability of data journalists to innovate is limited (Shuling Zhang and Feng 2019).

2.5 Exploring the Possibilities of AI in Journalism

Artificial intelligence is a fast-expanding area of computer science and in recent years, has become increasingly intertwined with journalism, offering a range of possibilities for the profession. The term AI was first coined in 1956 by John McCarthy at Stanford University, where he argued that machines could be programmed to perform tasks that would require human intelligence if done by a person. For this reason, he gave the name “artificial intelligence” to this new field explaining that “it is the science and engineering of making intelligent machines, especially intelligent computer programs” (McCarthy 2004, p. 1).

An AI system is able to memorize information and use it to make decisions based on the changes in its environment (Alpaydin 2020), therefore it can be used to automate mundane tasks. Also, AI-based technologies are extremely versatile and can be applied in a broad range of disciplines, with dozens of methods providing different capabilities (Hansen et al. 2017). The desire for quicker and more accurate decision-making as well as the automation of complicated processes has led to the application of AI in a wide range of fields, from physics and agriculture to healthcare and finance. Furthermore, the research in the field

has advanced significantly over the past decades, with a massive influx of new technologies and applications ranging from computer vision, language and speech processing, robotics, knowledge representation and reasoning, problem-solving, machine learning, expert systems, man-machine interaction, and artificial life (Aarts and Encarnaç o 2006). However, these topics have been further refined and developed over the years.

One of the most fundamental uses of AI in journalism is automation. This can be used in journalism to automate tasks such as fact-checking, researching, and data entry, freeing up time for journalists to focus on more complex tasks (S. Cohen et al. 2011). In order to give journalists real-time information on the stories they are working on, automated tools can also be employed to scan news wires, social media, and other sources. In addition, automation can be utilized to generate reports and visualizations, allowing journalists to quickly and easily analyze data, and produce stories from complex datasets (Ausserhofer et al. 2020). Moreover, AI-based tools can assist journalists, by transforming scanned documents into usable formats, transcribing audio, and video recordings, and monitoring social media posts to uncover trends that may be complementary to a story that a journalist is working on. For instance, *Bloomberg*, *ABC News*, *the New York Times* and *ESPN* have developed AI-powered tools, specifically to transcribe audio files into text and generate summaries of news articles.

In newsrooms, machine learning (ML) has been increasingly utilized over the past years. ML is a subfield of artificial intelligence that enables systems to learn from data, recognize patterns, and generate knowledge autonomously, without human intervention or programmed guidance (Broussard et al. 2019). Alpaydin (2020) defines ML as mathematical models that, based on the training data or past experience, can make predictions in the future (predictive models) or extract insights from the data (descriptive models). Moreover, ML is mainly dependent on algorithms, which are sets of dynamic instructions that, when followed, convert the input to output (Alpaydin 2020) and ultimately lead to the desired results. In particular, ML develops classification or regression algorithms through the training set and evaluates their performance through the test set (J. Zhang et al. 2020). Finally, a very promising field of machine learning is deep learning which uses neural networks to learn from data sets and identify patterns within them (J. Zhang et al. 2020). Those types of methods can be used in

journalism to make predictions about future events, such as election results, or to identify sources or information related to a particular story.

In general, machine learning methods and tools can improve the data-driven journalism process in newsrooms. More specifically, machine learning algorithms have the ability to classify data into meaningful categories, helping data journalists to spot patterns and discover stories that otherwise would have been overlooked, for example, to detect fraud and other illegal activities by analyzing financial transactions and other data (Flew et al. 2012; S. Cohen et al. 2011). Additionally, machine learning can be used to help journalists identify and classify topics, images, people, and events, to predict the likely impact of a news story on public opinion or engagement on social media. Also, these methods are currently used in newsrooms to detect breaking news and newsworthy events, leverage users' preferences to produce personalized content and detect fake or suspicious content and online bots.

The story “We Trained A Computer To Search For Hidden Spy Planes. This Is What It Found¹⁷” by BuzzFeed News shows the great potential of machine learning algorithms and how they can be used to identify patterns in public databases. The journalists wanted to find out if secret spy planes fly over the U.S. and for what reason, so they trained a random forest classifier on data from the open flight-tracking service “Flightradar24”. By giving to the model information on how spy planes fly from FBI and the Department of Homeland Security, such as moving in circles, and skydiving, they could easily divide all the planes into two categories the surveillance aircrafts and the normal ones. The journalists were able to detect surveillance aircrafts flying over U.S. cities and uncover their origins, publishing two stories about secret operations like spying on drug cartels and terrorism.

Additionally, “Machine Bias¹⁸” by ProPublica investigated the effects of employing an algorithm to forecast future criminal recidivism rates for white and black criminals. The journalists were able to discover the variables that created the racial gap by reverse engineering the rating system findings and comparing the results to the actual facts. The investigation's conclusions showed that, despite the fact that this was not true in practice, the algorithm

¹⁷<https://www.buzzfeednews.com/article/peteraldhous/hidden-spy-planes>

¹⁸<https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>

gave black criminals a lower likelihood than white offenders of committing a crime again in the next two years. This mismatch raises serious questions about the accuracy of algorithms in foretelling criminal behavior, particularly when race is taken into account. Furthermore, it suggests that the algorithm was biased toward black offenders because it ignored other significant characteristics such as age, gender, and the type of crime committed.

The conclusions of the article are particularly alarming since they imply that while algorithms are widely utilized in the United States, they may not be sufficiently trained to successfully predict criminal behavior, particularly where race is an input variable. This research highlights the potential for algorithms to further marginalize minority populations in a criminal justice system that is already significantly prejudiced in favor of white perpetrators. Therefore, it is crucial that the algorithms used for predicting criminal recidivism rates receive thorough testing to guarantee they are unbiased and accurately represent the true likelihood of criminal behavior. As a result, data journalism now requires a new skill called algorithmic accountability: the ability to evaluate and explain such algorithms. A skill that some computational journalists have already started to develop, but one that will probably require many more years of practice and trial-and-error before becoming a common expertise amongst journalists and technologists in general.

Text Analysis is used in the context of computational journalism, as it provides journalists with powerful methods for extracting and analyzing the textual data they need to effectively report the news, such as political speeches, or tweets about a particular issue. Similarly, it is able to map the use of specific words or phrases, name entities, hate speech, and sentiment, and uncover trends or insights that might otherwise be difficult to find (Stray 2019). For the journalists to use these techniques themselves, programming knowledge is required to organize and structure text-based data and assign tags or labels to specific words or phrases, which can then be used to group and compare texts. A more advanced way of analyzing text is with Natural language processing (NLP), a subdomain of artificial intelligence that enables machines to understand, interpret and manipulate human language, allowing them to process large amounts of unstructured data quickly and accurately, automatically generate content such as captions and summaries and more (Liddy 2001).

Furthermore, text analysis is a great way for newsrooms to identify the sentiment of text, such as determining if an article is positive or negative in tone for gauging public opinion on a particular issue. The article “Trump Sounds a Different Tone in First Address to Congress” by the New York Times is one of the first examples of how sentiment analysis can be used to gain insight into the emotional content of presidential speeches. The data journalists used the NRC EMOLEX emotional lexicons (Saif M Mohammad 2017) to examine former President Trump’s speeches and create a visualization of the results. The visualization revealed that Trump appeared to be less angry when he addressed the State of the Union compared to other speeches. This news story highlights the potential of sentiment analysis as a tool for understanding the sentiment of political discourse (Stray 2019). By applying sentiment analysis to political texts, journalists and researchers can gain a better understanding of the emotions expressed in political rhetoric and track changes in sentiment over time, which can provide insights into political trends and developments.

Robot journalism is about the incorporation of these technologies that have enabled journalists to automate tasks such as textual data collection, and machine-aided writing (Carlson and Lewis 2015). NLP models are used by news organizations for summarizing large documents, while automated writing assistants generate stories from raw data sources, such as weather or earthquake APIs, helping journalists to quickly produce relevant stories, save time, and produce more content at a lower cost (Marconi 2020). Other NLP tools such as Newswhip¹⁹ and Awario²⁰ monitor stories as they evolve over time, allowing news desks to track changes in public opinion on a particular issue, receive customized alerts when a story breaks online, allowing journalists to stay on top of emerging trends and to get ahead of the competition. Additionally, NLP tools can classify news articles by theme or geographic location, quickly comb through large data sets to uncover correlations between events or identify potential leads that can be used to further research a story increasing the speed and efficiency with which journalists can work (Hansen et al. 2017; Flew et al. 2012).

Similar to robot journalism, augmented journalism, as proposed by Marconi (2020), is the

¹⁹<https://www.newswhip.com/spike-real-time-media-monitoring/>

²⁰<https://awario.com/news-monitoring/>

integration of AI-assisted tools into reporting, research, writing, and editing. This concept of “augmentation” implies that complex and time-consuming tasks for humans can be facilitated and improved through machines. In an optimal augmented newsroom, the computational journalists, utilize data analysis methods for their research, while automation editors combine computer and journalism skills to guarantee the editorial reliability of automated processes. AI ethics editors are responsible for examining training data, managing biased algorithms, and interpreting the results. Newsroom tool managers assess the effectiveness and efficiency of AI technologies and conduct training for reporters. In this regard, news providers acquire skills of acting as amateur data scientists (Hansen et al. 2017). Additionally, the concept of the “exo journalist” proposed by Tejedor and Vila (2021) suggests that automation assists journalists in achieving their goals and improving their work quality, thus forming an alliance between AI and journalism.

The use of Natural Language Processing is common in many news organizations. In an article for the Associated Press, the computational journalist Jonathan Stray applied text analysis techniques, TF-IDF method, a bag-of-words model, and cosine similarity metric, to thousands of private security documents in Iraq (Stray 2010a). Similar techniques with a focus on investigating the misconduct of police officers were used by the newspaper Newsday in their Pulitzer Prize winner data story “Police misconduct on Long Island²¹”. Moreover, the news website Vox used NLP to analyze and identify the main issues addressed in Obama’s political speeches, in the story “Barack Obama’s crucial insight into the crisis of American democracy²²”. The Digital news outlet, Quartz used language analyses in their story, “Analysis of 141 hours of cable news reveals how mass killers are really portrayed²³” to examine cases of media bias during the coverage of the 2017 mass shooting in Las Vegas. The findings showed that the perpetrators were given specific terms depending on their nationality or race (Y. Zhou 2017).

²¹<https://www.pulitzer.org/finalists/newsday-0>

²²<https://www.vox.com/policy-and-politics/2018/9/7/17832266/obama-speech-illinois-transcript-trump-republicans>

²³<https://qz.com/1099083/analysis-of-141-hours-of-cable-news-reveals-how-mass-killers-are-really-portrayed>

The Atlanta Journal-Constitution's "Doctors and Sex Abuse²⁴" project conducted a survey to examine how doctors maintain their licenses after being disciplined. Danny Robbins, the data journalist, scraped website content to gather pertinent complaints and applied machine learning algorithms to analyze a large number of disciplinary documents. In contrast, the Wall Street Journal has taken a more business-oriented approach by using an algorithm to predict which readers will subscribe to the newspaper, thus increasing its revenue. In addition, the Wall Street Journal utilizes Narrativa's Natural Language Generation (NLG) services to produce informative material related to the stock market and producer and consumer price indices (Sánchez 2021).

The use of automated methods has been extended to the field of disinformation, as fake news is becoming increasingly widespread on news agencies and social media platforms. The term "fake news" was popularized during the 2016 US Presidential elections, and refers to the deliberate spreading of false information for political or financial gain, or in an attempt to cause public harm (Shu et al. 2017). This has encouraged computer scientists to develop and use artificial intelligence techniques to combat the issue, such as Neural Networks and other deep learning models to verify sources and detect deepfakes.

Although applications of artificial intelligence in journalism raise some questions about journalists and the potential for their replacement by automation, it is clear that raw data cannot become news content without human interpretation. Thus, AI methods are being used to expand the range of work, creating tasks such as data production and classification, knowledge management, parametrization, and configuration. Moreover, news verification is one of the primary responsibilities of reporters, so checking the validity of the output of AI tools is necessary (Hansen et al. 2017). To ensure accuracy, automated tools and systems always require human oversight. The Quakebot²⁵ of the newspaper Los Angeles Times is an example of this; it erroneously published a piece of earthquake news off the coast of California.

Large news organizations, such as the Washington Post, BBC News Labs, the Wall Street

²⁴<https://doctors.ajc.com/>

²⁵<https://laist.com/news/quakebot-error>

Journal, and Quartz AI Studio, are leading the way in technological modernization. The Washington Post's Research and Development lab is teaching reporters computational techniques to use in election coverage. BBC News Labs is promoting semi-automatic journalism and providing resources to convert audio files into video. The Wall Street Journal's Research and Development team is focusing on natural language processing (NLP) techniques and other methods such as verifying deep-fakes and algorithmic transparency reporting. Additionally, Quartz AI Studio is helping journalists apply machine learning (ML) algorithms in order to develop new methods of storytelling (Marconi 2020).

AI-powered tools have changed the media landscape by adding new opportunities to the ecosystem. Examples of these tools are systems that identify breaking news, convert written text to audio files, automatically produce articles from data, reveal hidden information in social media posts, predict user subscriptions, and detect untrustworthy and fabricated content. As helpful as these tools may be for journalism, they can still pose a threat to the profession itself. In the future, news organizations may hire fewer professionals or rely so heavily on Artificial Intelligence that they make serious mistakes in reporting the news. Therefore, publishers should bare in mind that these systems need constant human oversight.

2.5.1 Existing Initiatives

With the introduction of ChatGPT²⁶, a software based on OpenAI's new language model that responds to queries by mimicking human conversation, the discussion about AI's use is becoming more prominent in newsrooms. While the free availability of this model may have democratized the use of AI, with everyone having full access to a state-of-the-art chatbot, it still raises concerns about what is to come in the journalism industry.

LSE Professor, Beckett (2022), analyzed in an article the ways in which smaller news outlets can benefit by implementing AI practices and automation, highlighting the importance of digital technologies in the modern era, and emphasizing the competitiveness prevailing in the media landscape. The article suggests that AI can enhance the work of human reporters

²⁶<https://openai.com/blog/chatgpt>

by producing a more complete journalistic result and also ensure the survival of the organization. Furthermore, a report on “AI in local news” by the Associated Press in the US demonstrated potential AI applications that can be incorporated into newswork (Aimee Rinehart 2022).

Exploring innovative ways to diversify revenue streams, building a loyal audience, and organizing a digital transformation in a strategic and sustainable manner are some of the current challenges that publishers try to overcome with AI. At the same time, ethical considerations like journalism independence and trust in media remain prevalent. AI methods can be utilized for a variety of tasks, from newsgathering and production to distribution, and they are particularly useful for finding significant features in complex data, and for identifying trends by grouping common data (Hansen et al. 2017). Some examples of how AI can be used in newsrooms, include the automatic production of content such as the automated publication of sports results, the personalization of user experience, the automated news gathering, such as the transcription of interviews and data journalism.

Incorporating AI is a trend that has become increasingly prevalent across all sectors of society (Tejedor and Vila 2021), with robot journalism being introduced the last decade (Carlson 2015). This trend is driven by the need to maintain a competitive advantage in an ever-evolving news landscape, which is being shaped and reshaped by the introduction of new digital technologies (Broussard et al. 2019). As a result, AI methods are being applied to newsroom processes in order to optimize them, saving human effort and time, while simultaneously improving performance. Specifically, the definition of the “robotic reporter” (Carlson 2015), is used to describe an algorithm that transforms data into natural language and publishes new stories without the assistance of journalists. Systems that enable the rapid generation of automated messages or alerts are often employed for texts that are directly derived from data, such as athletic results and stock market news, or for catastrophic phenomena like earthquakes or wildfires.

Early adopters of AI practices in every day news production were Bloomberg, BBC, Washington Post and Forbes (Thompson 2016). The financial news organization Bloomberg, used

such algorithms to turn financial reports into news and ushering in a new era for the sector. Likewise, programmers at the Washington Post developed a reporting bot, Heliograph²⁷, to cover major events. The program was initially used during the 2016 Rio Olympics to help journalists broadcast the medal ceremonies. Subsequently, the system's capabilities were expanded to cover the 2016 elections with over 500 articles being published on the election day. Interestingly, the texts of the election results were automatically converted into podcasts, which were spoken by a voice assistant and were delivered to listeners according to their location. Furthermore, the every day use of this automated storytelling tool is to use natural language modeling to write news stories for local and international sports, such as high school football²⁸. The bot produces continuous updates from numerical data, which are then posted on social media platforms through spoken commands from Amazon Alexa. Furthermore, in 2017, the Washington Post implemented ModBot²⁹ to reduce the time-consuming process of reading comments. ModBot reads, checks, and filters comments for offensive language.

Le Monde used algorithms to gather data of the municipalities during the 2015 election, enabling bots to produce articles on the results of each municipality in the newspaper's editorial style. The Financial Times' AI tool "He Said She Said" is based on text analysis, and is designed to detect gender bias in articles and promote more feminine views (Waterson 2018). Similarly, Forbes' system "Bertie"³⁰ uses AI to select the headlines and images for its stories, in order to facilitate the storytelling process (Willens 2019). Also, Bloomberg's innovation lab called BHIVE has launched "The Bulletin"³¹ a function to its mobile application, which is apt to generate the summary of the most important news articles leveraging NLP and semantic analysis. Finally, in the United Kingdom the Press Association's news service, RADAR³², which inside the newsroom is also called: "reporters and data and robots", uses

²⁷<https://www.washingtonpost.com/pr/2020/10/13/washington-post-debut-ai-powered-audio-updates-2020-election-results/>

²⁸<https://www.washingtonpost.com/pr/wp/2017/09/01/the-washington-post-leverages-heliograf-to-cover-high-school-football/>

²⁹<https://www.washingtonpost.com/pr/wp/2017/06/22/the-washington-post-leverages-artificial-intelligence-in-comment-moderation/>

³⁰<https://bertie.forbes.com/>

³¹<https://www.bloombergmedia.com/press/bloomberg-medias-innovation-lab-launches-bulletin/>

³²<https://pa.media/radar/>

open data related to crime, health, education, and transport, map them to the local authority or postcode level, and produces a news story. The PA Media Group owns several popular news outlets, with the editor-in-chief of the company, confirming that AI writes 50K individual local news articles every three months³³.

By utilizing AI-assisted tools, journalism can move away from the traditional linear model and move towards a more personalized experience for each individual reader. The use of AI-assisted tools such as text analysis, sentiment analysis, document classification, and breaking news detection can help shape content to fit the interests and individual needs of the audience (Hansen et al. 2017). Journalists can produce content that is more pertinent, interesting, and valuable to their audience by utilizing this data. For instance, a news website can tailor the content to a reader's interests using data on their browsing history, geography, and demographics to improve reader engagement and their reading experience. Therefore, news organizations are collaborating more and more with tech firms to create their own AI-assisted tools. The first to use such technologies in the media industry was BBC. The tool "MyBBC"³⁴, collects data on users' behaviour and interests to then suggest relevant news.

Social media platforms, such as Facebook and Twitter, serve as sources of data, with news organizations exploiting the information contained in the posts for detecting breaking news and newsworthy stories (Stray 2019). Journalists at Reuters automatically find interesting trends on social media platforms, through the AI-powered tool News Tracer³⁵. News Tracer uses trained models for tweet classification and clustering (Stray 2019). This system is especially useful in the case of breaking news, such as natural phenomena or murderous acts in public, as internet users participate by posting relevant content. A similar tool is used by Radio France giving the first update on the Brussels bombing in 2016 and the terrorist attack in Nice (Marconi 2020).

Graphext³⁶, an AI-based tool, has been used by El País newspaper to identify the presence of

³³<https://www.newsrewired.com/2018/11/07/press-associations-news-service-radar-has-written-50000-individual-local-news-stories-in-three-months-with-ai-technology/>

³⁴<https://www.bbc.co.uk/blogs/aboutthebbc/entries/46a896ea-e587-4c63-ae7e-9781bca58dd3>

³⁵<https://www.reutersagency.com/en/reuters-community/reuters-news-tracer-filtering-through-the-noise-of-social-media/>

³⁶<https://www.graphext.com/>

politicians on social media by mapping their posts. Stray (2019) note that sentiment analysis of social media content is often used by news providers to gauge public opinion. The Laboratory for Social Machines at the MIT Media Lab analyzed the Twitter followers of Donald Trump and Hillary Clinton in the 2016 U.S. presidential election (Maffei 2016). Supervised learning classifiers were employed to categorize the Twitter users according to their political view, thus demonstrating the divide among supporters.

European projects have thus far aimed to use data research, visualisation, clear reporting, contextualised story timelines, audience engagement, and interaction to understand complex topics. The Digicom project (2016-2019), an initiative of Eurostat's CROS portal³⁷, generated results into creating "new, innovative dissemination products, tools and services for European statistics", as seen in the "Eurostat: Statistics in Support of Data Journalism" workshop in 2019. Two consortia, the European Data News Hub and the European Data Journalism Network³⁸ (EDJNet), also focus on the use of data for covering European affairs. Journalism++³⁹, a member of EDJNet, is one of the few international teams of data journalists specialising in data analysis, data-driven storytelling, and newsroom programming. The European Journalism Center's DataJournalism.com, supported by Google News Initiative, provides journalists with free resources, materials, online courses, and forums. Furthermore, the London School of Economics and Google News Initiative launched the project JournalismAI⁴⁰ in 2019 to help news organisations use Artificial Intelligence responsibly. Those involved map the problems they would like to solve and work together to create prototypes and reports.

AI and data analytics can be used in innovative ways to build an audience. For example, the Globe and Mail, one of Canada's most acclaimed national newspapers, won the "2020 Award for Technical Innovation" in the Service of Digital Journalism for developing a machine learning product called Sophi.io⁴¹. This product uses machine learning to determine when to prompt readers to pay for a subscription to the website, increasing subscriptions

³⁷<https://ec.europa.eu/eurostat/cros/system/files/datajournalistoutreachstrategy.pdf>

³⁸<https://www.europeandatajournalism.eu/>

³⁹<https://jplusplus.org/en/>

⁴⁰<https://www.lse.ac.uk/media-and-communications/polis/JournalismAI>

⁴¹<https://www.sophi.io/>

and providing a better user experience. Additionally, when data analytics tools are combined with qualitative analysis and editorial expertise, they can offer an insightful view of audience engagement and lead to change (Kalsnes and Krumsvik 2019).

Tools such as Ophan⁴² and Lantern⁴³ provide journalists from the Guardian and Financial Times with user-friendly metrics and graphics that can help them to develop an audience engagement strategy. Furthermore, AI technology can be employed to create and disseminate personalized content tailored towards each reader's interests and concerns, such as recommending articles (Hansen et al. 2017). BBC News and Reuters' news apps offer examples of this application, where users can interact with a personalized interface.

The field of artificial intelligence is fast developing and brings with it broad consequences with the need for a common regulatory framework becoming imperative. As newsrooms are starting to understand its potential, it is critical for journalists and publishers to think about the moral aspects of using it while also ensuring the security and precision of AI systems.

2.5.2 Algorithmic Accountability

Although the use of artificial intelligence in the media industry has the potential to simplify and automate news production, it also raises concerns about the output's reliability and accuracy, as well as possible implications for the roles and responsibilities of journalists (Diakopoulos 2015). Since these technologies are still rather new, implementing them could result in ethical ramifications and data privacy issues. In order for these technologies to be utilized for computational journalism, a number of additional technical and practical issues must be resolved. For instance, the outcomes of a machine learning model could yield bias that misrepresents reality, which is in contrast to a news organization's primary purpose of representing an accurate and honest report on the facts. In order to use AI in their reporting ethically and responsibly, journalists must always be aware of the danger of bias in artificial intelligence systems.

⁴²<https://www.theguardian.com/info/2021/jul/12/how-we-backfilled-the-guardians-in-house-analytics-tool-to-provide-greater-journalistic-insight>

⁴³<https://www.niemanlab.org/2016/03/the-ft-is-launching-a-new-analytics-tool-to-make-metrics-more-understandable-for-its-newsroom/>

The potential for AI systems to produce fake or misleading news, the danger of biased or unfair reporting, and the loss of human expertise and judgment in the news production process are just a few of the dangers that AI poses to the media industry. Strong ethical norms and constant human oversight are required to guarantee that AI is being used in a responsible and accountable manner. This may entail taking steps to make sure AI systems are transparent and comprehensible, subject to human review and evaluation, and trained on a variety of representative data sets (Doran et al. 2017; Shah et al. 2015; Cranor 2008). The ideal of algorithmic accountability refers to the process of monitoring, evaluating, and regulating algorithms to ensure they are ethical and do not discriminate or otherwise harm certain people or groups (Shah et al. 2015). It involves designing systems that are transparent and accountable for the decisions made by algorithms and the impacts those decisions have on people. This includes creating mechanisms for people to understand how algorithms arrive at the results, providing people with recourse when algorithms make mistakes, and establishing processes for holding algorithm creators and users accountable for any harm caused by those algorithms (Ananny and Crawford 2018).

Algorithms are becoming increasingly influential in business and government decisions, even though their power is often hard to discern. Algorithmic accountability is important because algorithms are increasingly being used to make important decisions that affect individuals and society, such as in the criminal justice system, hiring and admissions processes, and content moderation on social media platforms. Ensuring accountability in the use of algorithms can help to prevent discrimination and other negative effects on individuals and society. Ananny and Crawford (2018) highlight the fact that there are numerous limitations to the application of transparency for understanding and governing algorithmic systems. Therefore, journalists are responsible of holding those in power accountable in cases of misuse, such as the Propublica's story "Machine Bias" on page 30. Diakopoulos (2015) introduced the term algorithmic accountability reporting in journalism referring to the process of understanding the power dynamics, biases, and effects of computational artifacts in society. He proposed a framework for analyzing algorithmic power and urges journalists to learn AI methods to be able to reverse engineering government algorithms.

Toward the idea that AI systems should be designed and used in a way that involves human input and decision-making at key points in the process, is the “Human-in-the-loop Artificial Intelligence” (HitAI) model (Zanzotto 2019). The human in the loop approach is a methodology for incorporating human input and oversight into artificial intelligence systems rather than relying solely on automated decision-making by the AI models (Berendt and Preibusch 2017). By bringing human knowledge and judgment into the process, this method is sometimes considered as a solution to solve some of the difficulties and issues related to the usage of AI, such as the possibility of prejudice or inaccuracies. This strategy can also assist in making sure that AI systems are utilized ethically and legally and that their use is transparent and accountable. A person must examine and approve the AI system’s final decisions or output in some situations such as data analysis or pattern identification. In other cases, the human in the loop approach may involve using AI to make initial decisions, but allowing for human intervention in cases where the AI system is uncertain or the consequences of a decision are particularly significant (Cranor 2008).

Another approach for addressing these challenges and ensuring the ethical use of AI in the media while minimizing the potential risks and negative impacts, is the Explainable AI, also known as XAI. This is a new branch of AI which tries to provide explanations for algorithmic decisions and actions in a way that is understandable to humans (Hoffman et al. 2018). In other words, all AI systems should be transparent and understandable, and enough information should be known about how they make decisions and what factors they take into account. Even though the aim of XAI is to increase accountability and trustworthiness, this is not attainable for all models; neural networks are the hardest to interpret. This entails ensuring that algorithms are open, impartial, and fair as well as that their application adheres to moral and legal principles.

Doran et al. (2017) examined various approaches to explainable AI across research fields, and considered three types of Explainable AI: total black box systems that have no insight into their algorithms; interpretable systems that can be explained using mathematics; and comprehensible systems that produce symbols that allow for further interpretation. Finally, they introduced a fourth concept: truly explainable systems, which are able to generate au-

tomated explanations without requiring a human to post-process them. Even though this field is currently developing, some examples of explainable AI techniques include using natural language explanations, and visualizations of the decision-making process. Both techniques are used in this dissertation. By using the aforementioned practices, the AI systems could be used in a responsible and ethical manner, and can help to avoid bias or unfairness in decision-making by AI systems.

Some experts have cautioned that with the extended use of AI in many fields policymakers should start weighing in (H. Shah 2018), with the US National Institute of Standards and Technology (NIST) initiative to examine the potential misuse and suggest standards for creating trustworthy AI. At least 17 U.S. states introduced general laws on AI in 2022 (*Legislation Related to Artificial Intelligence 2022*), and a number of other laws that forbid discrimination based on protected characteristics, like age, gender, and race, may also apply to AI systems that make decisions that have an impact on people. Additionally, regulations that safeguard individuals' privacy and personal data, such as the California Consumer Privacy Act (CCPA) in the United States and the General Data Protection Regulation (GDPR) in the European Union, might be applicable to AI systems that gather or handle personal data.

A collection of recommendations for the ethical application of artificial intelligence in the European Union has also been produced by the European Commission. The moral principles that ought to direct the development and use of AI are generally referred to as the moral standards of employing AI. Fairness, transparency, accountability, privacy, and non-discrimination are some important ethical factors to take into account when using AI.

Following those considerations, the Commission's guidelines⁴⁴ are meant to serve as a foundation for ethical development and deployment of AI in the EU. The "Association for Computing Machinery" (ACM)⁴⁵ a professional association for computer scientists, has also offered advice on the moral use of AI. In general, the code of ethics revolves around the following principles:

Respect for human autonomy: AI systems should respect individuals' autonomy and

⁴⁴<https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>

⁴⁵<https://www.acm.org/code-of-ethics>

should not be used to undermine or subvert their rights or interests.

Respect the privacy of others: Computing professionals should respect the privacy of individuals and should protect personal information from unauthorized access or disclosure.

Prevention of harm: AI algorithms should be designed and used in a way that minimizes the potential for harm to individuals or society.

Fairness and non-discrimination: AI models should not make decisions that are based on discriminatory criteria, and should be designed and evaluated to minimize the potential for bias.

Transparency and explainability: AI systems should be transparent and explainable , and should provide individuals with access to information about how decisions are made.

These guidelines offer general direction for the moral use of AI and can be applied as a framework to evaluate the ethical implications of various AI applications.

Chapter 3

Exploring the Impact of Social Media

3.1 Adapting to the Digital Age

Having explained the foundations of journalism and its prevalence as the Fourth Estate, it is expected to continue with what is considered from many to be the “Fifth Estate”, the active networked users on social media (Dutton 2009). The very architecture of those platforms promotes participation, discussion, content generation and dissemination by offering new opportunities to citizens to express themselves, share ideas, report wrong doings and act as watchdogs, like journalists do. Moreover, the free use of online tools and mobile apps for almost professional video production, paired with the vast information flow and the acceleration of good content thanks to their recommendation algorithms, provides citizens with a real chance to reshape journalism. The media ecosystem has been converted from one-way communication into a participatory process between news producers and the public (Hermida 2012), changing the concept of a centralized news market to a conversation (Alejandro 2010, p. 5).

One of the many definitions of social media, presents them as digital applications that were the foundations of Web 2.0 and provided internet users with the means to generate and exchange content A. M. Kaplan and Haenlein (2010, p. 61), therefore decentralizing journalistic authority and changing the rules of the media game. Akhgar et al. (2017) characterizes

social media as a collection of platforms with responsive functions used for open dialogue, exchange of viewpoints, and encouraging users' to form relationships with each other. All these affordances along with the produced rivers of user-generated-content inevitably invited users to join forces with professional reporters in the so-called "Crowdsourcing in Journalism" (Lamprou et al. 2021). As the result of this disruption, journalism took a turn to openness and became more participatory and dynamic by giving in to the two-way conversation resulting to what is called social media journalism (Kuyucu 2020). As usually happens with everything new, journalists and publishers used social media in a way that reinforces the traditional techniques of their craft and upholds the values of the profession thus normalizing the transition (Djerf-Pierre et al. 2016; Hermida 2012).

Not able to change the sequence of events regarding the digital transformation of the media industry and in fear of succumbing to "digital darwinism" (Schwartz 2002), publishers implemented social networking technologies as effective means for daily journalistic practice, marketing reasons and profit-generation. Undoubtedly, the media sector is undergoing a mediatization process (Papathanassopoulos et al. 2021; Papathanassopoulos 2002), incorporating the social media logic (Hedman 2020) in practices such as gathering, reporting and disseminating news. Except for a vast pool of available, and costless textual and visual information to select from, journalists use social media to discover breaking news such as eyewitness footage which have proven to be particularly effective in crisis situations. Using user-generated content, tweets and other information directly from social media adds an extra level of information to the news stories (Kuyucu 2020). Another important aspect is the "digital first" strategy followed by newsrooms that creates an opportunity for the online reader to be informed faster, since news stories are breaking online at least 20 minutes ahead of traditional media (Akhgar et al. 2017).

Furthermore, newsrooms publish content relevant to the online trending topics and customize their coverage to match the given standards of each medium. For instance pictures with 1:1 aspect ratio, and short explainable videos are produced to be used as stories or reels on Instagram, but short texts on Twitter often promoting interaction with the audience. Also, the hashtag feature is very important for journalists as they can easily identify major topics

of discussion, and prevailing trends on a daily basis. Many newsrooms have automated the process of publishing on social media (Marconi 2020) in which they both inform and promote content. A usual tactic for media organisations is to use URLs in their posts, and profile pages that lead to the actual news piece to increase the traffic and make money (Hermida 2012).

A Swedish study of newsrooms that are very active in social media showed that the primary purpose was to gain information, get ideas for stories, observe interesting conversations, and keep track of trends (Hedman and Djerf-Pierre 2013). They found that the views of the readers and listeners were seen as important. A 2012 study of Greek media (Serafeim 2012) examined freedom of expression and found that social media had a significant effect on independent and democratic production and consumption of information, leading to more freedom of the press and expression. Furthermore, the research showed that the content on social media was different from what appeared in print, as journalists shared stories without being driven by political, commercial or other factors.

3.1.1 From Imagined Audience to Micro-targeting

The digitalization and datafication of aspects of everyday life and forms of organisation, have advanced new ways for the media industry to relate to their audiences by leveraging social media and their “like economies” (Gerlitz and Helmond 2013). It used to be the circulation of a newspaper or the TV ratings that determined the size of its audience and the success of a media company, although nowadays the intangible digital essence of online journalism demands different measurements such as the audience engagement with the news product itself. That drove journalists to worry about how well will their reporting perform online and alter their news pieces to help social media algorithms to give prominence to their content amidst an overloaded information environment. These conscious efforts to capture the attention of the public (Tandoc and Vos 2016) made journalists concentrate on who they ultimately want to approach with their writing (Nelson 2020). Several scholars have argued for a new media regime that is more inclusive, diverse, and open to collaboration and interaction with its audience (McCollough et al. 2017; Kiesow 2015; Lewis and

Usher 2014; Serafeim 2012; Davis Mersey et al. 2010).

Accordingly, reporters' gatekeeping function, in the sense that they select the news pieces to be published thus setting the agenda of the public debate, has been modified (Munger 2020; Molyneux 2015). With the intervention of the public in the news making process, the gates are no longer managed by journalists. The latter acquire the role of content aggregators by attaching links and making reposts.

The social media users play a key role in news agencies' routines as their feedback is omnipresent (Kramp and Loosen 2018). In this way, news producers in order to be informed of current developments, monitor and scan the social media content which is constantly updated, thus forming a new genre of reporting (Kuyucu 2020). In keeping with this sense, some journalists mainly publish user generated content, and are responsible only for editing the news as they navigate the internet and distinguish the newsworthy events from an extensive volume of data (Hermida 2012). Bruns (2018) describes this process for news coverage as "gatewatching".

What contributed the most to the emergence of engagement as a structural element of the journalistic experience is the disruption of the traditional journalistic model and news consumption patterns. Before the advent of social media and its subsequent platform logic, a traditional news consumer would have a regular relationship with the news outlet in order to read/listen/watch the news. The classical model proposed by Janowitz (1975), known as "gatekeeping", which would define the very procedure of selection and publication of news, predominantly assumes that "he who controls the flow of information can influence social reality" (Shoemaker and T. Vos 2009; Janowitz 1975). So far its regulatory role establishes journalism as the filter between vast informational flows and the public (Benson 2008), adopting an asymmetric relationship between journalists (Deuze 2003) as the authorized "gatekeepers" who decide which events will pass through the information gates and the public. In short, the newsrooms used to shape the news content based on the assumption that what they find interesting is what the public would -and should- find interesting (Harcup and O'neill 2017; Gans 2004; Harcup and O'neill 2001).

Social media has largely displaced this one-way perception in news circulation by replacing it with the concept of sharing (Noguera-Vivo 2018). In this environment the news consumer is no longer the final recipient of the message but a node in a wider network (Carlson 2016). A change that suggests that users play a very important role in the dissemination of news and therefore in the final configuration of the audience of each medium, deciding for themselves which content to share with their peers on their various platforms. Bruns (2018, p. 17) coined for the readers the term “gatewatchers”, engaging in identifying the interesting material that comes out of the gates and in sharing particular bits of available information, thus positioning users as a secondary level of observation, evaluation and news distribution (J. Singer 2014) in a network pipeline that eventually augments and expands the original publication. Therefore, engagement, as the new concept that enters the structure of management and distribution of news is placed somewhere between the creation and reception of news.

Publishers have always wanted to know more about what type of content people consume. For years they have been relying on focus groups, ratings, paid surveys, and circulation numbers to arrive at conclusions about audience attitudes and behavior related to their products. Decisions for future articles were made by editors based on vague assumptions of what the imaginative audience would have wanted to read and journalists used to write their stories having a certain persona in their minds usually created by the managers (Darnton 1975). More specifically, the concept of the public has taken on various dimensions and several roles in media history but also in communication research and journalism. The mass audience, news consumers, citizens with the right to information, and prosumers are some of the conceptualizations of the public at every stage of the evolution of information technologies and media culture. Before the web, estimating how many people read a given article out of a newspaper or a magazine or what the readers thought of it was hard to assess and only possible on a small scale. Technological developments and the plethora of available data have revolutionized the marketing and advertising industry, providing publishers with monitoring tools to track key performance indicators (KPIs) and sometimes even give access to micro-level data about individual user’s behavior. These kind of metrics which directly target the identity and preferences of each user are used in a range of marketing strategies

from recommender systems optimisation to microtargeting and according to (Manovich 2018) have dramatically exceeded simple quantification. Moreover, the debate around data processing and user profiling is increasingly being dominated by the obscured nature of algorithmically-based behavioral targeting, which is being presented as a continuation of traditional audience targeting techniques (Bolin and Andersson Schwarz 2015). Such tools have become a crucial element of modern newsrooms, which, for instance, display live ratings of the performance of individual stories on screens in their offices (Diakopoulos 2019).

In particular, social media play an increasing role in news dissemination, as do search engines and news aggregators (Eliza Shearer 2019; Newman et al. 2019). This means that more and more people rely on secret news feed algorithms and various other ranking programs for their daily news update. Therefore, positive feedback from online readers is coveted more than ever before, since social interactions such as likes, shares, and comments can determine how many people have a chance to come across the news content on social media (Nechushtai and Lewis 2019; Thorson and Wells 2016a). Thus, for many digital organizations, it is an important goal to gain meaningful interactions from the users, so the filtering algorithms that decide the ordering of the posts will have their content appear higher in the news feeds of potential readers.

Although, this shift in gatekeeping created a different kind of problem for news organisations related this time to their audience loyalty. The use of social media for finding news, the so-called “News-Finds-Me” (NFM) attitude, has led some users to believe that there is no reason to actively follow a news outlet to stay informed about the recent developments since social media will ultimately bring the latest news to their own digital “doorstep” (Gil de Zúñiga and Diehl 2019; Weeks and Holbert 2013). This tension of a user relying on social media to provide them with the daily news updates, is associated with lower levels of political knowledge and participation and has a negative impact in citizenship (Gil de Zúñiga and Diehl 2019; Shehata and Strömbäck 2021).

3.1.2 Definition of Engagement

News organizations are increasingly dependent on engagement metrics (Ferrer-Conill and Tandoc Jr 2018; Coddington 2015) as the system to measure their audience and evaluate their approach to an increasingly complex and competitive marketplace. For some, engagement means more online interaction with the public, while for others it promises a more interactive and perhaps rewarding approach to reporting and journalism itself (Lewis et al. 2014). Even though the theoretical starting point of engagement could be traced to the concept of interaction, several scholars have raised the issue of how vague the very definition of the term is (Ferrer-Conill and Tandoc Jr 2018). For example, McMillan (2005) distinguished interaction into three broad categories, namely between humans, humans to machines, and humans to content. Moreover, he argued that interaction consists of features which are the very elements that make the communication environment interactive (technologies, platforms, etc.), the processes which are the ways of performing the act of interaction itself, and the perceptions that concern the public's belief in the interactive value of the provided technical affordances and the way they perceive and evaluate their degree of interactivity.

From a more broad perspective, scholars agree that with audience engagement we mean the consumer's experiences on a cognitive and emotional level derived from their interaction with the media elements (Broersma 2019), emphasizing a subjective relationship with the media. However, (Hill 2018) describes engagement as an inclusive term that represents how the audience experiences news content and engage or actively participates on social media. Although this holistic and more audience-centric approach covers a wide range of engagement experiences, it is somewhat limited concerning the journalistic product. Ksiazek et al. (2016) places engagement on a continuum, from mere exposure to a news story to interaction, while McMillan (2005) stressed the importance to consider the technical and behavioral dimension of audience engagement. Finally, other scholars promote the idea of engagement as a radical way of reconnecting journalists with the community, decentralizing the role of the journalist in gatekeeping the news and reorienting it to collaborate with the public (Lawrence et al. 2018; Guzman 2016).

Ørmen (2015, p. 25) argues that the audience engagement with the news includes the attention users pay to a news story and how they interact with it, therefore in this context, engagement describes audience's consumption and participation. Nelson (2021) defined engagement as the very means to urge the public into consuming and participating in the news, while Guzman (2016) notes that from a journalistic point, engagement is the confirmation to journalists that their work is relevant and important to their audience. Ha et al. (2018) proposes various levels of engagement based on the user's effort to acquire and utilize news content either for personal or social purposes, while Sang et al. (2020, p. 468), talks for interaction instead of engagement and defines it as "different ways in which news consumers interact with the news content and/or each other".

Nelson (2021) proposes a broad categorization of reception-centric and production-centric engagement. The "reception-centric" orientation focuses mainly on the way the public receives the news and considers their involvement with it in terms of minutes spent, shares, likes and so on. This is more from a marketing perspective that treats news as a commodity and readers as customers. The "production-centric" involvement concerns the characteristics of the news item (topic, news values, text features). In short, the "production-centric" approach refers to the ways in which journalists engage with the public, while the "reception-centric" approach refers to the ways in which the public responds to the news.

Engaging with news on social platforms takes various forms, such as liking, commenting and sharing (D. H. Kim et al. 2021). In this framework, Choi (2016) proposes a distinction between engaging behaviors of internalization on an individual level and the engaging behaviors of externalization on a social level. As a key element of news engagement, internalization refers to the reception of news by consuming content that appears on a user's timeline, feed, or even when it is delivered in a more personalized manner (Choi 2016, p. 819). Externalization explains a more social attitude toward news content and is related to the dissemination of news to others (Choi 2016). However, this distinction does not touch the underlying epistemological problem that engagement is perceived primarily as behavioral, and its detection and quantification metrics favor behavior over emotion.

Opting to address the issue, Steensen et al. (2020) disguises felt engagement from “practical engagement” (Lawrence et al. 2018) which relates to the behavior of the user, and methodically categorizes it as technical-behavioral, emotional, normative and spatio-temporal. Sora Park et al. (2021) also suggests three dimensions to the public’s engagement with news on social platforms: the first is the technological affordances (S. K. Evans et al. 2017) which affect the level of the public’s exposure to the news, their consumption and interaction with the content (Sang et al. 2020; Halpern and J. Gibbs 2013) and the second is the human factor, namely personality traits (Dafonte-Gómez 2018) or personal political preferences (Arendt et al. 2016) that affect the reception, perception and sharing of news. Selective exposure theory, for instance, assumes that people find it difficult to consume or interact with diverse and mixed content and news (Shin and Thorson 2017; Stroud 2011; Stroud 2008). The third dimension is social endorsements, and they involve a regulatory aspect that connects the technological affordances with the human factor. The pervasive and ubiquitous opportunities offered by social media of accessing and engaging with the news create an environment where social endorsements or signaling behaviors play a very important role in the sharing and consuming news stories (Park et al. 2020; Sang et al. 2020). Specifically, extending the users’ role as “gatewatchers”, social endorsements, such as likes, shares, most popular labels, and lists of recommendations, function as heuristic signals of newsworthy content, subsequently other users will share it as well (Anspach 2017; Messing and Westwood 2014).

As it seems engagement metrics are useful not only for the content creators but for the audience as well. A study on Facebook usage (Bolin and Velkova 2020) showed that when a user is confronted with posts with no representational metrics, namely all the numbers a user sees on post, they were troubled and confused since they felt uncertain on how to engage with the news items. Therefore, it is evident that the technological affordances of views, likes, shares and comments, construct the community and sociality while defining the content’s worth. Eventually, this study showcased that eliminating technical-behavioural aspects of online content seriously affected the other three categories of engagement.

In short, the very design of the social media’s technological and social affordances, influences how a user will engage with the news. This is confirmed by the research of (Masip et

al. 2021) on news sharing on Whatsapp which showed that what determines the decision of users to share a news item is the identity of the sender and the level of social endorsement of the content. Similarly, Johannesson and Knudsen (2021, p. 17) found that people tend to share “content they agree with politically coming from friends and acquaintances who they also agree with politically” thus contributing to ideologically homogenized news feeds or the so-called filter bubbles. Another study, focused on the social interaction of South Korean news producers with their audience on Twitter and how the political inclination can be related to their Twitter activity and behavior (N. Y. Lee et al. 2016) The results showed that more than half of the reporters posted content on Twitter in an attempt to interact with other users. In terms of political ideology, the survey data showed that journalists working for liberal newspapers talked more about public issues than those working for conservatives. This was explained by the fact that the liberal reporters observed that users with similar political views were active on Twitter.

3.1.3 Popularity Metrics

Journalism is an emotional industry that raises fear or anger, inspires joy or affection and evokes feelings of sadness or disgust (Wahl-Jorgensen 2019; Steensen 2017; Peters 2011). The public engagement with the news is not novel as people have been arguing about the election results or which team won the championship in social gatherings for many years. It's the dynamics of subjectivity, the idea of a connection between the particular and the general and the identification with something bigger (Steensen 2017), that is often perceived as engagement, with emotion being considered as its most intense and visible manifestation. Through their engagement with journalistic activity, people may worry, argue passionately, be inspired and even mobilize an entire dialogue or an online campaign. Although, when the public engages with what Picone et al. (2019) refer to as “small acts of engagement,” the level of engagement is limited since it involves only limited activation, such as liking, commenting, and sharing.

Some of the popularity metrics used by the media organizations to distinguish valuable content except for the traditional methods (ratings, circulation) are based on user online behav-

ior and include shares, favorites, likes, “most emailed”, “trending”, “top ranked” and so forth. These metrics along with certain exposure measurements such as page views and time spent on an article, can provide an overview of the perceived value of a news story. Even though clicks are a sign that a news article caught the user’s attention, they may not provide too valuable actual feedback on the user’s perceptions, but rather be a function of other reasons, such as the items’ position in the design of the site (Robinson n.d.). However, when a user consciously decides to hit “like” or “share”, that reflects a more personal attachment to and connection with the content and can serve as a valuable indication of the perceived quality of the article. After all, people also use these actions to construct a personal image, because what they share reflects on who they are or want to be (Trilling et al. 2017).

The predominant form of engagement is the like button (Gonzalez 2015), it is as thumbs-up icon, indicating the user’s intention to approve, agree or support a post, and it is “the most easily recognizable icon on Facebook” (*Logo vs icon* 2016). Commenting is a deeper form of interacting with the news and another way of measuring user engagement that exceeds the simplicity of the like button or the emotionally expressive act of sharing. The participation of a reader in the comments below a news item presupposes a more cognitive intensive involvement and signals their intent to express opinions and feelings and to participate in the discussion (J. B. Singer et al. 2011). In the framework of an interactive digital democracy (Dahlberg 2011), there is a great prospect for public dialogue that can be activated and emerge through this particular form of user communication.

Drawing on various studies that used the number of comments as a predictor for digital content popularity, from YouTube videos (Szabo and Huberman 2010) to movie revenues (Asur and Huberman 2010), Tatar et al. (2011) analyzed 338,394 articles and 2.6 million comments from the French news site 20minutes. They correlated the number of comments with the topic of the articles and the time of publication and they found that opinion articles and domestic news received the largest number of comments whereas international and economic news received the smallest. Also, news items that were posted between 06.00 to 11.00 in the morning were most commented compared with news items that were posted from midnight to 06.00 am. Tsagkias et al. (2009) coded the metadata of the articles (e.g., author name, date,

time, summary), their degree of reposting, name entities and keywords, and found positive correlations between these features and users' participation in the comment section.

Commenting is a powerful social interaction in the digital news space concerning a wide number of users, however, it involves greater personal risk than just hitting the like button or sharing a news item (Almgren and Olsson 2015). Comments expose the author in the public view (Gummerus et al. 2012) which is why many users are reluctant to write their opinion underneath the controversial news items, in fear they may damage their online identity, status or their relationships. Empirical studies have shown that when a user publishes or shares a post reveals a lot about their personality (Dominick 1999) and by these online acts they construct their personal identity (Pempek et al. 2009). It goes without saying that the representation of the user through posts must be a positive one (Zhao et al. 2008). This was also shown by Tenenboim and A. A. Cohen (2015) who, analyzed the clicks and the comments in 15,431 articles of a well-known Israeli newspaper and found that, while the news with a greater degree of deviance (strange news or just entertaining) had more clicks, political and social news and articles with a strong element of conflict attracted most comments. Commenting under a deviant story can make the person look just as radical, extreme and even ridiculous as the story might be, while commenting under a news item of importance to the general public is more socially acceptable.

If engagement had a hierarchy, sharing would be the pinnacle as it is a micro-activity vital to consuming news content (Hermida 2012). Online audience engagement has disrupted the enduring unilateral flow of information from the media but also the power structures of information and the dynamics between publishers and their readers (Wendelin et al. 2017). The act of sharing news on social media facilitates the selective proliferation of news and encourages discussions within users' networks and their engagement with current events (Oeldorf-Hirsch and Sundar 2015) while it is the safest way for the users to participate in the media ecosystem compared to commenting that might jeopardise their online identity or provoke attacks from others (Almgren and Olsson 2015; Tierney 2013). Unlike the journalist who performs a professional role and publishes something following a general orientation of interest obeying to the regulatory role of journalism, for a user the decision to share an

article is, as mentioned, a manifestation of his online identity, a projection of self (Giddens 1984).

3.1.4 Shareability

Meanwhile, a user able to share the news of their liking inevitably becomes a potential gatekeeper who disseminates information on their social media profiles and blogs under the criterion of “shareworthiness” instead of newsworthiness. The possibility of a news piece going “viral” depends on audience engagement, therefore every single user on social media contributes to its “virality” (Thorson and Wells 2016b). Several studies have focused on the reasons behind viral news, drawing from the Uses and Gratifications theory (Ruggiero 2000) searching for motivations like information acquisition, social status, fun, expression (Berger 2014; C. S. Lee et al. 2011), personal meaning (Weeks and Holbert 2013, p. 215) or consistency with predisposed attitudes that results in “selective sharing” (Shin and Thorson 2017; Barnidge 2015). Uses and gratifications can explain various patterns of behavior that determine what kind of needs people seek to satisfy, what type of content they consume and how they are affected by engaging with online news (McLeod 2000). Following this approach, recent studies emphasized personality traits as predictors of social media behavior (Gil de Zúñiga et al. 2017), highlighting people’s tendency to consume content based on confirmation bias or selective exposure (Cappella et al. 2015; Messing and Westwood 2014; Knobloch-Westerwick and Kleinman 2012) that prompts them to share news content with which they agree.

In general, people with a lot of followers tend to share posts that seem to them as interesting, valuable, useful or emotional, however in the case of a controversial post they do not want to risk their following and skip sharing (Mitchelstein and Boczkowski 2010). Going back to news values, Harcup and O’neill (2017) revisited their famous list and added “shareability”, referring to shareable news that are funny or provoke anger, making shareability is an inherent quality of news. Similar to this, Hurcombe et al. (2021), highlights inter alia humor, culture and intelligence as the necessary features for a shareable article.

Moreover, every social media network has its own “platform vernaculars” (M. Gibbs et al. 2015), the “platform-centric” aesthetics and mechanics by which the journalistic product is published and consumed, including funny micro-symbols in the flow of text, such as memes, GIFs (Milner 2018; Miltner and Highfield 2017) and emojis. These elements compose what Highfield and Leaver (2016) calls the “irreverent internet”, a space where the engagement with social issues transforms into a joke often framed in an easy to digest and informal manner. These affective gestures (Papacharissi 2015), contradict the apathetic tone of conventional journalism, and redefine how people relate to the news, however the slogan ‘bad news equals good news’ for publishers existed before social media came to picture. Bad policies attract more attention and their shareability is a given for the news outlets, a fact that is supported by several scholars (Ørmen 2019; Harcup and O’neill 2017; Soroka and McAdams 2015). In a recent study on political news that investigated a corpus of 5 major Mexican newspapers and their shares on Facebook, researchers confirmed a “negativity bias” as negative political news were most shared but they also identified a “sadness bias” as the sad emoji reaction accompanying a news item was a strong positive prediction for its sharing (León and Trilling 2021). On the contrary, if a news item had more love reactions, that was a negative sharing predictor. To explain virality and emotional contagion, Dafonte-Gómez (2018) investigated what motivated users to share the news, and found that both negative and positive arousal influences shareability. In a similar vein, Guerini and Staiano (2015) investigated the relations between emotions and news virality and discovered that dominance was a defining predictor of sharing content as it made readers feel more in control (Russell 1980).

Looking for the role of emotionality to news virality as well as the type of political news that tend to be most shared in social media, Hasell (2021) analyzed 300K tweets from 22 news organizations and found that hyperpartisan articles were shared more frequently and also employed more emotional language compared to non-partisan articles, which is why they tend to appear “amplified” on social media. Similarly, the study by Hopp et al. (2021) has shown that people on the far ends of both the conservative and progressive spectrum tend to share false or highly biased news and draw their material from unreliable sources mainly

due to a lack of trust in established news organizations.

Another issue related to audience engagement in social media and sharing political news is the users' political positioning, the partisanship of the news stories and how they influence public opinion and participation. A study on news values and sharing articles about politics in 3 Korean newspapers from different ideologies, showed that content sharing was influenced by the political attitude of the reader and their different ways of viewing society (Sora Park et al. 2021). Looking at why people share political news in social media, D. H. Kim et al. (2021) found that users will share political information primarily to inform, socialize and interact with other users or to criticize opposing views but not to express personal opinions or to promote their interests. Apart from regular reading and political behavior, Sturm Wilkerson et al. (2021) explored the phenomenon of partisan media, which often dominate public discourse and social media. Analyzing the posts on Facebook pages of partisan news websites, they found that political posts on the right pages gained more likes and were shared more than those on the left pages, while articles that provoked or contained anger were more often shared on left pages than on right pages.

3.1.5 News Values and Engagement

The idea that news comes with some inherent qualities that renders it publishable dates back to 1965 and the works of Östgaard (1965) and Galtung and Ruge (1965), although the very definition of “news values” is attributed to Lippmann (1946). Notions like sensationalism, proximity, relevance, unambiguity and facticity were presented as factors that explain what makes a story worthy of publishing. These factors were later extended by Galtung and Ruge (1965) to an empirical set of 12 values that journalists rely on to determine if an event is ultimately newsworthy. News values evolve over time, therefore subsequent studies developed novel news values such as the unexpected/surprise value (Staab 1990), deviance and social significance (Shoemaker and A. A. Cohen 2012) or the impact/magnitude value (O’neill and Harcup 2009). News values within deviance include novelty, oddity, controversy, conflict, and sensationalism, while examples of news values within the social significance dimension include importance and impact (Shoemaker and A. A. Cohen 2012).

Except for journalists also the audience judges the news based on news values (Eilders 2006), a fact that led communication scholars to hypothesize that the same values could predict virality on social media. Various studies have investigated the contribution of specific news values to audience engagement and the speed with which news spreads. P. Weber (2014) found that news values influence readers' engagement and interaction in comments, while Schaudt and Carpenter (2009) found that the news values that attract online readers are: proximity (how close to home an event occurs) and conflict, while temporality (the timeliness and duration of the event) and prominence (news involving prominent persons) are the least attractive.

Other studies showed that news of high social significance (political events, welfare state issues, etc.) have more comments than news with high deviance, such as novelty, controversy, and sensationalism (Ziegele et al. 2014). In the contrary, Rayson (2017) reports that the most successful articles on Facebook were characterized by the news values of novelty and surprise, and that political posts have more comments and shares than likes. Sora Park et al. (2021) analyzed the performance of 2.5K articles from 3 Korean newspapers on Facebook based on the model on deviance and social significance by Shoemaker and A. A. Cohen (2012). They argued that both dimensions were predictive factors for engagement, with deviance best predicting likes and social significance comments. However, they found large differences in the intention to share texts, as news of social significance was shared considerably more frequently than deviant news.

Berger and Milkman (2012) analyzed 7K news stories from the New York Times and found a high degree of sharing useful news that competed with the sharing of popular news items with a strong element of surprise. They also highlighted that great emotional arousal causes news to become viral, regardless of the emotions it inspires (admiration, joy, anger or anxiety). Moreover, sharing on Facebook has been related to the news values of eliteness (Caple and Bednarek 2016), human interest, conflict and controversy García-Perdomo et al. (2018). Switching newsworthiness with shareworthiness, Trilling et al. (2017) analyzed 132,682 articles from 6 major Dutch news sites and their circulation on both Facebook and Twitter to test whether news sharing on social media could be predicted using news values. They

found that texts with high geographical proximity, positivity and the exclusives had the highest degree of sharing, contrary to news with the values of conflict and human interest. Another finding was that news about a personal story was more shareable on Facebook than on Twitter. The exclusivity of news stories was a positive sharing factor on both platforms, exactly like an “exclusive” story is an important factor of journalistic success.

Another way of approaching the shareability of news is the information utility model which has three dimensions, a) the perceived importance of the challenges or satisfactions an event may bring b) the perceived probability that these challenges or satisfactions will be realized and c) their perceived immediacy (Knobloch et al. 2004). For example, a news story that informs about which roads will be closed due to construction can be perceived by readers as important, as it has a high probability of affecting their lives (if the roads are part of their daily route) and it may occur promptly. This was also shown by two experiments conducted by Bobkowski (2015) that positively correlated the perceived utility of two articles with the intention of the participants to share them.

3.2 The Power of Emotions

In newsrooms across the world, emotions play an important role in how stories are pitched, how they're reported and how they're ultimately presented to readers, viewers and listeners. For journalists, understanding and utilizing emotions is essential to creating engaging, impactful and high-quality journalism. While emotions have always been a part of journalism, the rise of social media and the 24-hour news cycle has made them even more important. With audiences now more fragmented than ever, journalists must be strategic in how they use emotions to capture attention and keep people engaged. One of the most effective ways to do this is through the use of empathy.

Empathy is the ability to understand and share the feelings of another. And it's a key ingredient in many of the most successful and impactful journalism stories. When journalists are able to tap into the empathy of their audience, they're able to create a connection that can be powerful and long lasting. The best journalism tells stories that make people feel some-

thing. They make us laugh, they make us cry, they make us angry, they make us think. And while not every story needs to be emotional, utilizing emotions judiciously.

The internet is full of information and audience engagement can make a news story stand out. As we have seen in the previous chapters, studies have revealed the importance of emotions expressed in news stories because a story that evokes an emotional response makes readers more invested in the narrative. While expressing emotions in news coverage may challenge the traditional norm of impartiality, several scholars have made a call for an “emotional turn” in journalism (Wahl-Jorgensen and Pantti 2021; Lecheler 2020), especially in light of the current hybrid media system (Chadwick 2017) and networked publics.

Emotional stories have a specific effect on the audience, and some scholars conceptualize emotions as an essential part of journalism (Wahl-Jorgensen 2020; Beckett 2015). Emotional responses in journalism encourage readers to attend to coverage and offer opportunities for journalists to express their opinions about political issues (F. Moore 2018). Moreover, emotions help people to judge the news (Gluck 2019), which may result in feeling disappointment, anger, frustration or sadness. It is important to be noted that emotions are experienced based on certain circumstances a person is in, therefore an emotional news story may instigate reactions, like anger or sadness in one individual and fear in another. This is because people react differently to the same emotional stimulus based on past experiences and future goals. Additionally, emotions are amongst the human processes that guide attention to information (De Sousa 1990) and act as an anchor for people to navigate through the plethora of signs in their everyday lives (Reddy 1997, p. 331:332).

While emotionality has always been intertwined in journalism, the aim of this chapter is to explore how media coverage uses emotion and how it influences journalistic quality and audience engagement.

3.2.1 Definition of Emotions

Emotions have been characterized as being the way in which a body understands affect, and they depend on what the individual thinks is happening to their body and can be shown to

the public by way of description, articulation, and circulation (Davidson and Milligan 2004). Emotions are not simply a human response to an internal or external stimulus; they arise from the interplay between a person and their environment (Davidson and Milligan 2004), and there are different types of feelings that arise from an event or experience. Therefore, the word emotion is not just a word used to describe anger or sadness, but rather refers to reactions and experiences that involve the body, mind, heart, and spirit (Freeden 2013), and influence the ways people experience and understand their surroundings.

Although there are different emotional definitions in the literature, here we take a broad definition of emotions. Emotions can be described in general as feelings (Davidson and Milligan 2004), but as Nummenmaa et al. (2014) suggest emotions provoke bodily reaction, which translates into changes in the skeletal, muscular, endocrine, and autonomic neural systems (Levenson 2003). These physiological responses stimulated by an impulse are referred to as affect (Massumi 1995), while emotions are subjective feelings of an individual that are not verbalized or otherwise made explicit (Massumi 1995, p.88). Therefore, affect is an emotional reaction that causes a body to act in a certain way (Deleuze and Guattari 1987), while emotion is the feeling that we get when something happens, which can be something both good and bad. In other words, affect is to cause a person to do something or have an effect on them; however, emotion is the feeling that they get when experiencing something. Those feelings may vary from person to person and they can be easily shared or exchanged through the use of non-verbal communication. Moreover, affect precedes emotion and directs the intensity with which emotions are experienced, since affect itself is characterized by its intensity (Deleuze and Guattari 1987; Massumi 2021).

Emotions can make the audience identify with the people presented on the news who are in a dramatic presentation, and cause responses to the audience through mental, cognitive, and physiological reactions. These can be felt physically, but also expressively and behaviorally; and they influence previous thoughts, predispose to certain responses, or shape future actions (Freeden 2013). It is believed that emotions guide attention to information relevant to the goal of the individual who is experiencing them, and therefore serve as a tool for social interaction. For instance, anger drives people into participating in diverse debates

not only with like-minded people but also with out-groups, although angry users tend to seek information that confirm their prior beliefs, resulting in “media diets” (Wollebæk et al. 2019). Additionally, cultural and social norms can function as “amplifiers or ‘brakes’ of emotions” in some contexts, but in others may function to moderate emotions (Elster 1999, p. 262). Furthermore, there is specific research to support that emotions are not at odds with logical thinking (Bandes and Salerno 2014), and as De Sousa (1990) explains purely rational or emotionless thinking is a myth. Emotions work together with reason and intuition to facilitate and evaluate information processing. Also, he argues that emotions are composed of both descriptive and subjective elements. The emotions involved in news stories are usually very specific, as they are induced by a concrete situation. This specificity is achieved through the use of specific vocabulary, as well as by describing emotion-conveying voice or facial expressions.

3.2.2 Emotions and Journalism

Emotions affect the way people interact with media and create an atmosphere of presentational style, which is a form of storytelling (Lecheler 2020; Beckett 2015; Peters 2011). Therefore, it is important to understand how emotions are used to craft an effective narrative. Dealing with emotions within the newsroom has been found to be tricky for news professionals because personal feelings are perceived as intimate experiences, and journalists, anchors, politicians and public figures have been criticized for showing them in public. Instead, journalists feel comfortable with “objectivity”, a quality that has been associated with news content since the early days of journalism, which includes the notion that reporting is about facts, not about opinion. However, emotions have a place in the world of journalism and are used to get attention and connect with the audience on an emotional level.

The perception of emotions and how they are shown in the news has changed over time. Researchers have found that there is an increase in the number of emotionally charged stories, along with a decrease in the amount of objective or documentary reporting because emotions are related to high sharing on social media (Kilgo et al. 2018a). Furthermore, peo-

ple respond to emotions rather than ideas or facts, as argued by Beckett (2015) who points out that in order to find pertinent news, people need the stories to be more human. Since online news are no longer just a channel for mass communication but also a social space where people share information with others around the world, when journalists give a voice to these feelings in news stories, it could motivate people to take action. Moreover, emotional personal stories encourage the reader to share their own emotions, which may result in them being more likely to share the story (Kilgo et al. 2018a) More. Therefore, in modern newsrooms, journalists tend to write emotional stories because they have a greater chance of success on social media than those that lack emotion (Peters 2011). It could also increase empathy and understanding among readers, which would make journalism more successful as a whole (Pantti and Wahl-Jorgensen 2021; Kotišová 2019; Wahl-Jorgensen 2019; Beckett 2015).

Emotions function as dramatic devices that enliven and crystallize abstract material and can evoke a mental, aesthetic, and moving effect on its consumer (Freeden 2013, p. 7), while concurrently acting as “intensifiers of the ideas, concepts and attitudes to which they attach”. Thus, intensity is an aesthetic consideration, a parameter that signals and designs an audience response to the content of discourse. The use of emotions in politics has been extensively studied by researchers. In his study, Richards (2009) refers to the “emotional public sphere”, as an inherent affective activity constantly present in political discourse and to “emotional governance” as a process where political actors strategically use emotions to elicit particular political responses from the public (Richards 2007). Moreover, emotions can mobilize people in political and social movements so that they can be a driving force in politics (Bas and Grabe 2016), while the use of emotions in political discourse has been shown to influence the public’s response to political candidates, issues, and parties (Oschatz et al. 2021; Bas and Grabe 2016). In other words, it has been shown that the use of emotions by politicians can influence people’s support or opposition to a candidate or issue. Moreover, social emotions are connected to political and civic participation (Demertzis 2013, p. 9) and play a significant role in shaping the political decisions Papacharissi (2015) thus are fundamental to the perpetuation of a stable political environment.

According to Demertzis (2020, p. 5), emotions are mediated culturally as they emerge and are experienced situationally and relationally, expressed according to social conventions and emotion structures that define their valence, arousal and intensity, whereas they are discursively spread within and through language games thus participating in the processes of identity and will formation. Although traditionally the locus of emotion is the individual, by placing the analysis of emotional phenomena at the micro level, Demertzis notes that the composition of emotion does not shrink in a simple biological process as individuals are integrated within social structures. Norms as macrostructures are instantiated through emotional practical representations (Demertzis 2013, p. 13). The very concept of emotion implies that socially constructed rules in fact dictate the very experience and externalization of emotions in certain situations. Thus, the macrostructure of reality as perceived through the micro-experience of emotion coalesces the “macrocosm and the microcosm as two pre-existing autonomous entities”.

Emotionality is a very important factor when it comes to media, especially with the rise of social media where emotional communication between users lies at the heart of those platforms. In social media posts, the author’s emotional state is often relayed, along with the reason for it. Furthermore, the interactivity of news media has made people from passive readers to active or as Papacharissi (2015) calls them “affective publics” that reside in a set of affective gestures, such as the reaction emoticons as well as sharing, liking, and commenting. Affective publics are defined as “a set of actors with a history of sustained, non-hierarchical participation in media use who interact using affective media, who pursue emotionally meaningful goals in the context of online news” (Papacharissi 2015). The technological affordances of the social platforms give way to a number of affective microactivities that transform news items lending them their dynamic participatory structure that is “affective in nature”. Therefore, a news story intended for social media is a priori emotional. There has been extensive research on how emotions influence the share of information, consumption, and engagement, and how the role of emotions in news has changed over time. In general, there is no way to control the emotions of individuals, therefore the only way to steer emotions while reporting news stories is through a reporter’s choice of words.

Research showed that positive emotions encourage political participation (Capelos 2013; Brader 2005; Johnson and Tversky 1983) and risk-taking (Gross et al. 2009), while negative emotions cause aversion (Capelos 2013). Furthermore, emotions have been linked to the ways people organize information, and how they recall facts about politics and make decisions (Kinder 2013; Brader 2011). Also, it is well established that emotional stories tend to be more popular than non-emotional stories on social media, and tweets with strong emotional content are disseminated both more and faster than those with a neutral tone (Jensen et al. 2013).

Clough (2008) has highlighted the emergence of an “affective society”, wherein emotions are a significant factor in news consumption, as well as in journalism and news production. This shift has been largely driven by the rise of social media, which has made it possible for people to share their emotional reactions to news stories with large groups of people. The disruption of the traditional one-way information stream that has turned each reader into a node of a network, has also led markets seeking to capitalize on user’s emotions in order to engage with their “private publics”. In addition, communication scholars have observed that actors behind disinformation have taken advantage of this dynamic, using collective emotions to dominate public discourse (Davies 2018).

A decade ago, Papacharissi and Fatima Oliveira (2012, p. 279) studied the effect of Twitter in the Egyptian revolution of 2011 and how events have been recounted during an oppressing political climate. In the Arab Spring narrative, the role of social media in spreading information was vital in building pressure and awareness of issues that were previously thought to be minor concerns. Social media provided human faces on a global scale to the voices of concerned citizens who were no longer able to tolerate the injustices they had experienced for so long. Studying the role played by Twitter in disseminating and connecting the crowds over this historic upheaval, the researchers introduced the notion of “affective news” as tweets increased the “drama” of instantaneity that pulled in both journalists and the audience. Facts, opinions, and emotions were intertwined in -occasionally- subjective accounts and interpretations of events in real time, while allowing people to create bonds online. The researchers also found that the same news item would be repeated multiple times without

changing the information but the emotional input.

In his work, Kramer et al. (2014) described emotion as social and contagious, an observation that aligns with social media users who are attracted to platforms that will allow them to connect and share their opinions with others, and often want to be heard and either post their own content online or express themselves in the comments sections of posts and news items with anger, admiration, solidarity and so on. Furthermore, when people are confronting stressful or emotional subjects, they will naturally want to discuss it, socialize and react, because doing so allows them to regulate their own emotions via social interactions (Zaki and Williams 2013). Moreover, engaging with online content allows users to demonstrate their thoughts or delve deeper into the topic discussed (Yoo et al. 2017) and it strengthens their link to news as well as the events described therein (Oeldorf-Hirsch and Sundar 2015).

New platforms often bring with them new ways of perceiving and existing in the social world. In their work, Beckett and Deuze (2016) emphasize the shift in thinking about emotionality in journalism and talk about an “affective media ecosystem” where people interact with the news, sharing their own opinions and experiences, and journalists connect with their audience on an emotional level. In these different online interactions and conversations, the news becomes a key component of personal and social connections, voices of people whose stories and opinions would not have been heard are published and sidelined or less explored topics are discussed and communities both online and offline are formed (Howard and Hussain 2013).

Papacharissi (2015) explored the role of emotions in people’s reactions to news coverage by analyzing the hashtag #ThisIsACoup in relation to Greece’s long negotiations with its Eurozone partners in the summer of 2015. She noted that the driving force behind the hashtag’s engagement was public sentiment, turning it into an “open signifier” filled with emotions. Hashtags frame events and give their users the power to negotiate their meaning in an emotional way. Examples of this include #BringBackOurGirls and #BlackLivesMatter, which have ideological purposes and serve as signifiers “open to definition, redefinition, and re-appropriation”, guiding random users and undefined crowds to coalesce into publics and

“affective publics” through mediated interactions. A prime example of this is the #metoo campaign, which raised awareness of sexual harassment and prompted global discussions about the underlying sexism of our society (Hensbergen 2017).

In today’s competitive media industry, it is not easy for a news outlet to differentiate their product as the market has never been more saturated with content. Publishers need to ensure that they’re keeping a close eye on the competition’s social media strategy in order to get ahead of them. The key to doing this successfully is by appealing to social media users’ emotions, and for that understanding the audience’s mood can be important. Competition among news outlets and the informational abundance drives a journalistic quest for engaging strategies (Wahl-Jorgensen and Pantti 2021; Beckett and Deuze 2016) that go beyond technical affordances and simple marketing. The very nature of technology, which may be characterized as distant and cold, requires emotional cues to attract the attention of the user (Beckett 2015) and secure their engagement: strong visual elements or controversial language are certain to catch a wandering user’s eye so they may perform their distributing role in the media ecosystem. They also highlighted the fact that reading the news moves toward a more personalized and mobile experience with emotions becoming even more important (Beckett and Deuze 2016). Specifically, these changes have led to the creation of “news as you want it” model and aims to actively stimulate the emotions of the readers. Finally, news content published on platforms like Instagram and Tik Tok is often characterized by its short, quick, and quite often rather sensationalist content to create certain moods and trigger emotional responses.

Important events, such as the Asian tsunami in 2004 and the Boston marathon bombings in 2013, were first reported by citizens and captured by their mobile phones. The user-generated content is important in the news industry as it is considered real-time and authentic, and has been shown to change the public’s awareness of certain matters. According to Allan and Thorsen (2009) this “citizen witnessing” does not subscribe to any journalistic norms, but rather to the “vernacular of their lived experience” (Wahl-Jorgensen 2020; Chouliaraki 2013). Furthermore, people read the news via their smartphones on social media which can make their access to information something even more intimate and specific

than it is with mass media (J.-H. Schmidt 2014).

One of the most important changes has been the way that information is now shared. Today, news can be shared instantaneously through social media and other online platforms. This has made it possible for journalists to reach a much wider audience with their work. Inevitably, a news outlet is faced with the dilemma of finding a formula that will make its product as appealing as possible to its online audience. One of the ways publishers implement, even in hard news, is to increase the degree of emotionality present in the news piece (Wahl-Jorgensen 2020). In a study by Kilgo et al. (2018a), they found that a story was more likely to be shared when it contained a significant amount of emotional information. However, they ascertained that despite sensationalism's intent to attract readers' attention and evoke emotions, it was not correlated with more interaction, which may indicate that the emotional portion of news does not always work as intended.

Social media platforms have become common forums for news consumption and journalists often use social networks to report from the field. Additionally, technological affordances lead journalists to report, discuss and negotiate events online subjecting their work to much more public scrutiny than ever before, but also promoting the engagement of journalists with their readers. Many communication scholars believe that emotions lie in the heart of journalism practice: they act as its "inspiration, creation, style, appeal and its resonance or impact" (Beckett 2015). Compelling narratives and vivid descriptions were always among the notable characteristics of good journalism, but the growth of online media has enabled journalistic storytelling to become more personal, intimate, and engaging (Lecheler 2020).

Objectivity has long been the ideal in journalism, and the complete opposite of emotionality (Maras 2013). Yet, as journalism has moved on to the digital era, is evident that the old principles are no longer enough and emotions have become a key element in journalism. Even though emotions have been always present in journalistic practice they were believed to compromise quality in journalism and were linked to tabloidization, vulgarization, oversimplification, voyeurism, commercialization and so on (Peters 2011; Pantti 2010; Harring-

ton 2008; Sparks 1998; Slattery and Hakanen 1994). Objective journalism is meant to be “eye versus” “I” (Hopper and Huxford 2015) (p. 33) which has led to the idea that news, as reported, is impartial, not biased, without any personal feelings influencing the reporting (Beckett 2015).

Tuchman (1972) and Schudson (2008) and Schudson (2001) provide valuable insights into how the objectivity norm has changed over time. In the 19th century, objectivity was seen as the to-go-to norm in journalism (Tuchman 1972). However, in the first decade of the 21st century, this scenario changed and social empathy became a more accepted form of understanding different people’s experiences (Schudson 2008). T. R. Schmidt (2021) analyzed how emotion charged news became an acceptable form of journalism in the United States with reporters including in their coverage their own experiences. A number of studies show that the current assumption of emotionality in news media has shifted away from the old opposition between emotions and reason (or objectivity) (Wahl-Jorgensen 2013; Pantti 2010; Stenvall 2008; Kitch 2003), to a new view that focuses on how they tend to complement each other (Hermida 2016).

According to many scholars, people are more likely to take action against social issues when they are informed and engaged with the problem. Noticeably, Hermida (2016) claims that in order for audiences begin to understand a problem or issue, they need to see its impact on the lives of real people. Thus, in order to change people’s individual and collective behavior, empathy for members of a designated group is likely to be necessary. As Hermida (2016, p. 53:54) demonstrates, social empathy is a process of comprehending and communicating the experiences of others and offers a more holistic understanding of how emotions impact the transmission of news. Some scholars note that using emotion in journalism may improve the audience reception of news stories (Lecheler 2020; Hermida 2016; Peters 2011), and the need for emotional news coverage has been explained by Beckett (2015) who insists that “all journalists are human”. Specifically, he explains that everyday routines in the newsroom from selecting a story or conducting interviews to setting an agenda are all driven by subjectivity. Also, thinking about how emotion helps to ensure the reception of the news, at least in some cases, is well worth thinking about the impact on journalism itself of new

technologies that bring us closer to the experiences of others. Thus digital journalism enforces different emotional responses allowing the audience to identify with characters or protagonists in situations described within journalistic narratives thereby increasing levels of rapport or trust in the media.

Despite the protests of some scholars and journalists, remaining impartial is often regarded as an insufficient response to challenges like sexism and racism, homophobia, marginalized groups, and social injustice in general. Studies of emotional content in news articles Baden et al. (2019) have found that the use of emotion in news media is associated with affective reactions, especially in terms of activism and social change. In reporting on sensitive issues such as those related to sex and gender or disasters and trauma (Jukes 2017; Kovach and Rosenstiel 2014; Joye 2013; Kitch 2003), the use of emotion has been shown to increase audience engagement and motivate people to act on behalf of change. The theory is that by hearing the voices and seeing images of those affected disaster victims will be more sympathetic to those who are hurt by such events. This can lead to a greater sense of empathy as they see the actual people suffering rather than just numbers in terms of victims affected or displaced. Additionally, this type of coverage can help rally support for aid and donations which have been found to be significantly lower without emotional coverage during disasters. It was found that the proportion of public donations to disaster victims increased when emotional elements were included in the story (Maier et al. 2017). This news coverage can also make people more likely to engage positively with another even if they do not share their opinions congruently.

Journalists have the power to influence individuals through the content that they write or by how they deliver it. Specifically, the way journalists use emotions is also significant when determining how a particular story will be received by an audience. This can have an effect on how people respond to the story. Baden et al. (2019) investigated the impact of disaster stories that are written either with the intent to stress the catastrophe of an event or ways to resolve it. Their findings suggest that people reacted to sad news with more negativity and had less motivation to do something about the issue, whereas when the story was framed as a solution that evoked positive emotions and they were more likely to take action.

Some scholars also suggest that focusing on emotions helps audiences better connect with journalists which does not thwart the rationales of truth and objectivity and could increase trust in the media (Pantti and Wahl-Jorgensen 2021; Orgeret 2020; Kotišová 2019; Wahl-Jorgensen 2019). This is due to the fact, that people intuitively comprehend the world through both their cognitive and emotional systems, and journalists can use emotionality and subjectivity but at the same time affirm the values of truthfulness and credibility of news reporting (Orgeret 2020) to ensure the journalistic information they provide is reliable and factual. As Pantti (2010, p. 179) explains “the main objective of emotional storytelling was to enhance the political and social knowledge of the audience, to facilitate the understanding of news”.

Journalistic reportage can make readers feel something, rather than merely describe it, and in today’s digital world emotionality and personal feeling are impacting what journalists report. News articles are often told from a first-person perspective which means readers would have a better understanding of the context and the source. This is an integral part of news reporting on disasters (Kovach and Rosenstiel 2014; Joye 2013; Kitch 2003), migration and terrorism (Kotišová 2017), humanitarian crises (Nikunen 2018; Chouliaraki 2006), and reporting on the elderly where reporters are becoming more social to their audiences without necessarily translating into tabloidization (Meijer 2001). For instance, BBC war correspondent, Martin Bell, reporting from the Balkans war proposed emotive language rather than cold objective facts, with his strategy being to showcase the tales of victims to shock and emotionally persuade people (F. Moore 2018). In this vein, Bell argued that the viewers need to be emotionally invested, and this can be achieved through a journalism of attachment. This is backed by the study of Koivunen et al. (2021), who analyzed 4K news stories published in four Finish news outlets and found that direct quotes included more emotional identifiers than non-quoted text confirming the regulation of emotionality as being outsourced into the words of others.

To detect how emotions are presented in news items and challenge widespread assumptions about objectivity and factuality, Stenvall (2008) opted for the strongholds of objective journalism and examined hard news from the Associated Press and Reuters. One obser-

vation was that journalists used distancing and detached expressions in hard news by using nominalized emotions (nouns designating emotions) in a way that the emotion could discursively function like participating in the described events but itself, to mask the subject who is feeling the emotion, and what has provoked it. The study further demonstrated how the text was made impersonal and intended to hide the journalistic voice in the background, making emotions appear as “free-floating entities in material processes” (Stenvall 2008, p. 1577). Stenvall (2008) argues that this linguistic detachment is also employed when journalists want to distance themselves when describing abnormal behavior of big masses in public demonstrations but lead the narrative on how the audiences would feel (Martin and White 2003, p. 63). Maybe those methods of emotional detachment makes it easy for readers to accept the news without questioning its accuracy or asking questions.

As emotions help memory generate paths and create connections between information, empathy may ultimately define moral decisions and judgment-making processes (Gluck 2019), while affecting the audience’s recall of the reported events (Mujica and Bachmann 2018). The representation of emotion is found everywhere in journalism, from photos, video footage, editorial cartoons, and infographics, to the tone of voice in the narration, the article and the headlines. Beckett (2015) suggests the “cycle of sensitive content creation” that happens online, where a. the story breaks b. the journalist produces a piece live, in real-time, that elicits emotion from the audience; c. the audience has a reaction; and finally, d. the emotion is transformed into sharing and commenting on the story. Beckett (2015) argues that emotion is interwoven into the journalism process, and the emotional impact of a story is directly related to the reader’s willingness to share, comment, and interact with the piece. According to other scholars, the increasing influence of online news and social media is having a significant impact on how future news is produced, thus it is important for news organizations to remain receptive to “emotional output” from the audience.

The ability of the journalist to produce reporting that arouses readers’ emotions through writing that is rich with emotional detail, has been described by Wahl-Jorgensen (2019) and Wahl-Jorgensen (2013) as a “strategic ritual of emotionality”. In her analysis of Pulitzer-winning news stories, Wahl-Jorgensen (2020) discusses how journalists use emotionality

in reporting to create understanding between news subjects, and audiences, by letting the protagonists to share their personal pain, fear, or anger, without misrepresenting the facts. Similarly, Rosas (2018) spoke with 33 journalists from five Spanish digital news outlets aiming to explore the motivations of reporters to employ an emotional approach or to follow strictly the objectivity norm in their stories. She found that the practices vary significantly depending on the individual, others prefer to remain detached and objective, while others reported having to go beyond journalistic objectivity and allow for a greater focus on human emotions.

This shift away from objectivity has been a large topic of study among academic and industry leaders, who explored the phenomenon of “intimacy” in news production (Steensen 2016; Coward 2013; Enli 2014). One study by Enli (2014) found that the inclusion of the personal opinions and subjective views of the journalist in their news coverage happens in order to justify the veracity of the narrator and the trustworthiness of the journalists themselves. A study conducted by Coward, 2013 suggests the personal involvement of journalists in the news story is more present on soft news where the journalist takes a more personal and close approach, almost “confessional”. In addition to the academic interest in the intimacy of journalism, research has shown the popularity of the “story personification” of news (Maier et al. 2017). What the protagonists did, and how they felt after the events took place with the inclusion of their personal perspectives is used to heighten the empathy and connection of the audience with the event, thereby giving them a more complete understanding of what was going on. This practice of focusing on human stories has been found to educate lower socio-economic classes and those who are less educated about the world around them (Bas and Grabe 2015) and motivate citizens to actively participate in the democratic process (Oschatz et al. 2021). Moreover, the stories reporting the misfortunes of an individual result in a higher desire for support (Maier et al. 2017) and have the power to strengthen communities. A vivid example of this can be found in the news media coverage of the events of 9-11 (Zelizer and Allan 2011, p. 9).

In the journalistic practice of covering events that might otherwise disturb the social order, such as wildfires, earthquakes, and floods, in order to create a resonance between the

journalist, who is not a part of the story they are telling, and the people they report on, emotionality is crucial. Reporters are expected to feel with their subject in order to be able to connect with readers. Huan (2017) found that the proximity between the journalist and reporting subject is often seen as a means to reinforce shared principles and bring social order. After examining Chinese and Australian news, the author concluded that the use of emotions in news stories could bond social empathy, thus producing a sense of unity. However, Australian journalism focused on how ordinary people feel, whereas Chinese journalism emphasized how the elite feel. By employing emotionality in journalism, the individual members of the audience naturally associate with those collective feelings in the same social context and are more likely to behave accordingly (Stenvall 2008; Kitch 2003).

The concept of shared emotions in journalism was researched by (Capelos et al. 2018) in the context of the Greek financial crisis. For the analysis, they used news stories, opinion articles, reader's comments, and so on, published between 2009-2012 in four elite UK newspapers, to observe how the reporting of emotions would illustrate the way they were experienced in those articles. The data individuated three clusters of emotional information in order to assess the picture they would portray a. individual emotions discussed on a personal level (quoting I, me or my), b. collective when using pronouns like we, us, or ours, and c. social when defining an external group by the usage of plural pronouns they or them. The study generated three key insights about the "emotional economy" of the crisis (Capelos et al. 2018, p. 18), first commentaries reflected the public's negative attitudes towards the EU, while editorials expressed lower amounts of hostility; second emotions related to aversion, like shame or fear, were reported as a shared emotion either collective or social, whereas anger was depicted as individual emotion; and finally, positive valence emotions like joy or hope were sparse in the dataset. The authors claim this overall dissatisfaction with the EU from the UK population maybe paved the way to the Brexit debate.

Pantti and Sumiala (2009) focused on how the media report matters related to death, pointing out that journalists go beyond the role of the objective and distanced observer. In regards to the shooting of the Dutch right-wing politician, Pim Fortuyn, or a bus crash with many casualties in Finland, the researchers highlighted that the media generated a "public

mourning ritual” that focused on social inclusiveness, thus providing a sort of a “psychological instruction” indicating the new ways of how media narrate emotions on a collective level to serve larger purposes. In the same vein, Chouliaraki (2008) , refers to this news as ecstatic news that highlights the suffering and helps to present a realistic and more complete account of the events with which the audiences will identify.

The role of emotions in news production has, however, also been discussed in relation to mediating social responses, and interpreting complex events. In a recent study based on UK news coverage of protests, Wahl-Jorgensen (2018) found that anger was used in news reporting as an explanatory framework to break down the reasons for this kind of collective action. Therefore, mediated anger appeared collective, political and ultimately discursively constructed through the narrative of journalists, representing the anger of the people they're writing about. Within this context, the ideology of the news organizations also played a role in conditioning the audiences' favorable inclinations, as right-leaning newspapers would report in favor of anti-immigration protests, whereas left-leaning newspapers would tend to show their identification with matters of social injustice. The mediating role of emotions was researched also by DeLuca et al. (2012) during the Occupy Wall Street movement, where the conservatively slanted blogs presented the protesters as being angry and threatening, in contrast to bloggers who favor a liberal point of view who saw the anger of protesters as a legitimate expression of their political opinions. As a consequence, these left-leaning bloggers would not see the protesters as dangerous because they are not violent, but rather are angry because they are justified. These two positions are certainly biased and imply that emotions strengthen the political positions of journalists, but they also underline how emotions may incite people to action, and thus stimulate the development of collective consciousness. In the same line of research, D. K. Brown et al. (2020) analyzed Facebook posts and found that mainstream media outlets tend to over-emphasize the emotional and dramatic aspects of a protest when reporting on them, without providing a more complete description of why people were for instance so angry or disappointed. Creating negatively charged narratives leads to public support for a movement decreasing, and the author noted that the media have the potential to exert a significant influence on public perspectives and allegiances, as

they are able to legitimize some protests and delegitimize others.

To speak the language of younger generations news organizations like BuzzFeed and Vice in the UK tend to adopt a more personal and emotional tone, referred to as a 'youthquake' (Sloam and Henn 2018) because it "shook" British politics. Dennis and Sampaio-Dias (2021) analyzed news articles from the two aforementioned digital news outlets published during the 2017 UK general election and found that the informalization of news (Wahl-Jorgensen 2020) was not simply a matter of approaching the young voters but also an inextricable part of the current journalistic practice of those organizations that could serve as a means of introduction to serious political stories. The researchers detected the use of colloquial language and the presence of a conversational tone suited for a young and digitally-savvy audience as well as the use of humor and inclusiveness (frequent use of the pronoun we in the texts). This could be seen as an attempt to reach out to the younger generation and convince them that the news stories they read are worth their attention and create a community of young voters.

Journalists who frame stories around emotions such as fear, anxiety, or anger do not just provide information; their stories also create a sense of urgency, relevance, and strong reactions. Additionally, it is an established fact that the reader is more likely to share a negative or sad story with others, thus amplifying its effect (Choi et al. 2021; León and Trilling 2021; Soroka and McAdams 2015). However, positivity has also been linked with equally strong reactions (Lecheler and Vreese 2013; González-Bailón et al. 2012). A study by De los Santos and Nabi (2019) on how emotions influence memory and interaction with a news story, showed that individuals exposed to anger or hope-framed news stories recorded more information in their memory, than those who read about fear. Similarly, a study by Kühne and Schemer (2015) used a story about a car accident framed either with anger or sadness and revealed that people exposed to the anger-framed article were associated with unfavorable and punishable reactions to the perpetrators, while sadness evoked the intention to help the victims.

It is apparent, that the perception of a news story depends on the emotional frame it is

framed under, although some frames are more emotional than others. For instance, the human-interest frame (Aarøe 2011) or the conflict frame (Bartholomé et al. 2015) are supposed to bring complicated sociopolitical issues on an individual level that audiences find more relatable and are typically emotional in nature. Valenzuela et al. (2017a) analyzed 37K articles from six online news providers in Chile on the impact of four general frames on the audiences' comprehension and engagement with the news on social media, namely conflict, economic impact, human factor, and morality. The results showed that journalists mainly use the human factor and conflict frames and that the readers were more likely to be engaged with a moral-framed article, and shared it more. The experiment of Xu et al. (2020) who analyzed 1K news stories from 17 hyperpartisan American sites on immigration, agrees with the previous study; The moral framing specifically the authority/respect frame along with its emotions was the most engaging, while fairness had the opposite results, and care although disliked by the readers, were widely shared. However, further research is needed to better understand why some news frames are more or less effective at influencing people's emotions.

This call to engage the reader emotionally in news, and the related call for audience engagement has driven many scholars to study the role of specific emotions in driving readers' engagement with different news stories on social media, with their findings suggesting that emotion plays a fundamental role in shaping news diffusion on social media. Bright (2016, p. 348) wrote extensively on the role of negative emotions in news diffusion, and discussed the need for the audience to make sense of disasters and feel for the victims thus sharing sad news. However, news consumption is evolving and increasingly occurs via social media platforms, where affective appeals related to the enjoyment of news enhance sharing on social media. The importance of enjoyment in the narrative structure in journalism is considered by Knobloch et al. (2004, p. 282) who claim that the inverted-pyramid format is not entertaining for the readers and that publishers should consider alternatives based on emotions. In addition, Dafonte-Gómez (2018) found a strong link between pleasant emotions and sharing of the news on social media and highlighted the importance of arousal either positive or negative. In other words, the more arousal in a story the more viral it will

be. D. K. Brown et al. (2020) researched the Ice Bucket Challenge on social media, a trend for supporting patients with amyotrophic lateral sclerosis back in 2014, and found that the intention to share a story that involved emotions was much greater than the intent to share a story with the news values of prominence like a celebrity's involvement or human-interest.

A number of studies show that online news has a significant impact on the formation of political attitudes and engagement in political activities (Hasell and Weeks 2016), by using intense feelings in journalism and over-dramatization (Robertson and Mourão 2020). According to Hasell (2021), anger, anxiety, and enthusiasm trigger more political information-sharing on social media that not only increases direct political participation (e.g., voting) but also indirectly influences people's attitudes and behavior resulting in greater political polarization among voters and the general public, hostile behaviors and distrust toward politicians and public officials (Garrett and Stroud 2014; Levendusky 2013). Although, Wojcieszak et al. (2016)] believe that passionate political discourse may have a positive effect on political engagement and lead to meaningful discussions. After exploring the sentiment of 140K tweets from 44 media outlets in the US using the Vader lexicon (Hutto and E. Gilbert 2014), Bellovary et al. (2021) discovered that negative emotions are the ones that invite engagement and the ones that motivate people to share political stories because they feel angry or hatred towards others (Sánchez Laws 2020, p. 11) and that is true for both news outlets with left- and right-leaning orientations.

Despite the fact that emotions are increasingly important for news consumption, researchers also point out that an overreliance upon emotional stories would reduce decision-making critical thinking from readers, while the inclusion of emotions in news content is often related to disinformation (Vosoughi et al. 2018) or hate speech in the comments sections. In fact, researchers documented a number of 'dark participation' incidents on news comment sections that reflect upsetting, unkind or negative expressions (Quandt 2018). Research has also been conducted that shows how disinformation actors use emotional language to achieve virality (Hsu et al. 2020; Vosoughi et al. 2018). Furthermore, specific emotions were used to detect fake news. In particular, Vosoughi et al. (2018) research on the role of emotion on Twitter has shown that fake information evokes fear, disgust and surprise in

the reply section, in contrast with non-hoax stories that activate emotions, including joy, sadness, trust, and anticipation.

The way in which emotions are embedded in journalistic writing is important because of the impact on the audience's engagement. Emotional writing does not only make content more relatable, but it also stimulates audience engagement. This can be seen in the cases of the *New York Times*, *ESPN*, and *USA Today*, which created advertising products to leverage emotionality for profit (Rick 2019). Except for being a factor of audience engagement and a revenue booster for the news industry, emotions are very important in social and humanistic studies. Many analytical instruments such as lexicons and linguistic models have been created in order to capture emotions in writing. In addition, researchers are also working on improving tools in order to detect emotions from video, images, and speech.

One of the most important and widely used tools in the journalistic field is the SemEval project in 2007 (Strapparava and Mihalcea 2007), aiming to measure the affectivity expressed in the text. More specifically, the developed task was to use English news Headlines and seek to determine whether they bring emotional responses to the audience. Similarly, Mohammad, based on established psychological models, such as Ekman (1999), and Plutchik (1980b), used human annotators to create the NRC EMOLEX, a set of words that signify basic emotions (anger, joy, fear, sadness, disgust, surprise, trust, anticipation) and build several lexicons to extract emotional properties from text such as joy and anger with a given extract (Saif M Mohammad 2017; S. Mohammad and Turney 2010). Furthermore, pleasure (valence), arousal and dominance are also used by researchers to explain what motivates users to engage with the news. Specifically, the model was introduced by Russell (1980) as the Valence-Arousal-Dominance (VAD). The model that has been coded into an emotional lexicon (S. Mohammad 2018) has three dimensions valence (pleasure), which goes from feeling happy to unhappy, arousal which measures how active someone feels going from sleep to excitement, and dominance which calculates how much in control a person feels ranging from submissive to dominant.

Other lexicons that are based on the annotator's input are the NRC Affect Intensity Lexi-

con (Saif M Mohammad 2017), ANEW (Bradley and Lang 1999), MPQA Subjectivity Lexicon (Wilson et al. 2005), VADER (Hutto and E. Gilbert 2014), the norms lexicon by Warriner et al. (2013), TextBlob (Loria 2018) and more. Moreover, relevant research aims to automatically develop computational tools for the identification and extraction of emotions in the written text, such as SentiWordNet (SWN) (Esuli and Sebastiani 2006), and the LIWC (Tausczik and Pennebaker 2010).

3.2.3 Emotions for Disinformation Detection

The intentional spread of false and concocted information serves many purposes such as financial and political interests, influencing public discourse against marginalized populations, has a negative impact on society and democracy (Marwick A. and Weigel 2021; Shu et al. 2017), and can expose the public to immediate danger. Examples of false stories that went viral on social media platforms like the “Pizzagate”, a conspiracy theory that threatened the lives of the employees of a pizzeria (Shu et al. 2020) and coronavirus-related false content that led people to drink toxic chemicals with at least 800 people dead and thousands hospitalized ¹, show that online virality can become dangerous. More specifically, previous research has found that social bots are crucial in the spread of misinformation (Shao et al. 2017) since search engines, social media platforms, and news aggregators use algorithms that control the information a user sees. For instance, algorithmic curation on Google can promote a greatly visited news article very high on the search results, thus improving the likelihood of it being shared, read, and emailed. Audience metrics such as page views, likes, shares, and so on, unquestionably influence the number of people who see a given article on their screen. Therefore, experts in disinformation and online radicalization take advantage of these known algorithmic vulnerabilities by creating fabricated accounts which generate fake traffic that results in virality (Shao et al. 2017). Virality in turn guarantees that disinformation, trolling rumors, and coordinated campaigns are rapidly propagated across the internet, and as Lotan (2014) highlights what we need is “algorithms that optimize for an informed public, rather than page views and traffic”. Nevertheless, after much debate about

¹<https://www.bbc.com/news/world-53755067>

the need for Facebook to change its algorithm to reduce filter bubbles, and the platform's avoidance of taking responsibility for the distribution of deceptive content on its News Feed, since mid-December 2016 it started to alter its algorithm to make misleading information to appear lower and Google followed with raising the fact-checked stories higher (Bakir and McStay 2018). However, the Covid-19 pandemic proved those measures were insufficient, while also highlighting the challenges that journalists face as they need to manually check countless requests of potentially deceptive information daily², without sometimes possessing the necessary skills, or having the resources, time, and expert personnel to fight disinformation (Bakir and McStay 2018).

The urgent need for disinformation detection led many scientific disciplines in the search for new effective ways to mitigate this problem with promising approaches coming from various fields. In line with this, this paper proposes a computational approach to detect potentially fake information, by identifying textual and nontextual characteristics of both fake and real news articles and then using machine learning algorithms for disinformation prediction. More precisely, we consider two sets of machine-readable features i) content-based, and ii) engagement-based, and we conduct our analysis in two distinct phases. In phase A, only content-based features are explored, while in phase B we add features that correspond to the users' interactions on Facebook and test them on a subset of the original fake and real news dataset.

3.2.4 Disinformation

Fake and manipulated information is circulated in all forms and platforms, unverified videos are shared on Facebook, rumors are being forwarded via messaging apps, while conspiracy theories are being shared by Twitter influencers, and these are only a few of the distribution patterns of disinformation. According to Tandoc Jr et al. (2018) the role of social media platforms is crucial to understand the current state of disinformation globally since Facebook and Twitter changed both the news distribution and the trust to traditional media outlets. As they vividly note “now, a tweet, which at most is 140 characters long, is considered

²<https://www.poynter.org/coronavirusfactsalliance/>

a piece of news, particularly if it comes from a person in authority” (Tandoc Jr et al. 2018). In this work, we consider real news as defined by Kovach and Rosenstiel (2014) to be “independent, reliable, accurate, and comprehensive information”, and “not include unverified facts”, thus disinformation campaigns threaten to curtail the actual purpose of journalism, which is “to provide citizens with the information they need to be free and self-governing” (Kovach and Rosenstiel 2014). In addition, to define fake news we use the description by the European Commission (Commission 2018) “disinformation is understood as verifiably false or misleading information that is created, presented, and disseminated for economic gain or to intentionally deceive the public, and may cause public harm”. Journalists and professional fact-checkers can determine the correctness of potential threats based on their expertise and the use of many digital tools designed to detect a plethora of manipulated elements inside a fake story. Finally, news verification can be a procedure done inside a news outlet that checks all the information before publication or it can be done after the piece is published or shared in social media networks.

The rise of disinformation has attracted strong interest from computer scientists who employ machine learning and other automated methods to help identify disinformation. Fake news detection in computer science is defined as the task of classifying news by its veracity (Olivieri et al. 2019) with many studies of this phenomenon aiming to extract useful linguistic and other types of features and then build effective models that can identify and predict fake news from real content. A useful overview of the computational methods used for automated disinformation detection (Conroy et al. 2015) separates two categories, notably machine learning research using linguistic cues, and network analysis using behavioral data. In this section, we will focus only on previous work around the former category, linguistic approaches.

The thought behind linguistic approaches for fake news detection based on content is to find deception elements which can lead to distinguishing the fakeness of news (Rubin et al. 2016). Rubin et al. (2016) built a Support Vector Machine (SVM) model to identify satire and humor articles. Their model performed with 87% accuracy and the results showed that the best predictive features were absurdity, and the use of grammar and punctuation. A sim-

ilar study by Horne and Adali (2017) compared real news against satire articles using also SVM with an accuracy of 91%, and found that headlines, complexity, and style of content are good predictors of satire news. However, when classifying real and fake news the accuracy dropped dramatically. Ahmed et al. (2017) experimented with n-grams and examined different feature extraction methods and multiple machine learning models, to find the best algorithm to classify disinformation. The results showed that overall linear-based classifiers are better than nonlinear ones, with the highest accuracy achieved by a Linear SVM. Furthermore, Shu et al. (2017) conducted a survey providing a comprehensive review of fake news detection on social media. They discussed existing fake news detection approaches from a data mining perspective, including feature extraction, model construction, and evaluation metrics.

The fast-paced creation and dissemination of disinformation have resulted in a new wave of data-driven tools developed to help filter and analyze this information with the potential to become key in detecting fake news. Although the identification of fake news and its creators is a complicated task; False information spreads more quickly than it is debunked by real news outlets, it is copied and pasted across platforms and to multiple websites, and often relies on fake accounts, purchased followers, and non-human entities to create the illusion of popularity and public trust. Coupled with artificial intelligence (deep fakes), the spread of social media, and the rapid information on the internet, disinformation is really difficult to distinguish from actual events or sources. Tools are needed to help filter and detect false information, but they can only be created with knowledge about why people consume and share fake news. The manipulation of emotions in news stories has been related to disinformation (Vosoughi et al. 2018), with many studies testing emotion lexicons in the detection of fake news.

The use of emotional language in fake news is a still a promising topic in the world of computational techniques, with many different studies testing how emotions such as anger and fear can be used to detect fake news. These studies infer emotions from the frequency of certain words and from the potential combinations of emotion-related words. The level of emotionality in a text is controlled by the presence or absence of certain words as well as

by their order and the frequency with which they are used. Research in the emotional language in fake news compared with non-fake news using an LSTM machine learning model showed that hoaxes contained more sympathetic words like “beautiful”, “friendly”, “soft-hearted” thus obtaining the reader’s attention. Additionally, articles related to propaganda elicited joy, fear, and calmness in an attempt to manipulate the public’s opinion while providing some sort of confidence (Ghanem et al. 2020). Using again deep learning models Ghanem et al. (2021) found that the existence of fear, sadness, and surprise emotions at the beginning of the article was important for the model to predict fake news. Towards the end of the article, on the other hand, they noticed that not only does negativity not exist, but joy and anticipation take its place. Another study by Vosoughi et al. (2018) found that fake articles usually provoked fear, disgust and surprise in the comments, a fact that was later proved by Giachanou et al. (2019) with the use of an LSTM model.

For the fake news corpuses, many researchers use ready-to-use datasets, such as BuzzFeedNews³, BuzzFace⁴, BS Detector⁵, CREDBANK⁶ and FacebookHoax⁷ Shu et al. (2020) and others construct their own using potentially false stories from websites marked as fake news by PolitiFact (Asubiaro and Rubin 2018; W. Y. Wang 2017). W. Y. Wang (2017) introduced LIAR, a benchmark dataset for fake news detection about politics created from manually labeled reports from Politifact.com. In this work, the authors used a Convolutional Neural Network and showed that the combination of meta-data with text improves disinformation detection. Asubiaro and Rubin (2018) downloaded fabricated articles from websites marked as fake news sources by PolitiFact.com and matched them with real news around the same political topics. Their computational content analysis showed that false political news articles tend to have fewer words and paragraphs than the real ones although the fabricated stories have lengthier paragraphs and include more profanity and affectivity. Finally, the titles of the fake stories are bigger and more emotional, including more punctuation marks, demonstratives, and fewer verifiable facts.

³<https://github.com/BuzzFeedNews/2016-10-facebook-fact-check/tree/master/data>

⁴<https://github.com/gstantia/BuzzFace>

⁵<https://github.com/bs-detector/bs-detector>

⁶<http://compsocial.github.io/CREDBANK-data/>

⁷<https://github.com/gabll/some-like-it-hoax>

Several studies related to fake news detection examined social media aiming to extract useful features and build effective models that can differentiate potentially fabricated stories over truthful news. The study of Tacchini et al. (2017) focused on whether a hoax post can be identified based on how many people “liked” it on Facebook. Using two different classification techniques, which both provided a performance of 99% accuracy, the research proved that hoax posts have, on average, more likes than non-hoax posts, indicating that the users’ interactions on news posts on social media platforms can be used to predict whether posts are hoaxes. Similarly, the study of Idrees et al. (2019) showed that the users’ reactions to Facebook news-related posts are an important factor for determining if they are fake or not. The authors proposed a model based on both users’ comments and expressed emotions (emoji) and suggested that a future Support Vector Machine approach would increase its accuracy. Finally, the work of Reis et al. (2019) examined features such as language use and source reliability, while also examining the social network structure. The authors studied the degree of users’ engagement and the temporal patterns and evaluated the discriminative power of the features using several classifiers with the best results obtained by a Random Forest and an XGBoost which both had an F1 score of 81%.

Chapter 4

The Role of Images in Journalism

4.1 The Importance of Visual Content

Journalism requires the use of images as they aid in grabbing the audience's attention and drawing the readers into the narrative. According to research presented in this chapter, social media users are more willing to share news articles with images than those without, because images can make a story more aesthetically appealing and easy to understand. Moreover, images can assist in illustrating the problems and communicating the emotion of a story. This is crucial in stories about humanitarian crises since "iconic" visuals can raise public awareness of the suffering of people impacted. However, it is crucial to make sure that the pictures utilized are respectful, follow journalistic guidelines, and do not take advantage of the subjects.

It is widely acknowledged that social media engagement is crucial to the success of any brand irrespective of the industry. Yet this means that in order to increase online audience loyalty, digital news outlets, newspapers, magazines, and publishing groups, in general, must pay close attention to visuals. However, there have long been concerns about how to use photos on social media to draw attention and boost audience interaction. And there is no clear answer: Is it the quality of the images, or the use of them, that generates interest?

Research on social media engagement mostly focuses on the written text of a news-related post, although one of the most important aspects of generating engagement with any news story is the use of relevant and compelling images. With the everyday use of mobile devices for news consumption and social media like Instagram, users appreciate more than ever the importance of using images to tell a story. Particularly, for social media users, the image is not only a way to describe something (although this is also a key part), but also a way to communicate with other people, for instance by sending a gif, a meme, an emoji, or a sticker (Highfield and Leaver 2016). Therefore, users do not only embrace images but create their own with smartphones. Inevitably, the proliferation of “selfie” images and snapshots on social media has led to the development of a culture that embraces visual content and supports it (Frosh 2015). The fact that images are more popular than text guides publishers in their decision to use more images in news stories. Studies reveal that the use of images in any social media post has shown a steady trend of increasing engagement over the last years (Stepaniuk 2015), with images receiving the most engagement of all types of posts on Facebook for 2021 ¹. It is important to note that while a majority of the research focuses on how images relate to authorship and text, there are ongoing discussions surrounding the idea that content and images are wholly separate entities (J. C. Alexander 2010). Some scholars believe that visuals have become the norm over text-based news articles, and their study is under-researched (Pearce et al. 2020; Highfield and Leaver 2016).

4.1.1 Pictures as the Main Way of Expression

Beginning in the 1990s, visual practices and visual analyses were taken seriously and played a significant role in academia as well as in society. Thus, the so-called “pictorial”, “visual”, or “iconic turn” became not only a theoretical paradigm but also a methodological practice in everyday life, with two main theorists Gottfried Boehm and W. J. T. Mitchell. This turn addresses the idea that images have become the main means of expression in our time (W. T. Mitchell 2005). Boehm explains that the shift from a predominantly verbal to a visual culture is part of a process of replacing talking about something to showing it (Boehm and W. J.

¹<https://locowise.com/blog/digital-2021-global-report-what-can-we-learn>

Mitchell 2009). While according to W. T. Mitchell (2005) and W. T. Mitchell (1995), this shift was oriented toward a symbiosis of images and text, recognising both the power of the word as well as that of the image. The last decade, cultural sociologists have been exploring the notion of an image existing by itself, separate from the discourse that is produced around them (J. C. Alexander 2010). In his research on images from art and history to online news, J. Alexander et al. (2012) focuses on how the image's existence can be understood in terms of iconic rituals and humanitarian visual icons.

Whereas traditional media had been predominantly text-based while becoming increasingly image-based as a result of developments such as television, satellite, and internet technology, modern media have changed their structure significantly toward an entanglement of image and text. News photographs that “tell a story” and get readers to feel more connected to the protagonists are widely deemed to be “iconic” in their power to rally public interest in a decisive manner (Hariman and Lucaites 2007) and due to their wide visibility and near “sacred” quality, they are also referred to as “secular” icons (J. Alexander et al. 2012; Maynard 1983). A lot of studies focused on the topic of immigration and particularly on the photographs of the dead body of Alan Kurdi on the shores of Turkey in 2015 (Mortensen et al. 2017; Devichand 2016; Vis and Goriunova 2015), that gained momentum in an unimaginable way and reached more than 20 million online users the first 12 hours it was published online (Vis and Goriunova 2015). The photos of the dead child changed the public discourse over the topic of immigration, created awareness about the Syrian war and drove people from all over the world to volunteer and donate to refugee agencies (Devichand 2016). The study of Mortensen et al. (2017) investigated how those editorial processing of an iconic yet disturbing set of news photographs was communicated. The researchers examining the Alan Kurdi case analyzed articles from Denmark, Canada and the UK news outlets, and highlighted that the narrative structure used by journalists to describe the images was “self-fulfilling”, meaning that they proclaimed them as iconic and newsworthy since the first day of publication based on their virality on social media.

Domke et al. (2002) study explored the potential of visual media to shape public opinion. Specifically, they conducted an experiment in which news coverage was altered by the in-

clusion of either influential or uninfluential images. The results of the study suggested that the inclusion of images had an effect on the way readers comprehended the information. Furthermore, the study pointed to the possibility that the media can use visual images to influence public opinion in subtle ways. This research has implications for understanding how the media can manipulate public opinion. In other words, it suggests that visual elements can be used to influence how audiences interpret news stories. This highlights the importance of visual media in contemporary journalism and reinforces the need for careful consideration of the images that accompany news stories.

4.2 Social Media and the Power of the Image

With the world's population sharing stories on Whatsapp, Facebook, Twitter and Instagram at an ever-increasing rate, it follows that news has also come to be delivered through a lens of visual imagery.

The increased use of smartphones further emphasizes the importance of images in journalism. Nowadays, it is very common for online readers to consume news on their mobile devices while they are engaged in another activity. As a result, in today's fast-paced environment, images can be more very effective in conveying information instead of lengthy articles. Additionally, it is especially important to note that with the avalanche of content on social media, a post needs to stand out from the rest in order to get readers' attention. Thus, the image should be carefully selected to create interest (Malthouse et al. 2013). News images are increasingly elemental in providing a visual scaffolding for news stories that online readers can access easily on their phones, whether directly or from a feed. The ease with which information can be understood and consumed on platforms like Instagram or Tik Tok, allows news stories to be easily interpreted with little need for readers to stop and digest complicated arguments or text-heavy articles.

The use of visual components in news articles is not new, but the current global media landscape provides outlets with a far greater possibility to harness user-generated content and for individuals to develop their own material and share it with other viewers via social media.

The internet provides media professionals with a greater range of visual content to choose from (or, at the very least, the capacity to find aesthetically intriguing content more easily), as well as the ability to curate and modify that content as they see appropriate. Along with the vast pool of visual elements to select from, visualisation is also playing a growing role in transmitting the information. When coupled with the more overt emphasis on visual images, the unquestionably positive role of imagery in journalism comes into focus, since images are an effective way of grabbing readers' attention and for driving engagement.

Research has focused on how people use pictures on social media for self-branding, communication, narration, and group identity (Mendelson and Papacharissi 2010). While these practices are used differently by different people, Van Dijck (2008, p. 57) argues that cameras serve as the preferred devices for expression of the online users with pictures acting as words that are disseminated and circulated to create bonds. More toward emotions, Palmer (2010) noted that platforms based on photographs serve as "emotional archives" for users and their online friends. One promising project researching visual content with computational methods is the "Selfiecity"² (Tifentale and Manovich 2015), which aims to draw insights related to facial expressions, pose and demographics depending on the city the selfie is taken, and offers vibrant infographics, picture maps, and explanations.

As mentioned in Chapter 3, M. Gibbs et al. (2015) refer to platform vernaculars to explain the interactions and the "way of understanding how communication practices emerge" in each social media network with images to play a decisive role in the "iconic turn" of those platforms. Similarly, Pearce et al. (2020) attempts to describe visual vernaculars by researching the use of images about climate change across five different social media platforms, namely Instagram, Twitter, Reddit, Tumblr, and Facebook. By using a multi-platform approach instead of just focusing on one platform the author finds the best practices for impactful image use on social media. Highfield and Leaver (2016) argue that multimodal analyses in social media are difficult compared to text because to analyze a picture one needs to intervene in several stages of the process. Additionally, Bossio (2021) chose to study how journalists are utilizing the visual culture of Instagram. The research was conducted through visual anal-

²<https://selfiecity.net/>

ysis and semi-structured interviews with journalists and the results showed that journalists tend to leverage Instagram's "popular" social practices in order to produce and distribute news, but also to generate "traffic" to the site.

4.2.1 What Makes an Image Newsworthy and Engaging?

In the 1980s, the visual nature of modern mass media prompted the reevaluation of traditional visual media across programs of art, communication, and journalism. This led some journalism schools to recast their photojournalism and publication graphics tracks into multidisciplinary visual communication curricula. Communication scholars argued that media like newspapers should be viewed as an inherently visual phenomenon (Barnhurst 1994). Despite this, most journalism and mass communication programs tended to regard photojournalism, publication graphics, and video production as technical support for written reports. However, some schools did begin to integrate visual communication concerns into their professional and research curricula, including concepts from visual anthropology, the history of film and photography, and so on (Griffin 2001).

Our brains process information faster than words (Thorpe et al. 1996), just in 13 milliseconds, as a team of neuroscientists from MIT proved (Potter et al. 2014), while pictures are easier to remember than words (Paivio 1991; Snodgrass et al. 1974). An essential part of journalistic quality is visual information which can contribute to a variety of aspects such as the interpretation of information. Vivid images also can provoke excitement, surprise, and intrigue in the human brain. Also, visualizations, are an effective way to convey information, since they can help illustrate concepts that are otherwise difficult to explain. As we have seen in chapter one by using charts and eye-catching graphics, data journalists can convey complex ideas in an easily understood format. For example, a chart can show how a number such as the unemployment rate has changed over time. Charts and graphs can help people see patterns or trends that might not be obvious without visual representation (Henshall and Ingram 1991). Additionally, Skyword³ research has revealed that articles with images garner 94% more views than those without. Thus, the choice of photo is of paramount importance

³<https://www.skyword.com/resources/>

to the post and article's success. Furthermore, Wojdyski (2015) found that the presence of interactive graphics in articles increased the interest of users who had no prior knowledge or interest in the subject but did not affect the intention to engage for users who had prior knowledge or interest. These results are consistent with the model of Liu and Shrum (2002) according to which interactivity encourages online engaging behaviors in users of low or non-existent prior interest.

The role of photos in journalism cannot be over-emphasized. Not only can photographs contribute to a better understanding of the event they accompany, but they also convey emotion and change attitudes. As discussed in the chapter on "Emotionality" images help readers become more engaged in the story and keep them focused on the most important points. Studies show that when an image accompanies a news story, readers are more likely to pay attention, remember details, and care about what they have read (Griffin 2001), while visual media have been known to maintain collective memory (Zelizer 1998). Moreover, Lagun and Lalmas (2016) investigated the reading patterns of users by analysing screen viewing data from Yahoo news, and identified four reading profiles (Bounce, Shallow, Deep, Complete) based on the amount of time spent on each part of a website. They found that the reading behaviour was closely related to the topic and keywords of an article, and that users often spend most of their time at the top of an article and then scroll back to the top before leaving. Additionally, they discovered that the presence of a photo or a video encouraged users to stay at the beginning of an article for twice as long.

In light of these factors, it is essential to delve deeper into the factors that define a newsworthy photograph in the context of photojournalism. More specifically, photojournalism is a form of journalism focused on the visual representation of events, places, and people, and as with any form of journalism, there are certain criteria that must be met to ensure the quality of the work. Pictures can sometimes tell the news just by themselves, with a caption to say who the people are and where the event is taking place. At other times, the picture may go with a story, to work as a team with the words. In either case, a news picture must always leave the reader knowing more than they did before. Although many online news stories use stock images or even irrelevant ones to attract the reader's attention, the ultimate purpose of

a picture is to carry information. Therefore, one of the criteria of a good picture is accuracy. A photojournalist must accurately represent the scenes and people they are photographing Henshall and Ingram (1991).

A good picture is also a matter of composition. A photojournalist must be able to capture a scene in an aesthetically pleasing way (Thomson and Greenwood 2017). It means having an understanding of design concepts like balance, contrast, and symmetry as well as the ability to apply the principles of light and shadow to produce a captivating image. The resolution of the photo is another crucial consideration because photos with low-quality or hard-to-distinguish objects receive fewer clicks and likes Tenk (2021) and Ding et al. (2019). However, sometimes publishers use user-generated content in cases of breaking news like the bombing in 2017 at the Manchester Arena⁴. Additionally, the editor highlights that the more striking the image and the stronger the message, the more likely readers will share the article. Finally, images that are surprising are also very shareable. Moreover, Lai (2011) investigated the influence of images and the media's choice of emotionally charged shots, particularly in light of the new technology that has enabled users to become potential photojournalists. Results from the study suggested that snapshots taken by amateur photographers are perceived as more genuine and reliable by viewers than those taken by professionals, with the credibility of the content being more valued than the quality of the photography.

Similar to news values a picture should be newsworthy. A photojournalist needs to be able to take a picture that the general audience will find interesting and relevant. This requires keeping up with current affairs as well as being able to spot a scene or moment that will engage viewers. Also, it is through the people that photographs can tell the story Henshall and Ingram (1991), therefore a picture must involve people, show the protagonists, their actions, what happened to them, and so on.

Also, for a photojournalist to tell the story, they must provide context to the scene, for example, a picture of a salesperson talking to a customer in front of a car is more relevant than sitting at their desk. Some theorists overemphasize images and underestimate the role

⁴<https://www.aljazeera.com/news/2017/5/23/witnesses-panic-after-deadly-manchester-arena-blast>

of captions and other linguistic devices that guide the reader on how to “read” the image. For instance, Henshall and Ingram (1991) state that a picture can often tell more about an event than words ever could. A photo can show what a place is like with its architecture, the colors or shapes of its houses, where the area is located, and how it changes over time. Similarly, the findings of Isola et al. (2013) provide valuable insights into the nature of human memory and how it affects our perception of the world. Their work highlights the fact that photographs featuring individuals tend to be more memorable than pictures of nature. Furthermore, their findings imply that images with high semantic qualities—such as people, interiors, foregrounds, and human-scale objects are more likely to be remembered than those with low semantic features. This has interesting implications for the field of visual communication and for designers who create visuals to communicate a message. By understanding the power of human memory, designers can create visuals that are more likely to be remembered and to effectively communicate their desired message.

Undeniably, an essential criterion of photojournalism is emotion. The photograph should capture a moment that is emotionally resonant and powerful. This includes being aware of the emotions of the subjects in the photograph, as well as being able to convey a message or story through the image. As the Reuters Pictures editor, Rickey Rogers said to Tenk (2021) “What makes a photograph successful is its ability to provoke an emotion or a reaction in the viewer”. Nor is this surprising: A picture with a face on it is capable of conveying far more emotion than any words (Bakhshi et al. 2014). A reporter can use words to tell the audience how something made someone upset or angry – but this emotion cannot simply be added to the story without any context. However, through a photojournalist’s lens, an audience can see the person experiencing this emotion. That is why photojournalism has been crucial to show the truth and reality of war (F. Moore 2018). Nevertheless, when it comes to journalistic articles, the use of emotional images can be seen as a lack of objectivity (Orgeret 2020).

Editors use pictures in their articles also to provide the readers with a more positive experience. For example, readers are likely to look at the pictures first as they scroll over an article, or do it simultaneously when scrolling through a digital multimedia story. Planer et al. (2022) analyzed journalistic quality in multimedia long form journalism from the World

Press Photo contest, like the Pulitzer Prize winner the Snow Fall: The Avalanche at Tunnel Creek by the New York Times (Branch 2012). The analysis was based on the quality criteria of multimediality, usability, linking, interactivity, gamification and so on, aiming at discovering the most important characteristics of multimedia narratives. The authors established that the most important features were multimediality and continuous text, which meant that photographs and video along with graphics, sound and animation need a flow of text to support them and make the storytelling more dramatic.

From an advertising perspective, Messaris (1992) explores the power of visual media to influence, and convince viewers, by simulating reality, showcasing evidence, and indirectly proposing solutions e.g. to sell a product by highlighting their uniqueness, beauty or effectiveness. Barry (1997) examined the influence of images shown in films, politics, and television, and stressed the danger of the emerging violence shown in the media. Furthermore, nowadays more than ever the observation of Griffin (2001) that images act as a “transnational cultural currency in the modern global media environment” is true, with Western standards being propagated worldwide via social media. Due to the prevalence of visual media in modern western culture, recent studies have shown that visual media affects choices regarding which information people choose to consume.

While images are a fundamental part of news, they are also part of everyday life or as Mirzoeff (1999, p. 1) argues “it is everyday life”. Consumers rely on images to experience the world and social media platforms provide the most efficient way of sharing these images. Media companies are aware of this new trend of communicating on social media through pictures. Thus, the use of images has become more powerful than the written word and news organizations now have to choose how to present stories with both strong visual and textual content. Additionally, news organisations are required to engage with visual elements on their websites or in mobile apps in order to attract a wide range of readers required for a meaningful return, therefore knowing what audiences want would help publishers to preemptively assign resources in accordance with the potential needs of future users (Tatar et al. 2014). Although, this requires that media professionals understand how to effectively exploit this new environment while also ensuring they meet the basic tenets of journalism, and that re-

porting is supposed to adhere to ethical and professional standards making the media's visual orientation a matter of concern. This is a matter of concern in the wake of the proliferation of news images and journalists' often less than scrupulous ways of deploying them. The way in which news photos are taken - or staged - from the users on social media can raise doubts about their validity and the way in which they are manipulated (either during or after taking) can undermine journalistic integrity when used without verification. On top of that, photographs must not be used as clickbait, should originate from reputable sources, and be presented in a way that they do not confuse and mislead the readers. Due to the popularity of social media and the habit of users sharing news articles with their peers, individuals are now exposed to images in a way that is very different from how they were in the past.

The topic of how visual media affect how we opt-in to participate in journalism is currently being explored, but the implications on content production are significant. If visual media influence our engagement with news then it would be advantageous for media professionals to consider a more effective strategy for storytelling, by considering the "power of images" and the impact that images have on our online behavior. Previous studies have suggested that the content of images is important in how people interact with news and have evaluated the influence of posts' visual content on engagement. The study by Bakhshi et al. (2014) on Instagram suggests that pictures with human faces get significantly more likes and comments, while the use of warm filters or color correction and contrast also enhances engagement. A digital marketing brand (Lowry 2013) computationally analyzed 8M pictures on Instagram and found that using light images, blue as dominant color, muted palette, and single dominant color, denote higher "likes". Similarly, Jaakonmäki et al. (2017) noted higher engagement in pictures with people or scenery.

The outcomes hierarchy model proposed by Lavidge and Steiner (1961) consists of three distinct stages in the purchase process: knowledge, emotion, and action. Consequently, when a user engages with an advertisement, they first form an opinion in the knowledge stage and, based on this initial opinion, will decide whether they like or dislike the product in the emotional stage before finally making a purchase in the action stage. Social media can

be seen as the first contact point between the consumer and the product or brand, thereby determining the consumer's subsequent actions. Thus, the marketing strategies of news outlets should pay attention to how they appear as a brand, particularly when considering the growing use of social media.

As we have discussed, news brands use a variety of tactics to influence their audience through social media, and one of the most powerful is the use of images, which can draw attention to stories and evoke a range of emotions, from shock and awe to sympathy and empathy. For news outlets to optimize social media use, develop an effective content strategy, and build an engaged community they first need to know their audience. Thomson and Greenwood (2017) used news posts on Instagram to determine the different types of Instagram users. The researchers categorized users into three groups, namely the “feature lovers”, who showed a preference for attractive and uplifting images featuring people, sunsets, and instances of travel, and a clear preference for still photos over video content. The “News Hounds”, were users that exhibited a distinct preference towards photographs pertaining to tragedy, politics, and global culture, and found watermarks to impede their levels of engagement. Their favorite kind of shots were wide-angle, context-rich group images. The last group was the “The Optimists”, who had a preference for content featuring visible facial features, and avoided images of violence. People in this group were mainly interested in engaging with inspiring, motivational, and encouraging material. In terms of general observations, users liked easy-to-understand images, such as portraits, and did not trust brand watermarks, perceiving the UGC more accurate.

4.2.2 Brand Related Images

Numerous scholars have studied the impact of images in advertising and consider them to be a fundamental component of any effective marketing strategy. Images can be utilized to build meaningful connections between an item or service and the desired outcome. Studies using eye-tracking devices have shown that the use of relevant images can increase viewers' engagement with an advertisement, making them more likely to remember the message and potentially make a purchase more than text Pieters and Wedel (2004). Research done

by (Goldfarb and Tucker 2011) also demonstrated that although consumers found intrusive banners annoying in fact they did remember the brand name and the message attached to it, even if they did not like the placement or feel it was in their face. Furthermore, the content of an advertisement often relies on the target audience's emotional responses more than objective messages, because emotional content is more likely to earn their attention and consideration (S. G. Moore and McFerran 2017; Berger and Milkman 2012). In addition to creating strong emotional connections, images can also be used to communicate complex messages in a single glance. By combining different visual elements, advertisers can communicate ideas that may be difficult to put into words. Through images, advertisers can create a sense of urgency or a feeling of exclusivity that can be difficult to achieve with words alone.

Also, images have the power to influence the viewers' perception of a product or service. By carefully selecting relevant images, advertisers can influence viewers to view a product or service in a certain light. Finn (1988) found that colorfulness increases a reader's attention, whereas Xiao and M. Ding (2014) revealed that photographs of people turned out to be more persuasive in print advertisements. On the other hand, Shunyuan Zhang et al. (2017) suggest that the quality and composition of an image are of utmost significance when it comes to product sales, such selecting a place to stay from Airbnb properties. Building on these findings, the study of Yiyi Li and Xie (2020) explores the impacts of picture characteristics such as colorfulness, the presence of human faces and facial expressions, and picture quality on user engagement in relation to attributes of the written text such as sentiment, topic, emojis and so on. The study used almost 19K advertising posts on Instagram and Twitter, deep learning models for image recognition namely the Google Cloud Vision API⁵ and zero-inflated negative binomial regression for the analysis. The results showed that consistently professional and high-quality pictures drove more engagement in both platforms, while attributes like color or human faces are dependent on the context, for instance, colorful images from airlines had higher engagement on Twitter in contrast with lower for SUV cars, and humans or relevant text accompanying the picture are good for Twitter advertisements

⁵<https://cloud.google.com/vision>

but not for Instagram.

4.2.3 Machine Learning Approaches

Machine learning is a popular approach for researching engagement with visual content on social media. Studies have been conducted on Twitter (Cappallo et al. 2015) and other image-based platforms like Flickr (Qian et al. 2017; McParlane et al. 2014; Khosla et al. 2014) or Instagram (Purba et al. 2021; Z. Zhang et al. 2018; Zohourian et al. 2018; Qian et al. 2017; Mazloom et al. 2016; Bakhshi et al. 2014). Researchers have developed both regression and classification models on Instagram pictures with high accuracy (Zohourian et al. 2018). Additionally, Support Vector Regression (SVR) has been used to show that quality, time of publication, and image type are the most informative features for predicting engagement (Purba et al. 2021). A similar approach (Khosla et al. 2014), incorporating the relationship between visual content and user engagement trained a machine learning model on 2.3M pictures from the picture based platform, Flickr, and concluded that content like human faces or busy pictures were the most engaging. However, the behavior of the user and their online volume of followers are also decisive of the performance of images on the platform.

Most of the aforementioned studies used computer vision technology either toward a content analysis approach where they examined the category of the picture e.g., human face, selfie, brand logo, the objects in it, car, boy etc., the scenery indoor, outdoor, the composition of the image, contrast, quality and so on. Other studies focused on context like time (McParlane et al. 2013) examining time-related trends, others like Dhar et al. (2011) proposed a crowdsourced model for evaluation of image aesthetics or mixed visual characteristics with textual, temporal and social in multi modal analyses (McParlane et al. 2014). Specifically, the study by Mazloom et al. (2016) utilized a SVR model that considers features like logos, sentiment, face, and product, and showed that the brand logo, emotions and filters are the best predictors for Instagram advertising posts. Another study McParlane et al. (2014) showed that more comments with contents of images (on Flickr was depended on the user's activity and network and the associated tags. The study achieved 76% accuracy by using a framework of 16 features namely textual, contextual, user and content features,

and highlights that the best predictions were the result of the combination of all the features together pointing to the conclusion that they are complementary to each other.

Chapter 5

Methodology

5.1 Computational Social Science in Journalism Studies

Computational social science (CSS) was introduced in 2009 Lazer et al. (2009), and represents a new frontier in the social sciences with its roots in traditional social sciences. CSS is a paradigm shift in social science that requires rigorous theorization, research design, statistical analysis, and results presentation (Peng et al. 2019). It has been actively adopted by a growing number of scholars in various fields as well as widely embraced by several influential journals¹. It has emerged as an interdisciplinary field because of its need to marry concepts from disparate fields such as computer science, sociology, economics, physiology, anthropology, and political science; it has the potential to become the optimal tool for private companies and government agencies (Lazer et al. 2009). The ongoing advancements in technology have created a multitude of opportunities for the field of social science to evolve, while the internet and technological advancements have made the retrieval of data for social science analysis much more efficient. Digital traces and user-generated content from social media, mobile phones, and the overall digitalization of most communications media, has led to a boom in data that are available for use. With the increasing availability of digital data, social scientists are able to access and analyse information in ways that were previously not possible (Peng et al. 2019). This access to data related to human interaction

¹<https://www.tandfonline.com/toc/hcms20/12/2-3>

has enabled the emergence of Computational Social Science, as a multidisciplinary domain that uses data processing techniques to explore the structure (see Figure 5.1, functions and behaviour of society (J. Zhang et al. 2020; Wallach 2018). These advances, combined with the efforts of computational social scientists, have provided a more detailed understanding of the complexity of social phenomena (Conte et al. 2012).

The application of computational methods in the social sciences has been steadily increasing in recent years. The use of automated methods has enabled scientists to conduct larger and more complex research projects, while also reducing the amount of time and effort needed to analyze data, while minimizing human error (Cioffi-Revilla 2014). Additionally, the development of more efficient algorithms has enabled researchers to uncover more intricate relationships between different variables and make more accurate predictions (Alpaydin 2020). As a result, computational methods have become an essential tool for social scientists and have had a profound impact on the way research is conducted in the field. This research highlighted the importance of focusing on the mental health of vulnerable populations during economic downturns. Moreover, CSS is used to measure the impact of public policies and governmental regulations on people's behavior (J. Zhang et al. 2020). By providing a more comprehensive view of the human society, computational social science can inform the decision-making process and potentially lead to more effective policies.

The emergence of computational methods has also provided social scientists with the ability to access and analyze large volumes of data from a variety of sources, such as online databases and social media platforms (Lazer et al. 2009). This is particularly useful in fields such as sociology, where the ability to collect and interpret data from large populations is essential for understanding social dynamics and behaviors. Additionally, the use of these automated algorithms has enabled researchers to make more informed decisions by taking into account data points that would otherwise be overlooked by manual analysis (Cioffi-Revilla 2014).

Moreover, computational social science (CSS) amplifies and enhances traditional analyses, while providing insights into the ethical implications of automated decisions (Wallach

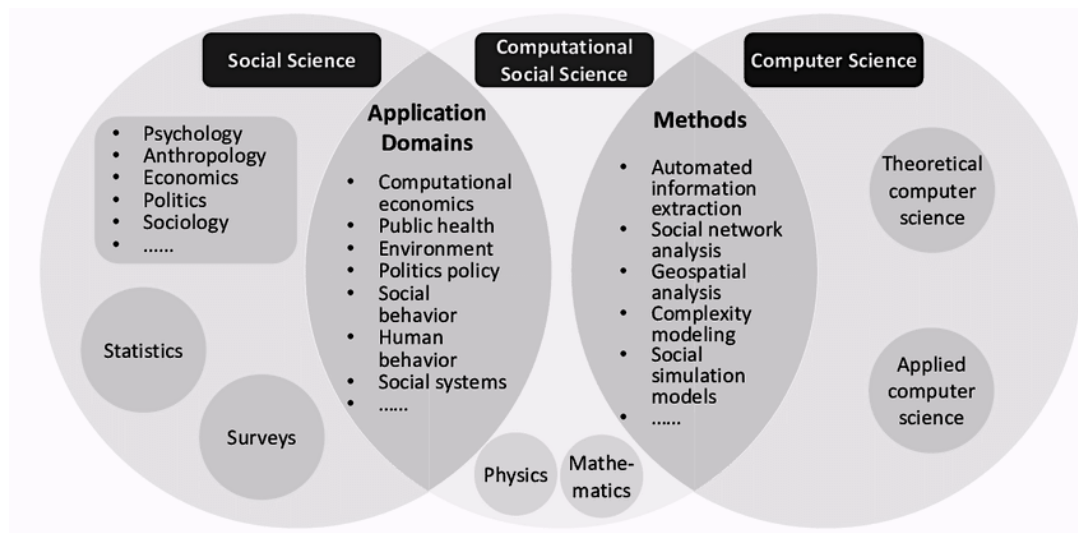


Figure 5.1: Interdisciplinary nature of CSS. Graphic by J. Zhang et al. (2020)

2018). Cioffi-Revilla (2014, p. 2) provides a more comprehensive definition, describing the field as “the interdisciplinary investigation of the social universe on many scales, ranging from individual actors to the largest groupings, through the medium of computation”. In this definition, the term “many scales” refers to many organizational, temporal, and spatial dimensions, while “computation” refers to a range of computer-based methodologies. The use of increasingly complicated data-driven techniques in CSS research to investigate the complexity of social systems includes network analysis, machine learning, and natural language processing.

The emergence of big data has led to the development of new methods of analysis in social sciences that necessitate the acquirement of new skills and investments on infrastructure (Atteveldt and Peng 2018; Jagadish et al. 2014). Large textual corpora and big data sets require more complex techniques than traditional statistics, thus require the use of programming languages such as Python and R, as well as computational applications to content analysis, network and topic models, and social simulations. De Mauro et al. (2016, p. 131) defines big data by taking into consideration its essential features, suggesting that it is the information asset characterized “by such a High Volume, Velocity and Variety to require specific Technology and Analytical Methods for its transformation into Value”. The terms “Volume”, “Velocity” and “Variety” describe the main features of information, “Technology” and “Analytical Methods” demonstrate the preconditions of handling the information and

the term “Value” summarizes the extraction of useful knowledge, exploited by societies and companies.

The use of big data has had a profound impact on many different scientific fields, such as economics, physical sciences, anthropology, sociology, retail, manufacturing, communication and journalism studies (Peng et al. 2019; Wallach 2018; Hesse et al. 2015; Jagadish et al. 2014). Big data has enabled researchers to gain deeper insights into these fields by providing them with more data points to analyze and draw conclusions from. Furthermore, the process of gathering, modeling, and analyzing massive data to generate insights has also led to the use of sophisticated and complex methods of analysis, such as machine learning and artificial intelligence (Atteveldt and Peng 2018). By leveraging these methods and technologies, social researchers can uncover hidden patterns in data, allowing them to gain a deeper understanding of the phenomena they are studying. However, it is important to note that computational social science research should be approached in a responsible and transparent manner, taking into consideration the potential implications of the results (Hesse et al. 2015). With that in mind, it is essential that computational social scientists develop a strong understanding of ethical principles when using people’s personal data.

The translation of big data to knowledge is important to the social and behavioural sciences in terms of the shift to data-intensive science, and the application of principles from the social sciences to the conduct of research. According to Hesse et al. (2015) the scientific enterprise needs to recalibrate to become more transparent, cumulative, integrative, cohesive, rapid, relevant and responsive.

5.1.1 Knowledge Discovery in Databases

Knowledge discovery in databases (KDD) is an interdisciplinary field of research that focuses on the development of methods and techniques to extract knowledge from large volumes of data and has been used in Computational Social Science extensively (Claudio Cioffi-Revilla 2017). Due to the availability of big datasets and the need to analyze and make sense of them, KDD has grown in significance in recent years. In order to turn raw data into valu-

able and actionable information that can be used to support decision-making, it employs a number of methodologies to find patterns, correlations, and trends in data.

KDD involves data preprocessing and data mining and consists of five stages (Figure 5.2), namely data selection, cleaning and preprocessing, transformation, data mining, and accurate interpretation of the mining results, which are fundamental to ensure that essential information is derived from the data (Han et al. 2011; Fayyad et al. 1996). The first step is to define the objective of the process and clarify the application domain and then use the research questions to specify data selection. After the data is selected the next stage is data cleaning and preprocessing followed by data reduction and feature selection. Precisely, the data preprocessing phase involves transforming raw data into a form that is suitable for data mining. This stage is composed of dimensionality reduction and other transformation methods such as managing missing values, and encoding categorical variables in order to detect important features. Afterward one of the more important steps is data mining, which is the “process that consists of applying data analysis and discovery algorithms that produce a particular enumeration of patterns (or models) over the data (Fayyad et al. 1996). It is conducted by using selected algorithms and methods such as classification rules or trees, regression, and clustering. The last stage is the interpretation and explanation of the mined patterns, which may be accompanied by visualizations and statistical tests before the final knowledge is produced by the whole process. According to Cioffi-Revilla (2014) the communication of the results is the final stage including any impacts identified.

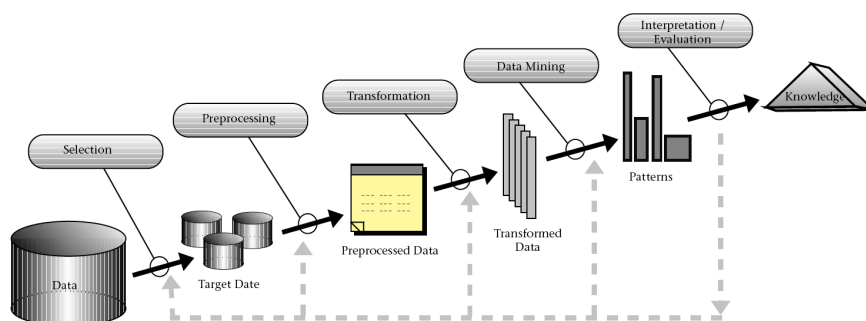


Figure 5.2: Overview of KDD Process, from the Fayyad et al. (1996)

KDD is an important tool for many business and scientific applications. KDD is a tool for

numerous commercial and academic uses, such as to identify customer categories, spot fraud, and predict customer behavior. It can also be used to reveal hidden relationships and patterns in data. KDD additionally has the potential to optimize data-driven decision-making and boost the precision of machine learning algorithms. In many industries, including health, marketing, finance, and computational linguistics, KDD is an important method for extracting knowledge from data. Its applications involve entity extraction from text, document categorization, document clustering, topic modeling, and more. Furthermore, KDD has become important in data mining algorithm design due to its well-known advantages. In order to obtain reliable results, the process should be conducted in an organized and structured way following the five stages previously mentioned. Therefore, the area of KDD is still challenging users because of its long time consuming effort, large amount of data, and often complex algorithms. The results obtained from KDD research have helped many organizations focus their resources on an effective data-driven decision process that could be used to optimize operational efficiency and promote business growth.

5.1.2 Cross Industry Standard Process for Data Mining

The Cross Industry Standard Process for Data Mining (CRISP-DM) is a framework developed to provide a structured approach to data mining projects. It was designed by Shearer (2000) to facilitate the efficient and effective analysis of large datasets to identify meaningful patterns and relationships. The CRISP-DM process, which can be seen in (Figure 5.3) consists of six separate phases: Business Understanding, Data Understanding, Data Preparation, Modeling, Evaluation, and Deployment. It offers a comprehensive set of principles for the end to end creation of data mining projects. Each phase contains particular tasks that must be finished in order for the project to be successful.

The Business Understanding phase involves defining the objectives of the project, identifying the data sources, and determining the success criteria. The Data Understanding phase involves exploring and visualizing the data, as well as identifying any potential problems and issues. The Data Preparation phase involves cleaning the data, transforming it, and selecting the appropriate data for modeling. The Modeling phase involves selecting and

applying the appropriate algorithms to build the model. The Evaluation phase involves assessing the model's performance, and the Deployment phase involves deploying the model and assessing its performance in the real world. By following the CRISP-DM process, data scientists can ensure that the project is executed properly and efficiently, and that the results are meaningful and actionable.

CRISP-DM provides a structured and systematic approach to data mining projects and remains the most widely used framework for executing data science projects (Saltz 2022) in many domains such as health (Rivo et al. 2012). Based on the CRISP-DM Guide²,

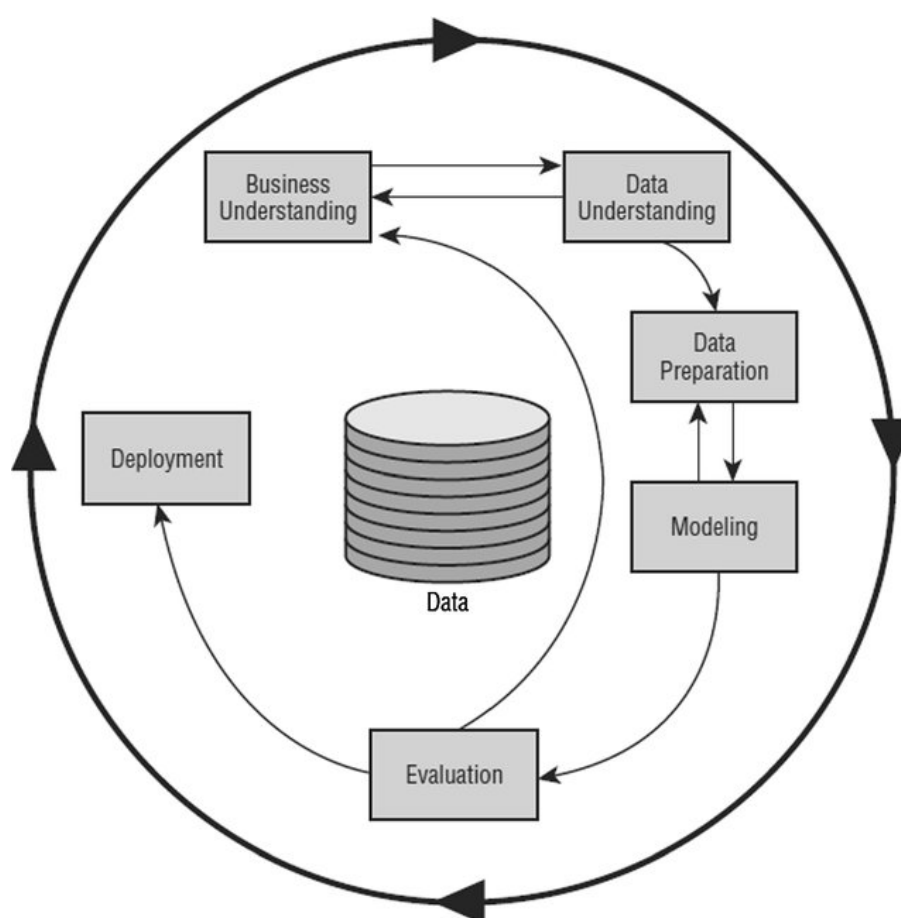


Figure 5.3: The life cycle of a data mining project, from the Chapman et al. (2000)

The order of phases in data mining is not fixed. As a result, it requires transitioning between different phases, and the result of each phase will direct where the next phase or task should be. Figure 5.3 illustrates the relationship between phases, as well as the cyclical nature of data mining, depicted by the outer circle. Even after a solution is delivered, data mining does

²<https://web.archive.org/web/20220401041957/https://www.the-modeling-agency.com/crisp-dm.pdf>

not end. The process can be continuous, as the newly acquired knowledge and deployed solution can create new business questions. Prior experiences can aid in the development of subsequent data mining processes. To delve deeper into the life cycle of a data mining project based on the CRISP-DM reference model (Chapman et al. 2000), a brief analysis of the six distinct phases follows.

Phase 1: Business Understanding: In the first phase, the goal is to understand the business objectives of the data mining project and to define business questions that can be answered by data analysis. This phase also involves understanding the context of the project, such as the available data sources, the target audience, the timeline for completion, and any constraints or risks associated with the project.

Phase 2: Data Understanding: After the goals of the project have been established, the next step is to understand the data that will be used for the analysis. This includes exploring the data to get an understanding of its structure, content, and quality. This phase also involves identifying any data gaps and determining how to fill them, spotting initial insights and intriguing subsets to form hypotheses about not (yet) obvious information.

Phase 3: Data Preparation: This phase involves preparing the data for analysis. This includes cleaning the data, transforming it into a usable format, dealing with any missing or incorrect values, and selecting the features.

Phase 4: Modeling: Modeling is the process of constructing models to address the questions that have been identified. Various modeling techniques are chosen to be implemented, and their parameters are adjusted to maximize their efficiency. There are usually multiple approaches for any given data mining problem. However, some methods necessitate certain data formats, so it might be necessary to return to the data preparation step.

Phase 5: Evaluation: In this phase, the models are evaluated to determine if they are achieving the desired results. This includes validating the models and assessing their accuracy, performance, and robustness regarding new data that include outliers and so on.

Phase 6: Deployment: The final phase of the CRISP-DM model is deployment. Even if the

model is made to gain knowledge of the data, the knowledge needs to be organized and delivered in a way that is useful to the customer. Depending on the requirements, the deployment phase can be as simple as generating a report or as complex as implementing an ongoing data mining process across an enterprise.

In conclusion, decision-makers must be given the capability to interpret and verify the results of analysis. It is not enough to provide just the results: one must also provide users with the ability to repeat the analysis with different assumptions, parameters, or datasets to support the human thought process and social circumstances (Jagadish et al. 2014).

5.1.3 Data Mining Methods in CSS

Computational social science involves a wide range of methodologies that can be used to answer specific research questions. The topics vary from simple computational uses to complex statistical analysis and are designed to work with numerical, textual, visual and audio data. As Tao et al. (2020) suggests CSS has the potential of becoming a continuous learning process when social phenomena are examined by data-driven methods such as machine learning. A large body of research in computational social science is focusing on three applications of CSS in the domain of communication research, namely researching user analytics, text mining, and experiment research (Peng et al. 2019). User analytics is the practice of understanding and analyzing behavioral patterns on social media platforms. This includes tracking a user's response to social stimuli such as posts and images, as well as uncovering psychological and social influences that shape user behavior. Also, in order to effectively analyze the large-scale and unstructured user-generated content on social media, text mining techniques must be employed, allowing for the discovery of underlying patterns and relations. Finally, moving away from the traditional self-report surveys, researchers can now measure the real behavioral changes of subjects (Atteveldt and Peng 2018), such as their interaction in response to experimental stimuli on social media, and design online experiments in a more dynamic and flexible way. Compared to offline experiments which are often done with University students, researchers can easily adjust the experimental treatments or add new ones during the experiments. Besides, online experiments are usually more eco-

nomical and practical than offline experiments (Salganik 2019), due to social media and crowdsourced platforms such as Amazon Mechanical Turk (Golder and Macy 2014).

In the Computational Social Science field the most extensively used machine learning methods are categorization, also known as classification, and clustering (Cioffi-Revilla 2014). Cluster analysis is an inductive form of unsupervised machine learning, which aims to group the data with common features. Clustering dimensions are extracted in numerous ways, such as dendrograms which present a hierarchical structure (Cioffi-Revilla 2014; Han et al. 2011). In classification, the goal is to predict for a new object the correct category (group) to which belongs, while in clustering, the model aims to predict for a set of objects (potentially belonging to different categories) their categories or groupings based on their similarity. It can be used to find patterns in data, to identify groups of objects that are similar to each other, and to detect outliers. Classification is a supervised machine learning form with the object of creating an array of categorized information as output. The data mining algorithms that are used for the information extraction are called classifiers, such as K-nearest neighbor and naive Bayes classifier, which will be explained in detail in the next section. The selection of the training data and the parameters as well as highlighting the feature importance are some of the human interventions in the process.

Additionally, a very commonly used model is the Bag-of-Words (BoW), which is simple to implement, fast, and efficient. This algorithm is a method for extracting features from text. The basic idea is to represent each text document as a vector of word counts, where each word in the document gets a count. The algorithm starts by tokenizing the text into words and building a vocabulary of all the unique words. Then, for each document, it counts the occurrences of each word in the vocabulary and creates a word count vector. This text representation method is beneficial for a variety of natural language processing problems, including text classification, clustering, and topic modeling. This form of analysis is considered semi-supervised learning as a manual annotation in training data is required (Cioffi-Revilla 2014).

Furthermore, natural language processing (NLP) is used for text analysis purposes and is

very commonly utilized in semantic analysis. NLP algorithms have the ability to extract key phrases, words, name entities, topics, and other important attributes from a text, allowing for a more effective way of exploring a document's content. NLP may also be used to identify word associations like antonyms, synonyms, and collocations. Furthermore, semantic analysis is a powerful technique for determining the meaning of a sentence and its components, whereas sentiment analysis can be used to identify the emotional context of textual data. This is accomplished by extracting the subjective information from the data, which can be used to infer the sentiment of the text (Chen 2018). Furthermore, similarity analysis can be employed to detect similarities between two written pieces. For instance, this can be used to compare and contrast the source material and help to identify any potential bias, or identify the author of a piece of writing (Cioffi-Revilla 2014). Finally, machine learning algorithms can be utilized to detect a text's sentiment and contribute to increasing sentiment analysis accuracy. Researchers can acquire a more comprehensive knowledge of the content of a corpus by combining the results of lexical and NLP analysis.

Oakes (2019) reported that statistics are employed in natural language processing for “detective work, especially in the context of large corpora. Specifically, one application of statistics in corpus linguistics is the analysis of writing style, with much of the effort being directed towards automated methods for “training” computers to identify certain features, such as sentence length, word placement, choice and frequency, syntactic analysis, and more complex statistical measures of style. This approach enables the attribution of texts to authors, the comparison of styles of different or the same authors, the establishment of the chronology of writings and many more applications.

Other applications of data mining algorithms include geographic data, network and sequence analysis, and anomaly detection. Chen (2018) suggests that geographic data can also be used in CSS for applications like satellite remote sensing. Spatial analysis encompasses techniques such as geocoding and geographic clustering aiming to reveal insights on spatial patterns. Such patterns are data about the social world, disasters, migration or social movements. Thus, spatial technology has enriched the meaning of the terms place and space (Chen 2018). Furthermore, network analysis is increasingly being used in the field of

Computational Social Science (CSS) to explore the various connections between actors (Tao et al. 2020). Through the visualization of networks consisting of nodes and relationships with a range of attributes, CSS researchers are able to gain insights into cognitive belief systems, political Twitter accounts, and online disinformation patterns, which could not be gleaned from more traditional methods or direct observation.

Sequence analysis seeks to offer insights into a given process considering the phase transitions. For instance, temporal patterns can be analyzed in opinion data, political events data, and financial data (Chen 2018). In intensity analysis, the extracted features of source data can be used as input for conducting subsequent analysis, as with complexity-theoretic models. For example, this form of analysis provides information to researchers about the risk of extreme events, the fragility of unstable conditions, or the early-warning indicators of impending abrupt change (Claudio Cioffi-Revilla 2017). In addition, anomaly detection analysis, as a data mining methodology, seeks to identify anomalies or modifications in data through the prerequisite that data present a significant degree of stability (Cioffi-Revilla 2014). This stability is measured through statistical measurements such as dispersion or central tendency, which should not exhibit drastic changes during the test stage. Conversely, sonification analysis focuses on sound as a source of discovering new knowledge.

5.1.4 Machine Learning Algorithms

Machine learning has grown as a field in computer science and has become a prominent research topic for academia and a new trend for various industries. More specifically, ML is a subfield of AI in which models learn from data in order to discover patterns and correlations, with the aim to predict, extract usable information, make decisions, and categorize data. Scholars use ML to automate many aspects of their research and obtain insights while in the industry the same models are used for example in predicting the stock market, analyzing consumer behavior and creating self-driving cars. Additionally, machine learning has become a valuable tool for newsrooms as discussed in Chapter 2.

In the CSS, machine learning algorithms are a popular approach for data analysis and pre-

diction. This method has a wide range of possible applications in social science research, from extracting insights from massive datasets to automating data processing and even forecasting future trends. The analysis of massive datasets was the first use of machine learning in the social sciences. Researchers are able to discover insights and patterns that would not be obvious using typical data analysis approaches by applying state-of-the-art machine learning algorithms. ML for example, may be used to detect audience segmentation or psychological trends in online user behavior. Another application is to automate data analysis and create estimates for future trends. For example, political scientists can apply machine learning to forecast election results or the impact of policy changes. This can assist policy-makers understand the potential consequences of their policies.

The distinction between qualitative and quantitative variables is critical for the application of machine learning methods. Qualitative, or categorical variables, are those that describe data that fit into specific classes, such as gender, eye color, level of education or profession. Quantitative variables, on the other hand, are measured on a numeric scale, for example temperature, prices, or a person's age, height, weight. Depending on the response, machine learning methods are then applied accordingly. Regression problems are used to address those with a quantitative response, while those with a qualitative response are defined as classification problems (James et al. 2013). Although logistic and multinomial regression models are also used for categorical response data (Liang et al. 2020; Zhou et al. 2019).

There are three main types of machine learning: supervised learning, unsupervised learning, and reinforcement learning. Supervised learning is using a known training dataset to train an algorithm using labeled input data, known as features, to generate predictions. The objective of the supervised learning technique is to identify correlations between the features and the output data by looking at similarities and discrepancies (Marconi 2020) to accurately predict the correct label. On the other hand, unsupervised learning does not have any labeled data given to the learning algorithm, leaving it to explore and find patterns and knowledge on its own. Lastly, reinforcement learning is based on feedback data. This means that the training data is supplied in the form of rewards and punishments (J. Zhang et al. 2020).

Supervised learning uses a training dataset in which the labels are “known” and the algorithm is trained to predict the right output label. The model is not free to create its own labels like in the case of unsupervised learning. Supervised models are widely used in classification and regression problems, such as recognizing objects in photographs, categorizing emails as spam or not, and forecasting market values. After training, the model is tested on unseen data to evaluate its accuracy. This process is also known as supervised learning because the algorithm is being “supervised” by the labeled data. Supervised learning can be used for both classification and regression problems. In a classification problem, the algorithm is trained to generate a mapping between the inputs and the output labels that are categorical, such as “spam” or “not spam”. For regression tasks, the algorithm is trained to provide a mapping between inputs and continuous outputs, such as developing predictive models for house prices or type 2 diabetes. Also, a supervised learning model is very dependent on the data on which it is trained, meaning that if the data is incorrectly labeled, biased, or imbalanced, the model will be unable to provide reliable predictions. So it is important to train the models on high-quality data. Furthermore, in supervised learning the algorithms’ results can be evaluated and investigated in depth. The Decision Tree Classifier, Support Vector Machines, and Logistic Regression are examples of supervised algorithms.

Unsupervised learning works with unlabeled data and searches for previously unknown patterns in a dataset without the use of labels or predefined categories. The way it draws conclusions is by using statistical approaches and in contrast to supervised machine learning, does not require any prior “known” information or labels to function. Instead, it uses algorithms to discover patterns, trends, and relationships in data that are not immediately obvious. It can be used to improve the understanding of complex data sets, including those with many features, or to identify clusters of similar data points. These types of models are usually applied to problems like anomaly detection, clustering, and automated visualization creation. It can also be used to make predictions about unseen data like text generation where it predicts the successive word in a given sentence. Unsupervised learning offers several advantages, including the capacity to make predictions without depending on labeled data, making it excellent for situations where labeled data is unavailable or difficult to get.

However, there are several disadvantages, such as being heavily reliant on data quality, and the findings might be difficult to comprehend and evaluate.

Reinforcement Learning works with agents acting in an environment to maximize a reward without depending on explicit instructions. More specifically, it identifies the optimal strategies to be applied via trial and error. The agent then learns from this experience and continues to take new actions in order to maximize the reward. The goal is that over time, the agent learns an optimal behavior, which is the combination of actions that maximizes the reward. From robotics to gaming, this approach has been utilized to tackle a wide range of challenges.

Deep learning is a famous subfield of ML, because of its potential to bring change to a variety of industries, including robotics, language translation, text generation, image recognition, and more. In short, deep learning uses big data to be trained and employs a collection of algorithms to replicate the activity of neurons in the human brain. The way this approach works is by giving a collection of data to the model and letting it learn how to recognize complex patterns and make precise predictions. Furthermore, deep learning algorithms use the errors that occur and correct them using a process known as “backpropagation” in order to get better outcomes. The neural network may then extrapolate the patterns and predict the outcomes of new data.

Moreover, deep learning has opened up new possibilities in many fields such as autonomous driving, machine translation, and text generation. Examples of easy-to-use pre-trained deep learning models include BERT and OPEN AI, which have been applied to a variety of NLP tasks, such as sentiment analysis, question-answering, and text summarization. Additionally, great progress has been achieved in the area of computer vision, and specifically in recognizing and categorizing objects in photographs. Convolutional neural networks (CNNs) are also commonly used, for instance, in Google Vision and Amazon Rekognition to identify objects in images by extracting characteristics from the input data. Similarly, these methods are being applied to medical imaging to find anomalies in medical images such as breast cancer.

5.2 Research Questions

This dissertation examines the concept of journalistic quality in digital journalism using quantitative and computational methodologies drawn from communication studies and computational social science. The first overarching research question is: How can AI provide a nuanced understanding of what constitutes quality journalism in the digital age? To answer this question, two further research questions need to be posed. Can subjectivity, emotionality, entertainment, and quality of language predict journalistic quality in online news? And what is the exact contribution of the quality criteria to quality prediction? For instance, one hypothesis is that lower subjectivity predicts high-quality news stories while emotional coverage contributes to stories being of low quality. Additionally, research questions are posed concerning audience engagement, such as: Can a machine learning model accurately predict social media engagement of news? And what images are significant for predicting the engagement or the quality of a news story? Finally, a subsequent question is: Is the quality predictive model able to detect online disinformation?

Therefore, the main argument put forward in this work is that computational social science can be used for answering those questions and help journalists to better understand how to reach and retain their audience and produce high-quality journalism that resonates with readers. The findings of this work, which used computational approaches to investigate the concept of journalistic quality in the digital era, as well as the characteristics that lead to audience engagement, may be of interest to journalists, academics, and news consumers.

5.3 Research Design

In order to take advantage of computational methods used in data science and artificial intelligence and benefit from them, the social science researchers must invest in skills (Wal-lach 2018). To make use of digital traces and “big” amounts of textual data as Shah et al. (2015) point out, scientists for one, need to develop advanced natural language processing skills. Atteveldt and Peng (2018) analyzed the benefits, challenges and drawbacks of using CSS in Communication studies, and concluded that sometimes the solution is to work to-

gether with computer scientists, computational linguists, and data analysts, but in any case the presence of a skilled programmer is necessary in order to collect, clean, analyze, and visualize data. However not all researchers have the resources to hire a programmer and the ability to supervise and motivate them. Therefore, as the researchers underlined in their study, to get the most out of computational methods, many young scholars invest in learning programming languages and how to use a variety of tools to perform data mining projects.

For this work the Python programming language was used and the CRISP-DM reference model (Chapman et al. 2000) was followed as a data mining framework for all the different studies presented in the following sections. As Jagadish et al. (2014) highlights data analysis pipelines are composed of more than just the analysis/modeling phase, with each stage of the data analysis life circle to pose certain challenges. Moreover, the data used for all the studies are news articles, therefore a special focus will be on techniques to handle and analyze textual data.

5.3.1 Business Understanding

Beginning with the first phase, for every study first the relevant scholarship is reviewed, the key gaps are identified and then the research questions are formed. Afterward, the theoretical framework is designed based on the theory and the research approach is discussed. Also, the technical aspects of the studies must be addressed here.

To start, a user-friendly and interactive environment was needed, since Python scripts are not easy for keeping notes, maintaining track of progress or collaborating and reproducing the results due to different installations and versions on one's computer. To address these issues, the Jupyter Notebook was used, which is a free, open source application that enables users to create and collaborate on interactive computational documents and is widely used by the scientific community (Pimentel et al. 2021). Also, designed to support reproducible research, Jupyter Notebook allows users to easily reproduce their results, as the environment stores all of their scripts. Thus users can create computational stories that combine code, results, and text and can be easily shared. Furthermore, this environment provides a graph-

ical interface to the Python language, and enables users to run their scripts in a secure and controlled environment.

The first step in using Jupyter notebook for Python is to install the necessary packages. The most popular packages for data mining tasks are Numpy, Pandas, and Scikit-learn. It is also important to install the Jupyter notebook package itself. Once these packages were installed on a server of the Faculty of Communication and Media Studies, the environment for the research was ready.

Data Acquisition

Web scraping is a technique used to extract data from websites. It involves using computer software or a program to request information from a website, and then extracting the information from the HTML of the website. It can be used to extract data from a single website or multiple websites, depending on the needs of the research. The data extracted can be structured or unstructured, for instance, an Application Programming Interface (API), which is a program that allows communication between two applications provides well-structured information, whereas a scraping program gathers raw information from the internet often containing HTML code that require a lot of cleaning, filtering and organizing. An example of unstructured textual data from scraping the *Sun* news outlet is shown in Figure 5.4.

Web scraping can be done either manually or with the aid of automated tools. Manual web scraping requires the researcher to write a program to request information from the website and then parse the HTML in order to extract the desired data. This method demands programming skills and is time-consuming, particularly when dealing with intricate and large websites. Automated web scraping tools and crawlers also exist, but they can only extract specific data formats such as URLs, or download information from straightforward sites, such as downloading tables from a Wikipedia page. For websites that contain intrusive banners, complex architecture, or scripts that prohibit web scraping, these tools do not work. In addition to the internet and online databases, web scraping can also be used to acquire data from text files, images, and videos.

author	Title	Noimages	body
By JILL ROBINSON	Flack cashes in despite all the flak	8	Fans slag off 'awkward' X Factor host Caroline...
By ELLIE FLYNN	India's biggest baby has been born – weighing ...	7	AN ENORMOUS 14.77lb baby boy has been born t...
By ED DYSON, Showbiz Reporter	Cowell says sorry for mocking Malik	7	X Factor boss Simon joked: 'Who misses Zayn?' ...
EXCLUSIVE by RUTH WARRANDER	WAGGRO	7	Charlie Adam's missus in Twitter spat with pai...
thesun	Wine advert banned for naughty 'taste the bush...	5	NaN
...
EXCLUSIVE by LEIGH HOLMWOOD	A happy Ender as Nat gets engaged	5	EXCLUSIVE: Soap star and cameraman engaged t...
thesun	Are aliens trying to communicate to us via thi...	5	NaN

Figure 5.4: Example of unstructured data from the *Sun*.

For the studies of this dissertation, all the data have been gathered through manual written code in Python programming language, except for the Facebook data that were directly downloaded from the CrowdTangle³ platform. For gathering news articles from websites like *the Sun*, or the *Washington Post* the Python library Beautiful Soup⁴ was implemented.

Beautiful Soup is a Python library designed for developers to create web scraping applications. It parses HTML and XML documents, and provides a wide range of methods to extract data from web pages, process HTML forms, and automate web interactions. With that library a researcher can also use XPath queries and CSS selectors, in order to quickly locate and extract specific data from a web page by looking at its source code.

For dynamic web pages like the Medium.com which is more complex and also has a pay-wall, a subscription was paid and the Selenium⁵ package was used to create an automated program that visits the website, fills out search forms, clicks buttons and links and navigates through different pages. Specifically, Selenium is a free and open source library for Python that is used for web scraping and automated testing of web applications since it can sim-

³<https://www.crowdtangle.com/>

⁴<https://beautiful-soup-4.readthedocs.io/en/latest/>

⁵<https://selenium-python.readthedocs.io/>

ulate user interactions. It provides features such as browser emulation, and supports for multiple browsers.

Finally, for news organizations that provide a developer's API, like the *New York Times* and the *Guardian* (see <https://developer.nytimes.com/> and <https://open-platform.theguardian.com/>), an application was built to access the news organization's content and transform the single format (JSON) file to a pandas dataframe to match the rest of the data. In these cases the data was already structured, but extra data, such as the images of the news stories could not be obtained, therefore the freedom of web scraping is limited.

5.3.2 Data Understanding

In this stage, the goal is to explore the data and get a better understanding of it in order to prepare the data for analysis. This stage involves examining if the data has the correct format for the analysis. When multiple formats of valuable information exist in the same column, it should be transformed into a form suitable for the models. Therefore, the right strategies must be designed and deployed to convert the multiple formats into a structured dataset.

5.3.3 Data Preparation

The process of extracting significant data and changing it from one format to another is known as data transformation. It is an essential aspect of the data pre-processing stage as it converts the data into a format suited for the target assignment. For instance, this method can be used in the medical domain to transform data from random texts about a patient into a structured format that can be used to train AI models for disease diagnosis (Ford et al. 2020). Other examples of data transformation are normalizing data, and transforming it into a format that can be understood by a model, for example change a variable from categorical "Yes", "No" to numerical 1, 0. The data preparation step also includes cleaning the data from noise, missing values, irrelevant information, and inconsistencies. It is necessary to check if the data is valid and correct and if needed, replace it.

In the case of analyzing textual data, namely news stories, preprocessing is an essential part of Natural Language Processing. It is the process that prepares the text for further analysis, such as topic modeling, sentiment analysis, and text classification. The following preprocessing steps were performed on all the datasets:

1. Removing punctuation from the text.
2. Converting all characters to lowercase.
3. Removing stop words that have little meaning (e.g. conjunctions, prepositions, etc.).
4. Converting dates to machine usable format.
5. Tokenizing the text into individual words.
6. Stemming or lemmatizing the words, when needed. For example some lexicons need the original form of a word to work properly.
7. Part-of-Speech Tagging. This step assigns a part-of-speech to each word in the article.
8. Named Entity Recognition: This step identifies entities in the text, such as people, organizations, and locations.

Feature Engineering

Feature engineering is the process of using domain knowledge about the data to create features for machine learning algorithms. In other words, it is a process of transforming raw data into features that better represent the underlying problem to the predictive models, resulting in improved model accuracy on unseen data. This process involves selecting, creating, and transforming features from data to make them suitable for modeling. This is a key step in machine learning, as it can have a significant impact on the performance of the models.

Various dictionaries are used in this phase as well. They are collections of words or phrases with associated numerical values that represent contextual sentiment, emotion or other characteristics of the text. They can be used to analyze text for sentiment (positive, neutral,

negative) or emotions (such as fear, anger, joy, and sadness), or to identify topics, entities and other words that belong in a list (e.g. animals, politicians, or celebrities). Other libraries that are commonly used are the readability and semantic complexity metrics (Flesch Reading Ease, Measure of Textual Lexical Diversity, Gunning Fog Index). After the features are created, certain methods for dimensionality reduction, like dendrograms and matrices can be used to find correlated variables and reduce the number of features. In a dendrogram, features are connected depending on their similarity; the more similar two features are, the sooner they are connected in the dendrogram and thus one feature can represent a set of similar other features. More information on the lexicons and Python libraries used for feature engineering can be found in Appedix C.

5.3.4 Modeling

The modeling phase consists of developing, testing, and improving a model in order to discover the optimal solution to a given problem. This usually entails choosing and fine-tuning a suitable algorithm, such as a decision tree or neural network, and then performing a number of trials to determine which model performs best. Numerous methods have been proposed in the literature to construct text classifiers, such as decision trees, naive-Bayes, neural networks, nearest neighbors, and support vector machines (Ikonomakis et al. 2005).

5.3.5 Supervised Machine Learning Algorithms

A classifier is a function that is used to assign a label or category to a given set of data points. Specifically, it uses the values of the features (independent variables) to predict the class that a data point belongs to (the dependent variable) (Pereira et al. 2009). In supervised machine learning, where the data is labeled, the classifier is trained to recognize the pattern in the data and predict the correct label for new data points. In a journalistic context, the emotion of anger, subjective expression, and readability of a news article could be used as features to predict whether it is real or disinformation. The class in this case would be the label of fake or true.

For example, an article is indicated by the row vector $\mathbf{x} = [x_1, \dots, x_v]$ and its corresponding class label is denoted as y . The task of a classifier is to learn the parameters from a set of training examples that have been reserved for this purpose. After the parameters have been learned, the classifier essentially forms a trained model that is able to predict the label for a given example x it has not seen before. In other words, given an example x , the classifier is a function f that predicts the label $\hat{y} = f(x)$ (Pereira et al. 2009).

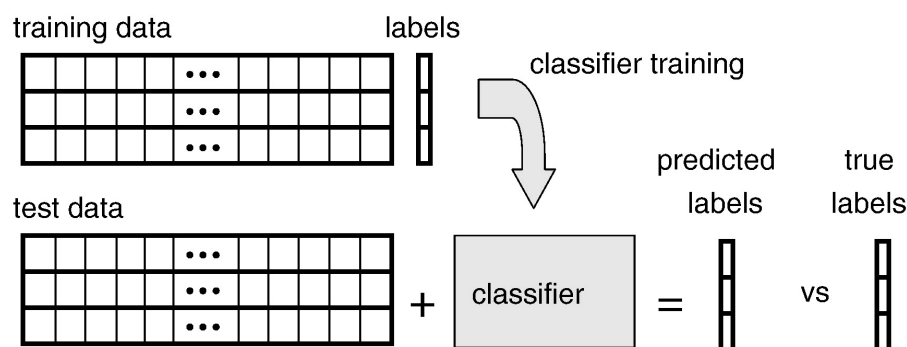


Figure 5.5: A classifier is trained by using examples that are labeled, and then used to predict labels for a test set that has not been seen before. The predicted labels are then compared with the true labels and the accuracy of the classifier can be determined by finding out the fraction of examples that the prediction was correct. Graphic by Pereira et al. (2009)

After being trained, the classifier can be utilized to see if the features used provide info about the class of the example. This connection is checked by utilizing the trained classifier on a distinct set of examples, the test data. Put simply, if the classifier really caught the relationship between features and classes, it should be able to anticipate the classes of examples it has not seen before. The common assumption for these models is that the training (and testing) examples are independently taken from an “example distribution”; when judging a classifier on a test set, an estimation of its performance on any test set from the same distribution is obtained. The training and test sets will be denoted by X_{train} and X_{test} , with respective matrices containing n_{train} and n_{test} examples as their rows, and example labels as column vectors y_{train} and y_{test} . The most widely used measure of how well a classifier performs on the test set is its accuracy. This is calculated as the percentage of examples in the test set with correct label predictions. A graphic representation is shown in Figure 5.5.

For the creation of the models the Scikit-learn (Pedregosa et al. 2011) machine learning library was used that includes both classifiers, and regressors along with many methods for

fitting, cross-validation, evaluation, and more. In the following, an overview of the classifiers that are used for the studies is presented along with their function as well as their advantages and disadvantages.

Naive Bayes Algorithm

Starting with the most simple classification algorithm, Naive Bayes is a probabilistic classifier, meaning that it uses the probabilities of each class to make decisions. It is a type of “naive” classifier because it assumes that the features of a data set are independent of one another, which is often not the case, although it provides very accurate results for text classification tasks, such as sentiment analysis and spam filtering. Naive Bayes works by creating a probability distribution for each class, based on the features of the data set. It then uses Bayes’ theorem to estimate the conditional probability of each class given the input data, and chooses the class with the highest probability. Naive Bayes classifier is simple and fast to train, and requires few parameters, while effective with limited training data (S.-B. Kim et al. 2002). On the other hand, it can be prone to errors when the data is highly correlated or there are many categories.

Support Vector Machine

Support Vector Machines (SVMs) are a strong and widely used machine learning tool for solving a wide range of classification and regression issues. SVM seeks the hyperplane in a high-dimensional space that best separates two classes of data points. The hyperplane is determined by maximizing the distance between the closest data points between the two classes, called the support vectors. Once the optimal hyperplane is determined, the SVM algorithm then uses the support vectors to classify the data points. The support vectors are the data points closest to the boundary of the classes, and they determine the decision boundary. The ability of SVM to handle both linear and nonlinear datasets, its good generalization performance, and its ability to learn from tiny datasets are its key benefits. Furthermore, SVMs have proven to be a useful tool for dealing with huge datasets and high-dimensionality challenges.

K-Nearest Neighbour Classifier

K-Nearest Neighbour Classifiers (KNNs) are a type of supervised learning algorithm that are used in pattern recognition and classification tasks. KNNs are based on a simple concept: the idea that “similar” data points should have a similar label. This means that when classifying a given data point, the label will be determined by comparing that data point to the closest “neighbors” in the training set. For calculating distances KNN uses a distance metric from the list of available metrics. There are a lot of different distance metrics available, but we focused on Euclidean distance function which is the most popular one among all of them as it is set default in the SKlearn KNN classifier library in Python. One of the easiest and most popular supervised learning algorithms is the KNN. It is particularly common in applications like picture identification and NLP tasks when the input is highly dimensional, namely facial recognition, spam filtering, and document categorization. The key benefit of a Nearest Neighbor Classifier is its ease of implementation and little need for data preparation alongside the fact that they can be accurate even when there is noise in the data (outliers). Also, these types of models are able to manage many features, are robust and have the ability to capture non-linear relationships in the data. While, the primary disadvantage of an KNN model is that it is computationally expensive since the computational complexity grows exponentially with the number of features. This is also known as the “curse of dimensionality”, indicating that the models become less efficient as the number of features increases.

Logistic Regression

Logistic regression, is a powerful supervised ML algorithm that uses a logistic function to estimate the likelihood of a discrete outcome given an input variable. More specifically, the model uses a sigmoid function to generate a probability score between 0 and 1 for each data point, and it is mostly used for binary classification with values like true/false, success/fail but can be used for multiclass classification as well.

The model is built on a linear equation that multiplies the input variables by a weight vector, a bias factor is included and a sigmoid function is employed to translate the model’s output into a probability score. A maximum likelihood estimation approach is employed to train

the model, that tries to discover the model parameters that maximize the likelihood of the observed data given the model. Logistic regression is a relatively simple model to implement, does not depend on a linear relationship between inputs and output variables and is limited to prediction tasks where the output is a binary or categorical value. Another advantage is the fact that it is resistant to overfitting, and because it can readily manage missing values, the model is suited for application with sparse data. Lastly, logistic regression is computationally efficient therefore it can produce predictions very quickly. However, the disadvantages of the model are that when the data is non-linear, it is susceptible to overfitting, and due to its assumption that the data is regularly distributed, it is sensitive to outliers.

Decision Tree Classifier

Decision tree classifiers are an important machine learning tool used to classify objects or predict outcomes. The model constructs a tree-like structure, with each branch representing a rule. At the root of the tree is the initial input data. As the tree is built, each branch is assigned a decision or rule that is based on the characteristics of the data. When the decision tree is constructed, it involves selecting the best attribute to split the data on, or to identify which variable will be used to make the rule. This is accomplished through the use of several approaches such as information gain, entropy, the Gini index, and chi-square. The chosen feature is then utilized to segment the data into subsets, which are subsequently partitioned further using the same or alternative features. The above process is repeated until each branch of the tree has only one conclusion. Then the decision tree is tested by using the constructed tree to classify new cases based on the input data. This is done by following the decisions or rules along the tree until a decision is made. The accuracy of the tree can then be determined by comparing predicted and actual outcomes with the advantage that the tree can be seen by the researcher allowing for a direct explanation. Finally, decision trees are not prone to overfitting, making them an excellent choice for a wide range of purposes.

Random Forest Classifier

The Random Forest algorithm combines multiple decision trees to create an ensemble that can make more accurate predictions than an individual decision tree, this is known as “ensemble learning” and improves significantly both the accuracy and stability of the model by reducing overfitting and variance. More specifically, this process is able to take into account the variance in the data in contrary to a single tree. Once a Random Forest model is trained, it feeds the new data through each of the individual decision trees and combines their predictions to produce a final output. The specific method used to combine the predictions of the individual trees depends on the type of problem being solved (classification or regression). The Random Forest classifier is easy to implement and is robust since it can deal with missing values and outliers in the data. Also, this model can work with a large number of features and identify the important ones to make accurate predictions. Despite its advantages, a random forest can be computationally expensive, as it needs to create multiple decision trees and combine them, while it is difficult to interpret the results of the models.

XGBoost Classifier

The tree-ensemble model is built using a gradient-based learning approach, where at each step, the algorithm constructs a decision tree and then refines the tree depending on the gradient of the loss function. The decision trees are then combined to produce an ensemble, which is subsequently used to make predictions. Furthermore, to increase the accuracy of the model it uses regularization and feature selection, along with a number of methods to decrease the complexity, increase the training speed and decrease overfitting, such as pruning, weight sharing, and early stopping. In terms of accuracy, speed, and scalability, XGBoost has been demonstrated to surpass many other methods.

5.3.6 Unsupervised Deep Learning for Text

Pre-trained Language Models

Another popular technique widely used for language processing are word embeddings, which are distributed representations of words in vector space, able to capture semantic information. Word embeddings, in particular, use unsupervised learning to generate a numerical representation of words from a big corpus. The fundamental idea behind these embeddings is to generate a low-dimensional vector representation of a word that may also collect additional data about its usage in a given text. The most popular word embedding models are Word2vec, GloVe, and fastText, and they are used for sentiment analysis, machine translation, and text classification or in more complex tasks such as question answering, dialogue systems, and machine reading comprehension.

Word2vec (Mikolov et al. 2013) is a shallow neural network that is trained on a large corpus of text to produce a vector representation of words. This model consists of two algorithms—the Continuous Bag of words (CBOW) and the Skip-gram. CBOW takes context words as input and predicts a center word, while skip-gram takes a center word as input and predicts context words. Both algorithms use a softmax activation function to produce a probability distribution over the output words. Additionally, fastText (Bojanowski et al. 2017) is an extension of the word2vec algorithm that uses sub-word information to generate word embeddings. This model takes into account the morphological features of words by breaking the words into n-grams. This model is especially useful for languages with rich morphological features. Lastly, GloVe is a global log-bilinear regression model that uses global co-occurrence statistics to learn the vector representations of words (Pennington et al. 2014).

Transformer-Based Models

Transformer models, which were introduced by Vaswani et al. (2017) in 2017, have revolutionized the field of NLP by taking into account the sequential nature of language. As discussed earlier, deep learning methods and particularly recurrent neural networks process words in a sentence one at a time in a sequential manner. In contrast, transformers

which are a novel type of architecture consider also the context of each word in a sentence simultaneously, rather than processing the words sequentially. This allows for more efficient and effective language processing, and has made it possible to build larger, more powerful language models.

The two most popular state-of-the-art NLP models, which have significantly outperformed previous models on a wide range of NLP tasks, are BERT, and OpenAI's Generative Pretrained Transformer (GPTn) (Aggarwal 2018). The GPT-n series of architectures is a family of large-scale language models developed by OpenAI⁶ (T. Brown et al. 2020). The "n" in GPT-n refers to the number of layers in the model, with each edition including more layers and more parameters than the previous one. The current GTP-3 edition was trained on 499 billion on tokens from Wikipedia, Books1, Books2, and WebText2 datasets, along with the Common Crawl dataset including 410B tokens.

Finally, BERT, which stands for "Bidirectional Encoder Representations from Transformers", is capable of learning contextual relations between words (or sub-words) in a text (Devlin et al. 2018). This model, which was trained on a large corpus, can be used to perform a variety of NLP tasks, such as predicting the next word in a sentence or generating summaries. Because of BERT's capacity to understand the context of words in a sentence, it is particularly well-suited for problems requiring a sophisticated comprehension of language, such as question answering or sentiment analysis.

Overall, word embeddings are now widely used in NLP tasks such as sentiment analysis, machine translation, and text classification and summarization. They have proven to be useful in capturing the semantic information of words, and have also been shown to improve the performance of most NLP tasks.

Cross-validation

When the training dataset is small, cross-validation is used in order to enable the model to be evaluated on a more extensive set of data than a single train/test split. Moreover, this

⁶<https://openai.com>

technique can be implemented to evaluate the performance of the same model on new data, or tune the hyperparameters of the model. It divides the training data into a number of folds, usually 5 or 10. The model is then trained on a different fold of the data each time, and the performance of the model is evaluated on the remaining folds. This process is repeated until each fold has been used as the evaluation set once. The performance of the model is then averaged across all of the folds, providing an estimate of the model's accuracy on unseen data (Pereira et al. 2009). This technique is more appropriate in algorithms that require the dataset to be divided into train and test set where there is a need to limit the effect of bias in a "random" selection of samples.

In its most extreme form, the "leave-one-out" cross-validation (Kohavi 1995) iteratively involves the entire dataset for training the model by leaving out one sample at a time, making predictions for the left out sample, and repeating the process until all samples in the dataset have been used for prediction. The model's accuracy is then calculated using the predictions made for each sample in the dataset. As the model is evaluated using all of the available data, this technique offers an unbiased evaluation of how well it works. Additionally, it can be used to compare different models and choose the best one for a certain problem, while it can assess the accuracy of models in high-dimensional datasets where the quantity of samples can be quite small for machine learning tasks. Overall, this estimate is typically more accurate than a single train/test split, because it uses more of the training data for evaluation.

5.3.7 Evaluation

A crucial stage in creating effective models is evaluating the outcomes of machine learning. The researcher must examine the performance and decide whether or not additional modifications are required. The most widely used metric for assessing a machine learning model is accuracy, which calculates the proportion of the correct predictions. Unfortunately, this is not always right since in the case of imbalanced datasets, where one class is significantly more common than the other, the model may obtain high accuracy by consistently predicting the majority class. Therefore, in order to create a robust model that can generalize to

unknown data, the training set needs to be balanced.

For the evaluation of a model various metrics can be used, such as precision, recall, or F1-score. Precision is defined as the percentage of the true positive predictions made by the model, while recall is the percentage of the actual positive cases that were correctly predicted. Ideally a model is considered good when it has both high precision and high recall, indicating that it is able to accurately identify the positive cases without generating too many false positives (Novaković et al. 2017). In addition, the F1-score which combines precision and recall is very often used to evaluate binary classification models, while other methods such as confusion matrix, and Receiver Operator Characteristic (ROC) curves, can also be used to evaluate a machine learning model (Davis and Goadrich 2006). Confusion matrix is used to compare the model's predictions with the actual labels of the examples in the test set, and ROC curves plot the true positive rate against the false positive rate, and are useful also for comparing the performance of different models.

Finally, a researcher should consider the context of the model when evaluating its performance. For instance a model might have a high accuracy, but its results may not be meaningful in terms of the context in which it is being used. Therefore, the model's performance should be evaluated in the context of the application in order to ensure that the model is providing useful results (Ribeiro et al. 2016). Model Interpretation allows for greater trust and accountability in the algorithmic decision-making process. That is the reason for the creation of Explainable AI, also known as XAI.

Traditionally, ML models make decisions based on complex mathematical models that may be difficult to comprehend. For a model to be interpretable there is a need for the generation of simple, human-readable explanations of the decisions. For instance, decision trees can be visualized, so a direct observation can provide insight into the model's process. This can be useful for understanding the structure of the model, as well as its strengths and weaknesses (Doran et al. 2017). For example, if the tree is too deep or has too many branches, it may be overfitting the data therefore it will not be able to generalize well. Furthermore, by observing the decision tree, one can gain insight into the relationships between the features and the

output, as well as the importance of each feature in the model. Also feature importance provides insight into the algorithm, by presenting the most important features based on the weights of each one. The features that turn out to be irrelevant for the prediction can be omitted.

However, neural networks are considered a “black box”, and people cannot always trust the results, since it is challenging to identify the exact decision-making process or to correct errors in the system. Nevertheless, attempts to make them more transparent and trustworthy for users are being explored with model interpretation libraries like SHAP⁷, ELI5⁸, Treeinterpreter⁹, Dtreviz¹⁰ and LIME¹¹.

5.3.8 Deployment

In a research setting the Deployment phase includes the organization of the results in an easy to understand and useful manner. A thorough analysis of the results, the interpretation of the AI models and a discussion of the value of each study are presented. In the next section, the methodology outlined here is followed for five different research studies.

The first study examined the use of AI to detect the perceived quality of stories on the blogging website *Medium.com*. The study revealed that machine learning models, such as decision trees, and random forest, can be used to accurately identify the success of an article in the platform based on author, style, content and context features.

The second study focused on using AI to determine the quality of news stories based on a theoretical framework derived from the literature. The analysis revealed that there are three level of journalistic quality, low, medium, and high with the classifier to predict best the low and high quality labels.

The third study looked at the use of machine learning algorithms to predict engagement on social media. The results showed that tree-based approaches, support vector machines and

⁷<https://shap.readthedocs.io/en/latest/>

⁸<https://eli5.readthedocs.io/en/latest/>

⁹<https://pypi.org/project/treeinterpreter/>

¹⁰<https://github.com/parrt/dtreviz>

¹¹<https://christophm.github.io/interpretable-ml-book/lime.html>

logistic regression, can be used to accurately predict audience engagement on Facebook.

The fourth study focused on the use of the visual elements of journalism to predict both the engagement and the quality of a news article. The featured image of each news story was first used for image recognition and then as an input to the classifiers. The findings showed that images alone can reach an accuracy of 0.70 (F1-score), but when combined with textual features they do not improve the performance of the models.

Finally, the fifth study focused on the application of AI for disinformation detection. The results showed that textual attributes can predict fake news with high accuracy. Furthermore, the accuracy of the models was improved when combined with additional engagement features from Facebook such as likes and comments.

In conclusion, this research indicates that AI models have the potential to be utilized for predicting a variety of outcomes, from social media engagement to journalistic quality. By understanding which characteristics of a news story are the most influential, news organizations can leverage this knowledge to optimize their content and better engage their audience.

Chapter 6

Studies

6.1 Study:1 Audience Engagement Metrics and Perceived Quality

The relationship between the perceived quality of a news article and a person's intention to share it has been studied, with a positive correlation being established (Ma et al. 2014). This suggests that higher quality articles have a greater likelihood of being shared by readers. However, it may be unwise to simply conclude that the most engaging articles are better than the least engaging and make business decisions based on this assumption. To come to more reliable conclusions, a more systematic approach is necessary. By analyzing on-line textual data, meaningful answers can be obtained to important questions such as what makes one article better than the others?

In an effort to answer questions of user engagement, many communication scientists have studied traffic analytics (e.g. reading time and visits) and user metrics from social media networks such as Facebook and Twitter. Different measures of interest, including attention, content popularity, and sub-dimensions of engagement, have been used to try to explain why users may prefer one news article over others (García-Perdomo et al. 2018; Y. Jin et al. 2017; Trilling et al. 2017; Valenzuela et al. 2017b; Arapakis et al. 2016). Additionally, it has been demonstrated that article features can be used to accurately predict the social pop-

ularity of news stories on Twitter before they are published, with social popularity being measured as the total number of times an article's link appears on Twitter (Bandari et al. 2012).

Scholars have also been attempting to use early measurements of an article's popularity to predict its potential success. For example, data scientists at the Washington Post have been able to anticipate the level of popularity after publication, in order to provide their readers with a better experience (Keneshloo et al. 2016). Popularity is defined as the number of page views that an article receives on the first day, and a system has been created which takes into account factors such as contemporariness, writing quality, and social media networks, during the first half hour after publication, to accurately predict the article's popularity. It has been suggested that these features can be linked to the *quality* of the text. For instance, D. Park et al. (2016) predicted the quality of online comments. For this, they used the selections that a newsrooms' editors made by picking "top comments" as a ground truth.

6.1.1 Study Overview

This study focuses on the features that have been previously linked to engagement with news articles, referred to as engagement metrics, engagement features, or features of engagement. It builds upon efforts that have successfully used different sets of features (García-Perdomo et al. 2018; Trilling et al. 2017; Valenzuela et al. 2017b), and taxonomies (Orellana-Rodriguez and Keane 2018), and proposes a framework with new features and extended feature dimensions. The results of this study were published at two different research papers (Sotirakou et al. 2019; Sotirakou et al. 2018).

For this study, 200,000 articles from *Medium.com*, a popular blog platform, were used. The number of "claps" (similar to likes) an article receives was taken as an indication of its perceived quality. Additionally, new features were generated to explore how certain content characteristics, such as the length of the article, affected reader satisfaction. Furthermore, features that evaluated the author's tone of speech and narrative style, as well as their level of social media network on the blogging platform, were created to measure how an author

affects people's attitudes towards a news story.

A preliminary evaluation of this work produced some interesting outcomes. The perceived quality of the articles may take different meanings depending on the category that they belong to. Although some generic, rather expected results were found, like the significant contribution that the number of followers may have on the popularity of an article (irrespective of the category), a more close investigation revealed the almost equal importance of other features and the verification of our proposed extended model. Each news category generates a different classification of features as the most important, which is able to influence their acceptance from the digital audience. As such, perceived quality should always be regarded taking into consideration the category that an article belongs to (or other contextual characteristics) since a possible generalization of the term could lead to incorrect conclusions and interpretations concerning the fundamental elements that constitute it as a term and what it represents. Furthermore, this study was able to extract several important rules based on these engagement features for each category, which, if followed, may guide the writing of the next article of the authors, and provide them with insights into the probability of popularity it may gain.

This study is structured as follows: Section 2, presents an extensive literature review around the topic of investigation as well as a comparison with related work indicating the main contributions. In section 3, the proposed-extended model is detailed, along with the main dimensions and employed features. Section 4 refers to the method of evaluation and the data-set used and section 5 discusses the two phases of the data analysis, the results in perspective, as well as the proposed rules for each category of news that might increase their perceived quality. Section 6, concludes this study and addresses future work.

To find if there are certain dimensions of online article's perceived quality, a trustworthy proxy is required that signals that when a news story is valuable to the audience. Thus, the blog platform *Medium.com*, which consists of a miscellaneous collection of articles, features an interesting combination of liking and sharing: It allows the user to "clap" for an article, which shows the author that their story was liked. Arguably, these claps can be a better proxy

for perceived quality than the mere number of views or the number of “shared” on other platforms is, as reading or sharing does not necessarily imply a positive value judgment. One can, for instance, read an article because they feel bored, or share it with a disapproving comment or to mock it.

Although, there is not, to our knowledge, previous research that explicitly supports the correlation of claps with perceived quality, given the fact that the claps on *Medium.com* are public and other users can see the articles one clapped for, claps can result not only in feedback to the author but also as reader’s recommendations to each another. Furthermore, the number of claps (it is possible to clap up to 50 times for the same story) is also used as an input for the personalization algorithm to determine which stories one’s followers will see, and hence it can be seen as an indirect way of sharing. The current approach argues that claps on *Medium.com* can be used as a proxy to measure the perceived quality of an article and until now they have been under-researched.

Next follows a review of the research on the digital transformation of the modern newsroom in terms of the evaluation of the quality of their products and how online users on social platforms and recommendation algorithms have altered the traditional gatekeeping process. The choice of *Medium.com* as an optimal case for the research on the topic of perceived quality is explained and the concept of quality is defined. Then the literature review continues with a brief overview of how the digital journalist tends to think like a marketer implicating in a way the digital audience’s feedback in the decision making of the news agenda. After that, a section about previous attempts to use audience engagement to measure quality through popularity is presented and the actors involved in online readership that play a decisive role in what content ends up being evaluated are mentioned. The backbone of the literature review is then used to help build the framework and the model for the predictions.

Measuring Perceived Quality

Quality is a fundamental concept in building a brand name and customer satisfaction; it serves as an advantage against competitors and in the case of quality journalism “provides

citizens with the information they need to be free and self governing” (Kovach and Rosenstiel 2014). Although the quality of journalism is very hard to define, scholars have used a variety of different methods to measure it in order to improve the quality of news stories. According to literature, there are three distinct ways to measure the quality of journalism: examine the attributes of respected news organizations, analyze the content-based features of the products of these organizations and investigate the engagement metrics to unravel which characteristics resonate with the audience (Lacy and Rosenstiel 2015). In this study, we focus on the third method and delve deeper into the audience’s perception of quality. As mentioned by Aaker (2009), perceived product quality is “the customer’s perception of the overall quality or superiority of the product or service for its intended purpose, relative to alternatives”. This work considers the standpoint of Aaker, building upon the assumption that online articles could be regarded as the content of the services that *Medium.com* offers to the digital audience, and investigates how their perceived quality is influenced by their features of engagement and the category (purpose) that they belong, in relation to the next important article close to it.

Recent studies have shown that decisions within news organizations (for instance, about which articles to promote) are increasingly driven by engagement metrics (Hagar and Diakopoulos 2019; Welbers et al. 2016; Tandoc and Thomas 2015; Lee et al. 2014). Engagement data are crucial when it comes to news placement online, with evidence suggesting that in some cases they are even more important than editorial opinions (Lee et al. 2014). Tandoc and Vos (2016) after observing and interviewing journalists in three online newsrooms for 150 hours, concluded among other things that journalists balance their editorial decisions between traditional journalistic norms and audience influence expressed through social media engagement and traffic metrics. The same study stressed also that the journalists, promoting their stories on social media, thinking like marketers and giving in to the market demand, can compromise editorial autonomy. Thus, it is of crucial importance to find out how to reconcile the high popularity of news content with the perceived quality as expressed by the user through certain engagement measures.

This study is focused in particular on the perceived quality of online news articles. A very

broad definition of news is being used here that does not limit the focus on hard news (such as politics or economy) but rather views journalism as “a serial presentation of information and conversation about public events, trends, and issues distributed through various media with the primary purpose of informing, entertaining, and connecting citizens in communities” (Lacy and Rosenstiel 2015). Therefore, in this study, the articles that appeared in the following categories on *Medium.com*: Lifestyle, Business, Technology, Sports, Health along the general News category are considered to belong in the broad spectrum of news articles.

6.1.2 Audience Engagement

Finding a news story interesting is a matter of personal taste, yet there is a taxonomy of news values that suggests that certain attributes influence audience selection as well as engagement with digital news stories (Harcup and O’neill 2017; Trilling et al. 2017; P. Weber 2014; Kepplinger 2008; Eilders 2006; Galtung and Ruge 1965). There is no full consensus on which features have the largest impact on the so-called “shareworthiness” (Trilling et al. 2017), though Chapter 3 mentions several which have been used to predict engagement, such as geographical proximity, conflict, human interest, negativity, positivity (Trilling et al. 2017), exclusiveness, deviance and relevance to the list (Valenzuela et al. 2017b; García-Perdomo et al. 2018; Kilgo et al. 2018b), sensasionalism (Kilgo et al. 2018b), a normative approach, diversity, relevance, ethics, impartiality, objectivity, and comprehensibility (Urban and Schweiger 2014).

Actors Involved in Online Readership

Search engines, social media networks like Facebook, news aggregators, and publishing platforms such as *Medium.com* use algorithms that control the information a user sees (Carlson 2018). Thus, algorithmic curation delivers alternative results to similar search terms according to the previously collected behavioral data on the given user who sought information (Pariser 2011; Nechushtai and Lewis 2019). For instance, every user that visits *Medium.com* has a different experience due to Medium’s recommendation system that delivers a different mix of stories to every visitor based on factors like their reading history, the

authors, and publications they follow, and so on. These “algorithmic gatekeepers” (P. Napoli 2015) are updated frequently so there is not an easy way for authors to figure out how the algorithm forms its decisions.

In addition, editorial decisions about the article position in a news website are crucial concerning what content users see and inevitably end up rating. Previous research has shown that article promotion on the homepage of a news outlet improves the likelihood of it being shared, read, emailed, and several case studies have shown that these decisions are heavily influenced by audience choices (Tandoc and Vos 2016; Tandoc and Thomas 2015). Similarly, the editors of *Medium.com* use specific algorithms on a daily basis to review and select the best stories to be distributed across the platform, the app, and the newsletter. The selection is based on high editorial standards and the guidelines for writing an article eligible for curation are available on the website (Medium.com n.d.).

Furthermore, users themselves curate their personal information environment (Thorson and Wells 2016a) to serve their individual goals, while the relationships between users that follow one another can shape readership as well. Similar to other social media networks *Medium.com* consists of a community of users that through their behavior can impact what stories other users see on their homepage (Thorson and Wells 2016a). As Singer (J. Singer 2014) showed in her work about user-generated visibility, nowadays the traditional editorial decisions about which stories are important for the audience to read, are followed by a secondary gatekeeping process in which the users themselves filter out valuable information by liking, upvoting, downvoting, sharing, emailing, and so on. Since on *Medium.com* there is no negative feedback button like for instance in Reddit (Reddit n.d.), the possibility of more balanced feedback that users can even out through time does not exist. Therefore, more followers imply more potential readers which can lead to more positive feedback.

The actors discussed above, unquestionably influence the number of people who see a given article on their screen, therefore it can also affect directly the perceived quality of this article.

Main Contribution in Relation to Previous Literature

Despite the advantages that “claps” offer compared to “shares” or “likes” for understanding the perceived quality of articles, they are rarely studied. To the best of our knowledge, no academic paper so far has studied how features that explain the perceived quality of articles on Facebook and Twitter can be used to explain “claps” on *Medium.com*. A reason for this may be that Twitter and Facebook as widely popular social networking sites draw most attention of both the public and researchers. However, as discussed above, the “clap” as a signal that both implicates approval (as a “like”) *and* disseminates the story to one’s social network (like a “share”) is worth studying in detail for those who are interested in the perceived quality of an article. This study fills this gap. More specifically, the main contribution of this work is three-fold:

1. It proposes a model of engagement metrics based on the literature review that can quantify the perceived quality of different types of articles.
2. It identifies how different classifications of engagement metrics contribute to the perceived quality of articles depending on the various categories of news. Also, the results show that an article’s genre affects the importance of features that reveal its perceived quality differently and should be treated as a situation-specific factor given the nature of each category. However, instead of merely focusing on identifying the role of different feature categories, this work goes a step further and explores the nature of their effects.
3. It recommends a set of rules, as guidelines to authors, based on the engagement metrics of the respective articles’ categories, for increasing the probability to gain popularity. While prior work on online popularity prediction on Twitter shows that there are ways to forecast a given article’s success before (Bandari et al. 2012) and after publication (Keneshloo et al. 2016), little work has examined how different attributes of an article before publication can influence its popularity. This study demonstrates the importance of specific features that when combined, can shape online popularity and thus the perceived quality of a given article.

For this study various characteristics of an online news article have been used, such as its sentiment, novelty, and tone (García-Perdomo et al. 2018; Trilling et al. 2017; Valenzuela et al. 2017b), along with the useful taxonomy of news tweets by Orellana-Rodriguez and Keane (Orellana-Rodriguez and Keane 2018), that distinguished three dimensions of a tweet that influence its dissemination: user, content, and context features. The current study considers these three Twitter related feature dimensions and adjusts them to news articles, while also extends them by additionally proposing writing “style” as a category within content features. Then this taxonomy is used to an empirical test. By considering such a broad set of possible features, the current study sheds light on which ones can shape the perceived quality of news articles and which do not “work” depending on the genre, and proposes a specific set of guidelines to authors. In doing so, this study will set a baseline and starting point for others who are interested in the prediction of “claps” or other online popularity metrics as well as for those who want to further refine the taxonomy developed in Orellana-Rodriguez and Keane (2018).

6.1.3 Model & Feature Extraction

The primary concern of this study is to create an inclusive model to measure and quantify the value of perceived quality across various article categories. The backbone of the model is structured using the overarching dimensions proposed by Orellana-Rodriguez and Keane (Orellana-Rodriguez and Keane 2018), namely (i) author, (ii) content, (iii) context. (Figure 6.1).

As discussed above, perceived quality may be considered as a widely used validation factor and can be measured using various voting mechanisms for understanding the discrete actions of users on articles like acceptability, shareability, etc. Additionally, voting mechanisms are a way for communities such as Digg.com and Instagram.com to define content popularity and have been proven effective also for social Q&A communities like Quora.com to filter important answers (Y. Jin et al. 2017). Similarly, on *Medium.com* users’ claps act as social signals to draw attention to influential people (Paul et al. 2012). Besides following an

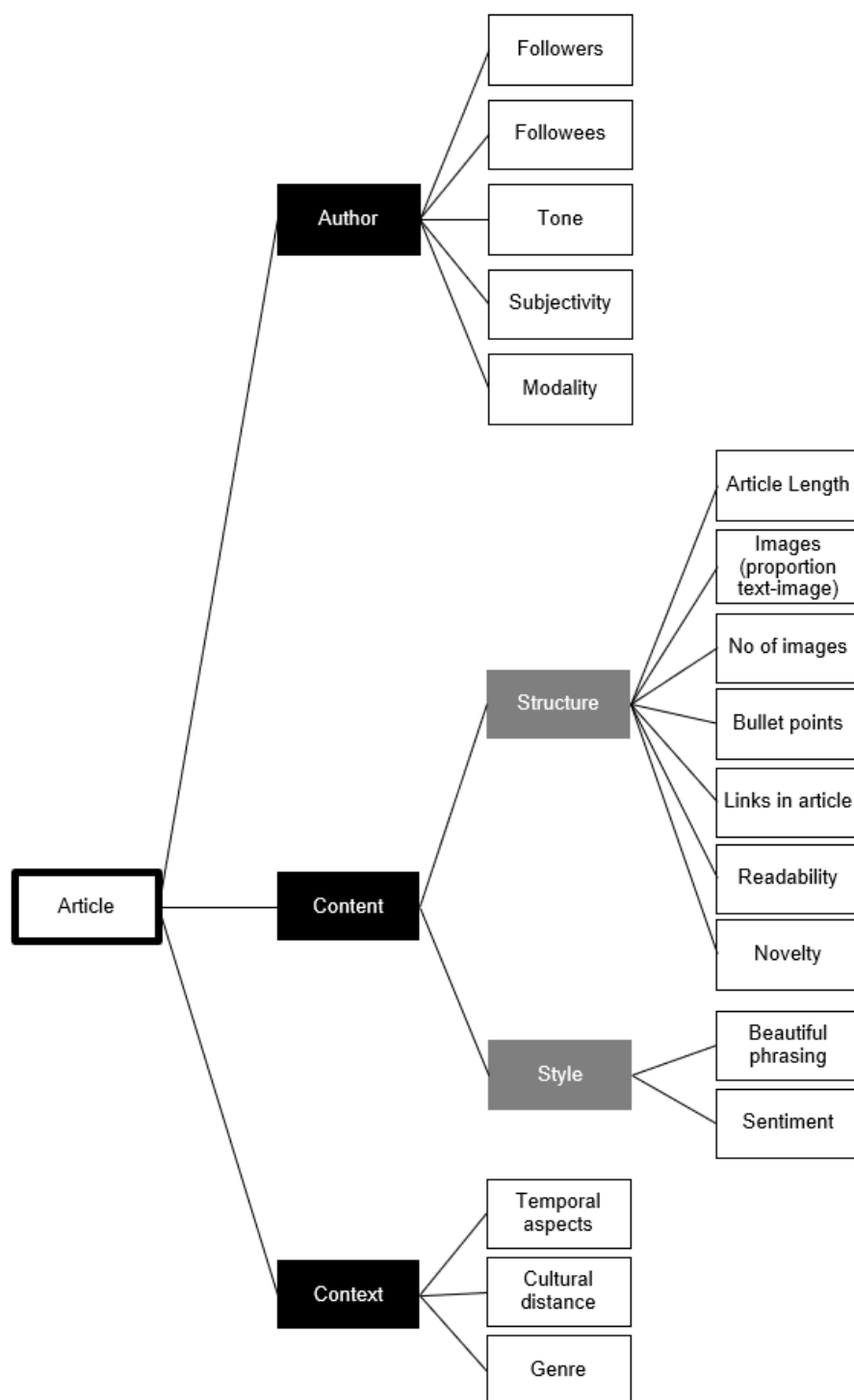


Figure 6.1: Features that affect the perceived quality of the readers

author, every content creator can follow other users to expand their network.

Hereafter, the main features of our model are explained in detail along with the rationale for their selection. Information about the text analysis techniques and specific lexicons used can be found at the Appendix C.

Author Features

Author's Popularity: The author's popularity which translates into the total number of followers on Medium can easily be calculated. Another way to determine an author's popularity is to look for the number of followees, as this can also contribute to build reputation and gain more votes.

Author's tone: According to marketing and communication specialists, the use of personal tone is a way for having a conversation with the audience, and expand one's brand, thus self-referential or self-reflexive messages are frequently used to attract the readers' attention (Nöth 2007). This study investigates the use of first person pronouns (I, you, we, etc.), reflexive pronouns (myself, yourself, etc.), and possessive pronouns (mine, yours, etc).

Subjectivity: In principle, journalistic coverage is meant to be neutral, but recent research shows that social media drive journalists towards more subjective ways of presenting information (Lasorsa et al. 2012). For the subjectivity feature we used the Python library TextBlob (Loria 2018).

Modality: Grammatical modality is signaled by grammatical moods that express a speaker's general intentions such as certainty, willingness, necessity and so on. To capture the degree of certainty of an author, by the modal words and expressions used in the text, the pattern.en (De Smedt et al. 2012) was implemented, which measures modality as a value between -1.0 and 1.0 , where values > 0.5 represent certainty.

Content features

In the model, the focus is on content features that are expected to influence the reader.

Structure

Article Length: The information depth is also depending on the average length of the articles. For instance, the New York Times prefers to focus on news depth rather than on news variation (Quandt 2008). For that feature, the so-called “reading time” is used provided by *Medium.com*.

Images: The use of salient pictures is a way to attract audience attention and achieve link clicks. Maintaining a good balance of text and images is an important factor for quality (Bogart 1989).

Bullet points: Listicles, meaning a mixture of “list” and “article” are very popular types of articles and according to research, they tend to be highly shareable (Okrent 2014).

Links: According to literature (Karlsson 2010) one of the signals of transparency in an article is the inclusion of external links to sources and documents, that could provide more verified sources, complementary information, or may present different angles.

Readability: Readability testing is a way to automatically check the clarity of writing. Good readability means a text is accessible and easy to understand, therefore, when an article is too difficult (low reading ease), respondents might just be unable to comprehend. In this study, the Flesch-Kincaid Grade Level (Kincaid et al. 1975) is used.

Novelty: If an article covers new, exclusive information that is not present elsewhere, it may be more interesting for readers and enhance engagement (Trilling et al. 2017). A three-day sliding window was used and the normalized (to account for title length) number of nouns (which was determined using spacy (Honnibal and Montani 2017)) that a headline shared with all concatenated headlines in the window to calculate a popularity score. A low popularity score can be interpreted as an indicator of high novelty (Trilling et al. 2017).

Style

Beautiful phrasing: High-quality articles are often written using beautiful language (Arapakis et al. 2016; Louis and Nenkova 2013), meaning more unusual phrasing and creative

words that lead to more positive feedback (Sotirakou et al. 2018). This study uses Term Frequency-Inverse Document Frequency (tf-idf) as a proxy because the use of rarer words can be seen as an indicator for the use of beautiful language.

Sentiment: Emotional text can spark feelings in the audience and trigger different kinds of emotional responses. Both positivity and negativity have been shown to influence news sharing (Sotirakou et al. 2018; Valenzuela et al. 2017b; Trilling et al. 2017). This study used Vader (Hutto and E. Gilbert 2014) package for sentiment and polarity calculation and NRC lexicon for the emotions extraction. The eight basic emotions of the lexicon were: joy, sadness, anger, fear, trust, surprise, disgust, and anticipation (Nissim and Patti 2017). For the measurement, the counts of emotional words were computed, each normalized by the total number of article words.

Title polarity: The headline of an article is a significant aspect of every story as modern news readers and social media users scan the headlines before they decide to click (Dor 2003; Scacco and Muddiman 2016). To measure polarity expressed in the headline, again, Vader Lexicon was used (Hutto and E. Gilbert 2014).

Context Features

Context features explore the conditions in which an article was posted and include temporal and locational features.

Genre: The topic of the article reflects certain characteristics of its nature and it has been greatly investigated in previous work (Louis and Nenkova 2013; Arapakis et al. 2016). The categories considered are: News, Technology, Business, Health, Sports and Lifestyle.

Temporal: The data set consists of articles from only one year, so only hour and type of day are used as temporal features.

Cultural distance: Previous research (Trilling et al. 2017) found that geographical proximity and stories that mention Western countries increase social media engagement. This study, uses references to geolocations from the texts using the Mordecai system (Halterman 2017), which extracts toponyms and returns their coordinates.

6.1.4 Method & Dataset

The analysis was conducted on a data set of a total of 200K articles taken from a single source, namely Medium.com, so that the extraction of the features would be consistent across all categories. The choice of the independent variables was made based on the characteristics that according to the studied literature seem to contribute to online content popularity. To quantify the perceived quality of the articles, this study opted for the official metric of this online community which is called claps, and serves as the dependent variable for our model. To determine if and to what extent the proposed characteristics of an article are related to its perceived quality a series of experiments took place. To this end, the model was trained using the aforementioned features and divided the prediction problem into two phases, the first was to predict the perceived quality of an article taking as input a mixed data set, while the input data for the second phase was each category separately.

Data

The data set was collected from the website Medium.com, which covers a range of topics from tech to politics to well-being, and for this study, articles from the following six categories News, Technology, Health, Business, Sports, and Lifestyle were retrieved. A Python program was written that scraped the website in January 2019 and collected articles published between September 2017 to September 2018. The total amount of articles before cleaning was 247,071 and had a large distribution of claps ranging from 0 to 292K. Before the analysis the data set had to be prepared, hence all the news articles were processed for stop words, nonstandard words and characters removal, stemming, and tokenization. In addition, regular expressions were used to remove non English stories, empty values, and HTML code from the text. After the preparation for the analysis 200,710 cases remained for inclusion in the experiments (Table 6.1). One of the issues that appear studying Table 6.1, is the fact that the vast majority of the articles have equal or near to zero claps. This causes the creation of imbalanced classes which is a common problem in machine learning classification, since there is a disproportionate ratio of observations in each class. Standard classification algorithms, that do not take into account class distribution, are overwhelmed

Table 6.1: Descriptive statistics

Category	No. Articles	Mean (claps)	Median (claps)	75 th percentile
News	17762	150	0	22
Technology	30883	236	2	59
Business	35768	65	0	14
Health	49072	39	0	3
Sports	12664	19	0	5
Lifestyle	54561	347	22	121

by the low-popularity class and they ignore and misclassify the minority of successful articles since there are not enough examples to recognize the patterns and the properties of the popular class. In this work, attention is taken using specific techniques and measures of quality to predict the rare but important class of successful articles.

For the extraction of the features, we used several Python libraries for text analysis, cleaning, filtering, counting words, and processing textual corpora. After determining the main categories of features for the model, the importance of each one was examined using a Random Forest Classifier. More specifically, the number of claps was used as the dependent variable and created decision trees able to predict whether the claps count of a given article will be high or low.

Articles with claps above the 99.8th percentile were considered outliers and were removed. Furthermore, the features were normalized and checked for feature correlation. (For normalization we used the scaling technique that is converting floating-point feature values from their natural range (for example, 100 to 900) into a standard range— in our case between 0 and 1). Also, a Principal Component Analysis (PCA) was conducted to search for redundant features. However, the analysis showed that the large set of components could not meaningfully be reduced into a smaller one without losing a lot of variance.

Data Analysis & Findings – in two Distinctive Evaluation Phases

The analysis was divided into two phases. In Phase A the whole data set of the articles was examined, to find the most important features and create a baseline for further analysis.

After that, the study was focused on every different news genre (see Table 6.1) trying to capture if there was any variation in the predictive power of the importance of each feature according to the genre the articles belonged to and in that case to what degree influences the perceived quality. In Phase B, a decision tree classifier was used to extract the rules of the best performing leaves that provided specific combinations of the engagement features. These features may increase the chances that a given article has to gain greater popularity and thus higher perceived quality when they reach certain values according to the specific leaf.

This study separates articles with low and high claps by converting the numerical variable that represents the total number of claps into a categorical one. In both phases, a binary classification task was created aiming at constructing a model able to give us the importance of each feature as well as an interpretation of the prediction to enhance our understanding of a successful article. In our experiments, we used 80% of the data set for the training set and the other 20% for testing purposes. We used the scikit-learn (Pedregosa et al. 2011) implementation of tree-based machine learning approaches (decision tree, random forest, xgboost) since tree models can be more interpretable than other complex models and provide more descriptive explanations of the effect of the features.

6.1.5 Phase A - Evaluating the Importance of the Proposed Model Engagement Metrics in Relation to the Perceived Quality

In the first phase of experiments, the original data set was randomly sampled and a new one was created that consisted of 60K articles, 10K from each category (News, Technology, Health, Business, Sports, Lifestyle). The original data set can be seen in Table 6.1. To predict the appreciation of the articles it is crucial to get the perceived quality rank right, particularly at the high ranks, by correctly distinguishing the most highly appreciated articles. Therefore, in the new data set, two buckets of articles were randomly sampled, the high and low with each consisting of 14K articles. More specifically, 14K articles that were placed in the top 25% of claps with the highest values were in the high-claps bucket, while 14K of 0 claps articles

Table 6.2: Permutation Importance for the top 10 features

Feature	Weight
Followed By	0.5810 ± 0.0240
Article length	0.0517 ± 0.0033
Following	0.0367 ± 0.0088
Flesch	0.0166 ± 0.0027
Polarity	0.0133 ± 0.0010
Tone of speech	0.0125 ± 0.0039
img/word	0.0087 ± 0.0026
No. Images	0.0081 ± 0.0020
Subjectivity	0.0062 ± 0.0044
Contains USA	0.0061 ± 0.0019

ended up to the low-claps bucket.

For prediction purposes, the Random Forest package from the Python library scikit-learn (Pedregosa et al. 2011) was used to create a classifier of the articles and to measure the importance of the predictor variables. The significance of a variable is calculated internally during the construction of the decision trees by checking the raise of prediction error when data for that variable is permuted while all others are left unaltered. The F-measure (F1) was adopted to evaluate the performance of the model. The best score of all between the three different classifiers was generated by the XGBoost algorithm with an F1-Score of 80%.

Afterward, for the interpretation of the random forest, the contribution of every feature on the prediction was measured. The variables were ordered from highest to lowest rated and the 10 most important out of 30 were kept. The ELI5 Python package for “Inspecting Black-Box Estimators” (Mikhail Korobov 2016), was used for the calculation of permutation importance (Table 6.2).

To get better insights into which features contribute to the high or low-claps bucket, the importance score of every bucket was examined separately, starting from the most important variables that have more predictive power over the high-claps bucket.

Figure 6.2 presents the importance of the variables of the selected dimensions. The dimension of author is the first with Authors’ popularity (Followers & Followees, Subjectivity, Tone, Modality), second is the Content dimension and specifically Structure (Article Length, Im-

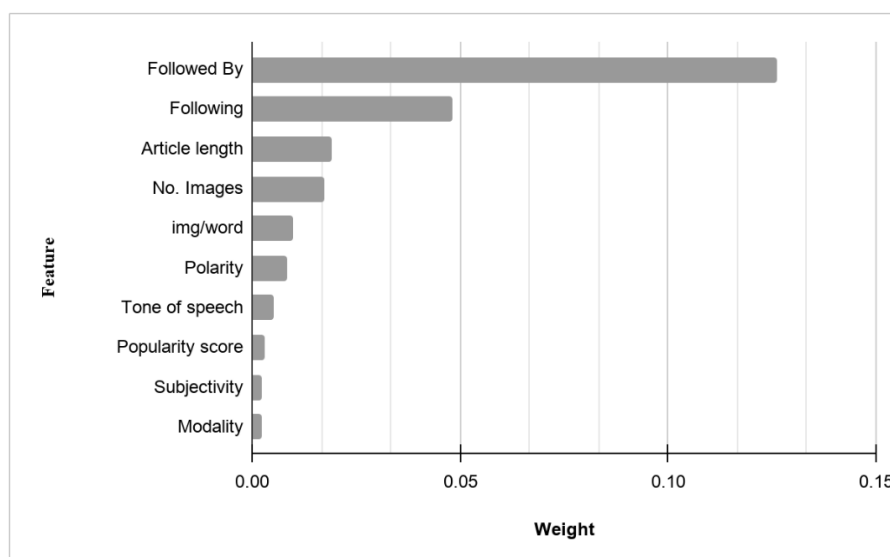


Figure 6.2: Feature importance score for high claps bucket

Table 6.3: Precision, Recall, and F1 scores

	precision	recall	f1-score
Baseline	0.8	0.8	0.8
News	0.85	0.79	0.82
Tech	0.8	0.85	0.83
Health	0.75	0.83	0.78
Lifestyle	0.88	0.85	0.86
Business	0.77	0.75	0.76
Sports	0.74	0.74	0.74

ages, Readability, Popularity score) and Style (Polarity).

In general, as shown in Table 6.2, the social network of an author, along with length, pictures and polarity are the most important characteristics that influence the perceived quality of an article. Interestingly, the results reported in Figure 6.2 give an idea about the predictive power of the features contributing to the high-claps bucket, which will be used as the baseline for the analysis of the different news genres that will follow.

Table 6.3 shows the precision, recall and F1 scores of each bucket of the seven data sets.

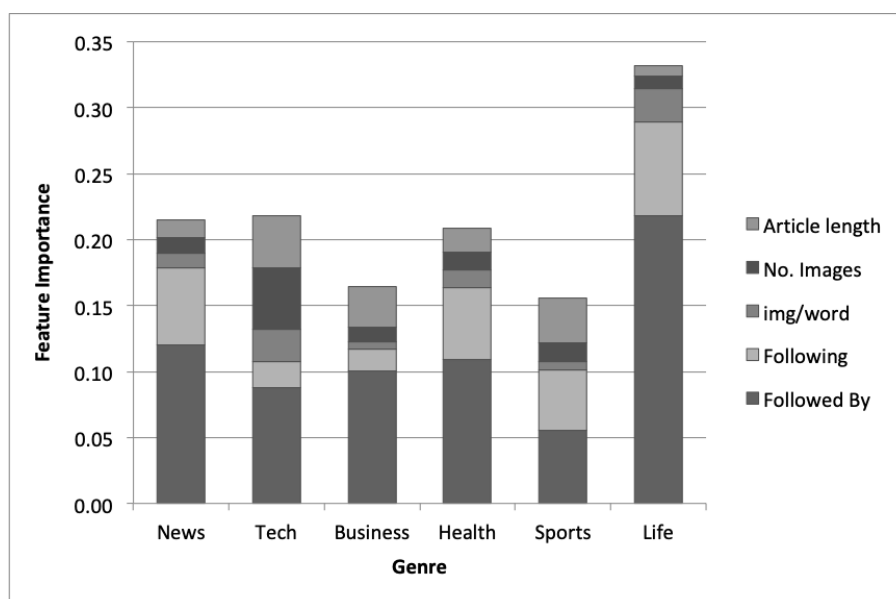


Figure 6.3: The importance measures of Top-5 features in six different article categories

Evaluating the Importance of Each Category's Engagement Metrics and Their Influence to the Perceived Quality

To check whether some features are more important depending on the category, different data sets were created for every category and the same process was applied to split the articles into two classes of the same proportion of texts, one bucket with the 25% higher clapped texts and the low-clapped bucket with the ones that got 0 claps. Despite the fact that all articles were taken from a single source, differences exist on the distribution of the claps in every category as shown in Figure 6.3.

The importance scores of the features revealed several interesting insights into the three dimensions of the model. The author dimension is the most influential regarding perceived quality, content follows with structure features being more important than style ones and the least significant is context.

Discussion of the Results

First, the most important feature is the number of followers an author has, which is the only feature that remains stable in all data sets. This is somewhat expected because the more followers the author has, the more claps their articles will receive. Similarly, the number of

Table 6.4: Comparison of Top-10 and Top-5 Articles Categories' Features Similarity

Article Types	Similarity (top-10)	Distance (top-10)	Similarity (top-5)	Distance (top-5)
<i>Baseline</i>	<i>1.000</i>	<i>0.000</i>	<i>1.000</i>	<i>0.000</i>
News	0.632	0.368	0.775	0.225
Technology	0.447	0.553	0.632	0.368
Business	0.707	0.293	0.775	0.225
Health	0.775	0.225	1.000	0.000
Sports	0.707	0.293	0.894	0.106
Lifestyle	0.632	0.368	0.775	0.225

people that the author follows influences the perceived quality of their article, meaning that the social network of the author has a positive effect. Second, the length of an article is associated with perceived quality. However, from the results of Phase A it cannot be interpreted if longer articles are correlated with better metrics or the opposite is true in this case.

Third, images seem to be of great significance in terms of positive feedback. Both the number of pictures and image-to-text ratio are among the top five important features in most data sets. Fourth, sentiment has an impact on the crowd's endorsements, since the polarity of a text, which translates to the sum of positive and negative scores, appears to be strongly correlated to perceived quality.

Interestingly, the findings suggest some differences depending on the category. Specifically, the existence of external links in the body of the article appears to be somewhat important in the case of Lifestyle, while beautiful phrasing is fairly relevant to Health and Sports categories. Meanwhile, readability (Flesch score) is associated with perceived quality when it comes to News articles. Notably, subjectivity is associated with high claps only in Business and Sports and tone only in the News data set. Likewise, the presence of bullet points in the text is quite significant in News, Technology, and Lifestyle, whereas grammatical modality is slightly correlated to the Sports category. Surprisingly, the only temporal feature that ended up having a small degree of significance is the hour, which came up in Health and Lifestyle data sets.

According to the results, articles that belonged to different categories presented some slight

differences regarding the importance of the predictive features. Additionally, one of the main objectives of this study was to understand the relationship between the various features and their role in the identification of the perceived quality of the articles. Accordingly, the features of each article's category were sorted based on their importance and the top-10 and top-5 ones were compared with the baseline in Figure 6.2. The main concern was to identify their similarity, and enrich the understanding regarding the influence that different categories have on the features that drive the value of the perceived quality of the articles. Calculating the cosine similarity and distance of each vector from the baseline (see Table 6.4) some interesting observations were drawn: (a) An article's category affects the importance of the engagement metrics that reveal its perceived quality differently, and (b) the importance of the engagement metrics is positively correlated with the similarity that each genre has with the baseline vector (the average similarity of vectors for top-10 is 65% while for the top-5 raised to 80% – with the category Health to have even a correlation of 1 with the baseline). Such results dictate that perceived quality could not be quantified using a predetermined or fixed classification of the engagement metrics across articles' categories, but rather should be treated as a situation-specific factor considering always the contextual characteristics of an article. The latter might be able to influence the importance and the hierarchy of the features, assigning to perceived quality a different semantic understanding of what it represents.

Given the aforementioned understanding, in the next section, the engagement metrics of each category are explored in more detail and a set of rules is generated for each one, expecting to estimate the perceived quality of articles in a specific category before their publication.

6.1.6 Phase B - Proposing a Set of Rules with Respect to the Engagement Metrics of each Category

In this phase, the creation of a set of rules for each category of articles based on the engagement metrics discovered in the previous phase was investigated. One of the main aims of

the study is to help authors during the writing process, by providing them with some initial insights regarding the possible popularity that their articles may have if written in a particular way. The proposed rules can point towards the direction of a high probability of gaining popularity depending on the category that an article belongs to.

Method & Rule Extraction

More specifically, for each category, a threshold value was set, that results from the classification of articles based on the number of claps they have received and the definition of the 10% of them with the most claps as high popularity and the remaining 90% as low popularity articles. The articles then were given as an input to the decision tree classifier of scikit-learn library (Pedregosa et al. 2011), the tree formation of which produced the rules for writing popular articles. These rules come from following the path from the root of the tree to every leaf node. Each node contains a numeric expression that involves an article feature, for instance, number of images greater than or equal to 2.5, so the final rule for each path leading to a leaf node contains a set of arithmetic expressions associated with the logical “AND”. At this point, it has to be mentioned, that binary decision trees introduce difficulties in rule extraction, due to the fact that they create many splits where attributes appear more than once in any path from the root to the leaf. Additionally, two more features were used in this phase, the (Followed By / Followers) ratio and the numbers of words used in the title.

Leaf Node Selection

From the above it can be concluded that each leaf node can lead to a rule, but the question is whether this rule is (a) important and (b) valid and robust. These two constraints can be satisfied, if the following conditions become true:

Importance: Since the purpose of this research is to find rules for popular articles, attention is focused on leaf nodes that improve the initial probability of high popularity, which is 10%. Thus, for a rule to be important, attention is turned to leaves with a probability of high popularity at least three times that of the original, i.e., greater than or equal to 30%.

Validity and Robustness: Machine learning algorithms, such as decision trees, use a large amount of input data for training and the rest for testing. So, there are two kinds of probability of classification, depending on the data set. For the leaf node selection, what matters is the predicted probability calculated by applying the algorithm to the train set, that is, the number of observations of high popularity class that has been “captured” by that leaf over the entire number of observations captured by that leaf (during training). However, the probability for the same leaf node may differ when the decision tree algorithm is applied to the test data. If these two probabilities vary widely, it is obvious that we have an unstable node, which will not lead to a reliable conclusion. The difference between these two probabilities is called *misclassification error* for the specific leaf and it is a measure of validity and robustness of the extracted rule. Leaf nodes with misclassification errors lower than 0.1 are considered appropriate for rule extraction. It was noticeable, that those leaf nodes having a large number of articles, i.e. more than 1% of the training set, appeared to be more stable having low misclassification error. Based on the above preconditions for leaf selection, the quality of the extraction process cannot be jeopardized and the results can be generalized with certainty.

According to phase A, to answer the questions about the nature of the effects of specific important features to high perceived quality, the genre of the given article should be considered. Therefore, every news category was examined separately, with the rules produced from the decision tree being different. Also, for each category, the mean misclassification error value was calculated for all the rules produced for the articles belonging to that genre.

Results - Rules

In this section, the rules that have been produced for each category are presented:

News Articles

RULE #1 - Probability for high popularity: 88,59%

– Followers 3206 - 102679

- No. Images > 3,5
- Tone of speech > 8,587%
- Subjectivity > 0.452%
- F/F Ratio 6.406 - 155,913

RULE #2 - Probability for high popularity: 100,00%

- No. Followers > 102679
- No. Words > 421.5

misclassification error (mean): 0.0199

Tech Articles

RULE #1 - Probability for high popularity: 32,99%

- Followers 41.5 - 509.5
- No. Words > 812.5
- F/F > 1.159 - 5.905
- Tone of speech > 15.053 & Contains USA <= 0.5

RULE #2 - Probability for high popularity: 31,90%

- Followers 509.5 - 5004.5
- No. Words 767.5 - 1193.5
- Bullets present > 0.5
- F/F <= 11.504
- Anticipation > 4.202

RULE #3 - Probability for high popularity: 54,77%

- Followers > 509.5

- No. Words 767.5 - 2914.5
- Bullets present > 0.5
- F/F > 11.504
- Tone of speech > 13.563

misclassification error (mean): 0.0223

Business Articles

RULE #1 - Probability for high popularity: 34,34%

- Followers 112.5 - 802.5
- F/F > 5.164
- No. Images > 1.5
- Links present <= 0.5
- Negativity > 0.033

RULE #2 - Probability for high popularity: 32,69%

- Followers 802.5 - 22747.5
- No. Words <= 722.5
- F/F 20.341 - 347.402
- Flesch <= 16.25

RULE #3 - Probability for high popularity: 39,61%

- Followed By 802.5 - 3790.5
- F/F > 2.711
- No. Words > 722.5
- No. Images <= 2.5

- Contains USA ≤ 0.5

RULE #4 - Probability for high popularity: 57,65%

- Followed By 802.5 - 3610.5
- F/F 2.711 - 337.208
- No. Words > 722.5
- No. Images > 2.5

misclassification error (mean): 0.0278

Health Articles

RULE #1 - Probability for high popularity: 37,73%

- Followed By 162.5 - 804.0
- No. Words > 812.5
- F/F > 1.123
- Negativity > 0.065

RULE #2 - Probability for high popularity: 41,91%

- Followed By 1872.5 - 10073.5
- No. Words 351.5 - 1140.0
- Beautiful phrasing ≤ 3.104
- Tone of speech > 13.052

RULE #3 - Probability for high popularity: 52,10%

- Followers 804.0 - 10073.5
- No. Words > 1140.0
- F/F 1.489 - 18.16

- Negativity ≤ 0.106

misclassification error (mean): 0.0226

Sports Articles

RULE #1 - Probability for high popularity: 34,61%

- Followers > 1955.5
- No. Words 261.0 - 1305.0
- Modality ≤ 0.661
- No. TitleWords > 5.5
- Bullets present ≤ 0.5
- Beautiful phrasing ≤ 4.083

misclassification error (mean): 0.0379

Lifestyle Articles

RULE #1 - Probability for high popularity: 33,85%

- Followers 826.5 - 4098.5
- No. TitleWords > 5.5
- No. Words > 1476.0
- Positivity > 0.123

RULE #2 - Probability for high popularity: 40,49%

- Followers 4098.5 - 7976.5
- No. Words > 754.5
- F/F $> 1.031 - 120.281$

RULE #3 - Probability for high popularity: 95,06%

- Followed By > 9047.0 - 38530.5
- No. Words > 550.5
- F/F <= 230.884
- Surprise <= 4.936
- Positivity > 0.097

RULE #4 - Probability for high popularity: 99,79%

- Followers 9047.0 - 38530.5
- No. Words > 550.5
- F/F > 230.893
- Sadness <= 9.0
- img/word > 0.001
- Beautiful phrasing > 0.757

misclassification error (mean): 0.0188

Discussion of the Results

As can be observed from the previous section, a sufficient number of rules (from one to four) were extracted for each article category. Those rules may be used as guidelines for the authors in the writing of their next articles; indicating the probability of popularity that they may gain if they are taken into consideration. For increasing the quality, a specific number of leaf nodes was exploited for structuring the rules, rather than the whole data set, which might bring some variation in the engagement features' importance. More specifically, only one feature was involved in all the rules, the number of user's Followers. Each rule aims at authors of different popularity, as it is expressed through Followers, Followees, or the (Followers / Followees) ratio. Regarding the rest of the features taken into account, the Author and Content dimensions were primarily involved.

For the News articles, the results revealed two rules, with a very high probability of gaining popularity. Authors who have a large number of followers, from 3206 to 102679, can write articles that are very likely to become popular if they add more than three images in their articles. Furthermore, they should be more objective, but not absolutely, while expressing a personal opinion is desirable. They also have to use personal pronouns (more than 8,6% of the total words) and be followed by a lot more people than they follow. The second rule for the News articles, that guarantees 100% success relates to authors with over 102679 followers and articles of at least 421 words. Probably, these are extremely popular accounts, producing material that always leads to the public's satisfaction.

Tech articles' authors who have not managed to gain a reputation on the platform can achieve high popularity if they have many more Followers than Followees, write long enough stories, use personal pronouns, and don't refer to the United States. If they have a smaller number of Followers, writing shorter stories is acceptable, but they have to use lists with bullets, and words that indicate anticipation (more than 4.2% of the total words). If they want to increase the probability of high popularity they can write medium length or extensive stories in which use list-based articles, personal pronouns, and over 14% emotionally charged words. It is very important to have many more Followers than Followees.

The first rule for Business articles is for authors with a few Followers. To produce articles with an increased probability of popularity, they must be followed by at least five times the number of people following them. Also, they have to use at least two images, more than 3.3% negatively charged words, and not to use external links. If they have gained a little more Followers (more than 802.5), there are three ways to increase the popularity of their stories. They can write short and easy to read stories, otherwise, they have to write more extensive stories, including more (or less) than three images depending on the existence of US toponyms.

Health articles' authors with a small number of Followers, or a ratio of Followers – Followees near to one, can simply increase the probability of high popularity by using negatively charged words to more than 6.5% of the total words. It seems like the negative news

in this area can win the audience. If authors are more popular, they could write an article ranging from 352 to 1140 words in which they do not use rare vocabulary to a great extent (at most 3.1% of total words) but use enough personal pronouns as a percentage over than 13%. Alternatively, they should spend time writing a lengthy read of over 1140 words in which the negatively charged words do not exceed the 10.6% of the total words. Also, they should have 1.5 to 18 times more Followers than Followees.

To write stories with an increased probability of gaining popularity, Sports authors have to be already popular and have more than 2000 Followers. In this case, they should write a text of 261 to 1305 words in which they express certainty, but not in an absolute way. The use of rare vocabulary should be limited, to absence of a list format from the text is desirable, and the title should not be short, as it should consist of at least six words.

Finally, for the Lifestyle articles, moderately popular authors should write an extensive story, longer than 1476 words, of which over 12.3% should be positively charged. Finally, they should give their story a title of at least six words. For writers who have gained a bit of a reputation, their chances of success are easily increased if they have more Followers than Followees and write over 755 words. In the field of Lifestyle, if authors have acquired a large number of Followers already, they can almost certainly write articles that will become popular if the article's length is bigger than 550 words. If the Followers - Followees ratio is less than 230, the use of words showing surprise should be limited (not exceeding 5% of total words), however, positive words should exceed 12.3%. Otherwise, sad words have to be less than the 9%, rare words more than 0.75%, and at least one image per thousand words should be included.

6.1.7 Conclusion & Future Work

This work primarily discusses the prediction task of the perceived quality of articles published on the blogging platform Medium.com, considering a framework with three dimensions, namely, author, content and context (see (Orellana-Rodriguez and Keane 2018)). The best predictor is consistently the number of followers, which is not surprising: if more peo-

ple see it, more people will clap. Yet, other important factors have been revealed that could be interpreted as indicators for the popularity of an online article. Moreover, this research is one of the first in the area, to our knowledge, that places systematic emphasis on the concept of perceived quality of various articles in different categories and explores a number of characteristics (engagement features) so to be able to forecast the article's popularity. To quantify the perceived quality of the articles, the number of "claps" was the dependent variable of the proposed model. Furthermore, to define if and to what extent the features contribute to the future online success of an article, tree-based machine learning algorithms have been employed for the prediction tasks. The findings support what has been showcased already in previous research, that indeed factors related to the author's reputation, along with content-based characteristics of an article mainly length and images, and sentiment expressed in the text can be predictors of a reader's perceived quality.

In addition, results have also shown that an article's category affects the importance of features that reveal its perceived quality differently and should be treated as a situation-specific factor given the nature of each category. The latter, constitutes a main contribution of this work, making an important step toward the understanding of the perceived quality of online articles.

Also, on online publishing platforms like Medium, the main goal is to keep users interested for as long as possible by supplying them with the most relevant content. To that end, the underlying personalization scheme leverages data from previous interactions with the individual user to tailor their experience and provide them with the perfect fit between readers' preferences and the article's actual attributes. Articles posted on Medium.com are recommended to readers with similar preferences which means that articles are mainly judged by the right community of users. The only limitation in the proposed approach originates from the unknown level of dissemination of each article based on the personalization strategy used. Each article receives a different degree of exposure and, as expected, more famous authors often reach a greater audience. As a result, the "Followers" feature is proved to be the most important factor for a successful story.

The current research delves deeper into the exact relationship between the engagement features and the popularity by examining the production of rules that improve online articles' probability of gaining high popularity on the Medium platform, based on the prediction. This could be considered as one more significant contribution of this study, where essentially, the proposed model's features have been further expanded and with the help of the machine learning algorithm, i.e., Decision Tree, a total of 17 rules have been formulated for the six categories of the articles. These rules can increase the probability of gaining high popularity from the initial 10% to a percent between 30% and 100%, by giving specific values to some of the features of the articles.

Of the 32 features that shaped the bottom-up model, 21 were those that participated in the rules that could affect the popularity of articles. Features concerning the popularity of the author, the writing style and the structure of the article were the ones that helped to gain high popularity. However, the feature that is involved in all the rules was again the number of followers, which due to its high importance, is a decisive feature of the rules. Observing the rules, it can be concluded that in most cases the author or the media outlet must have more followers than followees. Regarding the remaining of the engagement features, a universal conclusion cannot be drawn since the rules differ significantly, especially when it comes to different types of articles. This also adds an important nuance to previous work, based on which one may have assumed that stylistic features, for instance, may have played a greater role than we could empirically confirm.

This work can inspire more discussion and research towards different approaches to identify characteristics of news articles that resonate with perceived quality. Based on this study, the future work could head to two directions. One is to try a top to bottom approach and use machine learning models to explore how and if features related to journalistic quality norms shape audience evaluations of an article's quality, or digital engagement metrics are only appropriate to capture user behaviors towards the quality of popular journalism. The second direction is to further substantiate the premise that claps are indeed a proxy for perceived quality. To that end, it would be worthwhile to design a study in which human annotators judge the quality of a large number of Medium.com articles without knowing the number

of claps, and then correlating their judgment with the number of claps. This approach will be useful also in answering to what degree do unique characteristics of the individuals such as income, education level, gender, and age influence their preferences and perceptions of news quality.

6.2 Study:2 Predicting the Quality of News Articles

This study examines quality indicators on the performance level, focusing on news content features that reflect the outputs of both highbrow news organizations and tabloids. It specifically asks what the defining features of high-quality journalistic output are and whether they can be used to classify journalistic texts into higher or lower quality categories. The study is based on the theory discussed in Section 2.2, Section 2.3 and Chapter 2 and the preliminary results of this work were presented at the “COMPUTATION + JOURNALISM SYMPOSIUM 2021: Data Journalism in an Expanded Field”. At this section, a framework to evaluate online news stories based on quality criteria derived from previous empirical work is presented and these criteria are operationalized into concrete computational measures. Then, a machine learning model is used to predict the quality of digital news stories and explainable AI methods are employed to elucidate the results. Doing so yields important and new insights into the actual features that contribute to high-quality journalistic output which, hopefully, will be valuable for scholars, practitioners, and possibly even for engaged audiences.

6.2.1 Toward a Model of Quality in Journalistic Texts

As outlined above, the primary concern of this contribution is to create an inclusive theoretical model the specific quality dimensions of which enable the measurement and quantification of journalistic quality on the level of an individual news story. Therefore, the backbone of the model is structured based on the overarching journalistic norms and best practices proposed by previous journalism scholarship reviewed in chapter 2 (McQuail 2015; Harcup 2015; Kovach and Rosenstiel 2014; Rosenstiel et al. 2007; McQuail 1992). We propose

four dimensions to be central to understanding quality journalism based on text features. The model is put to the test through the development of a supervised machine learning algorithm. So in what follows we specify the dimensions' features that can subsequently be transferred into concrete measurements.

Impartiality

Although the objectivity norm serves as “the cornerstone principle” (Muñoz-Torres 2012) of journalism for more than a century it is one of the most controversial concepts that divides both communication scholars and practitioners throughout the years. Journalists must convey news in a neutral manner however, according to literature, digital journalists nowadays are less committed to the idea of neutrality (Steensen 2016; Lasorsa et al. 2012) and “instead of an objective, neutral tone, a journalist’s own voice can shine through” (Ruotsalainen 2018, p. 19). According to McQuail (1992, p. 481:482) impartiality dictates “a combination of balance (equal or proportional time/space/emphasis) as between opposing interpretations, points of view or versions of events, and neutrality in presentation” (McQuail 1992). In essence, the concept of impartiality can be broken down into and operationalized by means of a few subdimensions. First, it appears that *subjectivity* is important, in that journalists who write news stories for tabloids tend to use more intimacy in their news writing, thus they are more likely to “adhere to an ideal of subjectivity in their journalism” (Steensen 2016, p. 118). Second, some researchers suggest that *self-disclosure* is often used by popular journalists when they communicate on social media (Steensen 2016; Lasorsa et al. 2012) seeking to connect with the reader on an intimate level, by for instance using the personal pronoun “I” and similar self-reflective pronouns. Both subjectivity and self-disclosure may be considered indicators of lower quality news then. The notion of diversity has been extensively used in content analysis for measuring quality over the years (Urban and Schweiger 2014; McQuail 2005; P. M. Napoli 2011; P. M. Napoli 1999) since citing multiple viewpoints has been related to impartial reporting. For instance, Carpenter (2008) examined the *diversity* of sources presented on newspaper websites compared to online citizen journalism webpages and found that news that appeared in online newspapers had twice the amount of

sources than citizen articles and cited more authoritative sources. Therefore, the number of sources cited in a news story might be a sign of its journalistic quality.

Language Quality

One consistent characteristic in many studies is the use of language and its relationship to journalistic quality. For instance, complex or abstract notions along with ambiguity and unanswered questions are impermissible in news stories (Harcup 2015, p. 145) consequently, a story with a good *readability score*, within the range of 50 - 60, is accessible and easy to understand, while a story with a high score might be dull or a low score, for instance under 30 might be incomprehensible (Flesch 1948). Additionally, the ratio of visual to textual content was also considered a significant quality indicator by (Bogart 1989). Furthermore, one of the steps of the “Magic Formula” (Rosenstiel et al. 2007, p. 116:117) to improve news and generate larger audiences was “Make Important Stories Longer”. The average length of articles can potentially provide insights into the information depth of a news story according to Quandt (2008, p. 724) who observed that quality newspapers like the *New York Times* and the *Times* “prefer to focus on news depth rather than on news variation The difference is true for the title since long headlines have been linked to low-quality and fake information in an attempt to lure readers away from the small and repetitive text in the article and only read the argument in the title instead (Horne and Adali 2017). Even though there is a general consensus on a *language of news* - a basic grammar of journalism - such as past tense, short paragraphs, active sentences, and so on, there are stylistic variations like the number of adjectives or the amount of color allowed into the news story (Harcup 2015, p. 144). For instance, tabloid newspapers are rarely afraid of using adjectives to describe situations or people (Harcup 2015, p. 149). Other characteristics related to the language used are the existence of numbers in a news article that could refer to concrete figures, and give credence to the story, while a published high-quality news article should not include typos and misspellings.

Entertainment

Previous studies on audience expectations of excellent journalism reveal that individual members of the audience enjoy journalism, emphasizing the “sense of delight” (Costera-Meijer 2012). In contrast, when it comes to journalists’ role conception, providing entertainment is not portrayed as an important goal of journalism (Van Der Wurff and Schoenbach 2014), and is rather connected to the concept of tabloidisation, that according to its critics sacrifices information for entertainment (G. Turner 2013). From a news values perspective, stories that may entertain the audience are easily getting picked by news desks according to Harcup (2015), who puts entertainment within a journalist’s scope of responsibilities. The following categories come from journalistic best practices (Harcup 2015, p. 118:120). First, if there is a sex angle to a story, a crime, or an unfolding drama it is regarded as more entertaining and therefore is more likely to be used. Second, animal stories are considered entertaining and, according to (Darnton 1975, p. 190), “go over very well with the city desk”. Third, stories about famous people and their actions are engaging because audiences may identify with media characters (J. Cohen 2001). Therefore, it is safe to assume that entertaining stories point to lower quality.

Emotionality

As thoroughly discussed in Chapter 3.2, the presentation of people and situations in an emotive, polished, and dramatic fashion in news stories with the intent to garner empathy (Graber 1994), shock the audience (V. Popović and P. Popović 2014), or arouse people’s curiosity has been characterized as sensationalism (Pantti 2010; Slattery and Hakanen 1994) and is generally viewed as unacceptable. Scholars have heavily criticized sensationalism as a way of manipulating people’s public opinion. Emotionality in literature is considered the opposite of objectivity, is connected to tabloidization, vulgarization, and commercialization, and poses a threat to serious quality journalism (Harrington 2008; Pantti 2010; Peters 2011; Sparks 1998; Agarwal and Barthel 2015). However other studies suggest that emotionality can be used in case of disaster coverage (Kovach and Rosenstiel 2014), can convey important, and comprehensive information and does not necessarily translate into tabloidi-

sation (Meijer 2003). For instance, Kilgo et al. (2018a) and her colleagues investigated 400 online-native news organizations in the United States and proved that hard news were also treated sensationally. Additionally, Wahl-Jorgensen (2013, p. 141) who examined Pulitzer Prize-winning news stories from 1995 to 2011 concludes by saying “emotional story-telling is a driving force behind award-winning journalism, with the aim of drawing the audience’s attention to complex topics of social and political import” although reporters never talk about their own emotions.

News outlets like the *New York Times*, *ESPN* and *USA Today*, produced ad products indicative of how emotionality in news stories can be leveraged for making profit (Rick 2019). The *New York Times*’s “Project Feels”, was developed by the Data Science Group at the Times, and harnesses the power of state-of-the-art machine learning models, to predict the combination of emotions each article might evoke to the reader. A list of moods such as inspired, amused, boredom, love and so on, was used by the volunteer annotators who created a corpus of 150K news articles for training the algorithms. Based on the model-predictions the Times created premium advertising spaces, targeting ads depending on the reader’s feelings producing strong revenue results (Rick 2019; Spangher 2018).

In this study, we propose that the journalistic quality question can be examined in great detail if it is treated as a classification problem by using automatic machine learning approaches to find hidden patterns that confirm or refute the dimensions of the theoretical framework, or even provide new insights on important indicators. In general, many machine learning algorithms, such as neural networks, are considered a black box, hence there are difficulties in explaining the underlying patterns in the data. Therefore, an important part of the study focuses on ways to interpret machine learning algorithms and especially tree-based (Chen 2018; T. Chen and Guestrin 2016) by analyzing their algorithmic mechanisms and drawing generic insights based on rule extraction and feature importance. Doing so yields conclusions not only about the usefulness of the entire model but more specifically on certain characteristics of journalistic texts that are more or less important in classifying news stories as high quality. Hence the overarching question is to what degree the four dimensions reviewed above contribute to good classification. Tentative expectations based

on the literature would be that higher degrees of emotionality and entertainment lead to a lower likelihood of high quality, while higher degrees of impartiality and language quality result in a higher probability of high quality and vice versa.

6.2.2 Methods and Data

The analysis was conducted on a corpus of 2K articles from 15 different English-language newspapers, seven highbrow newspapers, namely the *Guardian*, *the Independent*, *Washington Post*, *Politico*, *the New York Times*, *CNN*, *Reuters*, and six tabloids the *Daily Mail*, *the Daily Mirror*, *the Sun*, *the Daily Star*, *the Daily Express*, and the *New York Post*. The dataset consists of news coverage during the whole year of 2019, covering news articles that appeared on the online version of the newspapers under two categories, namely news and politics. The reasons for selecting these two categories are first the fact that stories on politics often appear in the main news category and second that these two categories exist in all 15 news outlets. This dataset was manually annotated by three journalism students enrolled in the master's program "Digital Media and Interactive Environments" at the National and Kapodistrian University of Athens. The annotators were first introduced to a guide created by the American Press Institute¹ based on the essential elements of journalism from the book of Kovach and Rosenstiel (2014) that is being used as a manual for several news outlets and then were instructed to evaluate whether the news stories were of quality or not, by answering a yes-or-no question. The dataset given to the annotators included only the title and the body of each article and was stripped from any keyword indicating its origin, such as the names of any media organizations, and call-to-action sentences like: "The Daily Star's FREE newsletter is spectacular! Sign up today for the best stories straight to your inbox". In general, the three annotators completely agreed on the quality evaluation only in 54.78% of the cases, resulting in a Cronbach Alpha Reliability Coefficient of 0.66.

The focus of the current study is on the article level, meaning that the centre of attention is the unique textual characteristics of an individual newspaper article, irrespective of the

¹<https://www.americanpressinstitute.org/journalism-essentials/what-is-journalism/elements-journalism/>

category. The total amount of articles before the cleaning was 2000, but 1935 cases remained for inclusion in the model building. For the preparation of the dataset, several Python scripts were used for processing stop words, and nonstandard words, removing NaN values, and HTML code from the texts.

In general, text mining involves the use of techniques to extract and analyze information from unstructured text data. One common approach is to use vectorizers, such as Bag of Words or TF-IDF, to convert text into numerical representations that can be used as features in classification algorithms. These vectorizers allow for analyzing the frequency of words or groups of words (such as bigrams or trigrams) in the text, and use this information to train classifiers that can predict the class or category of a given piece of text. Text vectorization is effective for basic text classification tasks, however, as the complexity and difficulty of classification tasks increases, more sophisticated feature engineering is necessary. Therefore, in this study, simple techniques for trivial representation of text were used as a baseline. For the creation of the sophisticated features that correspond to the theoretical framework various text analysis techniques, Python packages, and natural language processing libraries were utilized. These included tokenization, stemming, lemmatization, and part-of-speech tagging, along with many lexicons which required the original form of words. In other cases, the original text was used instead of the preprocessed text to identify adjectives or assess the readability of a text.

The dependent variable (categorical) was determined by the annotators, therefore the articles with at least two positive answers to the quality question were scored with a value of 1 and the remaining with a value of 0. The independent variables of the model were chosen and operationalized based on the dimensions of the theoretical framework that according to the studied literature can define journalistic quality (see Table 6.5 for details on their measurement). Given that our goal is the investigation of general features across a large corpus, computer-assisted text analysis methodology is used through concept measurement.

Table 6.5: Creation of the features.

Dimension	Measurement
Impartiality	
Subjectivity	Subjectivity Lexicon (Wilson et al. 2005) which includes a list of subjectivity clues is used to calculate weak and strong subjectivity in texts. The score is the total number of subjectivity clues divided by the number of words in the text.
Self disclosure	Total count of the self-reflective pronouns divided by the number of words in the text.
Diversity of sources	Total count of the unique sources presented in a news story.
Language Quality	
Readability of text	Flesch Reading Ease Score (Flesch 1948) with a score ranging from 0 to 100; the higher score the easiest the text is to read. For the calculation the py-readability-metrics package ² was used.
Article Length	Total count of words except for stopwords.
Adjectives	Total count of adjectives.
Typographical errors	Autocorrect Python ³ library corrects mistakes in the text. Then the pairs of sequences are compared to calculate the difference.
Numbers	Total count of numbers
Images	The ratio of illustrations to text (Bogart 1989).

Continued on next page

²<https://pypi.org/project/py-readability-metrics/>

³Autocorrect Python Library: <https://pypi.org/project/autocorrect/>

Table 6.5 – continued from previous page

Dimension	Measurement
Headline words	The total number of words in the title. (Horne and Adali 2017).
Entertainment	We created four lists: a) sensual words, b) animals, c) crime, and d) celebrities. For the latter we included all the names of the 100 most influential people in the world that appeared in the TIME magazine since 2004. All the counts were used as separate features.
Emotionality	
Emotional head- lines	Textblob ⁴ sentiment analysis library works well with short text and can detect clickbait titles that have extremely negative or positive words, ranging from -1 to 1 respectively to detect negative or positive titles.
Emotions	NRC Affect Intensity Lexicon (Saif M. Mohammad 2018) includes 6k words annotated with an intensity label for each emotion using crowdsourcing that provides binary categorization for emotion The AIL measures intensity scores for four basic emotions: anger, fear, sadness, joy, based on theories of emotion (Plutchik 1980a).

To test our general assumption that certain characteristics of a news story are associated with quality, we used the corpus of newspapers as an input for the machine learning model that used the proposed features for training, to predict the quality of an article. To shed light on the hidden decision mechanisms, after the initial modeling follows a series of anal-

⁴<https://textblob.readthedocs.io/en/dev/>

yses to interpret the machine learning models in order to examine the importance of each dimension in greater detail.

6.2.3 Analysis

For the evaluation of the importance of the proposed dimensions concerning journalistic quality, a binary classification task was applied to build a model that could provide us with the importance of each feature along with an explanation of the system's predictions to improve our understanding of what comprises a high-quality news story. For building the machine learning algorithms we employed multiple classification models from the scikit-learn Python library (Pedregosa et al. 2011). Specifically, we used seven different approaches, namely the Naive Bayes, Logistic Regression, Support Vector Machine, K-Nearest Neighbors, Decision tree, Random forest, and XGBoost classifier, which uses gradient boosting to build and add new trees to the prior model for better prediction. For the experiments, 80% of the news articles were used for training reasons and the other 20% for testing; for the evaluation of the three different classification methods, the weighted average F-measure (F1) was implemented, which is a mean of precision and recall.

Initially, for the baseline models three techniques were used, namely the bag of word (bigrams and trigrams), the TF-IDF, and the BERT model. More specifically, the bag of words model was used, which is a technique to extract features from text data. It involves representing a piece of text as a numerical vector, where each element in the vector corresponds to a specific word in a dictionary, and the value of the element represents the frequency of that word in the text. Usually, bigrams and trigrams which are sequences of two and three words, respectively, that occur together within a document, are used as features in text classification tasks. Similarly, the TF-IDF (Term Frequency-Inverse Document Frequency) which is a statistical measure used to evaluate the importance of a word in a document within a collection of documents was used as a feature representation method. After the vectors for each news story were created using the different methods, they were used as an input for the machine learning algorithms. Additionally, the deep learning model BERT, which is trained on a large, unannotated corpus of text and can learn contextual relation-

Table 6.6: Baseline Models

	Full dataset			100% Agreement		
	TF-IDF	2-gram	3-gram	TF-IDF	2-gram	3-gram
Logistic Regression	0.78	0.77	0.69	0.9	0.87	0.76
Naive Bayes	0.72	0.71	0.7	0.78	0.77	0.75
Support Vector Machine	0.77	0.77	0.69	0.89	0.88	0.78
K-Nearest Neighbors	0.65	0.65	0.65	0.73	0.73	0.73
Decision tree	0.63	0.66	0.65	0.82	0.77	0.74
Random Forest	0.78	0.74	0.69	0.89	0.84	0.74
XGBoost	0.7	0.76	0.7	0.88	0.86	0.82

ships between words, was used. However, the highest accuracy it achieved was an F1-score of 0.66. It is possible that the dataset was too small for this type of model.

The results presented in Table 6.6, demonstrate that the baseline models perform well on the given task. However, the interpretability of these models is limited. For example, when examining the most important words for the tf-idf models, the terms “president”, “government”, and “health” are among the top words. While these words may be relevant to the task at hand, they do not necessarily provide insight into what makes a high-quality article. Similarly, the bag of words (BOW), and BERT models do not offer any clear insights into the characteristics that define a high-quality article. This highlights the importance of considering both performance and interpretability when evaluating NLP models.

In regards to the dimensions of the framework, after the application of the different classification methods described above, the Random Forest classifier proved to perform best with an F1-Score of 80%. The other models yielded F1-Scores ranging from 72% to 79%. The accuracy of the models signifies that the quality framework is effective, but the algorithms misclassify a fair number of samples. However, given the fact that the human annotators did not present a unified perception of the quality as well, and in 45.22% of the cases there was a two out of three agreement on the data samples, there is a chance that the misclassification corresponds to the medium quality of the training data (Song et al. 2020). Therefore, we decided to run the model again, this time using only the news stories with a 100% agreement. The smaller dataset consists of 1060 news stories, 583 low-quality and 477 high-quality. This time, the results of all algorithms, as shown in Table 6.7, were significantly improved, with

the best model being the Random Forest Classifier. Hence our preferred classifier was able to classify correctly 93% of news stories into the high or low-quality category based on the features that were identified.

Table 6.7: Accuracy scores of the different classification algorithms

Model	Full dataset	100% Agreement
Logistic Regression	0.79	0.91
Naive Bayes	0.76	0.88
Support Vector Machine	0.78	0.92
K-Nearest Neighbors	0.75	0.86
Decision tree	0.72	0.83
Random forest	0.80	0.93
XGBoost	0.78	0.92

Understanding Model Predictions

Although four models achieved over 90% accuracy on the given task, it was decided to focus on tree models, namely Random Forest and XGBoost because, according to previous work (Lundberg et al. 2020; Lundberg and S.-I. Lee 2017; T. Chen and Guestrin 2016), these types of classifiers can be explained by providing in-depth interpretations of the model predictions. Therefore the ELI5 Python library (Mikhail Korobov 2016) for inspecting black box estimators was used to examine the unique contribution of each variable to the model's predictions. When a random forest classifier constructs its decision trees the weight of every feature is determined internally, thus the significance of the predictor variables is automatically measured.

In general, as shown in Figure 6.4, the level of linguistic difficulty, and the length are the most influential for the quality of a news story. This is consistent with previous work (Tolochko and Boomgaarden 2018) which showed that quality news outlets have a higher level of difficulty than political blogs and tabloid news, due to their news stories demonstrating a more complex syntactic structure. Interestingly, features related to the ratio of illustrations to the words, self-reflective pronouns, and the diversity of the sources are also important predictors. Finally, the adjectives and numbers in the body of an article, followed by how many words exist in the title, along with the emotions of trust and joy, are of lower significance.

The weights of the features show what matters most for the classifications and seem to relate well with the theoretical dimensions of quality as discussed previously. However, from Figure 6.4 we cannot be sure which characteristics act positively or negatively to distinguish the high-quality articles from the low-ones.

The dimension of Language Quality most significantly contributed to correct classifications and in particular, the readability, length, presence of adjectives, numbers, and the ratio of images to the text were important features. The Impartiality dimension then appears to be the second most important with the diversity of sources, and first-person pronouns while features from both the Emotionality and Entertainment dimensions appear in the top ten. This already is very insightful in terms of how the algorithm classifies texts. We still need, however, a more meaningful analysis of how the classifier arrived at the predictions and the exact role played by the more important predictors in the model.

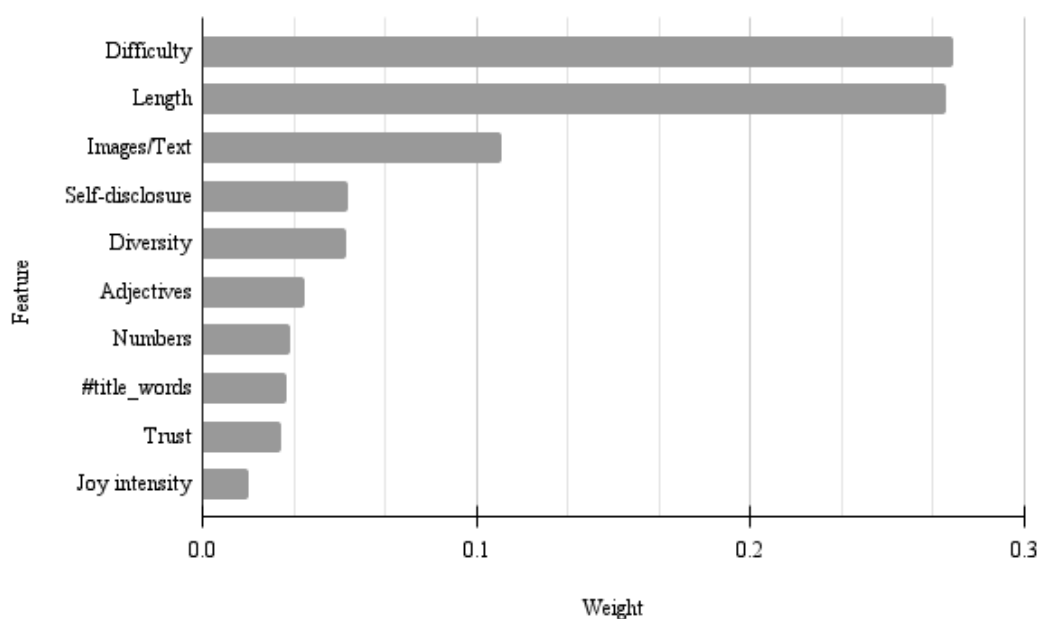


Figure 6.4: Feature Importance from the Random Forest Classifier

The results from the Random Forest Classifier on the whole dataset, with an F-1 of 0.80 are similar to the ones presented in Figure B.5. The graph of the feature importance of the whole dataset can be found in Appendix. Furthermore, the confusion matrix of the XGBoost Classifier predictions is shown in Figure 6.5.

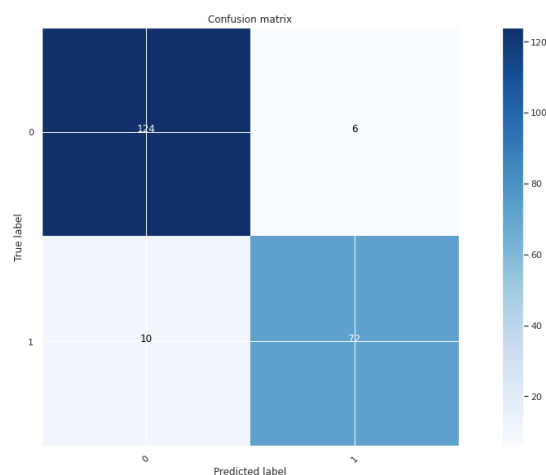


Figure 6.5: Confusion Matrix of the XGBoost Classifier for the small dataset

Rule Extraction

After examining the importance of different features, to strengthen our understanding of the predictions we investigate the graphic representation of a single decision tree from the XGBoost classifier with an F1-Score of 92% to propose a set of rules with respect to the journalistic quality. The problem with collecting all the rules from a decision tree is that as the tree grows deeper it becomes more complex and hard to comprehend. Besides, the constraints are conjunctive along one path and different paths may provide contradictory information and require extra effort for post-processing the initial set of rules. To avoid these issues we apply measures of confidence and support and opt for the “most important” rules. For this, we used the dtreeviz⁵ Python library, which shows how the algorithm arrives at a certain prediction by illustrating the decisions made in the feature-space splits of the leaf nodes (Parr and Howard 2018). Based on the above preconditions for node selection, the quality of the extraction process cannot be jeopardized and the results can be generalized with certainty. As shown in Figure 6.6, the tree begins from its root with the most important feature “Length” acting as the parent node creating a left child when the answer is less than the split value (wedge in the X-axis) and a right child for greater or equal. The histograms depict each variable’s feature space distribution, the Y-axis is the sum of the samples from the two categories that are stacked together, while the node’s size is in proportion to the total number of

⁵<https://github.com/parr/dtreeviz>

samples in each leaf. Leading toward the bottom, the nodes become purer, as the decision splits create finer regions until the samples are completely separated into two buckets, the low-quality (0) and high-quality (1). For rule extraction, only the four large final nodes will be taken into consideration, visualized as pie charts for better comprehension.

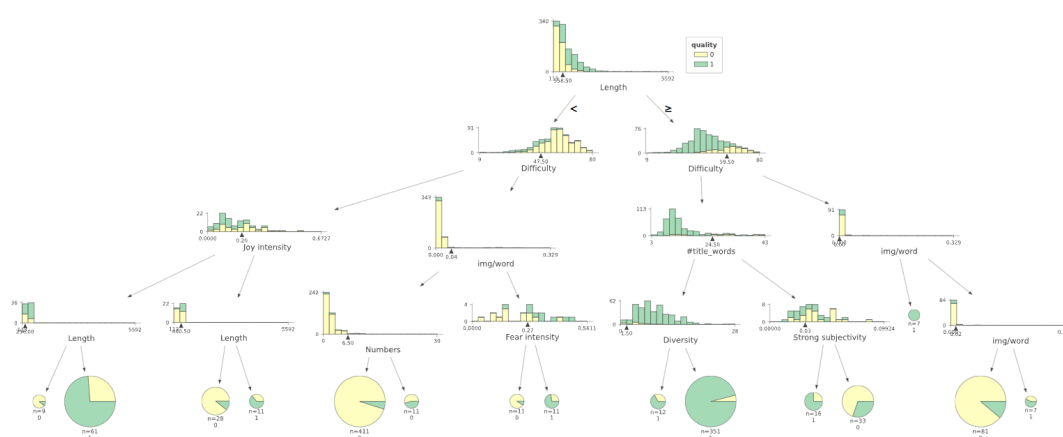


Figure 6.6: Visualization of one decision tree of the XGBoost classifier

The rules extracted from the model may be used as guidelines for journalists and editors alike for crafting their news stories, indicating the probability of an article being of high quality if the rules are taken into consideration. In the light of AI explainability, the results reveal four main rules able to discriminate between news stories based on seven important sub-dimensions of the framework. Derived from these rules shown in Table 6.8, a high-quality article is longer than 558 words, it is not very easy to read, its headline includes no more than 24 words, and cites at least two named sources in the story. According to the second rule (Table 6.9, if a news story is less than 558 words, is fairly difficult to read, includes less than 20% words expressing joy, and has a minimum length of 256 words might also fall into the high-quality category.

Furthermore, we used the ELI5 library to extract from the model the explanations of two random news stories that were correctly classified by the Random Forest algorithm. A graphic representation of the framework can be seen in Figure 6.7, in which the quality score on every particular dimension is shown for the two news articles.

Table 6.8: The rules extracted from the large final leaves of the XGBoost classifier for the High Quality class

1st Rule for	Explanation High-quality (class 1)
Length ≥ 558.5	Minimum length of 558 words
Difficulty < 58.5	Text with less than a readability index of 60 is somewhat difficult to read
Title Words < 24.5	Less than 25 words in the headline
Diversity ≥ 1.5	Two or more named sources in the story
2nd Rule for High-quality (class 1):	
Length < 558.5	Less than 558 words in the text
Difficulty < 47.5	Text with a readability index under 40 is fairly difficult to read
Joy intensity < 0.2	Less than 20% words expressing joy
Length ≥ 256	No less than 256 words in the text

Table 6.9: The rules extracted from the large final leaves of the XGBoost classifier for the Low Quality class

1st Rule	Explanation for Low-quality (class 0)
Length < 558.5	Minimum length of 557 words
Difficulty ≥ 47.5	Plain English, text is somewhat easy to read
Images/Text < 0.04	Less than 4 illustrations per 100 words
Numbers < 6.5	Less than 7 numbers
2nd Rule for Low-quality (class 0):	
Length ≥ 558.5	Minimum length of 558 words
Difficulty ≥ 58.5	Text with more than a readability index of 60 is considered easy to read
Images/Text < 0.02	Less than 2 illustrations per 100 words

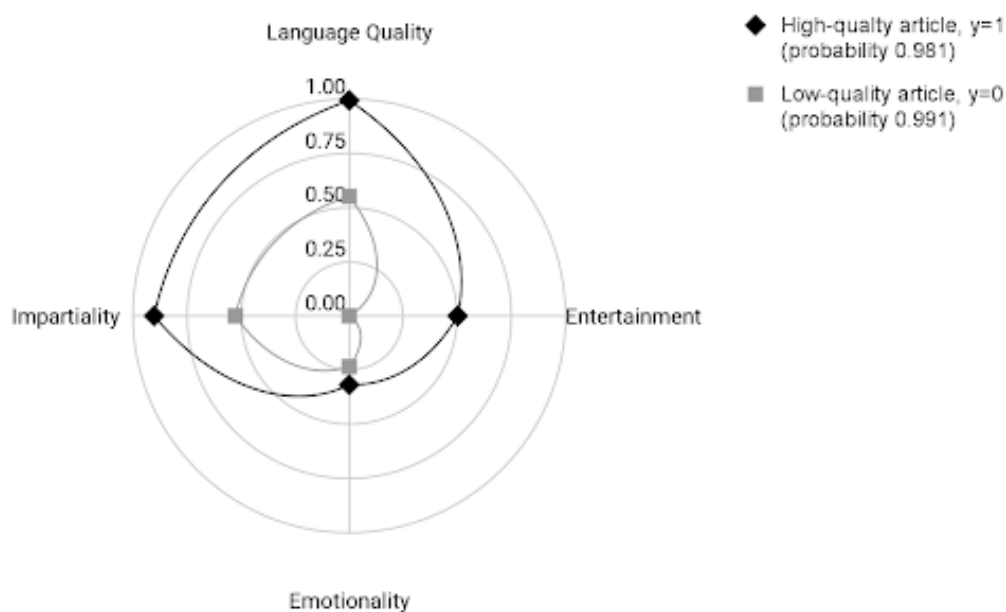


Figure 6.7: Visualization Example of Quality Dimensions using the Theoretical Framework

6.2.4 Discussion & Future work

The question of news quality was usually answered using human judgment (Arapakis et al. 2016; McQuail 2015; Kovach and Rosenstiel 2014; Urban and Schweiger 2014; P. M. Napoli 2011; Carpenter 2008; Rosenstiel et al. 2007), which has yielded important insights to the literature. This study takes another strategy and proposes the test of a theoretical model that considers the importance of textual features for the identification of quality in journalism. It does so by drawing on a broad corpus of news articles and an inclusive list of computationally measurable text features. Accordingly, our study employed content-based metrics drawn from theoretical constructs in an attempt to model quality in digital journalism with high accuracy. To this end, we created a newspaper corpus that was used to train a novel classification task. A machine learning investigation of four different dimensions - *Impartiality*, *Language Quality*, *Entertainment*, and *Emotionality* - to human annotated digital articles from traditional high-quality media and tabloids provided us with a robust understanding of quality elements.

Previous work methodologically close to this study has been conducted investigating how content-related features of news stories correlate with good quality journalism (Arapakis et

al. 2016; Tang et al. 2003), newsworthiness (Trilling et al. 2017), popularity (Keneshloo et al. 2016; Arapakis et al. 2014), and journalistic values (Choi et al. 2021) showing that automated metrics can measure rather well various characteristics of the text like their readability or diverse vocabulary. However, the approach of traditional news content analysis translated into machine-coded features will not replace the human judgment necessary to analyze thoroughly all the elements of journalistic quality, it can free up time for the researchers to focus on more sophisticated tasks (Flaounas et al. 2013).

The results of this study support a number of quality characteristics dictated by the theoretical foundations of the notion of journalistic quality. More specifically, the dimension of Language Quality almost monopolizes the quality forecast, with readability, length, and image to text ratio metrics to largely determine which news stories will follow the path toward the high- or the low-quality class. This is much in line with previous work in journalism studies that highlight the role of good readability, fluency, and rich vocabulary (Tolochko and Boomgaarden 2018; Arapakis et al. 2016; McQuail 2015; Tang et al. 2003) along with providing lengthy articles (Rosenstiel et al. 2007; Tang et al. 2003). Also, the numbers, the adjectives, and how long the headline is, were distinguished as having some discriminative power. Moreover, the Impartiality dimension is highlighted with the diversity of sources (Carpenter 2008; Choi et al. 2021) and self-disclosure being ranked among the top ten important features of the model. This is of no surprise since the essential attributes of Impartiality in the context of quality journalism were recognized early on by many scholars (McQuail 1992; Westerståhl 1983; Tuchman 1972) and is one of the quality criteria widely agreed upon by scholars, professionals, and audiences alike (Urban and Schweiger 2014; Rosenstiel et al. 2007; Gladney 1996). Additionally, the dimension of Emotionality (Peters 2011; Pantti 2010; Harrington 2008; Sparks 1998; Slattery and Hakanen 1994), and particularly words associated with trust and joy, was identified by the model as being significant in determining quality. Interestingly enough, the dimension often related to low-quality news, namely Entertainment (Van Der Wurff and Schoenbach 2014; G. Turner 2013) is not present in the top-quality predictors. Finally, our “recipe” for writing a high-quality news article is aligned with the “magic formula” (Rosenstiel et al. 2007), especially concerning, quoting

various sources, and writing lengthy articles.

Nonetheless, the findings presented here must be understood in conjunction with some limitations. A shortcoming of this work that we would like to acknowledge is the fact that the most important features of the model are rather language-related and thus it is unclear how this would travel to a multilingual or non-English setting. This limitation can be addressed by training the model with a more diverse corpus that incorporates other languages in order for the model to generalize better and consequently provide a more comprehensive view of quality identifiers outside the English-language boundaries. Further, we focus on articles that appeared under politics and the generic category news of the selected publications which is mostly hard news, therefore the results might not correspond to soft news like lifestyle and sports, since quality indicators are category-specific (Sotirakou et al. 2019). Also, we experimented with two datasets, a large dataset that included three levels of quality, namely low, medium and high, and a smaller one that contained only the high- and low-quality articles. We decided to focus on the smaller dataset and discard the medium-quality stories to obtain more accurate results from the machine learning model. Despite this limitation, this study proves for the first time to our knowledge the usefulness of explainable AI in disentangling some of the most significant quality criteria in journalism studies, and it facilitates a refined understanding of the relationship between these criteria and high-quality news stories.

6.3 Study: 3 An Analysis of News Engagement using AI

The digitalization and datafication of everyday life and forms of organization have brought about a new way for the media industry to interact with their audience through the use of social media and its “like economy” (Gerlitz and Helmond 2013). Before, the size and success of a media company were determined by the circulation of a newspaper or the TV ratings. However, in the digital age, the intangible aspect of online journalism necessitates different metrics to measure its success, such as the level of audience engagement with the news product. This has caused journalists to become more conscious of how their reporting will

be received online, often tweaking their stories in order to gain prominence in the crowded digital environment. This increased focus on who the target audience is for the news (Nelson 2020; Tandoc and Vos 2016) is indicative of the shift in the media industry towards social media, as journalists strive to capture the attention of the public.

The current study focuses on the characteristics of a news story before its publication and proposes a framework that can predict engagement on Facebook. The aim is not only to build a machine learning model for prediction, but to use explainable AI methods to disentangle the results and recommend best practices to journalists. More specifically, the study suggests treating engagement prediction as a classification problem and implementing machine learning algorithms to explore the data and identify patterns that may not be recognizable by humans. In this study, the focus is on tree-based algorithms, due to their ability to provide insightful visualizations for rule extraction and feature importance (T. Chen and Guestrin 2016). Drawing generic conclusions from a model with good accuracy can potentially assist journalists and publishers in producing more relevant and customizable content for social media platforms. To do this, the study will first present a framework based on the literature with specific dimensions to quantify audience engagement with news stories on Facebook. Second, it will check the accuracy of the proposed model and discover the most important features. Finally, it will recommend some guidelines for journalists on how they can combine certain elements when writing their article to increase the chances of success on social media. This work was presented at the 73rd Annual ICA Conference 2023, International Communication Association, at the Computational Methods theme.

6.3.1 Methods and Dataset

To investigate the influence of certain inherent attributes of an article on its engagement on Facebook, a corpus of online news stories was used as input for the machine learning algorithms. Then, four tree-based models were created, each one predicting a different engagement metric (Likes, Shares, Comments, and the number of Total Interactions). In order to create the dependent variable, the data was separated into two buckets. The low engagement bucket consisted of articles belonging to the 5th percentile and had a value of 0, and

the high engagement bucket was comprised of articles belonging to the 95th percentile and was labeled as 1. After searching for the model with the best accuracy, a series of explainability methods were used to acquire a set of rules able to predict audience engagement. The independent variables quantified some of the engagement and quality criteria from the literature that could capture audience engagement with news on Facebook. These criteria were operationalized using the Python programming language, especially with text analysis and Natural Language Processing methods. To gain deeper insights, the really popular articles were compared to the completely irrelevant ones in terms of engagement. The results could enable news organizations to craft and optimize content based on their audience preferences and interests.

The corpus was comprised of news stories that appeared under the general category “News” or “Politics” in the digital version of nine English-language newspapers, five highbrow newspapers (*the Guardian, the New York Times, the Independent, Washington Post and Politico*) and four tabloids (*the Daily Mail, the Daily Mirror, the Sun, and the Daily Star*) during 2019. Initially, the corpus contained 112,707 news stories, but in order to ensure an equal sample size across both elite newspapers and tabloids, a sample of 73,036 news articles remained. To gather additional data, the CrowdTangle platform was used to download the corresponding Facebook metrics for each news story and integrate that information into the dataset. Before running the analysis, the body of the news article was preprocessed, tokenized and lemmatized for lexicon implementation, empty values and stopwords were also removed from the dataset.

The smaller version of the framework of the study of quality with the inclusion of the VAD lexicon, was used as the basis for the development of features for this study that were used to test for audience engagement to news stories on Facebook, as shown in Table 6.10. To enhance the study, the VAD lexicon was employed to measure the emotional state of a person (Russell 1980). The model was coded into an emotional lexicon by S. Mohammad (2018) which has three dimensions: valence (pleasure), ranging from feeling happy to unhappy; arousal, which measures the level of activity a person feels, from sleep to excitement; and dominance, which gauges how much in control a person feels, ranging from submissive to

dominant. This lexicon has been used in past studies to explain user motivations for engaging with the news (Dafonte-Gómez 2018; Guerini and Staiano 2015).

Table 6.10: Creation of the features

Dimension	Measurement
Impartiality	
Subjectivity	With the Subjectivity Lexicon (Wilson et al. 2005) that consists of a set of subjectivity clues, three features were created for title, Facebook headline, and the body of the article.
Diversity of sources	Total count of the unique sources presented in a news story.
Language Quality	
Readability of text	Flesch Reading Ease Score (Flesch 1948) with a score ranging from 0 to 100; the higher score the easiest the text is to read. For the calculation the py-readability-metrics package ⁶ was used.
Article & Headline Length	Total count of words except for stopwords.
Entertainment	
	We created four lists: a) sensual words, b) animals, c) crime, and d) celebrities. For the latter we included all the names of the 100 most influential people in the world that appeared in the TIME magazine since 2004. The lexicons produced four separate features.
Emotionality	
Emotional headlines	Textblob ⁷ package to detect very negative or positive titles.
Emotions	NRC EmoLex was used (Saif M Mohammad 2017) to capture the following emotions: Anger, Fear, Sadness, Joy, VAD for Valence, Arousal, and Dominance (Saif M. Mohammad 2018).

6.3.2 Analysis and Results

The analysis was based on tree models because according to previous work (Lundberg et al. 2020; T. Chen and Guestrin 2016) these types of classifiers can be explained by providing in-depth interpretations of the model predictions. Therefore, three different models from the scikit-learn Python library, were implemented for a binary classification task, namely a Decision tree, a Random Forest, and an XGBoost. From the dataset, 80% was used for training and the remaining 20% for testing, while the F-measure (F1) was the preferred accuracy

⁶<https://pypi.org/project/py-readability-metrics/>

⁷<https://textblob.readthedocs.io/en/dev/>

method.

Three experiments were conducted using three classifiers to predict Likes, Shares, and Comments. The model with the highest accuracy was XGBoost, with an F1-Score of 91% for predicting Likes (see Table 6.11). Permutation importance was obtained using the ELI5 Python library for “Inspecting Black-Box Estimators” (Mikhail Korobov 2016) to identify the most important features. The Language Quality dimension has been identified as the most important one, with the number of words in the title being the highest predictor with regards to the three different engagement metrics. Additionally, the difficulty and length features appear highly significant in the permutation importance table. All four dimensions of the framework were found to be significant for the model. Impartiality with diversity and subjectivity in the body of the article and the Facebook headline were found to influence prediction, while Emotionality was revealed to have anticipation, dominance, arousal, and valence as important features. Further, the Entertainment dimension showed that famous people hold importance. The significance of some features was found to change depending on the target variable, which is in line with the literature (Sora Park et al. 2021). Specifically, the number of celebrities referred to in the text was found to best predict Comments and Likes, while it did not influence Shares. Meanwhile, emotion of dominance was found to be crucial for Likes, fear was only relevant for Likes, and positivity for Comments.

Table 6.11: Permutation Importance for the top ten features

Comments		Shares		Likes	
F-1 Score:0.864		F-1 Score: 0.847		F-1 Score: 0.913	
Number of samples:7161		Number of samples:7199		Number of samples:7249	
Weight	Feature	Weight	Feature	Weight	Feature
0.0951 ± 0.0104	Title words	0.1011 ± 0.0155	Title words	0.0934 ± 0.0120	Title words
0.0262 ± 0.0135	Difficulty	0.0549 ± 0.0066	Difficulty	0.0888 ± 0.0138	Length
0.0255 ± 0.0087	No Celebs	0.0372 ± 0.0100	Length	0.0302 ± 0.0092	Difficulty
0.0161 ± 0.0058	Length	0.0147 ± 0.0055	Anticipation	0.0234 ± 0.0062	Dominance
0.0124 ± 0.0054	Diversity	0.0146 ± 0.0088	Diversity	0.0199 ± 0.0047	FB Subjectivity
0.0085 ± 0.0054	FB Subjectivity	0.0134 ± 0.0074	Arousal	0.0194 ± 0.0114	No Celebs
0.0069 ± 0.0091	Anticipation	0.0132 ± 0.0073	Subjectivity	0.0157 ± 0.0032	Anticipation
0.0056 ± 0.0032	Dominance	0.0105 ± 0.0035	Valence	0.0092 ± 0.0019	Arousal
0.0029 ± 0.0030	Arousal	0.0066 ± 0.0035	FB Subjectivity	0.0077 ± 0.0027	Diversity
0.0018 ± 0.0019	Positivity	0.0065 ± 0.0048	Dominance	0.0051 ± 0.0042	Fear

To go a step further and have a closer look at the specific contribution of every feature on the classifiers decisions, a single tree was visualized from the XGBoost model with an accuracy of 0.884 (F-1 score) using this time the number of “Total Interactions” as the target variable, and then relied on the bigger and most important final nodes to draw our conclusions. That is because as the tree develops, it includes more splits and since the constraints are conjunctive in each path there is a chance of different branches presenting contradicting insights. For identifying the significant final nodes the dtreeviz⁸ Python library was implemented, which produces a representation of the feature-space splits of the leafs (Terence Parr 2018), thus the outcomes can be generalizable. In Figure 6.8, the root of the tree begins with the length feature which is the parent node with a child on the left side representing that the answer to the split value is “less than” and at the right side for “greater than or equal to”. Moreover, as it leads to the bottom the leaves become purer and finally divide the samples in the low- or high-engagement bucket.

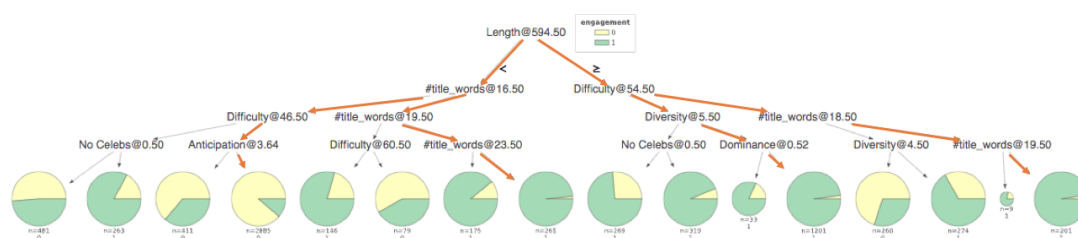


Figure 6.8: Visualization of one decision tree of the XGBoost classifier for “Total Interactions” prediction.

A leaf is considered important if it includes a high number of observations and the impurity or uncertainty in this group of observations is low. As it seems, four paths will be considered marked with orange that lead towards the most important final leaves visualized as large pie charts, and draw a set of rules that can lead to better engagement rates if followed by professional journalists before publishing their story on Facebook. Thus a news story with a high probability of becoming engaging consists of at least 595 words, has an easily comprehensible structure, includes a headline with at least 20 words, and features at least six sources. Additionally, readers should be made to feel in control of the story. If the article is shorter than 595 words, then the headline should contain more than 24 words in order to increase

⁸<https://pypi.org/project/dtreeviz/>

the probability of higher engagement.

Two randomly selected examples of news stories that were correctly classified by the model were visualized in order to ascertain the specific features that led to the prediction. From the graphical representation, we can inspect the feature space distribution that is shown as a histogram, while the variables are stacked on top of each other with the Y-axis to represent the sum of the cases under consideration. It was observed that four features were particularly influential in determining which article ended up in the low-engagement bucket.

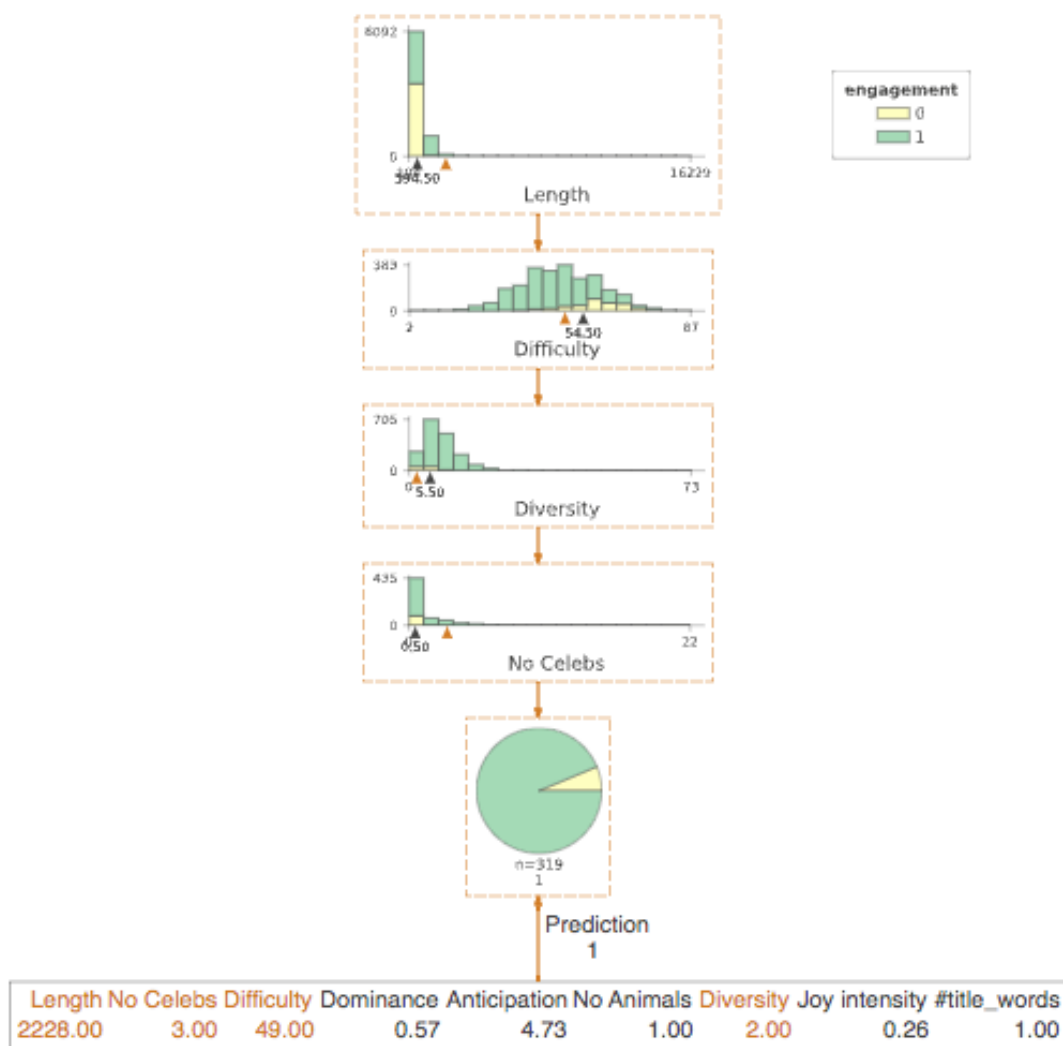


Figure 6.9: Contributions to the low- and high-engagement bucket for “Total Interactions” prediction

Specifically, articles in the low-engagement bucket had a word count of 158, fifteen words in the headline, were relatively uncomplicated, and generated a sense of expectation. Conversely, the articles with high engagement featured a more extensive amount of content,

were more difficult to comprehend, cited two named entities, and referred to three famous people.

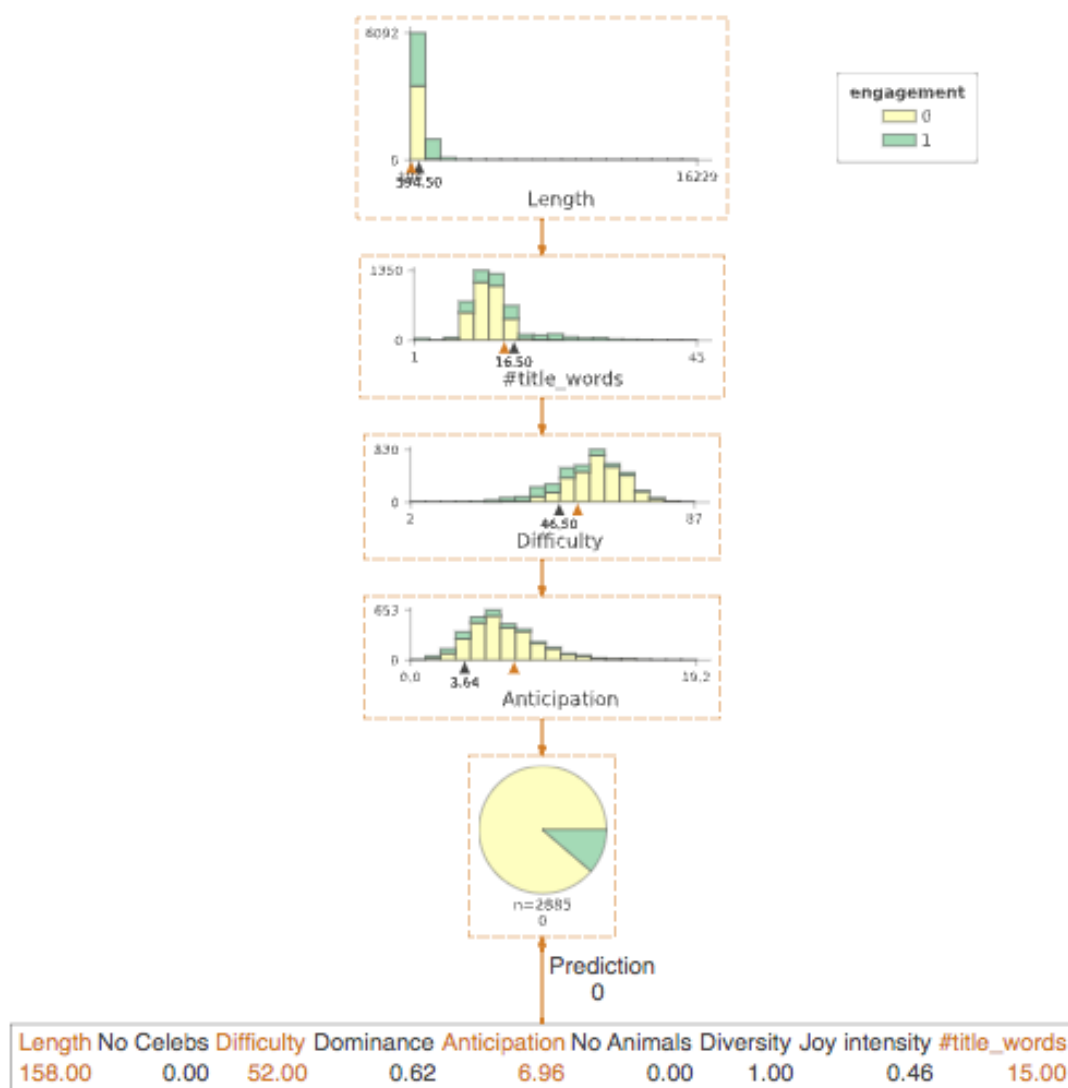


Figure 6.10: Contributions to the low- and high-engagement bucket for "Total Interactions" prediction

6.3.3 Discussion and Conclusion

This study proposes a machine learning approach to both predict and disentangle audience engagement with the news by testing a theoretical framework derived from the literature. The model focuses solely on textual characteristics at the article level to classify articles into two categories: low- and high-engagement. To measure this, a corpus of news articles from elite newspapers and tabloids from the news and politics category is used. Additionally, the four dimensions of the framework - Language Quality, Impartiality, Entertainment, and

Emotionality - are operationalized with computational measures based on the text of the headline, the body of the article, and the headline written on Facebook when the story was published.

The proposed model is able to accurately predict audience engagement with news content across all metrics, including Likes, Comments, Shares, and Total Interactions. Likes prediction was the most successful. In addition to building a robust classification algorithm, we also wanted to develop a better understanding of the factors that influence engagement on social media. To do this, we used explainable AI methods such as permutation importance and tree visualizations. These methods allowed us to draw conclusions about the characteristics of news that are most likely to lead to engagement on social media. Our findings align with many previous studies on the characteristics of news content that are most likely to predict audience engagement. Specifically, the dimension of Language Quality, as measured by factors such as the number of words in the headline, can increase readers' interest and encourage further engagement with the article confirming the work of (Kuiken et al. 2017). Length and readability were found to be important factors, with lengthy and comprehensible articles being more likely to engage readers, supporting many previous studies about the discriminative power of lengthy news stories (Sotirakou et al. 2019; Rosenstiel et al. 2007; Tang et al. 2003), and linguistic difficulty (Sotirakou et al. 2019; Tolochko and Boomgaarden 2018; Arapakis et al. 2016; McQuail 2015; Tang et al. 2003).

Furthermore, the model distinguished the dimension of Subjectivity as meaningful, with the diversity of sources used in a story being influential. This finding suggests that readers value stories that include multiple perspectives, which is consistent with the studies of Choi et al. (2021) and Urban and Schweiger (2014). Moreover, the subjectivity present both in the Facebook headline and in the article was found significant, attesting to the work of Kuiken et al. (2017) who correlated clickbait headlines with subjectivity, and other studies which noted that publishers massively use clickbait tactics and witty titles for marketing and advertising reasons (Lamot et al. 2022; Rony et al. 2017; Scacco and Muddiman 2016). From the Entertainment dimension, the presence of celebrities in a story was also found to be important, confirming the work of J. Cohen (2001) who found that the actions of famous people can en-

gage audiences. In terms of the Emotionality dimension, anticipation and dominance were found to be the most powerful factors, followed by arousal, pleasure, positivity, and fear. Positive dominance has previously been identified as a shareability factor due to the fact that makes people feel “in control” (Guerini and Staiano 2015), while overall emotional content has been shown to be related to engagement with news content (Sotirakou et al. 2019; Sotirakou et al. 2018; Dafonte-Gómez 2018; Scacco and Muddiman 2016; Berger and Milkman 2012). Lastly, the subjectivity of a story was found to be an important factor, in line with previous research that has identified a relationship between subjectivity and engagement, including the use of clickbait headlines.

The contribution of every dimension is already revealing for the way the predictive algorithm works, however, deeper insight into the model’s predictions was obtained in order to identify the unique influence of each feature and specific guidelines for journalists were obtained. According to the machine learning algorithm which predicted the total interactions on Facebook, if a news story has more than 595 words, is comprehensible, has a long headline, cites more than five named entities in the text, and evokes in the readers the feeling of control over their emotions, it has a high probability of becoming engaging. For shorter articles, the headline must be even longer to have better chances for achieving higher engagement. This work is not without limitations. The theoretical framework was tested only on news stories posted on Facebook, therefore different guidelines could apply to Instagram and Twitter. Furthermore, a qualitative approach to the model results could provide greater insights into the topics discussed, language use, context, and so on. Finally, for future work, the role of images used on Facebook posts and their impact on engagement is intended to be explored.

6.4 Study:4 The Impact of Images on News Quality and Engagement

This is a follow up study that investigates the impact of images on both the quality of news content and its engagement, with a focus on both highbrow and tabloid news organizations.

The aim is to identify the defining characteristics of images that accompanies high-quality news and determine whether it is possible to classify news articles based on images. To do this, a framework for evaluating images based on quality criteria derived from the principles of photojournalism, marketing, advertisement is developed and operationalized using computational measures. A machine learning model is then used to predict the quality of digital news stories, with explainable AI methods employed to shed light on the results. The findings of this study offer new insights into the image features that contribute to high-quality news and audience engagement on social media. This study was accepted for presentation and publication at the Future Technologies Conference (FTC) 2023, (Sotirakou et al. 2023).

Colors

Colors play a significant role in the engagement of advertisements. Different colors elicit various emotions and sensations, and they may be utilized to create a mood that motivates viewers to interact with a commercial. Warm hues, such as red and orange, may trigger feelings of energy joy, and enthusiasm, which can be utilized to entice visitors, whereas cooler hues like blue and green, are able to induce feelings of harmony and trust, which can be used to create a more comfortable environment. In addition, colors can also be used to create a sense of familiarity and connection between the viewer and the product or service being advertised (Zailskaitė-Jakštė et al. 2017). For example, companies often use the same colors in their branding throughout different advertisements, which can help viewers to recognize and remember the product or service (Ghaderi 2017).

Ultimately, the use of colors in advertisements can be a powerful tool for influencing engagement and sales since the first impression about a product happens incredibly fast and is by 62-90% based on colors (Lindgaard et al. 2006). By strategically selecting and combining colors, advertisers can create an atmosphere that encourages viewers to engage with their message and take action.

Image Aesthetics

Symmetry and contrast are two powerful visual elements that can be used to create an effective social media post. Symmetry is when elements are arranged in a balanced and harmo-

nious way, creating a pleasing visual aesthetic. For instance, a picture with symmetry has a self-similarity in at least one axis (Treder 2010). Symmetry can be used to make a photograph look organized and professional. It can also create a sense of balance and structure, which can make the image look more appealing and inviting. Symmetry can also be used to draw the viewer's eye to the most important elements of the branded content, such as the product or service being advertised. research has shown that people tend to appreciate symmetry on social media (Redies et al. 2020; Bertamini et al. 2019).

Contrast is when the light and dark elements of a photograph are different in brightness or luminance, often making one element stand out from the others. Contrast can be used to create an eye-catching effect, as it can make one element stand out against its surroundings e.g. by providing high contrast and color against its dark background (Kim and Lakshmanan 2015). It can also be used to make an image more striking, memorable, and visually pleasing (Pedersen et al. 2010). When used together, symmetry and contrast can create an engaging social media post (Kostyk and Huhmann 2021).

Human Faces

Visuals that include people can be especially effective, as they allow viewers to connect with the story that is being told. Studies have shown that human faces are the most powerful form of non-verbal communication, and that beautiful faces are more effective in terms of consumer engagement (Ding et al. 2019; Bakhshi et al. 2014). Also, portraits and direct gaze tend to engage viewers to a greater extent and can persuade them to take action (Lipovsky 2016). Additionally, Frankowska-Takhari et al. (2017) conducted a study on the selection and use of images in online journalism. Semi-structured interviews and observation of eight image professionals working in the field were used to gather data. Thematic analysis of the results revealed a pattern of characteristics associated with visually appealing images, such as depicting only one main object, preferably a person, and the shot, gaze, framing, positioning, color saturation, background type, resolution, and sharpness. The findings may be useful in improving the effectiveness of image retrieval. Finally, posters should strive to create content that features attractive faces with an emphasis on facial expressions in order

to maximize user engagement on social media platforms. Visuals featuring people can be especially powerful because they allow viewers to engage with the story being presented and are more likely to be remembered than images with items or landscapes.

Emotions

Emotional images have become one of the most potent marketing techniques, since they may elicit strong emotions in viewers and assist to leave a lasting impression. Furthermore, emotional images can be used to draw attention to a specific message, creating a stronger impact than words alone. When used effectively, emotional images can help to capture the attention of the viewer and make them more likely to remember the advertisement (Bakalash and Riemer 2013). From pleasure and excitement to fear and sadness, emotional visuals may be utilized to generate a variety of emotions. People who view this kind of pictures are more inclined to identify with the individual who uploaded them and the emotion they portray. Additionally, viewers might be more inclined to react to a photo's positive or negative emotion by liking or commenting on it. (She et al. 2021; Tafesse and Wien 2017). Different emotions can be used depending on the campaigns, for example, a happy image might be used to promote a product that brings joy, while a sad image might be used to draw attention to a social cause. However, it is also important to consider the target audience and the context in which the image will be seen, as the emotion may be interpreted differently depending on the emotional state of the viewer (Klassen et al. 2018).

The purpose of this study is to investigate how visuals affect the appeal of news items on social media platforms. Specifically, the aim is for journalists to be better equipped to choose which pictures to include in their articles by knowing the connection between image qualities and these results. Two specific hypotheses are put forth: It is hypothesized that certain attributes of a photograph are associated with the journalistic quality and social media engagement. To test this hypothesis, machine learning models will be used that incorporate image features as input variables. The second hypothesis is that the attributes of an image can improve the accuracy of models that use only textual features to predict journalistic quality and social media engagement. Therefore, in order to test this hypothesis, the per-

formance of the model that takes into account only textual features will be compared to the one of model that uses both textual and image features. The study builds upon the work presented in Study 6.2 and Study 6.3, in which only textual features were used.

6.4.1 Methods & Data

The annotated dataset from the study 6.2 was used here as the source of the images. Specifically, the link to the Facebook post that the website CrowdTangle provided served as means to download all the images. Then the Amazon Rekognition API⁹ was used, that allows for advanced image analysis including the detection and analysis of scenes, objects, faces, emotions of depicted faces, and famous people. Afterwards, based on the image filename, each image was matched with other features in the dataset for the subsequent analysis steps. The original dataset had 1935 news stories, but some of the article's featured images did not return results from the image recognition tools, therefore only 1451 were left for inclusion in the models.

The image used was the “featured image” of the news story, that was also the picture on the Facebook post. This is typically the image that is used to illustrate the story on the homepage of the news organization, and it is the first image that readers see when they access the story. The featured image is usually chosen to be visually striking and attention-grabbing as described in Chapter 4, and it is often intended to convey the main theme or message of the story. The featured picture is frequently an integral component of the narrative itself in addition to being used on the site and in social media feeds. To grab readers' attention and establish the tone for the remainder of the piece, it is often put at the beginning of the article or in the first few paragraphs.

The characteristics that were retrieved from the photographs with the help of the Amazon Rekognition service and other Python libraries were selected based on how well they served the objectives of the study and the theoretical framework that was being employed. These characteristics included the image's color, contrast, quality, inclusion of text or graphics, and

⁹<https://aws.amazon.com/rekognition/image-features/>

placement (inside or outside). The number of faces, their gender, ages, and if any of them were recognized as belonging to famous individuals were also retrieved from the photos. The emotions of the depicted faces were also included in the dataset. Lastly, the overall number of images in the body of the news story was also taken into consideration. This range of features (Table 6.12) allowed for a comprehensive analysis of the images, providing insight into various aspects of their content and context.

Table 6.12: The creation of the image-based features

Features	Measurement
Colors	ColorThief ¹⁰ Python library includes the following colors: grey, brown, red, blue, purple, white, green, pink, yellow, no color
Resolution	Pillow-Image ¹¹ package provides information about the weight and height of the images.
High Contrast	Scikit-image ¹² can capture the high and low contrast.
Text	Amazon Rekognition (AR) ¹³ detects text on the image.
Number of Images	How many images exist in the article.
Indoor/Outdoor	AR returns if the scene is outside or inside.
Animal	AR, whether there is an animal depicted.
Architecture	AR, depicted buildings.
Graphics	AR returns if there is a map, a diagram, or a plot.
Vehicle	AR provides information about pictures related to transportation.
Bonfire	AR
Nature	AR
Sports	AR
Food	AR, lunch, food, meal, dish and so on.
Crime	AR, pictures related to police, officer, blood, ambulance.
Faces	AR, whether there is a face depicted.
Number of faces	AR
Child	AR
Female/Male	AR
Portrait	AR
Celebrities	AR provides whether is a famous person and how many in total.
Smile	AR
Emotions	AR provides the following facial expression detection: calm, happy, disgusted, fear, sad, surprised, confused, angry.

To test the relationship between certain characteristics of a photograph and journalistic quality and audience engagement, two target variables were used in the machine learning

¹⁰<https://lokeshdhakar.com/projects/color-thief/>

¹¹<https://pillow.readthedocs.io/en/stable/reference/Image.html>

¹²<https://scikit-image.org/docs/stable/api/skimimage.exposure.html>

¹³<https://aws.amazon.com/rekognition/custom-labels-features/>

Table 6.13: The accuracy of the models

	Quality			Engagement		
	Images	Text	Combined	Images	Text	Combined
Logistic Regression	0.66	0.76	0.76	0.54	0.59	0.62
Naive Bayes	0.62	0.71	0.72	0.42	0.63	0.61
Support Vector Machine	0.64	0.78	0.77	0.53	0.6	0.63
K-Nearest Neighbors	0.65	0.73	0.72	0.67	0.58	0.58
Decision tree	0.57	0.73	0.72	0.63	0.64	0.61
Random Forest	0.65	0.76	0.76	0.67	0.64	0.69
XGBoost	0.68	0.76	0.76	0.66	0.68	0.69

models: quality (as annotated by humans) and Facebook likes. First the proposed image features were used to predict the quality of an article and engagement on Facebook. Then, the textual features from the previous study were used on their own, and finally the image and textual features were combined. By separating the features it was allowed to determine the extent to which different types of features contribute to predicting these variables.

6.4.2 Analysis & Results

To evaluate the importance of the proposed framework in relation to images, we applied a binary classification task and built a model using multiple classification algorithms from the scikit-learn Python library (Pedregosa et al. 2011). The algorithms used in this study include the Naive Bayes, Logistic Regression, Support Vector Machine, K-Nearest Neighbors, Decision Tree, Random Forest, and XGBoost classifier. The model was trained using 80% of the news articles and reserved the remaining 20% for testing. To evaluate the performance the weighted average F-measure (F1) was employed.

From Table 6.13, it can be seen that the image-based features were able to predict both the quality and the likes on Facebook of a news story, confirming the first hypothesis that the featured image influences those outcomes. The XGBoost classifier was found to be the best performing, with an F-1 score of 0.68 for quality prediction, and the Random Forest Classifier had a score of 0.67 for engagement prediction. Although the results are not particularly strong, it must be acknowledged that classifying news stories based solely on the featured image of an article is a very challenging task.

In addition, it appears that the combination of textual and image features did not benefit the model, which primarily used text-based variables. Therefore, the second hypothesis of this study, that the combination of textual and image-based attributes would result in better accuracy, is refuted by the findings.

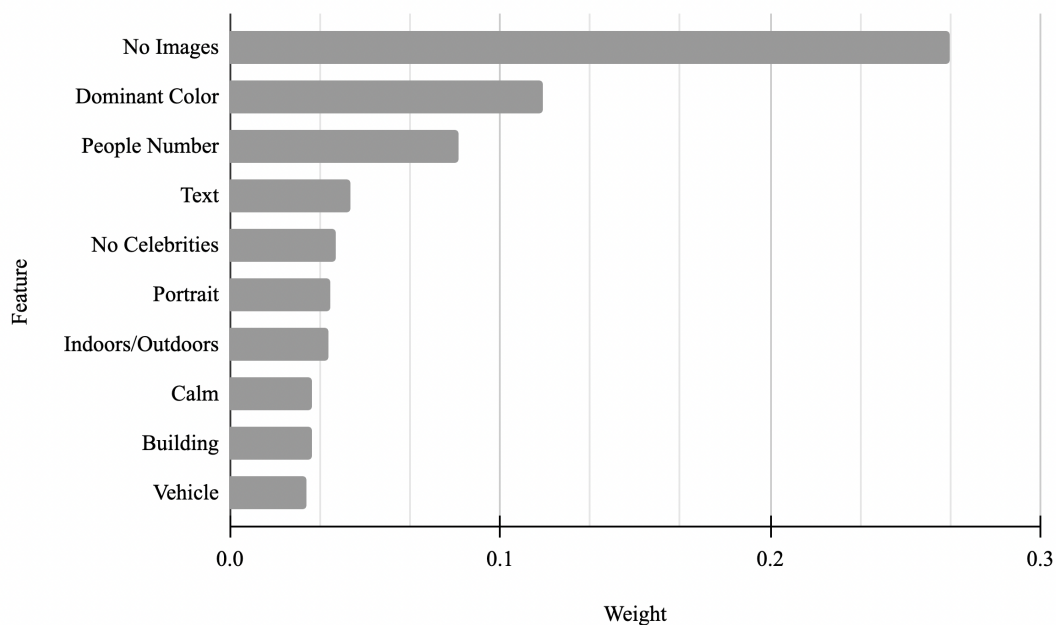


Figure 6.11: Visualization of the feature importance of the Random Forest Classifier for “Likes” prediction.

To have a better understanding of the most important characteristics the feature importance charts can be seen in Figure 6.12 and Figure 6.11.

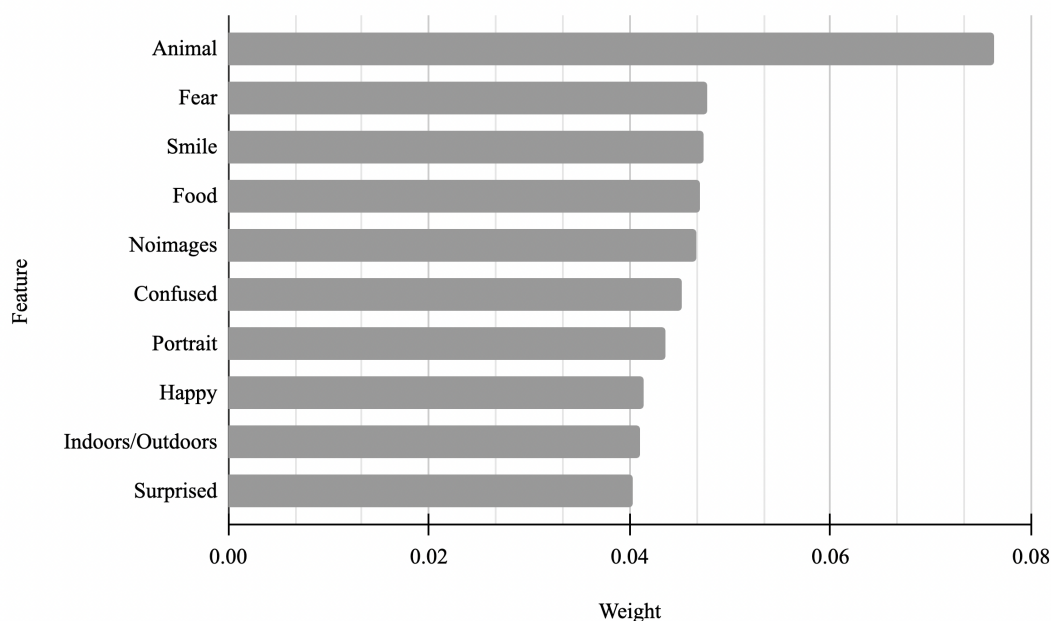


Figure 6.12: Visualization of the feature importance of the Random Forest Classifier for Quality prediction.

An interesting finding of this study was that the features that appeared in the top ten for feature importance were different between engagement and quality prediction. The amount of photographs in the narrative, the prominent color, the number of individuals, and the number of famous persons were the greatest indicators of high engagement on Facebook. The existence of an animal in the image was found to be highly significant for quality prediction, followed by facial expressions like a smile and the emotions of fear, confusion, happiness, and surprise. Moreover, the amount of photographs in the narrative, the presence of food, whether the image was a portrait, and whether the shot was indoors or outdoors were all important factors in the prediction.

6.4.3 Discussion

This study explored the impact of images on both the quality of news content and its engagement, with a focus on both highbrow and tabloid news organizations. The findings showed that in fact images are very important to both engagement and quality and the constructed framework managed to capture some of the defining characteristics of featured images in news stories. The models were able to classify the news articles based solely on

images, achieving nearly 70% accuracy. Also the most important characteristics of the images that contribute to the predictions were identified. While the first hypothesis, that certain photographic attributes are associated with both journalistic quality and social media engagement, was confirmed, the second hypothesis, that combining image-based features with textual features would improve model accuracy, was not confirmed. Textual features appear to outperform image-based characteristics, although additional research is needed to determine how pictures and text interact with one another.

Overall, the two models were different in regards to the weights they attributed to the features that contributed to the predictions. Specifically, the number of images in the story was found to be very important in predicting both quality and engagement, but the remaining features were found to differ significantly between the models, with the engagement model placing more importance on color and people, while the quality model placing more emphasis on animals and expressions. Although the predictive models were not particularly strong, it is a challenging task to classify news stories based solely on the featured image of an article. To our knowledge, this is the first study to predict the quality and audience engagement of news stories based only on the images. This work demonstrates the potential of images to predict the performance of news articles and it identifies the most influential characteristics that could be useful for editors, who are constantly looking for ways to publish both engaging and high quality news. Knowing which attributes of an image are significant towards these goals can help them in their decision-making process.

6.5 Study:5 News Quality and Engagement for Fake News Detection

Recognizing the global implications of online disinformation, this study proposes that detection of such campaigns can be effectively approached as a classification problem, utilizing explainable machine learning models. With this in mind, a model was constructed that focuses on textual attributes and user interactions on Facebook to detect deceptive content in both news articles and Facebook news-related posts. The purpose of the model is to (1)

assess the effectiveness of the quality framework presented in the second study, (2) extend it with more features and (3) draw conclusions regarding which factors predict fake news and why certain characteristics of news articles are more important in their classification as fake. Answers to these questions are essential in the battle against disinformation campaigns. To achieve this, the model exploits the power of language, emotions, and engagement features, and additionally leverages the predictive power of user interactions on Facebook.

Model & Feature Extraction

The main purpose of this study is to create an inclusive model able to accurately predict and identify disinformation campaigns on both news articles and Facebook news-related posts. The backbone of the model is structured based on an extensive review of previous studies in both communication and computational linguistics addressed particularly in section 3.2.4 of Chapter 3. In the light of the literature, the following types of features were identified:

Content-based features

Linguistic: The length of an article and the length of its headline have been identified as strong predictors of potentially false content (Horne and Adali 2017; Asubiaro and Rubin 2018). The use of capitalized words in the body and title of stories (Bradshaw et al. 2020), along with certain Part-of-Speech (POS) tags such as nouns, demonstratives, personal pronouns, and adverbs (Horne and Adali 2017; Asubiaro and Rubin 2018), have been employed to detect deceptive content. Additionally, complexity measures like the level of lexical diversity and readability have been used in prior studies, with lower levels of complexity indicating fake content (Horne and Adali 2017). Furthermore, a high number of swear words can increase the probability of an article being false (Asubiaro and Rubin 2018).

Emotional: The link between emotionality and disinformation has been investigated in several studies (Freelon and Lokot 2020; Horne and Adali 2017), which have found that false stories tend to contain more negativity than real news (Horne and Adali 2017), and that provocative content on social media is more likely to express anger in an attempt to exasperate the audience (Freelon and Lokot 2020). This study focuses on two distinct aspects of

emotionality: i) the emotions actually expressed in the text, measured using intensity scores for anger, fear, sadness, and joy based on the theories of basic emotions (Plutchik 1980a), and ii) the overall affect, which includes the level of valence, arousal, and dominance as described by Russel (Russell 2003). The difference between emotion and affect is explained by (Shouse 2005), which defines emotion as the demonstration of a feeling, whereas affect is the intensity of the non-conscious response of the body to an experience.

Engagement Features

Research has shown that the number of Facebook likes and user reactions to posts can indicate the presence of hoaxes and disinformation (Reis et al. 2019; Idrees et al. 2019; Tacchini et al. 2017). Thus, in this section the main features of our model will be discussed, along with the reasoning behind their selection.

Content-based features

Linguistic

Body length: The text size in characters. Real news articles are significantly longer than fake news articles (Horne and Adali 2017).

Title length: The title size in characters. The total number of words in fake news titles is higher than in real news titles (Horne and Adali 2017).

Capital letters: In fake news articles are used more capitalized words (Horne and Adali 2017).

Parts of speech: The identification of words as nouns, verbs, adverbs, adjectives, pronouns, prepositions, conjunctions, etc. The study (Rashkin et al. 2017) showed that words used to exaggerate, such as superlatives, and modal adverbs are indicative of fake news. However, a survey by Mahyoob et al. (2020) indicated that trustworthy news writers tend to use more personal pronouns, proper nouns, adverbs, numbers (Rashkin et al. 2017) and name entities (Rubin et al. 2016).

Noun/verb: The ratio of nouns to verbs in all words of the text (Marquardt 2019).

Lexical Diversity: Refers to the ratio of different unique words in a text (Horne and Adali 2017).

Readability: The Flesch readability score indicates how easy it is for someone to read a particular text, with high readability levels associated with real news (Pérez-Rosas et al. 2017).

Profanity: The number of swear words is a feature of fake news (Horne and Adali 2017).

Title and body similarity: The relevance of content between the title and the main body, clickbait headlines are often different from the main story (Tromble 2019).

Subjectivity: The quality of news is characterized by the personal author's tone, and personal opinions expressed in a text (Reinemann et al. 2012). Specifically, we measured the degree of weak or strong subjectivity using the MPQA Subjectivity Lexicon (Wilson et al. 2005).

Emotional

Emotions: For the emotion extraction, the NRC Affect Intensity Lexicon (NRC-AIL) was used that identifies the existence of four basic emotions, anger, fear, joy, and sadness (Saif M Mohammad 2017).

Affect: The NRC VAD Lexicon was used which identifies the sentiments of valence, arousal, and dominance (S. Mohammad 2018).

Bullet points: Listicles, meaning a mixture of “list” and “article” are very popular types of articles and according to research, they tend to be highly shareable (Okrent 2014).

Engagement

Likes: The number of likes of the post.

Love: Represents more appreciation than liking and expresses more empathy.

Wow: Indicates a surprising feeling that the post expresses something unexpected.

Haha: Represents a funny reaction, the post causes real laughter or an ironic expression.

Sad: Shows sadness about the post's content also is a sign of refusal (Idrees et al. 2019).

Angry: Represents the disliking of the post.

Shares: The number of shares may be related to news content truthfulness. (Granik and Mesyura 2017).

Comments: The total number of comments.

Total interactions: The total number of all interactions.

Overperforming Score: The overperforming metric is calculated automatically by Crowdtangle¹⁴ based on the performance of similar posts from the same page in similar time-frames.

6.5.1 Method & Dataset

For this study, a dataset of news articles was collected from both trustworthy and unreliable English-language websites using the Python programming language. The dataset consists of a total of 23,420 articles, both real and fake, published online during the years 2019 and 2020, covering a variety of genres. This study focuses only on the article level, and therefore characteristics such as the overall likes or followers of a Facebook page and other contextual attributes like the genre were not taken into consideration. To construct the dataset, the method of (Asubiaro and Rubin 2018) was followed, and 12,420 articles were retrieved from three widely-acknowledged fake news websites, which are listed in many disinformation indexes such as PolitiFact's fake news websites dataset¹⁵ and Wikipedia's list of fake news websites¹⁶. These websites are: *dailysurge.com*, *dcgazette.com*, and *newspunch.com*.

For this study, a total of 11,000 real news articles were collected from reliable sources, namely *the New York Times*, *Business Insider*, *Buzzfeed*, *New Yorker*, *Politico*, and *Washingtonpost*. The articles included the full text, title, date, author, and web address (URL). The dependent variable was set to 1 for all stories scraped from fake websites, and 0 for all truthful articles.

¹⁴<https://help.crowdtangle.com/en/articles/3213537-crowdtangle-codebook>

¹⁵<https://www.politifact.com/article/2017/apr/20/politifacts-guide-fake-news-websites-and-what-they/>

¹⁶https://en.wikipedia.org/wiki/List_of_fake_news_websites

The articles were then subjected to stop-word removal, NaN value treatment, stemming, tokenization, and lemmatization. Articles with less than 1K characters in the main body were also excluded from the study, as many fake stories were very short. After data cleaning, 19,340 qualified cases remained for model construction.

Data on engagement was collected from Facebook using CrowdTangle, a platform run by Facebook that gives access to statistics about public pages and groups. The analytics for all the articles in the dataset were searched using their headlines or URL; however, not all of the articles were present on Facebook, so out of the original dataset (from dailysurge.com and newspunch.com), only 4822 fake news articles had matching Facebook metrics. To balance the dataset, we added analytics for the same number of real articles, resulting in a total of 9,644 articles to use for the model.

For feature engineering, various Python libraries were used, such as the py-readability-metrics¹⁷ package and the Natural Language Toolkit¹⁸ (NLTK), to perform basic text analysis and filtering. After the features in each category (content-based, engagement-based) were created, redundant features were identified using a correlation matrix, and those with a correlation higher than 0.7 were removed from the data. Additionally, several similar features were removed using clustering techniques (for details see Appendix B). For the model, the Decision Tree and the Random Forest classifier from the Scikit-learn Python library¹⁹ were used, and their results were compared to determine the one with the highest prediction accuracy. Subsequently, the importance of each feature in this fake news classification problem was determined.

6.5.2 Data Analysis & Findings – in two Distinctive Phases

The data analysis for this study was conducted in two phases based on the two different datasets. In Phase A, the original dataset was used to evaluate the importance of only the content-based features, namely the linguistic and emotional features. Then, in Phase B, a

¹⁷<https://pypi.org/project/py-readability-metrics/>

¹⁸<https://www.nltk.org/>

¹⁹<https://scikit-learn.org/stable/>

subset of the dataset that included Facebook activity (engagement features) was used twice. First, only the engagement features were used as predictor variables, and then all the features were used. The goal at this stage was to add the predictive power of the engagement features and assess their effects on the accuracy scores. Additionally, the overall goal of the analysis is to explore the different sets of features in order to understand what elements of a story increase its probability of being fake. Therefore, models that provide in-depth explanations of the classifier's predictions, such as tree-based models (Lundberg et al. 2020), were used. For all experiments, 70% of the stories were used for training and the remaining 30% for testing, and three classification methods were used for evaluating the model: F-measure (F1), precision, and recall.

Phase A - Evaluating the content-based features

For phase A of the experiment, the original dataset (fake and real articles) was used to discover the most significant content-based features that can classify an article before publication, meaning that engagement features were not being considered at this stage. The two different classification methods were applied, and the algorithm with the highest accuracy was the Random Forest classifier with an F1-Score of 91%. The main interest lies in the feature importance of the classifier that will shed light into what matters most as the model constructs its decision trees, therefore except for calculating the contribution of every feature on the prediction (see Fig. 6.13), the ELI5²⁰ Python package for “Inspecting Black-Box Estimators” was used to measure the permutation importance (Table 6.14).

Figure 6.13 shows the importance of the content-based features. The category of linguistic features is the most significant with capital letters in the body of the article, POS tags (nouns, adpositions, particles), lexical diversity, headline length, article length, and weak subjectivity to be amongst the top-ten important predictors. From the emotional features, arousal is the only significant attribute for detecting false content.

²⁰<https://eli5.readthedocs.io/en/latest/overview.html>

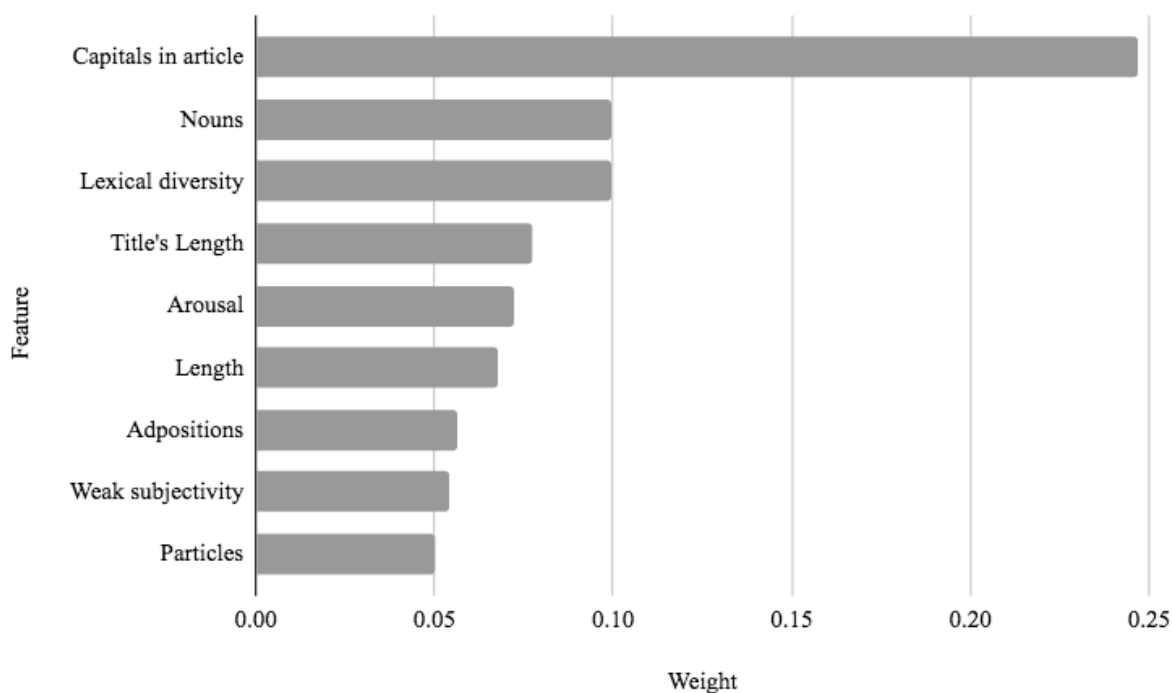


Figure 6.13: Feature importance score for the content-based features

Phase B - Combining the Content-based Features with the Engagement Features

The objective of this phase is to examine if the combination of the textual characteristics of an article (content-based features), together with audience metrics (engagement features), provides better accuracy in distinguishing the fake from real news. In this stage, a smaller dataset was used that includes the engagement features, and the models were executed twice; first, the performance results based only on the engagement features were examined, and then all the features were combined. The results of the two phases are presented in Table 6.15. When the model was run the first time using only the engagement features, the random forest correctly classified 95.8% of news-related posts into either the fake or real class, demonstrating that the model performs well based on users' interactions with the Facebook platform even without any textual features such as headline length or lexical diversity. Additionally, the total number of Facebook users who "liked" the post was the most important feature, followed by the overperforming score, which is calculated by CrowdTangle based on the performance of similar posts from the same page in similar timeframes.

It is observed that the combination of content-based and engagement features proved to

Table 6.14: Permutation Importance for the top 10 combined features

Feature	Weight
Capitals in article	0.0758 ± 0.0095
Likes	0.0732 ± 0.0092
Title's Length	0.0595 ± 0.0038
Numbers	0.0190 ± 0.0045
Overperforming	0.0128 ± 0.0053
Arousal	0.0088 ± 0.0023
Nouns	0.0063 ± 0.0027
Comments	0.0049 ± 0.0033
Readability Score	0.0040 ± 0.0006
Strong Subjectivity	0.0040 ± 0.0016

Table 6.15: Accuracy of Machine Learning Classifiers

Features	Measures in %	Machine Learning Classifiers	
		Decision Tree	Random Forest
News Content Features	Accuracy	84.1	91.0
Engagement Features	Accuracy	94.4	95.8
News Content Features + Engagement Features	Accuracy	94.9	98.0

have greater predictive power compared to any single group of features. In the top 3 of the permutation importance table (see Table 6.14), the number of capital letters in the body of the article is the most significant, with the significance of this feature remaining stable in both datasets. The second most important feature is the number of likes, followed by the length of the headline, which was also very important in phase A. Additionally, POS tags such as numbers and nouns are significant predictors, while the overperforming score is the fifth most significant characteristic. Similar to phase A, arousal is the only emotional feature that contributes to the prediction, while the total number of comments a news post received, the readability score, and the expressed subjectivity are of lower importance.

6.5.3 Discussion of the Results

In general, as shown in Table 6.14, the content-based features, particularly the linguistic ones, are the most informative for distinguishing real from fake news articles. These results are consistent with previous studies (Marquardt 2019; Horne and Adali 2017; Rashkin et al. 2017; Pérez-Rosas et al. 2017) that have found that textual attributes can accurately pre-

dict the probability of a news item being deceptive. The second most important category of features is the engagement features, with the number of likes being the best predictor. Interestingly, the emotional category of features is third, with only arousal being significant for the prediction. The weights of the features indicate what factors are most important for the classifications, and they align well with the proposed categories of features.

Overall, the findings of this study support several features that have been identified in other studies with similar methodologies. For example, features related to words in capital letters were highlighted in the study of Horne and Adali (2017), along with the headline length and the article length, which was also found to be significant in the work of Marquardt (Marquardt 2019). Facebook likes are essential for the model's predictions, and they have also been identified as a key factor in distinguishing hoax posts (Reis et al. 2019; Tacchini et al. 2017; Kuiken et al. 2017). Additionally, the use of audience reactions on the Facebook platform has been shown to provide patterns that can indicate disinformation (Idrees et al. 2019). Furthermore, the results show that the syntax of fake news articles is very significant, and this is one of the features that has been recognized by many researchers in the past. Specifically, false stories tend to include more adverbs (Asubiaro and Rubin 2018; Rashkin et al. 2017; Horne and Adali 2017), fewer nouns (Marquardt 2019; Horne and Adali 2017), more personal pronouns (Asubiaro and Rubin 2018; Rashkin et al. 2017; Pérez-Rosas et al. 2017), fewer numbers (Rashkin et al. 2017), and more demonstratives (Asubiaro and Rubin 2018). Additionally, lexical diversity and subjectivity were found to be significant in phase A, in line with previous findings that false stories have less lexical complexity and more self-referential words (Horne and Adali 2017). On the other hand, characteristics that are often associated with disinformation, such as profanity, negative sentiment (Marquardt 2019), and anger (Freelon and Lokot 2020), were not identified by the model as significant indicators of falsity.

6.5.4 Conclusion & Future Work

Many studies that focus on disinformation in news articles and social media treat fake news detection as a text classification problem, and they aim to extract features and build effective

models that can predict false stories (Conroy et al. 2015). Accordingly, this study employed content-based and engagement features drawn from previous theoretical constructs in an attempt to model online disinformation campaigns and cast light on its significant identifiers. To this end, two datasets were created, one that included real and fake news and a subset of the original that contained the audience's interactions to the same articles posted on Facebook. Then a number of experiments were performed comparing the different sets of features and two tree-based classifiers. The findings revealed that the content-based features such as Capitals in the article, Headline Length, POS tags, and the engagement feature of Facebook Likes were the most important predictors of deceptive online stories. The results provided us with insights of fake news attributes useful in the light of combating disinformation, in terms of proposing a machine learning approach to automatically detect false stories and of pointing to certain telling characteristics of these falsehoods that could be incorporated in media literacy education programs to bolster resilience against this devastating phenomenon.

However, the results of this study are based on a set of assumptions producing the following limitations. First of all, the dataset was built based on the fundamental assumption that all the articles from the sources listed as fake news websites by Politifact are 100% fake. Undoubtedly, there are better ways of constructing a fake news corpus such as asking fact-checkers to verify the potentially deceptive stories before incorporating them into the dataset or opting for a human-in-the-loop approach where the model would not rely so heavily on artificial intelligence but include more sophisticated human judgment. Except for the dependent variable of the model not being the optimal one, there is the limitation of the English language thus it is uncertain how the model would behave with datasets in other languages. Based on the current study, future work could use a more diverse dataset and design a study in which human fact-checkers define false stories based on certain features and their respective significance and then correlate their judgment with the feature importances of the model, or focus on rule extraction and investigate more closely the effect of each feature on disinformation detection.

*The results of this study were presented at 3rd Multidisciplinary International Symposium

on Disinformation in Open Online Media (Sotirakou et al. 2021).

Chapter 7

Conclusion and Future Work

7.1 Conclusion

Journalism has seen a transformation as a result of technical advancements with new approaches, novel concepts and creative ideas beginning to appear as a result of these shifts. Considering this, it's crucial to take into account journalistic quality in the digital era, as these advancements had a significant impact on how news is consumed. Using my own experience as a journalist, the objective was to investigate how AI could potentially be used as a tool to assist newsrooms produce better journalism. Through the years, I observed a shift towards more inclusive, informal, and intimate news coverage, which is discussed in chapter 3 and chapter 4. Despite the pervasive notion that they function in isolation, this dissertation explored the subtle nuances between the seemingly opposing ideas of rationality and irrationality, objectivity and subjectivity, neutrality and sensationalism, and elite newspapers versus tabloids. This work breaks down each news article into its components and analyzes these relationships through an AI lens, in order to provide a deeper understanding and appreciation of the subject.

My research extends beyond traditional theoretical studies by incorporating a technological viewpoint, which can be essential for creating more inclusive conversations when discussing the digital transformation of news media in the future. By collecting and analyzing

data from various sources, such as blogs, elite newspapers, tabloids, and Facebook, different patterns were identified that shed light into the contemporary news landscape. To refer to some of the research questions discussed in the introduction, the first question was about how artificial intelligence can provide a nuanced understanding of what constitutes quality journalism in the digital era.

However, in a news ecosystem led by digital technologies and social media platforms one cannot explore quality without taking under consideration the factors that contribute to audience engagement. Thus, the aim of the first study was to investigate the concept of quality in online news by examining how readers perceive good articles on the blogging platform *Medium.com*. The scope of the first study was broad, therefore considered articles from many categories, namely, news, technology, business, health, sports and lifestyle. The objective was to predict the perceived quality of an article and find writing rules that may be category-specific. These guidelines for different categories can be beneficial for writers since they allow them to consider the expectations and conventions of the genre for which they are writing and craft their narrative in a way that is more likely to be well-received by readers. Moreover, insights like these can be useful to editors in order to better analyze and choose articles for publication, as well as for writers to improve the quality of their work.

Specifically, the study used the number of “claps” an article received as a proxy to measure its perceived quality and classified articles as either successful or unsuccessful based on this measure. The appreciation of online audiences, in the form of clapping, could be translated as the willingness of a reader to recommend a certain story to other users, therefore signaling to the platform’s algorithm to promote this article to more readers. For the analysis, a framework with three dimensions (author, content, and context) was designed, and the number of “claps” served as the dependent variable for a series of classification tasks. Out of the 32 features that shaped the bottom-up model, 21 were identified by the model as significant. The best model was the XGBoost Classifier, with the accuracy (F1-score) to range between 0.74 and 0.88 depending on the news category. The best predictor consistently across the six categories was the number of users who follow the author of the article, which is not surprising, if more people see it, more people will clap. Yet, the decision trees for each

model were visualized and investigated in great detail to reveal other important factors for each category, some of which will be listed here.

Based on the rules extracted from the models, for a news story to have a high probability of gaining popularity the author must have more than 3K followers, include more than three images in their articles, and use first-person pronouns. Tech writers should consider incorporating longer stories into their work, using first-person pronouns, and using bullet points to organize their ideas. It can also be effective to express emotions, such as anticipation, in their writing. For Business articles, it is recommended that the authors have at least five times the number of followers as the number of people they follow. In addition, it is suggested to include at least two images and to avoid using external links. Finally, the article should not contain more than 3.3% negatively charged words. Health articles written by authors with a small number of followers or a close ratio of followers to followees may be able to increase their popularity by using negatively charged words in more than 6.5% of the total words. It seems like the negative news in this area can win the audience. More rules involve writing short and easy to understand pieces, use many personal pronouns and so on. In order for sports authors to be successful on the platform, it is essential for them to already have a well-established following of more than 2K followers. Their text must be between 261 and 1305 words and convey confidence, but not in an arrogant manner. Also, rare words and lists should be used sparingly and the title of the story should consist of at least six words. Finally, for Lifestyle articles, authors should focus on writing an extensive story that is longer than 1,476 words and has a positive sentiment of over 12.3% of total words. Also, the title should be at least six words in length. If the article's author has acquired a large number of followers, the article should be bigger than 550 words to become popular. Additionally, more rules were extracted depending on the number of followers the author has.

This study provided further insight into the concept of perceived quality of various articles in different categories. By emphasizing each category's respective qualities, it has demonstrated that the criteria used in judging an article's quality are not universal in nature; instead, they differ significantly based on the category the article belongs to. Not only does this study support the criteria outlined in previous literature, but it also provides for the first

time to our knowledge, a more comprehensive look at the rules of quality assessment specific to each article type. Therefore, this detailed examination of perceived quality aims to make a valuable addition to the field, even though the distinction between perceived quality and popularity is not outright clear at these early stages of my research, and further research is warranted to expand upon these findings.

To investigate further the notion of quality in digital journalism, the second study developed a theoretical model to measure the quality of journalistic texts based on norms and best practices in journalism reviewed in chapter 2, and specifically in section 2.2, and section 2.3. The study proposed four dimensions to be central to understanding quality journalism based on text features, and it specifically posed two more research questions: Can the dimensions of the model, namely Impartiality, Emotionality, Entertainment, and Language Quality predict the quality of online news? Additionally, what is the specific contribution of these quality criteria to the quality prediction? The study also proposed a hypothesis that lower levels of subjectivity and entertainment will predict higher quality, while emotional coverage and poor language use will predict lower quality. The second study was more focused in its scope compared to the first study; it focused on news stories from the “news” and “politics” categories from elite newspapers and tabloids and involved journalism students for evaluating the quality of 2K articles. Furthermore, the quality dimensions were operationalized into concrete computational measures, and a binary classification task was applied to build a model that could provide the importance of each feature along with an explanation of the system’s predictions to improve the understanding of what comprises a high-quality news story. The model was put to the test through the development of a seven different approaches, namely the Naive Bayes, Logistic Regression, Support Vector Machine, K-Nearest Neighbors, Decision tree, Random forest, and XGBoost classifier. Also the performance of the framework was compared with both more simple approaches like the bag of words, and TF-IDF vectorizers, and more advanced like the deep learning language model, BERT.

The best model was found to be the XGBoost classifier, which had an F1-Score of 92%. To better understand the predictions, the importance of different features was examined and

a graphic representation of a single decision tree from the best classifier was visualised to propose a set of rules that could be used as guidelines for journalists and editors to craft high-quality news stories. In the light of AI explainability, the results revealed rules able to discriminate between news stories based on seven important sub-dimensions of the framework. Based on these rules, a story with a high probability of being of good quality should be longer than 558 words, not very easy to read, have a headline with no more than 24 words, and cite at least two named sources. A second rule states that a news story that is less than 558 words, fairly difficult to read, has less than 20% words expressing joy, and has a minimum length of 256 words may also be considered high quality.

This study found that several quality characteristics identified in the theoretical foundations of journalistic quality were supported by the results. Specifically, the Language Quality dimension was found to be a major predictor of quality, with metrics such as readability, length, and image-to-text ratio largely determining whether a news story would be classified as high or low quality. This aligned with previous research in journalism that emphasizes the importance of good readability, fluency, and rich vocabulary (Tolochko and Boomgaarden 2018; Arapakis et al. 2016; McQuail 2015; Tang et al. 2003), as well as the provision of lengthy articles (Rosenstiel et al. 2007; Tang et al. 2003). Additionally, the use of numbers, adjectives, and headline length were found to have some discriminating power. Moreover, the Impartiality dimension was a key part of the model, illustrated by the importance of features like the diversity of sources (Carpenter 2008; Choi et al. 2021) and self-disclosure in its top ten features. This is unsurprising, as the essential qualities of Impartiality in the realm of quality journalism have been long recognised by scholars, experts and readers alike (Urban and Schweiger 2014; Rosenstiel et al. 2007; Gladney 1996; McQuail 1992; Westerståhl 1983; Tuchman 1972). Additionally, the dimension of Emotionality (Peters 2011; Pantti 2010; Harrington 2008; Sparks 1998; Slattery and Hakanen 1994), and specifically words related to emotions of trust and joy were identified by the model as significant predictors of quality. Interestingly, the Entertainment dimension (G. Turner 2013; Van Der Wurff and Schoenbach 2014), which is often linked to low-quality news, is absent from the top-ranking predictors. Finally, our “recipe” for writing a high-quality news article is in line with the “magic formula”

(Rosenstiel et al. 2007) when it comes to quoting multiple sources and producing longer articles.

The findings of this study provide new insights into the specific features that contribute to high-quality journalism, which may be useful for scholars, practitioners, and potentially even engaged audiences. This study is also the first to demonstrate the usefulness of explainable AI in identifying key quality criteria in journalism and enhancing our understanding of the relationship between these criteria and high-quality news stories.

Moving on to the second research question of the dissertation, which focused on audience engagement with news content, the third study aimed to determine if the level of engagement of a particular news article can be predicted based only on its individual characteristics. It trained multiple classifiers on 73K news from both elite newspapers and tabloids, with a focus on general news and politics like in the previous study. For this study, Facebook data were also used to calculate the engagement feature, with the most accurate model to be able to classify correctly 88% if a news story will be engaging or not. This study proposed a machine learning approach to both predict and disentangle audience engagement with the news by testing the theoretical framework derived from the literature on journalistic quality. The model focused solely on textual characteristics at the article level to classify articles into two categories: low- and high-engagement. Additionally, the four dimensions of the framework - Language Quality, Impartiality, Entertainment, and Emotionality - were operationalized with computational measures based on the text of the headline, the body of the article, and the headline written on Facebook when the story was published.

The proposed model effectively predicted audience engagement with news content across all metrics, including Likes, Comments, Shares, and Total Interactions, with the Likes prediction being the most successful. To further explore the factors that influence engagement on social media, sophisticated methods from the field of Explainable AI such as permutation importance and tree visualizations were used to identify the most significant features and their relationships with engagement. These findings corroborated previous studies that suggest that Language Quality, as measured by factors such as the number of words in the

headline, can increase readers' interest and encourage further engagement with the article (Kuiken et al. 2017). According to the model, length and readability are important factors in creating engaging news stories, while the diversity of sources used in a story was also significant. This finding suggests that readers value stories that include multiple perspectives, which is consistent with the studies of Choi et al. (2021) and Urban and Schweiger (2014). The presence of celebrities in a story was also found to be a key factor, confirming the findings of J. Cohen (2001), who suggested that the actions of famous people can engage audiences. Moreover, the subjectivity of the headline was identified as an important factor, in agreement with previous research that has established a relationship between subjectivity and engagement, including the use of clickbait headlines. Further, the expressed emotions of anticipation, dominance, arousal, pleasure, positivity and fear were found to be the among the important factors as well. Positive dominance has already been acknowledged as a shareability factor because it makes people feel "in control" (Guerini and Staiano 2015). In terms of rules, to ensure a higher probability of engagement on Facebook, a news story should be a minimum of 595 words and feature a headline of 20 or more words. Also, there should be at least six sources included, and readers should feel in control. In case the article is shorter than 595 words, the headline should be longer than 24 words.

Research on social media engagement has typically focused on the written text of a news-related post, but the use of relevant and compelling images is also crucial for generating engagement. The study of visual practices and visual analyses, known as the "pictorial", "visual", or "iconic turn", has gained importance in academia and society due to the increasing importance of images as a means of communication. Cultural sociologists have also studied the independent existence of images separate from the discourse surrounding them (J. C. Alexander 2010). While much research in journalism focuses on the relationship between images and text, some scholars believe that visuals have become more prevalent than text-based news articles, but their study is under-researched (Pearce et al. 2020; Highfield and Leaver 2016). Therefore, it is important to also investigate the role of featured images and their influence on both the quality and engagement of a news story, in addition to the textual and contextual attributes that have been previously examined.

The fourth study, is a follow up study that investigated the impact of images on both the quality of news content and its engagement, with a focus on both highbrow and tabloid news organizations. The aim was to identify the defining characteristics of images that accompanies high-quality news and determine whether it is possible to classify news articles based on images. To do this, a framework for evaluating images based on quality criteria derived from the principles of photojournalism, marketing, advertisement was developed and operationalized using computational measures. Two hypotheses were proposed for this study. The first hypothesis was that certain attributes of a photograph, including color, aesthetics, emotions, and human faces, are associated with both journalistic quality and social media engagement. The second hypothesis was that combining image-based features with previously used textual features would improve the accuracy of models.

For this study, the annotated dataset from the study 6.2 was used. Every image was downloaded and processed through automated image recognition tools and methods that label objects, faces, emotions of depicted faces, celebrities and more. Different machine learning models were then used to predict the quality and the engagement of digital news stories, first only using image-based features and then combined with textual. The models were able to classify the news articles based solely on images, achieving nearly 70% accuracy. The results revealed that the total number of images in the story was one of the most important characteristics for the models, but the remaining features differed significantly between the two models, with the engagement model emphasizing color and people and the quality model emphasizing animals and expressions. The first hypothesis, that certain characteristics of a photograph are related to both journalistic quality and social media engagement, was confirmed; however, the second hypothesis, that combining image-based features with textual features would improve model accuracy, was not supported. Textual features appear to outperform image-based features, although more research is needed to determine how images and text interact with one another.

This classification task was highly challenging, and to the best of our knowledge, it is the first time that the quality and audience engagement of news stories has been predicted solely based on the images. This study presented the potential of image-based features in predict-

ing the performance of news stories on social media and highlights the influence of these characteristics, which could potentially be used in the media industry, by news editors, who are constantly looking for ways to produce content that is both engaging and of high quality. Lastly, my approach has demonstrated that the “iconic turn” addressed in detail in Chapter 4, can be examined by using a computational perspective to identify, quantify, and predict the quality of news stories and their performance on social media. This research could be further extended by the exploration of other types of multimedia content, more platforms, and different news genres.

Finally, the last study of the dissertation, addressed the fourth research question and examined the use of the suggested quality model to detect online disinformation. Drawing from prior research from Communication and Computational Linguistics, reviewed in Chapter 3 and 3.2.4, it has been proposed that a classification problem can be leveraged to detect disinformation campaigns. To that end, a model was developed that focused on text-based features, as well as user interactions on the Facebook platform. Aiming to identify false news, the model analyzed language, emotions, and engagement features and used them to predict deceptive content in both news articles and Facebook news-related posts. The goal of the suggested model was twofold: first to extend the quality framework from the second study and evaluate its effectiveness and second determine which factors predict fake news, and why certain characteristics of news articles are more important in their classification as fake.

This study collected a dataset of 23,420 news articles from reputable and untrustworthy English-language websites using Python. The articles were published in 2019 and 2020 and span a variety of genres. The 12,420 potentially misleading stories were retrieved from three widely-acknowledged fake news websites which are listed on disinformation indexes like PolitiFact’s fake news websites dataset. The remaining 11K was collected from established news organizations like the *Politico*, *Buzzfeed*, and *New Yorker*. The findings revealed that 91% of the articles could be correctly classified using the Random Forest model with linguistic features of capital letters, POS tags (nouns, adpositions, particles), lexical diversity, headline length, article length, and weak subjectivity being amongst the top-ten important

predictors. Arousal was the only significant attribute for detecting false content from the emotional features.

Facebook engagement was also proven to be an excellent predictor for disinformation. When the model was run using only the engagement features, the random forest correctly classified 95.8% of news-related posts into either the fake or real class, demonstrating that the model performs well based on users' interactions with the Facebook platform even without any textual features. The two most important features of the model were the total number of Facebook users who liked the post and the overperforming score, which is calculated by CrowdTangle based on the performance of similar posts from the same page in similar timeframes. Combining content-based and engagement features improved the model significantly, yielding an F1 score of 0.98. This research revealed useful information about attributes of fake news which may be applied to machine learning models to automatically detect false stories and be incorporated into media literacy education programs to help fight against this destructive phenomenon.

The use of artificial intelligence in a range of contexts relating to the news was investigated in the studies described in Chapter 6. The main argument developed was that those computational social science methods can be used to answer questions related to the study of digital journalism and help journalists better understand how to reach and retain their audience, as well as produce high-quality journalism that resonates with their readers. The findings of this dissertation may be of interest to journalists, communication scholars, and news readers alike, and may also influence the news industry to employ these computational tools to investigate the concept of journalistic excellence in the digital era, as well as the aspects that contribute to audience engagement.

7.1.1 Discussion and Limitations

Artificial intelligence and algorithmic curation on social media can threaten the status of a professional journalist and can potentially blur the lines between journalism, user generated content, sponsored content and robot journalism. Technology giants can also pose

challenges to the survival of traditional media organizations forcing, for instance, local newspapers to offer digital news. News organizations too, are businesses with interests and owners, therefore journalism has become a collective enterprise with the editorial department working together with public relations, marketing and business departments to better serve their customers. Nevertheless, the essence of journalism remains the same through the years, unwavering in the face of evolving technologies. Journalism is an institution, steadfast in its commitment to uncover the truth, expose wrongdoings and inform the public, regardless of the different methods and tools it incorporates in its practice over time. Data journalism techniques, machine learning methods, and generative language models can significantly assist journalists in investigations and various other reporting tasks, while also potentially boosting the revenue and visibility of news. However, they remain tools in service of journalism.

My research sought to identify the factors that contribute to good journalism and explore how they can be integrated into algorithms that could enhance journalistic quality, as well as to investigate the use of AI to increase audience engagement and provide revenue for journalism organizations. Through this work, it was shown that machine learning algorithms are quickly becoming popular technologies in the media industry and the results of this research could provide the foundation for new, innovative solutions to the issues brought about by the contemporary media environment.

Summarizing the results of the five studies, the first study demonstrated that machine learning models can accurately predict the success of articles on the blogging platform *Medium.com* based on various features related to the author, style, content, and context and provided optimized writing guidelines for authors depending on the news genre. The second study revealed that AI can be used to determine the quality of news stories with high accuracy based on a theoretical framework that operationalized quality criteria derived from the literature. The results of this study confirmed the importance of journalistic norms like objectivity vs. subjectivity and facts vs. emotions; however, their role in the model was not as decisive as initially anticipated. Ultimately, it was concluded that although these polarity paradigms often appear rigid and absolute, in reality, the boundaries between them are

often highly ambiguous and subjective. The study found that, in addition to high- and low-quality news, there was a third category of medium quality news. As the first study also highlighted, it seems that there is no single answer that applies to all scenarios and oftentimes the exact same event can differ dramatically depending on which news outlet reported on it, without necessarily compromising its quality.

The third study showed that tree-based approaches, support vector machines, and logistic regression can accurately predict audience engagement on Facebook based on the content of the news stories from both high-brow newspapers and tabloid press. The fourth study found that images can contribute to the prediction of both engagement and quality in news articles, but their performance was not improved when combined with the textual features from the quality framework. Finally, the fifth study showed the same framework can be effectively applied to predict fake news, with the accuracy of the models improving when combined with additional engagement features from Facebook.

In the quest for the magic formula to craft good and engaging news stories, the findings heavily support the related literature showing that principles such as impartiality and objectivity remain significant in today's hybrid media environment. The depth of a story, its diversity as determined by the named entities mentioned in the text, the emotions expressed, the readability and the images used, can predict the quality of a news piece. Meanwhile, characteristics like the subjectivity and length of the headline can affect its reach on social media. More specifically, to achieve optimal results, a news story should be lengthy, employ a rich vocabulary, incorporate between two and six named sources, feature at least three images highlighting people's faces, make the reader feel "in control" and have a descriptive headline. In general, an unknown author needs to have at least 3K followers to succeed otherwise the curation algorithm will not prioritize the story. Moreover, the choice of emotions should vary depending on the news category; for instance, it is advisable to use more negative words for covering health-related issues, and to employ a positive sentiment in lifestyle and business articles. News headlines especially on Facebook can drive engagement if they are subjective. Additionally, a story that contains celebrities can engage audiences. Overall, these findings have the potential to make significant contributions to both academia and

the media industry and provide a roadmap for further development and implementation of AI in journalism.

The findings must be understood in conjunction with some limitations. AI-driven analysis is heavily reliant on data, which can be sometimes limited or even biased. To overcome this and be able to generalize the results, thousands of articles were scraped from different news outlets, namely *Guardian*, *the Independent*, *Washington Post*, *Politico*, *the New York Times*, *CNN*, *Reuters*, *Daily Mail*, *the Daily Mirror*, *the Sun*, *the Daily Star*, *the Daily Express*, *New York Post*, *Business Insider*, *Buzzfeed*, *New Yorker*, the blog *Medium.com* and three fake news sites. Nonetheless, to use “claps” or “likes” as a proxy to measure what the audience thinks of a news story or why they engage with it is only one way to see it. The results of the studies that do not use human input, such as the second study, need to be further validated by individual members of the audience or experts, to provide a more holistic understanding of news consumption and quality perception.

Furthermore, there has been a set of assumptions that produce the following limitations. First of all, the dataset used for the disinformation detection was built based on the fundamental assumption that all the articles from the sources listed as fake news websites by Politifact are 100% fake. Undoubtedly, there are better ways of constructing a fake news corpus such as asking fact-checkers to verify the potentially deceptive stories before incorporating them into the dataset or opting for a human-in-the-loop approach where the model would not rely so heavily on artificial intelligence but include more sophisticated human judgment. Similarly, the binary classification task for *Medium.com* and for Facebook engagement was based on extremely successful and unsuccessful news stories, ignoring those in the middle. This was done in order to create as pure as possible classes for the model to identify the significant features that contribute to success, so that the model could accurately predict the probability of a news story being engaging, for example. The future work will compare the model's important features and their respective significance to human judgment.

Another shortcoming, is that AI-driven analysis can be difficult to interpret and explain, as it

is based on algorithms that are often not transparent. In this case, all the models were thoroughly examined using various Python libraries including ELI5, LIME, SHAP, Dtreviz, and TreeInterpreter. It is important to note, however, that all insights refer to probabilities and are not absolute. For example, a different combination of features may make a news piece highly engaging. Additionally, a limitation of this work that must be acknowledged is the fact that the most important features of the model are rather language-related and thus it is unclear how this would travel to a multilingual or non-English setting. This limitation can be addressed in the future by training the model with a more diverse corpus that incorporates other languages in order for the model to generalize better and consequently provide a more comprehensive view of quality identifiers outside the English-language boundaries. Finally, the theoretical framework was tested only on news stories posted on Facebook, therefore different guidelines could apply to Instagram and Twitter.

7.2 Knowledge Transfer and Future Work

Technologies based on AI are rapidly being employed in newsrooms, and it is important to understand how they might be used to improve the quality and impact of journalism. This dissertation examined some of the potential provided by artificial intelligence and how they can be applied to the study of journalism. Ultimately, the use of machine learning models uncovered previously unseen patterns and relationships of the subject at hand. The pioneering work presented in this dissertation, envisages to be applied to future investigations in a wider array of perspectives, allowing for a greater appreciation of the complexities of news consumption. One key aspect of this work is its potential to make significant contributions to both the academic and practical spheres, by transferring the newfound knowledge to research projects. Therefore, the findings of the dissertation are currently used to the ongoing project “IQ Journalism”, that received a million euros in funding for industry application. The final product will be an intelligent agent that provides real-time guidance to journalists on how to craft compelling news pieces. This writing assistant will help news organizations create high-quality, profitable journalism that resonates with their audience.

The proposed intelligent advisor will leverage the cutting-edge technologies demonstrated here, to process Greek news articles and suggest edits to journalists and editors.

One of the objectives of the project is to both advance the theoretical foundations and practical applications of the research. To achieve this, the theoretical frameworks on measuring quality, engagement, and the influence of images, were confirmed through a qualitative study in the context of IQ Journalism. Specifically, 20 experts, academics, researchers, and media professionals were interviewed to confirm or refute the theoretical frameworks used for the studies. The majority of respondents attributed the quality of a news story to the information quality, the diversity of sources in the text, and the pluralism of opinions or information. Additionally, the proper use of language that respects spelling, syntactic, and grammatical rules was deemed an important aspect of good journalism.

According to the respondents, journalists should use comprehensible and simple language, while avoiding colorism, subjective bias, and do not let their emotions to impact the tone and timbre of their reporting. Moreover, the headline was identified as a crucial element that captures the audience's attention; it should not be misleading, but rather comprehensive and representative of the content. In addition, most experts advised against the use of capital letters in the title and body of the articles, considering it to be "unacceptable" and "offensive". Furthermore, they highlighted the usefulness of bullet points, while rejecting the overuse of punctuation. Audiovisual content, such as eye-catching photographs, infographics, videos, and hyperlinks to additional sources, were all considered as important components of quality journalism. Finally, experts argued about the right length of an online news piece, with some suggesting that it should be concise and clear, while others stated that the size might vary based on the genre of the article. The experts' responses to the writing of a journalistic article will be used as the training features for the machine learning model. The study was carried out by a team of academics and researchers at NKUA (Poulakidakos et al. 2023).

In regards to practical applications, the machine learning models will be trained on millions of Greek news stories from news outlets, blogs, and analytics from social media provided to

the University by the industry partner in the project. Thousands of news stories will be manually annotated, and various dictionaries will be either translated into Greek or created for the purpose of this project. The user interface of the intelligent agent will be designed with a user-centred approach and will provide suggestions for content customization for publication on a news website, blog, or social media. It will advise the journalist on how to tailor the language, tone, emotions, sources and so on, and also suggest what kind of multimedia to be included or removed to better match the reader's preferences. This AI system will be the first to be created and implemented in Greece with a hope to help revolutionize the media industry in the country. The anticipated impact of this project is to increase automation and efficiency in media production processes, reduce costs, and optimize the utilization of resources for media companies. It will also acquaint media professionals with the power of AI in the newsroom, provide accurate audience insights, and ultimately facilitate the production of high quality and engaging news content that yields optimal returns.

References

- Aaker, David A (2009). *Managing brand equity*. The Free Press: Simon and Schuster.
- Aarøe, Lene (2011). “Investigating frame strength: The case of episodic and thematic frames”. In: *Political communication* 28.2, pp. 207–226.
- Aarts, Emile HL and Encarnação, José Luis (2006). *True visions: The emergence of ambient intelligence*. Springer Science & Business Media.
- Agarwal, Sheetal D and Barthel, Michael L (2015). “The friendly barbarians: Professional norms and work routines of online journalists in the United States”. In: *Journalism* 16.3, pp. 376–391.
- Aggarwal, Charu C (2018). *Machine Learning for Text*. Springer.
- Ahmed, Hadeer, Traore, Issa, and Saad, Sherif (2017). “Detection of online fake news using n-gram analysis and machine learning techniques”. In: *International conference on intelligent, secure, and dependable systems in distributed and cloud environments*. Springer, pp. 127–138.
- Aimee Rinehart, Ernest Kung (2022). *Artificial Intelligence in Local News A survey of US newsrooms’ AI readiness*. URL: https://www.ap.org/assets/files/ap_local_news_ai_report_march_2022.pdf.
- Akhgar, Babak, Staniforth, Andrew, and Waddington, David (2017). *Application of social media in crisis management*. Springer.

- Alejandro, Jennifer (2010). "Journalism in the age of social media". In: *Reuters Institute Fellowship Paper* 5.1-47, p. 1.
- Alexander, Jeffrey, Bartmanski, Dominik, Giesen, Bernhard, et al. (2012). *Iconic power: materiality and meaning in social life*. Springer.
- Alexander, Jeffrey C (2010). "Iconic consciousness: The material feeling of meaning". In: *Thesis Eleven* 103.1, pp. 10–25.
- Allan, Stuart and Thorsen, Einar (2009). *Citizen journalism: Global perspectives*. Vol. 1. Peter Lang.
- Almgren, Susanne M and Olsson, Tobias (2015). "Let's get them involved'... to some extent: Analyzing online news participation". In: *Social media+ society* 1.2, p. 2056305115621934.
- Alpaydin, Ethem (2020). *Introduction to machine learning*. MIT press.
- Ananny, Mike and Crawford, Kate (2018). "Seeing without knowing: Limitations of the transparency ideal and its application to algorithmic accountability". In: *new media & society* 20.3, pp. 973–989.
- Anderson, Christopher W (2013). "Towards a sociology of computational and algorithmic journalism". In: *New media & society* 15.7, pp. 1005–1021.
- Anspach, Nicolas M (2017). "The new personal influence: How our Facebook friends influence the news we read". In: *Political communication* 34.4, pp. 590–606.
- Arapakis, Ioannis, Cambazoglu, B Barla, and Lalmas, Mounia (2014). "On the feasibility of predicting news popularity at cold start". In: *International Conference on Social Informatics*. Springer, pp. 290–299.
- Arapakis, Ioannis, Peleja, Filipa, Berkant, Barla, and Magalhaes, Joao (2016). "Linguistic benchmarks of online news article quality". In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1893–1902.

- Arendt, Florian, Steindl, Nina, and Kümpel, Anna (2016). "Implicit and explicit attitudes as predictors of gatekeeping, selective exposure, and news sharing: Testing a general model of media-related selection". In: *Journal of Communication* 66.5, pp. 717–740.
- Asubiaro, Toluwase Victor and Rubin, Victoria L (2018). "Comparing features of fabricated and legitimate political news in digital environments (2016-2017)". In: *Proceedings of the Association for Information Science and Technology* 55.1, pp. 747–750.
- Asur, Sitaram and Huberman, Bernardo A (2010). "Predicting the future with social media". In: *2010 IEEE/WIC/ACM international conference on web intelligence and intelligent agent technology*. Vol. 1. IEEE, pp. 492–499.
- Atteveldt, Wouter van and Peng, Tai-Quan (2018). "When communication meets computation: Opportunities, challenges, and pitfalls in computational communication science". In: *Communication Methods and Measures* 12.2-3, pp. 81–92.
- Ausserhofer, Julian, Gutounig, Robert, Oppermann, Michael, Matiasek, Sarah, and Goldgruber, Eva (2020). "The datafication of data journalism scholarship: Focal points, methods, and research propositions for the investigation of data-intensive newswork". In: *Journalism* 21.7, pp. 950–973.
- Baack, Stefan (2015). "Datafication and empowerment: How the open data movement rearticulates notions of democracy, participation, and journalism". In: *Big Data & Society* 2.2.
- Baden, Denise, McIntyre, Karen, and Homberg, Fabian (2019). "The impact of constructive news on affective and behavioural responses". In: *Journalism Studies* 20.13, pp. 1940–1959.
- Bakalash, Tomer and Riemer, Hila (2013). "Exploring ad-elicited emotional arousal and memory for the ad using fMRI". In: *Journal of Advertising* 42.4, pp. 275–291.
- Bakhshi, Saeideh, Shamma, David A, and Gilbert, Eric (2014). "Faces engage us: Photos with faces attract more likes and comments on instagram". In: *Proceedings of the SIGCHI conference on human factors in computing systems*, pp. 965–974.

- Bakir, Vian and McStay, Andrew (2018). "Fake news and the economy of emotions: Problems, causes, solutions". In: *Digital journalism* 6.2, pp. 154–175.
- Bandari, Roja, Asur, Sitaram, and Huberman, Bernardo A (2012). "The pulse of news in social media: Forecasting popularity". In: *Sixth International AAAI Conference on Weblogs and Social Media*.
- Bandes, Susan A and Salerno, Jessica M (2014). "Emotion, proof and prejudice: The cognitive science of gruesome photos and victim impact statements". In: *Ariz. St. LJ* 46, p. 1003.
- Barnhurst, Kevin G (1994). "Seeing the Newspaper. New York: St". In: *Martin's*.
- Barnidge, Matthew (2015). "The role of news in promoting political disagreement on social media". In: *Computers in Human Behavior* 52, pp. 211–218.
- Barry, Ann Marie Seward (1997). *Visual intelligence: Perception, image, and manipulation in visual communication*. SUNY Press.
- Bartholomé, Guus, Lecheler, Sophie, and Vreese, Claes de (2015). "Manufacturing conflict? How journalists intervene in the conflict frame building process". In: *The International Journal of Press/Politics* 20.4, pp. 438–457.
- Bas, Ozen and Grabe, Maria Elizabeth (2015). "Emotion-provoking personalization of news: Informing citizens and closing the knowledge gap?" In: *Communication Research* 42.2, pp. 159–185.
- (2016). "Personalized news and participatory intent: How emotional displays of everyday citizens promote political involvement". In: *American Behavioral Scientist* 60.14, pp. 1719–1736.
- Beckett, Charlie (2015). *How Journalism is Turning Emotional and What That Might Mean for News.*, Retrieved 11.11.2022. URL: <https://blogs.lse.ac.uk/polis/2015/09/10/how-journalism-is-turning-emotional-and-what-that-might-mean-for-news/>.

- Beckett, Charlie (2022). *Closing the AI gap: How small news publishers can benefit*. LSE. URL: <https://blogs.lse.ac.uk/polis/2022/07/18/closing-the-ai-gap/>.
- Beckett, Charlie and Deuze, Mark (2016). "On the role of emotion in the future of journalism". In: *Social media+ society* 2.3, p. 2056305116662395.
- Bell, Emily (2017). *Technology company? Publisher? The lines can no longer be blurred*. <https://www.theguardian.com/media/2017/apr/02/facebook-google-youtube-inappropriate-advertising-fake-news>, Retrieved 9-11-2020.
- Bellovary, Andrea K, Young, Nathaniel A, and Goldenberg, Amit (2021). "Left-and right-leaning news organizations use negative emotional content and elicit user engagement similarly". In: *Affective science* 2.4, pp. 391–396.
- Benson, Rodney (2008). "Journalism: normative theories". In: *The international encyclopedia of communication* 6, pp. 2591–2597.
- Berendt, Bettina and Preibusch, Sören (2017). "Toward accountable discrimination-aware data mining: the Importance of keeping the human in the loop—and under the looking glass". In: *Big data* 5.2, pp. 135–152.
- Berger, Jonah (2014). "Word of mouth and interpersonal communication: A review and directions for future research". In: *Journal of consumer psychology* 24.4, pp. 586–607.
- Berger, Jonah and Milkman, Katherine L (2012). "What makes online content viral?" In: *Journal of marketing research* 49.2, pp. 192–205.
- Bertamini, Marco, Rampone, Giulia, Makin, Alexis DJ, and Jessop, Andrew (2019). "Symmetry preference in shapes, faces, flowers and landscapes". In: *PeerJ* 7, e7078.
- Blumler, Jay G (2019). "Uses and Gratifications research". In: *The International Encyclopedia of Journalism Studies*, pp. 1–8.

- Bobkowski, Piotr S (2015). "Sharing the news: Effects of informational utility and opinion leadership on online news sharing". In: *Journalism & Mass Communication Quarterly* 92.2, pp. 320–345.
- Boehm, Gottfried and Mitchell, William JT (2009). "Pictorial versus iconic turn: Two letters". In:
- Bogart, Leo (1989). *Press and public: Who reads what, when, where, and why in American newspapers*. Psychology Press.
- Bojanowski, Piotr, Grave, Edouard, Joulin, Armand, and Mikolov, Tomas (2017). "Enriching Word Vectors with Subword Information". In: *Transactions of the Association for Computational Linguistics* 5, pp. 135–146.
- Bolin, Göran and Andersson Schwarz, Jonas (2015). "Heuristics of the algorithm: Big Data, user interpretation and institutional translation". In: *Big Data & Society* 2.2, p. 2053951715608406.
- Bolin, Göran and Velkova, Julia (2020). "Audience-metric continuity? Approaching the meaning of measurement in the digital everyday". In: *Media, Culture & Society* 42.7-8, pp. 1193–1209.
- Borges-Rey, Eddy (2016). "Unravelling data journalism: A study of data journalism practice in British newsrooms". In: *Journalism Practice* 10.7, pp. 833–843.
- Bossio, Diana (2021). "Journalists on Instagram: Presenting professional identity and role on image-focused social media". In: *Journalism Practice*, pp. 1–17.
- Boyles, Jan Lauren and Meyer, Eric (2016). "Letting the data speak: Role perceptions of data journalists in fostering democratic conversation". In: *Digital Journalism* 4.7, pp. 944–954.
- Brader, Ted (2005). "Striking a responsive chord: How political ads motivate and persuade voters by appealing to emotions". In: *American Journal of Political Science* 49.2, pp. 388–405.

- Brader, Ted (2011). "The political relevance of emotions: "Reassessing" revisited". In: *Political Psychology* 32.2, pp. 337–346.
- Bradley, Margaret M and Lang, Peter J (1999). *Affective norms for English words (ANEW): Instruction manual and affective ratings*. Tech. rep. Technical report C-1, the center for research in psychophysiology.
- Bradshaw, Samantha, Howard, Philip N, Kollanyi, Bence, and Neudert, Lisa-Maria (2020). "Sourcing and automation of political news and information over social media in the United States, 2016-2018". In: *Political Communication* 37.2, pp. 173–193.
- Branch, John (2012). "Snow Fall: The Avalanche at Tunnel Creek. The New York Times." In: URL: <http://www.nytimes.com/projects/2012/snow-fall/#/?part=tunnel-creek>.
- Bright, Jonathan (2016). "The social news gap: How news reading and news sharing diverge". In: *Journal of communication* 66.3, pp. 343–365.
- Broersma, Marcel (2019). "Audience engagement". In: *The International Encyclopedia of Journalism Studies*, pp. 1–6.
- Broussard, Meredith, Diakopoulos, Nicholas, Guzman, Andrea L, Abebe, Rediet, Dupagne, Michel, and Chuan, Ching-Hua (2019). "Artificial intelligence and journalism". In: *Journalism & mass communication quarterly* 96.3, pp. 673–695.
- Brown, Danielle K, Lough, Kyser, and Riedl, Martin J (2020). "Emotional appeals and news values as factors of shareworthiness in Ice Bucket Challenge coverage". In: *Digital Journalism* 8.2, pp. 267–286.
- Brown, Tom, Mann, Benjamin, Ryder, Nick, Subbiah, Melanie, Kaplan, Jared D, Dhariwal, Prfulla, Neelakantan, Arvind, Shyam, Pranav, Sastry, Girish, Askell, Amanda, et al. (2020). "Language models are few-shot learners". In: *Advances in neural information processing systems* 33, pp. 1877–1901.

- Brugger, Jérôme, Fraefel, Marianne, Riedl, Reinhard, Fehr, Hansjakob, Schönebeck, Daniel, and Weissbrod, Christoph Stähli (2016). “Current barriers to open government data use and visualization by political intermediaries”. In: *2016 Conference for E-Democracy and Open Government (CeDEM)*. IEEE, pp. 219–229.
- Bruns, Axel (2018). *Gatewatching and news curation: Journalism, social media, and the public sphere (Digital Formations, Volume 113)*. Peter Lang Publishing.
- Capelos, Tereza (2013). “Understanding anxiety and aversion: The origins and consequences of affectivity in political campaigns”. In: *Emotions in politics*. Springer, pp. 39–59.
- Capelos, Tereza, Exadaktylos, Theofanis, Chrona, Stavroula, and Pouloupoulou, Maria (2018). “News Media and the Emotional Public Sphere, The Emotional Economy of the European Financial Crisis in the UK Press”. In: *International Journal of Communication* 12, p. 26.
- Caple, Helen and Bednarek, Monika (2016). “Rethinking news values: What a discursive approach can tell us about the construction of news discourse and news photography”. In: *Journalism* 17.4, pp. 435–455.
- Cappallo, Spencer, Mensink, Thomas, and Snoek, Cees GM (2015). “Latent factors of visual popularity prediction”. In: *Proceedings of the 5th ACM on International Conference on Multimedia Retrieval*, pp. 195–202.
- Cappella, Joseph N, Kim, Hyun Suk, and Albarracín, Dolores (2015). “Selection and transmission processes for information in the emerging media environment: Psychological motives and message characteristics”. In: *Media psychology* 18.3, pp. 396–424.
- Carlson, Matt (2015). “The robotic reporter: Automated journalism and the redefinition of labor, compositional forms, and journalistic authority”. In: *Digital journalism* 3.3, pp. 416–431.
- (2016). “Embedded links, embedded meanings: Social media commentary and news sharing as mundane media criticism”. In: *Journalism studies* 17.7, pp. 915–924.

- Carlson, Matt (2018). "Facebook in the news: Social media, journalism, and public responsibility following the 2016 trending topics controversy". In: *Digital journalism* 6.1, pp. 4–20.
- Carlson, Matt and Lewis, Seth C (2015). *Boundaries of journalism: Professionalism, practices and participation*. Routledge.
- Carpenter, Serena (2008). "Source diversity in US online citizen journalism and online newspaper articles". In: *International Symposium on Online Journalism*. Vol. 4.
- Cellan-Jones, Rory (2020). *Google to pay for 'high quality' news in three countries*. <https://www.bbc.com/news/technology-53176945>, Retrieved 10.12.2020.
- Chadwick, Andrew (2017). *The hybrid media system: Politics and power*. Oxford University Press.
- Chapman, Pete, Clinton, Julian, Kerber, Randy, Khabaza, Thomas, Reinartz, Thomas, Shearer, Colin, Wirth, Rudiger, et al. (2000). "CRISP-DM 1.0: Step-by-step data mining guide". In: *SPSS inc* 9.13, pp. 1–73.
- Chen (2018). *Big data in computational social science and humanities*. Springer.
- Chen, Tianqi and Guestrin, Carlos (2016). "Xgboost: A scalable tree boosting system". In: *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794.
- Choi, Jihyang (2016). "News internalizing and externalizing: The dimensions of news sharing on online social networking sites". In: *Journalism & mass communication quarterly* 93.4, pp. 816–835.
- Choi, Jihyang, Lee, Sang Yup, and Ji, Sung Wook (2021). "Engagement in emotional news on social media: intensity and type of emotions". In: *Journalism & Mass Communication Quarterly* 98.4, pp. 1017–1040.
- Chouliaraki, Lilie (2006). *The spectatorship of suffering*. Sage.

- (2008). “The mediation of suffering and the vision of a cosmopolitan public”. In: *Television & new media* 9.5, pp. 371–391.
- (2013). “Ordinary witnessing in post-television news: Towards a new moral imagination”. In: *Self-Mediation*. Routledge, pp. 121–136.
- Cioffi-Revilla, C (2014). *Introduction to computational social science. Texts in computer science*.
- Cioffi-Revilla, Claudio (2017). “Computation and social science”. In: *Introduction to computational social science*. Springer, pp. 35–102.
- Clough, Patricia T (2008). “The affective turn: Political economy, biomedicine and bodies”. In: *Theory, culture & society* 25.1, pp. 1–22.
- Coddington, Mark (2015). “Clarifying journalism’s quantitative turn: A typology for evaluating data journalism, computational journalism, and computer-assisted reporting”. In: *Digital journalism* 3.3, pp. 331–348.
- Cohen, Jonathan (2001). “Defining identification: A theoretical look at the identification of audiences with media characters”. In: *Mass communication & society* 4.3, pp. 245–264.
- Cohen, Sarah, Hamilton, James T, and Turner, Fred (2011). “Computational journalism”. In: *Communications of the ACM* 54.10, pp. 66–71.
- Commission, European (2018). *Joint Communication to the European Parliament, the European council, the European economic and social committee and the committee of the regions: Action Plan against Disinformation*.
- Conroy, Nadia K, Rubin, Victoria L, and Chen, Yimin (2015). “Automatic deception detection: Methods for finding fake news”. In: *Proceedings of the Association for Information Science and Technology* 52.1, pp. 1–4.
- Conte, Rosaria, Gilbert, Nigel, Bonelli, Giulia, Cioffi-Revilla, Claudio, Deffuant, Guillaume, Kertesz, Janos, Loreto, Vittorio, Moat, Suzy, Nadal, J-P, Sanchez, Anxo, et al. (2012). “Man-

- ifesto of computational social science”. In: *The European Physical Journal Special Topics* 214.1, pp. 325–346.
- Costera-Meijer, IC (2012). “Valuable journalism: The search for quality from the vantage point of the user”. In: *Journalism* 14, p. 1.
- Coward, Rosalind (2013). *Speaking personally: The rise of subjective and confessional journalism*. Bloomsbury Publishing.
- Cranor, Lorrie F (2008). “A framework for reasoning about the human in the loop”. In:
- Dafonte-Gómez, Alberto (2018). “News Media and the Emotional Public Sphere| Audiences as Medium: Motivations and Emotions in News Sharing”. In: *International journal of communication* 12, p. 20.
- Dahlberg, Lincoln (2011). “Re-constructing digital democracy: An outline of four ‘positions’”. In: *New media & society* 13.6, pp. 855–872.
- Darnton, Robert (1975). “Writing news and telling stories”. In: *Daedalus*, pp. 175–194.
- Davidson, Joyce and Milligan, Christine (2004). *Embodying emotion sensing space: introducing emotional geographies*.
- Davies, W. (2018). *How Feelings Took over the World.*, *The Guardian* 2018, Retrieved 11.11.2022. URL: <https://www.theguardian.com/books/2018/sep/08/high-anxiety-how-feelings-took-over-the-world>.
- Davis, Jesse and Goadrich, Mark (2006). “The relationship between Precision-Recall and ROC curves”. In: *Proceedings of the 23rd international conference on Machine learning*, pp. 233–240.
- Davis Mersey, Rachel, Malthouse, Edward C, and Calder, Bobby J (2010). “Engagement with online media”. In: *Journal of Media Business Studies* 7.2, pp. 39–56.

- De los Santos, Theresa M and Nabi, Robin L (2019). “Emotionally charged: Exploring the role of emotion in online news information seeking and processing”. In: *Journal of Broadcasting & Electronic Media* 63.1, pp. 39–58.
- De Mauro, Andrea, Greco, Marco, and Grimaldi, Michele (2016). “A formal definition of Big Data based on its essential features”. In: *Library Review*.
- De Smedt, Tom, Daelemans, W, and Smedt, Tom De (2012). “Pattern for Python”. In: *The Journal of Machine Learning Research* 13, pp. 2063–2067.
- De Sousa, Ronald (1990). *The rationality of emotion*. Mit Press.
- Deleuze, Gilles and Guattari, Felix (1987). *A thousand plateaus* (B. Massumi, Trans.)
- DeLuca, Kevin M, Lawson, Sean, and Sun, Ye (2012). “Occupy Wall Street on the public screens of social media: The many framings of the birth of a protest movement”. In: *Communication, Culture & Critique* 5.4, pp. 483–509.
- Demertzis, Nicolas (2013). “Emotions in politics”. In: *The affect dimension in political tension*. Nueva York: Palgrave Macmillan.
- (2020). *The political sociology of emotions: Essays on trauma and resentment*. Routledge.
- Dennis, James and Sampaio-Dias, Susana (2021). ““Tell the Story as You’d Tell It to Your Friends in a Pub”: Emotional Storytelling in Election Reporting by BuzzFeed News and Vice News”. In: *Journalism Studies* 22.12, pp. 1608–1626.
- Deuze, Mark (2003). “The web and its journalisms: considering the consequences of different types of newsmedia online”. In: *New media & society* 5.2, pp. 203–230.
- (2004). “Journalism studies beyond media: On ideology and identity”. In: *Ecquid Novi: African Journalism Studies* 25.2, pp. 275–293.
- (2005). “What is journalism? Professional identity and ideology of journalists reconsidered”. In: *Journalism* 6.4, pp. 442–464.

- Deuze, Mark (2019). “What journalism is (not)”. In: *Social Media+ Society* 5.3, p. 2056305119857202.
- Deuze, Mark and Witschge, Tamara (2020). *Beyond journalism*. John Wiley & Sons.
- Devichand, Mukul (2016). “Did Alan Kurdi’s Death Change Anything?” In: *BBC News* 2.
- Devlin, Jacob, Chang, Ming-Wei, Lee, Kenton, and Toutanova, Kristina (2018). “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *arXiv preprint arXiv:1810.04805*.
- Dhar, Sagnik, Ordonez, Vicente, and Berg, Tamara L (2011). “High level describable attributes for predicting aesthetics and interestingness”. In: *CVPR 2011*. IEEE, pp. 1657–1664.
- Diakopoulos, Nicholas (2015). “Algorithmic accountability: Journalistic investigation of computational power structures”. In: *Digital journalism* 3.3, pp. 398–415.
- (2019). *Automating the news: How algorithms are rewriting the media*. Cambridge, MA: Harvard University Press.
- Ding, Wang, Ronggang, and Wang, Shiqi (2019). “Social media popularity prediction: A multiple feature fusion approach with deep neural networks”. In: *Proceedings of the 27th ACM International Conference on Multimedia*, pp. 2682–2686.
- Djerf-Pierre, Monika, Ghersetti, Marina, and Hedman, Ulrika (2016). “Appropriating Social Media: The changing uses of social media among journalists across time”. In: *Digital Journalism* 4.7, pp. 849–860.
- Dominick, Joseph R (1999). “Who do you think you are? Personal home pages and self-presentation on the World Wide Web”. In: *Journalism & Mass Communication Quarterly* 76.4, pp. 646–658.
- Domke, David, Perlmutter, David, and Spratt, Meg (2002). “The primes of our times? An examination of the ‘power’ of visual images”. In: *Journalism* 3.2, pp. 131–159.

- Dor, Daniel (2003). "On newspaper headlines as relevance optimizers". In: *Journal of Pragmatics* 35.5, pp. 695–721.
- Doran, Derek, Schulz, Sarah, and Besold, Tarek R (2017). "What does explainable AI really mean? A new conceptualization of perspectives". In: *arXiv preprint arXiv:1710.00794*.
- Drozdiak, Natalia (2020). *Google to Pay Publishers Over 1 Billion dollars for News Content*.
<https://www.bloomberg.com/news/articles/2020-10-01/google-to-pay-publishers-over-1-billion-to-license-news-content>, Retrieved 10.12.2020.
- Dutton, William H (2009). "The fifth estate emerging through the network of networks". In: *Prometheus* 27.1, pp. 1–15.
- Eilders, Christiane (2006). "News factors and news decisions. Theoretical and methodological advances in Germany". In: *Communications* 31.1, pp. 5–24.
- Ekman, Paul (1999). "Basic emotions". In: *Handbook of cognition and emotion* 98.45-60, p. 16.
- Eliza Shearer, Elizabeth Grieco (2019). "Americans Are Wary of the Role Social Media Sites Play in Delivering the News". In: *Pew Research Center, Journalism & Media*. Accessed 13-March-2020. URL: <https://www.journalism.org/2019/10/02/americans-are-wary-of-the-role-social-media-sites-play-in-delivering-the-news>.
- Elster, Jon (1999). *Alchemies of the Mind: Rationality and the Emotions*. Cambridge University Press.
- Enli, Gunn (2014). *Mediated authenticity*. Peter Lang Incorporated.
- Esuli, Andrea and Sebastiani, Fabrizio (2006). "Sentiwordnet: A publicly available lexical resource for opinion mining". In: *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*.

- Evans, Sandra K, Pearce, Katy E, Vitak, Jessica, and Treem, Jeffrey W (2017). “Explicating affordances: A conceptual framework for understanding affordances in communication research”. In: *Journal of Computer-Mediated Communication* 22.1, pp. 35–52.
- Evans, Taneth Autumn (2019). *Conscious commissioning and what (exactly) makes our readers tick*. <https://medium.com/digital-times/conscious-commissioning-and-what-exactly-makes-our-readers-tick-802cd4f5c868>, Retrieved 2020-02-16.
- Fayyad, Usama, Piatetsky-Shapiro, Gregory, and Smyth, Padhraic (1996). “From data mining to knowledge discovery in databases”. In: *AI magazine* 17.3, pp. 37–37.
- Felle, Tom (2016). “Digital watchdogs? Data reporting and the news media’s traditional ‘fourth estate’ function”. In: *Journalism* 17.1, pp. 85–96.
- Ferrer-Conill, Raul and Tandoc Jr, Edson C (2018). “The audience-oriented editor: Making sense of the audience in the newsroom”. In: *Digital Journalism* 6.4, pp. 436–453.
- Finn, Adam (1988). “Print ad recognition readership scores: An information processing perspective”. In: *Journal of Marketing Research* 25.2, pp. 168–177.
- Flaounas, Ilias, Ali, Omar, Lansdall-Welfare, Thomas, De Bie, Tijl, Mosdell, Nick, Lewis, Justin, and Cristianini, Nello (2013). “Research methods in the age of digital journalism: Massive-scale automated analysis of news-content—topics, style and gender”. In: *Digital journalism* 1.1, pp. 102–116.
- Flesch, Rudolph (1948). “A new readability yardstick.” In: *Journal of applied psychology* 32.3, p. 221.
- Flew, Terry, Spurgeon, Christina, Daniel, Anna, and Swift, Adam (2012). “The promise of computational journalism”. In: *Journalism practice* 6.2, pp. 157–171.
- Ford, Elizabeth, Oswald, Malcolm, Hassan, Lamiece, Bozentko, Kyle, Nenadic, Goran, and Cassell, Jackie (2020). “Should free-text data in electronic medical records be shared for research? A citizens’ jury study in the UK”. In: *Journal of medical ethics* 46.6, pp. 367–377.

- Frankowska-Takhari, Sylwia, MacFarlane, Andrew, Göker, Ayşe, and Stumpf, Simone (2017). "Selecting and tailoring of images for visual impact in online journalism". In:
- Freeden, Michael (2013). *Emotions, ideology and politics*.
- Freelon, Deen and Lokot, Tetyana (2020). "Russian Twitter disinformation campaigns reach across the American political spectrum". In: *Misinformation Review*.
- Frosh, Paul (2015). "Selfies| The gestural image: The selfie, photography theory, and kines-
thetic sociability". In: *International journal of communication* 9, p. 22.
- Galtung, Johan and Ruge, Mari Holmboe (1965). "The structure of foreign news: The presen-
tation of the Congo, Cuba and Cyprus crises in four Norwegian newspapers". In: *Journal
of peace research* 2.1, pp. 64–90.
- Gans, Herbert J (2004). *Deciding what's news: A study of CBS evening news, NBC nightly news,
Newsweek, and Time*. Northwestern University Press.
- García-Perdomo, Víctor, Salaverría, Ramón, Kilgo, Danielle K, and Harlow, Summer (2018).
"To share or not to share: The influence of news values and popular social media content
in the United States, Brzil, and Argentina". In: *Journalism Studies* 19.8, pp. 1180–1201.
ISSN: 1461-670X. DOI: 10.1080/1461670X.2016.1265896.
- Garrett, R Kelly and Stroud, Natalie Jomini (2014). "Partisan paths to exposure diversity: Dif-
ferences in pro-and counterattitudinal news consumption". In: *Journal of Communica-
tion* 64.4, pp. 680–701.
- Gerlitz, Carolin and Helmond, Anne (2013). "The like economy: Social buttons and the data-
intensive web". In: *New media & society* 15.8, pp. 1348–1365.
- Ghaderi, Mohammad (2017). "Preference disaggregation: Towards an integrated frame-
work". In: *Available at SSRN 2973415*.

- Ghanem, Bilal, Ponzetto, Simone Paolo, Rosso, Paolo, and Rangel, Francisco (2021). “Fake-flow: Fake news detection by modeling the flow of affective information”. In: *arXiv preprint arXiv:2101.09810*.
- Ghanem, Bilal, Rosso, Paolo, and Rangel, Francisco (2020). “An emotional analysis of false information in social media and news articles”. In: *ACM Transactions on Internet Technology (TOIT)* 20.2, pp. 1–18.
- Giachanou, Anastasia, Rosso, Paolo, and Crestani, Fabio (2019). “Leveraging emotional signals for credibility detection”. In: *Proceedings of the 42nd international ACM SIGIR conference on research and development in information retrieval*, pp. 877–880.
- Gibbs, Martin, Meese, James, Arnold, Michael, Nansen, Bjorn, and Carter, Marcus (2015). “# Funeral and Instagram: Death, social media, and platform vernacular”. In: *Information, communication & society* 18.3, pp. 255–268.
- Giddens, Anthony (1984). *The constitution of society: Outline of the theory of structuration*. Univ of California Press.
- Gil de Zúñiga, Homero and Diehl, Trevor (2019). “News finds me perception and democracy: Effects on political knowledge, political interest, and voting”. In: *New media & society* 21.6, pp. 1253–1271.
- Gil de Zúñiga, Homero, Diehl, Trevor, Huber, Brigitte, and Liu, James (2017). “Personality traits and social media use in 20 countries: How personality relates to frequency of social media use, social media news use, and social media use for social interaction”. In: *Cyberpsychology, Behavior, and Social Networking* 20.9, pp. 540–552.
- Gladney, George Albert (1996). “How editors and readers rank and rate the importance of eighteen traditional standards of newspaper excellence”. In: *Journalism & Mass Communication Quarterly* 73.2, pp. 319–331.

- Gluck, Antje (2019). "Should Journalists Be More Emotionally Literate?" In: URL: <https://en.ejo.ch/ethics-quality/should-journalists-be-more-emotionally-literate>.
- Golder, Scott A and Macy, Michael W (2014). "Digital footprints: Opportunities and challenges for online social research". In: *Annual Review of Sociology* 40.1, pp. 129–152.
- Goldfarb, Avi and Tucker, Catherine (2011). "Online display advertising: Targeting and obtrusiveness". In: *Marketing Science* 30.3, pp. 389–404.
- Gonzalez, Robbie (2015). "The surprisingly complex design of Facebook's new emoji". In: *Wired*. Available at: <http://www.wired.com/2015/10/facebook-reactions-design/> (accessed 16 December 2015).
- González-Bailón, Sandra, Banchs, Rafael E, and Kaltenbrunner, Andreas (2012). "Emotions, public opinion, and US presidential approval rates: A 5-year analysis of online political discussions". In: *Human Communication Research* 38.2, pp. 121–143.
- Graber, Doris A (1994). "The infotainment quotient in routine television news: A director's perspective". In: *Discourse & Society* 5.4, pp. 483–508.
- Granik, Mykhailo and Mesyura, Volodymyr (2017). "Fake news detection using naive Bayes classifier". In: *2017 IEEE First Ukraine Conference on Electrical and Computer Engineering (UKRCON)*. IEEE, pp. 900–903.
- Gray, Jonathan, Chambers, Lucy, and Bounegru, Liliana (2012). *The data journalism handbook: How journalists can use data to improve the news.* " O'Reilly Media, Inc."
- Griffin, Michael (2001). "Camera as witness, image as sign: The study of visual communication in communication research". In: *Annals of the International Communication Association* 24.1, pp. 433–463.

- Gross, Kimberly, Brewer, Paul R, and Aday, Sean (2009). "Confidence in government and emotional responses to terrorism after September 11, 2001". In: *American Politics Research* 37.1, pp. 107–128.
- Guerini, Marco and Staiano, Jacopo (2015). "Deep feelings: A massive cross-lingual study on the relation between emotions and virality". In: *Proceedings of the 24th International conference on world wide web*, pp. 299–305.
- Gummerus, Johanna, Liljander, Veronica, Weman, Emil, and Pihlström, Minna (2012). "Customer engagement in a Facebook brand community". In: *Management Research Review*.
- Guzman, Monica (2016). "The best ways to build audience and relevance by listening to and engaging your community". In: *Retrieved from American Press Institute: <https://www.americanpressinstitute.org/wp-content/uploads/2016/05/How-to-build-audiences-by-engaging-your-community.pdf>*.
- Ha, Louisa, Xu, Ying, Yang, Chen, Wang, Fang, Yang, Liu, Abuljadail, Mohammad, Hu, Xiao, Jiang, Weiwei, and Gabay, Itay (2018). "Decline in news content engagement or news medium engagement? A longitudinal analysis of news engagement since the rise of social and mobile media 2009–2012". In: *Journalism* 19.5, pp. 718–739.
- Hagar, Nick and Diakopoulos, Nicholas (2019). "Optimizing Content with A/B Headline Testing: Changing Newsroom Practices". In: *Media and Communication* 7.1, p. 117. ISSN: 2183-2439.
- Hallin, Daniel C (1992). "The passing of the "high modernism" of American journalism". In: *Journal of Communication* 42.3, pp. 14–25.
- Halpern, Daniel and Gibbs, Jennifer (2013). "Social media as a catalyst for online deliberation? Exploring the affordances of Facebook and YouTube for political expression". In: *Computers in Human Behavior* 29.3, pp. 1159–1168.
- Halterman, Andrew (Jan. 2017). "Mordecai: Full Text Geoparsing and Event Geocoding". In: *The Journal of Open Source Software* 2. DOI: 10.21105/joss.00091.

- Hamilton, James T and Turner, Fred (2009). "Accountability through algorithm: Developing the field of computational journalism". In: *Report from the Center for Advanced Study in the Behavioral Sciences, Summer Workshop*, pp. 27–41.
- Han, Jiawei, Pei, Jian, and Kamber, Micheline (2011). *Data mining: concepts and techniques*. Elsevier.
- Hanitzsch, Thomas, Hanusch, Folker, Mellado, Claudia, Anikina, Maria, Berganza, Rosa, Cangoz, Incilay, Coman, Mihai, Hamada, Basyouni, Elena Hernández, María, Karadjov, Christopher D, et al. (2011). "Mapping journalism cultures across nations: A comparative study of 18 countries". In: *Journalism Studies* 12.3, pp. 273–293.
- Hanitzsch, Thomas and Vos (2017). "Journalistic roles and the struggle over institutional identity: The discursive constitution of journalism". In: *Communication Theory* 27.2, pp. 115–135.
- Hannaford, Liz (2015). "Computational journalism in the UK newsroom". In: *Journalism education* 4.
- Hansen, Mark, Roca-Sales, Meritxell, Keegan, Jonathan M, and King, George (2017). "Artificial intelligence: Practice and implications for journalism". In:
- Hanusch, Folker (2017). "Web analytics and the functional differentiation of journalism cultures: Individual, organizational and platform-specific influences on newswork". In: *Information, Communication & Society* 20.10, pp. 1571–1586.
- Harcup, Tony (2015). *Journalism: principles and practice*. Sage.
- Harcup, Tony and O’neill, Deirdre (2001). "What is news? Galtung and Ruge revisited". In: *Journalism studies* 2.2, pp. 261–280.
- Harcup, Tony and O’neill, Deirdre (2017). "What is news? News values revisited (again)". In: *Journalism studies* 18.12, pp. 1470–1488.

- Hariman, Robert and Lucaites, John Louis (2007). *No caption needed: Iconic photographs, public culture, and liberal democracy*. University of Chicago Press.
- Harrington, Stephen (2008). "Popular news in the 21st century Time for a new critical approach?" In: *Journalism* 9.3, pp. 266–284.
- Hasell, Ariel (2021). "Shared emotion: The social amplification of partisan news on Twitter". In: *Digital Journalism* 9.8, pp. 1085–1102.
- Hasell, Ariel and Weeks, Brian E (2016). "Partisan provocation: The role of partisan news use and emotional responses in political information sharing in social media". In: *Human Communication Research* 42.4, pp. 641–661.
- Hedman, Ulrika (2020). "Making the most of Twitter: How technological affordances influence Swedish journalists' self-branding". In: *Journalism* 21.5, pp. 670–687.
- Hedman, Ulrika and Djerf-Pierre, Monika (2013). "The social journalist: Embracing the social media life or creating a new digital divide?" In: *Digital journalism* 1.3, pp. 368–385.
- Hensbergen, Van (2017). *The world speaks the language of men, but after metoo women must find their voice. The Conversation, October 25, Retrieved 11.11.2022*. URL: <https://theconversation.com/the-world-speaks-the-language-of-men-but-after-metoo-women-must-find-their-voice-86107>.
- Henshall, Peter and Ingram, David (1991). *The news manual: A training book for journalists*. Poroman Press.
- Hermida, Alfred (2012). "Social journalism: Exploring how social media is shaping journalism". In: *The handbook of global online journalism* 12, pp. 309–328.
- (2016). *Tell everyone: Why we share and why it matters*. Anchor Canada.
- Hermida, Alfred and Young, Mary Lynn (2017). "Finding the data unicorn: A hierarchy of hybridity in data and computational journalism". In: *Digital Journalism* 5.2, pp. 159–176.

- (2019). *Data journalism and the regeneration of news*. Routledge.
- Hesse, Bradford W, Moser, Richard P, and Riley, William T (2015). “From big data to knowledge in the social sciences”. In: *The Annals of the American Academy of Political and Social Science* 659.1, pp. 16–32.
- Highfield, Tim and Leaver, Tama (2016). “Instagrammatics and digital methods: Studying visual social media, from selfies and GIFs to memes and emoji”. In: *Communication research and practice* 2.1, pp. 47–62.
- Hill, Annette (2018). *Media experiences: engaging with drama and reality television*. Routledge.
- Hoffman, Robert R, Mueller, Shane T, Klein, Gary, and Litman, Jordan (2018). “Metrics for explainable AI: Challenges and prospects”. In: *arXiv preprint arXiv:1812.04608*.
- Honnibal, Matthew and Montani, Ines (2017). “spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing”. In:
- Hopp, Frederic R, Fisher, Jacob T, Cornell, Devin, Huskey, Richard, and Weber, René (2021). “The extended Moral Foundations Dictionary (eMFD): Development and applications of a crowd-sourced approach to extracting moral intuitions from text”. In: *Behavior research methods* 53.1, pp. 232–246.
- Hopper, K Megan and Huxford, John E (2015). “Gathering emotion: examining newspaper journalists’ engagement in emotional labor”. In: *Journal of Media Practice* 16.1, pp. 25–41.
- Horne, Benjamin and Adali, Sibel (2017). “This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news”. In: *Proceedings of the International AAAI Conference on Web and Social Media*. Vol. 11.
- Howard, Philip N and Hussain, Muzammil M (2013). *Democracy’s fourth wave?: digital media and the Arab Spring*. Oxford University Press.

- Hsu, Chin-Chia, Ajorlou, Amir, and Jadbabaie, Ali (2020). “News sharing, persuasion, and spread of misinformation on social networks”. In: *Persuasion, and Spread of Misinformation on Social Networks (July 1, 2020)*.
- Huan, Changpeng (2017). “The strategic ritual of emotionality in Chinese and Australian hard news: A corpus-based study”. In: *Critical Discourse Studies* 14.5, pp. 461–479.
- Hurcombe, Edward, Burgess, Jean, and Harrington, Stephen (2021). “What’s newsworthy about ‘social news’? Characteristics and potential of an emerging genre”. In: *Journalism* 22.2, pp. 378–394.
- Hutto, Clayton J and Gilbert, Eric (2014). “Vader: A parsimonious rule-based model for sentiment analysis of social media text”. In: *Eighth international AAI conference on weblogs and social media*.
- ICIJ (2017). *Paradise Papers*. ICIJ. URL: <https://www.icij.org/investigations/paradise-papers/>.
- Idrees, Amira M, Alsheref, Fahad Kamal, and ElSeddawy, Ahmed I (2019). “A Proposed Model for Detecting Facebook News’ Credibility”. In: *International Journal of Advanced Computer Science and Applications (IJACSA)* 10.7, pp. 311–316.
- Ikonomakis, M, Kotsiantis, Sotiris, and Tampakas, V (2005). “Text classification using machine learning techniques.” In: *WSEAS transactions on computers* 4.8, pp. 966–974.
- Isola, Phillip, Xiao, Jianxiong, Parikh, Devi, Torralba, Antonio, and Oliva, Aude (2013). “What makes a photograph memorable?” In: *IEEE transactions on pattern analysis and machine intelligence* 36.7, pp. 1469–1482.
- Jaakonmäki, Roope, Müller, Oliver, and Vom Brocke, Jan (2017). “The impact of content, context, and creator on user engagement in social media marketing”. In:

- Jagadish, Hosagrahar V, Gehrke, Johannes, Labrinidis, Alexandros, Papakonstantinou, Yannis, Patel, Jignesh M, Ramakrishnan, Raghu, and Shahabi, Cyrus (2014). "Big data and its technical challenges". In: *Communications of the ACM* 57.7, pp. 86–94.
- James, Gareth, Witten, Daniela, Hastie, Trevor, and Tibshirani, Robert (2013). *An introduction to statistical learning*. Vol. 112. Springer.
- Janowitz, Morris (1975). "Professional models in journalism: The gatekeeper and the advocate". In: *Journalism quarterly* 52.4, pp. 618–626.
- Jensen, Matthew L, Averbek, Joshua M, Zhang, Zhu, and Wright, Kevin B (2013). "Credibility of anonymous online product reviews: A language expectancy perspective". In: *Journal of Management Information Systems* 30.1, pp. 293–324.
- Jin, Yue, Huang, Jinghua, and Wang, Xinyao (2017). "What Influences Content Popularity? An Empirical Investigation of Voting in Social Q&A Communities." In: *PACIS*, p. 161.
- Johannesson, Mikael Poul and Knudsen, Erik (2021). "Disentangling the influence of recommender attributes and news-story attributes: A conjoint experiment on exposure and sharing decisions on social networking sites". In: *Digital Journalism* 9.8, pp. 1141–1161.
- Johnson, Eric J and Tversky, Amos (1983). "Affect, generalization, and the perception of risk." In: *Journal of personality and social psychology* 45.1, p. 20.
- Joye, Stijn (2013). "Pantti, M., Wahl-Jorgensen, K., & Cottle, S.(2012). Disasters and the media. New York, NY: Peter Lang Publishing. 235 pp." In: *Communications: The European Journal of Communication Research* 38.1, pp. 122–124.
- Jukes, Stephen (2017). "Affective journalism—uncovering the affective dimension of practice in the coverage of traumatic news". PhD thesis. Goldsmiths, University of London.
- Kalsnes, Bente and Krumsvik, Arne H (2019). "Building trust: media executives' perceptions of readers' trust". In: *Journal of Media Business Studies* 16.4, pp. 295–306.

- Kaplan, Andreas M and Haenlein, Michael (2010). "Users of the world, unite! The challenges and opportunities of Social Media". In: *Business horizons* 53.1, pp. 59–68.
- Karlsen, Joakim and Stavelin, Eirik (2014). "Computational journalism in Norwegian newsrooms". In: *Journalism practice* 8.1, pp. 34–48.
- Karlsson, Michael (2010). "Rituals of transparency: Evaluating online news outlets' uses of transparency rituals in the United States, United Kingdom and Sweden". In: *Journalism studies* 11.4, pp. 535–545.
- Katz, Elihu, Blumler, Jay G, and Gurevitch, Michael (1973). "Uses and gratifications research". In: *The public opinion quarterly* 37.4, pp. 509–523.
- Keneshloo, Yaser, Wang, Shuguang, Han, Eui-Hong, and Ramakrishnan, Naren (2016). "Predicting the popularity of news articles". In: *Proceedings of the 2016 SIAM International Conference on Data Mining*. SIAM, pp. 441–449.
- Kennedy, Helen, Poell, Thomas, and Dijck, José van (2015). "Introduction: Special issue on Data and agency". In: *Data & Society* 2.2, pp. 1–7.
- Kepplinger, Hans Mathias (2008). "News values". In: *The international encyclopedia of communication*.
- Khosla, Aditya, Das Sarma, Atish, and Hamid, Raffay (2014). "What makes an image popular?" In: *Proceedings of the 23rd international conference on World wide web*, pp. 867–876.
- Kiesow, Damon (2015). "The readers we ignore and the news they want". In: *Kiesow 8.0 Blog*.
- Kilgo, Danielle K, Harlow, Summer, García-Perdomo, Víctor, and Salaverría, Ramón (2018a). "A new sensation? An international exploration of sensationalism and social media recommendations in online news publications". In: *Journalism* 19.11, pp. 1497–1516.
- (2018b). "A new sensation? An international exploration of sensationalism and social media recommendations in online news publications". In: *Journalism: Theory, Practice & Criticism* 19 (11), pp. 1497–1516. ISSN: 1464-8849. DOI: 10.1177/1464884916683549.

- Kim, Dam Hee, Jones-Jang, Mo, and Kenski, Kate (2021). "Why do people share political information on social media?" In: *Digital Journalism* 9.8, pp. 1123–1140.
- Kim and Lakshmanan, Arun (2015). "How kinetic property shapes novelty perceptions". In: *Journal of Marketing* 79.6, pp. 94–111.
- Kim, Sang-Bum, Rim, Hae-Chang, Yook, DongSuk, and Lim, Heui-Seok (2002). "Effective methods for improving naive bayes text classifiers". In: *Pacific rim international conference on artificial intelligence*. Springer, pp. 414–423.
- Kincaid, Peter, Fishburne, Robert, Rogers, Richard, and Chissom, Brad (1975). *Derivation of new readability formulas (automated readability index, fog count and flesch reading ease formula) for navy enlisted personnel*. Tech. rep. Naval Technical Training Command Millington TN Research Branch.
- Kinder, Donald R (2013). "Reason and emotion in American political life". In: *Beliefs, reasoning, and decision making*. Psychology Press, pp. 287–324.
- Kitch, Carolyn (2003). "' Mourning in America': ritual, redemption, and recovery in news narrative after September 11". In: *Journalism Studies* 4.2, pp. 213–224.
- Klassen, Karen Michelle, Borleis, Emily S, Brennan, Linda, Reid, Mike, McCaffrey, Tracy A, and Lim, Megan SC (2018). "What people "like": Analysis of social media strategies used by food industry brands, lifestyle brands, and health promotion organizations on Facebook and Instagram". In: *Journal of medical Internet research* 20.6, e10227.
- Knobloch, Silvia, Patzig, Grit, Mende, Anna-Maria, and Hastall, Matthias (2004). "Affective news: Effects of discourse structure in narratives on suspense, curiosity, and enjoyment while reading news and novels". In: *Communication Research* 31.3, pp. 259–287.
- Knobloch-Westerwick, Silvia and Kleinman, Steven B (2012). "Preelection selective exposure: Confirmation bias versus informational utility". In: *Communication research* 39.2, pp. 170–193.

- Kohavi, Ron (1995). "A study of cross-validation and bootstrap for accuracy estimation and model selection". In: *International Joint Conference on Artificial Intelligence (IJCAI)*, pp. 1137–1143.
- Koivunen, Anu, Kanner, Antti, Janicki, Maciej, Harju, Auli, Hokkanen, Julius, and Mäkelä, Eetu (2021). "Emotive, evaluative, epistemic: a linguistic analysis of affectivity in news journalism". In: *Journalism* 22.5, pp. 1190–1206.
- Kostyk, Alena and Huhmann, Bruce A (2021). "Perfect social media image posts: Symmetry and contrast influence consumer response". In: *European Journal of Marketing*.
- Kotišová, Johana (2017). "Cynicism ex machina: The emotionality of reporting the 'refugee crisis' and Paris terrorist attacks in Czech Television". In: *European Journal of Communication* 32.3, pp. 242–256.
- (2019). *Crisis Reporters, Emotions, and Technology: An Ethnography*. Springer Nature.
- Kovach, Bill and Rosenstiel, Tom (2014). *The elements of journalism: What newspeople should know and the public should expect*. Three Rivers Press (CA).
- Kramer, Adam DI, Guillory, Jamie E, and Hancock, Jeffrey T (2014). "Experimental evidence of massive-scale emotional contagion through social networks". In: *Proceedings of the National Academy of Sciences* 111.24, pp. 8788–8790.
- Kramp, Leif and Loosen, Wiebke (2018). "The transformation of journalism: From changing newsroom cultures to a new communicative orientation?" In: *Communicative Figurations*. Palgrave Macmillan, Cham, pp. 205–239.
- Ksiazek, Thomas B, Peer, Limor, and Lessard, Kevin (2016). "User engagement with online news: Conceptualizing interactivity and exploring the relationship between online news videos and user comments". In: *New media & society* 18.3, pp. 502–520.

- Kühne, Rinaldo and Schemer, Christian (2015). "The emotional effects of news frames on information processing and opinion formation". In: *Communication Research* 42.3, pp. 387–407.
- Kuiken, Jeffrey, Schuth, Anne, Spitters, Martijn, and Marx, Maarten (2017). "Effective headlines of newspaper articles in a digital environment". In: *Digital Journalism* 5.10, pp. 1300–1314.
- Kuyucu, Mihalis Michael (2020). "Social Media and Journalism". In: *Academic Studies*, p. 72.
- Lacy, Stephen (2000). "Commitment of financial resources as a measure of quality". In: *Measuring media content, quality, and diversity*, pp. 25–50.
- Lacy, Stephen and Rosenstiel, Tom (2015). *Defining and measuring quality journalism*. Rutgers School of Communication and Information.
- Lagun, Dmitry and Lalmas, Mounia (2016). "Understanding user attention and engagement in online news reading". In: *Proceedings of the ninth ACM international conference on web search and data mining*, pp. 113–122.
- Lai, Susan (2011). "Iconic images and citizen journalism". In: *Proceedings of the 2011 iConference*, pp. 702–703.
- Lamot, Kenza, Kreutz, Tim, and Opgenhaffen, Michaël (2022). "'We Rewrote This Title': How News Headlines Are Remediated on Facebook and How This Affects Engagement". In: *Social Media+ Society* 8.3, p. 20563051221114827.
- Lamprou, Evangelos, Antonopoulos, Nikos, Anomeritou, Iouliani, and Apostolou, Chrysoula (2021). "Characteristics of fake news and misinformation in greece: the rise of new crowdsourcing-based journalistic fact-checking models". In: *Journalism and Media* 2.3, pp. 417–439.

- Lasorsa, Dominic L, Lewis, Seth C, and Holton, Avery E (2012). "Normalizing Twitter: Journalism practice in an emerging communication space". In: *Journalism studies* 13.1, pp. 19–36.
- Lavidge, Robert J and Steiner, Gary A (1961). "A model for predictive measurements of advertising effectiveness". In: *Journal of marketing* 25.6, pp. 59–62.
- Lawrence, Regina G, Radcliffe, Damian, and Schmidt, Thomas R (2018). "Practicing engagement: Participatory journalism in the Web 2.0 era". In: *Journalism Practice* 12.10, pp. 1220–1240.
- Lazer, David, Pentland, Alex, Adamic, Lada, Aral, Sinan, Barabási, Albert-László, Brewer, Devon, Christakis, Nicholas, Contractor, Noshir, Fowler, James, Gutmann, Myron, et al. (2009). "Computational social science". In: *Science* 323.5915, pp. 721–723.
- Lecheler, Sophie (2020). "The emotional turn in journalism needs to be about audience perceptions: Commentary-virtual special issue on the emotional turn". In: *Digital journalism* 8.2, pp. 287–291.
- Lecheler, Sophie and Vreese, Claes H de (2013). "What a difference a day makes? The effects of repetitive and competitive news framing over time". In: *Communication research* 40.2, pp. 147–175.
- Lee, Chei Sian, Ma, Long, and Goh, Dion Hoe-Lian (2011). "Why do people share news in social media?" In: *International Conference on Active Media Technology*. Springer, pp. 129–140.
- Lee, Lewis, Seth, and Powers, Matthew (2014). "Audience clicks and news placement: A study of time-lagged influence in online journalism". In: *Communication Research* 41.4, pp. 505–530.
- Lee, Na Yeon, Kim, Yonghwan, and Kim, Jiwon (2016). "Tweeting public affairs or personal affairs? Journalists' tweets, interactivity, and ideology". In: *Journalism* 17.7, pp. 845–864.

- Legislation Related to Artificial Intelligence* (2022). National Conference of State Legislatures. URL: <https://www.ncsl.org/research/telecommunications-and-information-technology/2020-legislation-related-to-artificial-intelligence.aspx> (visited on 09/21/2022).
- León, Ernesto de and Trilling, Damian (2021). “A sadness bias in political news sharing? The role of discrete emotions in the engagement and dissemination of political news on Facebook”. In: *Social media+ society* 7.4, p. 20563051211059710.
- Levendusky, Matthew (2013). “Partisan media exposure and attitudes toward the opposition”. In: *Political communication* 30.4, pp. 565–581.
- Levenson, Robert W (2003). “Blood, sweat, and fears: The autonomic architecture of emotion”. In: *Annals of the New York Academy of Sciences* 1000.1, pp. 348–366.
- Lewis, Seth C, Holton, Avery E, and Coddington, Mark (2014). “Reciprocal journalism: A concept of mutual exchange between journalists and audiences”. In: *Journalism Practice* 8.2, pp. 229–241.
- Lewis, Seth C and Usher, Nikki (2013). “Open source and journalism: Toward new frameworks for imagining news innovation”. In: *Media, culture & society* 35.5, pp. 602–619.
- (2014). “Code, collaboration, and the future of journalism: A case study of the Hackers/Hackers global network”. In: *Digital journalism* 2.3, pp. 383–393.
- Li, Yiyi and Xie, Ying (2020). “Is a picture worth a thousand words? An empirical study of image content and social media engagement”. In: *Journal of Marketing Research* 57.1, pp. 1–19.
- Liang, Jiaqi, Bi, Guoshu, and Zhan, Cheng (2020). “Multinomial and ordinal Logistic regression analyses with multi-categorical variables using R”. In: *Annals of translational medicine* 8.16.
- Liddy, Elizabeth D (2001). “Natural language processing”. In:

- Lindgaard, Gitte, Fernandes, Gary, Dudek, Cathy, and Brown, Judith (2006). "Attention web designers: You have 50 milliseconds to make a good first impression!" In: *Behaviour & information technology* 25.2, pp. 115–126.
- Lipovsky, Caroline (2016). "Negotiating solidarity with potential donors: a study of the images in fundraising letters by not-for-profit organizations". In: *Functional Linguistics* 3.1, pp. 1–18.
- Lippmann, Walter (1946). "Public opinion (Vol. 1)". In: *Transaction Publishers*. doi 10, pp. 14847–000.
- Liu, Yuping and Shrum, Lawrence J (2002). "What is interactivity and is it always such a good thing? Implications of definition, person, and situation for the influence of interactivity on advertising effectiveness". In: *Journal of advertising* 31.4, pp. 53–64.
- Logo vs icon* (2016). Medium.com. URL: <https://subsign.medium.com/logo-vs-icon-c2ad6eaab980> (visited on 06/04/2019).
- Loria, Steven (2018). "textblob Documentation". In: *Release 0.15 2*.
- Lotan, G (2014). "Networked audiences: attention and data-informed". In: *The New Ethics of Journalism: Principles for the 21st Century*, pp. 105–122.
- Louis, Annie and Nenkova, Ani (2013). "What makes writing great? First experiments on article quality prediction in the science journalism domain". In: *Transactions of the Association for Computational Linguistics* 1, pp. 341–352.
- Lowry, B (2013). *6 Image qualities which may drive more likes on Instagram*. URL: <https://www.digitalinformationworld.com/2013/11/how-to-get-more-likes-on-instagram.html>.
- Lundberg, Scott M, Erion, Gabriel, Chen, Hugh, DeGrave, Alex, Prutkin, Jordan M, Nair, Bala, Katz, Ronit, Himmelfarb, Jonathan, Bansal, Nisha, and Lee, Su-In (2020). "From local ex-

- planations to global understanding with explainable AI for trees”. In: *Nature machine intelligence* 2.1, pp. 56–67.
- Lundberg, Scott M and Lee, Su-In (2017). “A unified approach to interpreting model predictions”. In: *Advances in neural information processing systems*, pp. 4765–4774.
- Ma, Long, Lee, Chei Sian, and Goh, Dion Hoe-Lian (2014). “Understanding news sharing in social media: An explanation from the diffusion of innovations theory”. In: *Online information review* 38.5, pp. 598–615.
- Maffei, Lucia (2016). *Robots will cover the Olympics for The Washington Post*. Techcrunch. URL: https://techcrunch.com/2016/08/05/robots-will-cover-the-olympics-for-the-washington-post/?guccounter=1&guce_referrer=aHR0cHM6Ly93d3cuZ29vZ2x1LmNvbS8&guce_referrer_sig=AQAAAIuj7pyyIFZiXuc_CVFEMuJjk6SjxeAyIvrWFaf3Ph2UIz1AzXzkazZ2y_22D5jXTLr1VR1XiHIuXwZ2aj-I6ZAWWJIB2W6d7r8y8Z6AODLD5V62IRGA0yg7KaaUL8H6v7XkT7akt7XWVvXmvZREgIVvRxxzBPN16o2hE_kCokgE.
- Mahyoob, Mohammad, Al-Garaady, Jeehaan, and Alrahaili, Musaad (2020). “Linguistic-Based Detection of Fake News in Social Media”. In: *Forthcoming, International Journal of English Linguistics* 11.1.
- Maier, Scott R, Slovic, Paul, and Mayorga, Marcus (2017). “Reader reaction to news of mass suffering: Assessing the influence of story form and emotional response”. In: *Journalism* 18.8, pp. 1011–1029.
- Malthouse, Edward C, Haenlein, Michael, Skiera, Bernd, Wege, Egbert, and Zhang, Michael (2013). “Managing customer relationships in the social media era: Introducing the social CRM house”. In: *Journal of interactive marketing* 27.4, pp. 270–280.
- Manovich, Lev (2018). “Digital traces in context| 100 billion data rows per second: Media analytics in the early 21st century”. In: *International journal of communication* 12, p. 16.
- Maras, Steven (2013). *Objectivity in journalism*. John Wiley & Sons.

- Marconi, F. (2020). *Newsmakers: artificial intelligence and the future of journalism*. Columbia University Press.
- Marquardt, Dorota (2019). “Linguistic Indicators in the Identification of Fake News”. In: *Mediatization Studies* 3, pp. 95–114.
- Martin, James R and White, Peter R (2003). *The language of evaluation*. Vol. 2. Springer.
- Marwick A. Kuo R., Cameron S. J. and Weigel, M. (2021). *Critical Disinformation Studies: A Syllabus*.
- Masip, Pere, Suau, Jaume, Ruiz-Caballero, Carles, Capilla, Pablo, and Zilles, Klaus (2021). “News Engagement on Closed Platforms. Human Factors and Technological Affordances Influencing Exposure to News on WhatsApp”. In: *Digital Journalism* 9.8, pp. 1062–1084.
- Massumi, Brian (1995). “The autonomy of affect”. In: *Cultural critique* 31, pp. 83–109.
- (2021). *Parables for the virtual: Movement, affect, sensation*. Duke University Press.
- Maynard, Patrick (1983). “The secular icon: Photography and the functions of images”. In: *The Journal of Aesthetics and Art Criticism* 42.2, pp. 155–169.
- Mazloom, Masoud, Rietveld, Robert, Rudinac, Stevan, Worrying, Marcel, and Van Dolen, Willemijn (2016). “Multimodal popularity prediction of brand-related social media posts”. In: *Proceedings of the 24th ACM international conference on Multimedia*, pp. 197–201.
- McCarthy, John (2004). “What is artificial intelligence”. In: URL: <http://www-formal.stanford.edu/jmc/whatisai.html>.
- McCollough, Kathleen, Crowell, Jessica K, and Napoli, Philip M (2017). “Portrait of the online local news audience”. In: *Digital Journalism* 5.1, pp. 100–118.
- McLeod, Jack M (2000). “Media and civic socialization of youth”. In: *Journal of Adolescent Health* 27.2, pp. 45–51.

- McMillan, Sally J (2005). "The researchers and the concept: Moving beyond a blind examination of interactivity". In: *Journal of Interactive Advertising* 5.2, pp. 1–4.
- McParlane, Philip J, Moshfeghi, Yashar, and Jose, Joemon M (2013). "On contextual photo tag recommendation". In: *Proceedings of the 36th international ACM SIGIR conference on research and development in information retrieval*, pp. 965–968.
- (2014). "'Nobody comes here anymore, it's too crowded"; Predicting Image Popularity on Flickr". In: *Proceedings of international conference on multimedia retrieval*, pp. 385–391.
- McQuail, Denis (1992). *Media performance: Mass communication and the public interest*. Sage.
- (1994). "The rise of media of mass communication". In: *Mass communication theory: An introduction*, pp. 1–29.
- (2005). *McQuail's mass communication theory*. Sage publications.
- (2015). "Media performance". In: *The International Encyclopedia of Political Communication*, pp. 1–9.
- Medium.com (n.d.). *Medium's Curation Guidelines: everything writers need to know*. Accessed: 2020-04-09. URL: <https://help.medium.com/hc/en-us/articles/360006362473-Medium-s-Curation-Guidelines-everything-writers-need-to-know>.
- Meijer, Irene Costera (2001). "The public quality of popular journalism: Developing a normative framework". In: *Journalism studies* 2.2, pp. 189–205.
- (2003). "What is quality television news? A plea for extending the professional repertoire of newsmakers". In: *Journalism studies* 4.1, pp. 15–29.
- Mendelson, Andrew L and Papacharissi, Zizi (2010). "Look at us: Collective narcissism in college student Facebook photo galleries". In: *A Networked Self*. Routledge, pp. 259–281.

- Messaris, Paul (1992). "Visual" manipulation": Visual means of affecting responses to images." In: *Communication*.
- Messing, Solomon and Westwood, Sean J (2014). "Selective exposure in the age of social media: Endorsements trump partisan source affiliation when selecting news online". In: *Communication research* 41.8, pp. 1042–1063.
- Mikhail Korobov, Konstantin Lopuhin (2016). "ELI5 is a Python package which helps to debug machine learning classifiers and explain their predictions". In: *To appear*.
- Mikolov, Tomas, Chen, Kai, Corrado, Greg, and Dean, Jeffrey (2013). "Efficient Estimation of Word Representations in Vector Space". In: *Proceedings of the International Conference on Learning Representations (ICLR)*, pp. 1–12.
- Milner, Ryan M (2018). *The world made meme: Public conversations and participatory media*. mit Press.
- Miltner, Kate M and Highfield, Tim (2017). "Never gonna GIF you up: Analyzing the cultural significance of the animated GIF". In: *Social Media+ Society* 3.3, p. 2056305117725223.
- Mirzoeff, Nicholas (1999). *An introduction to visual culture*. Psychology press.
- Mitchell, WJ Thomas (1995). *Picture theory: Essays on verbal and visual representation*. University of Chicago Press.
- (2005). *What do pictures want?: The lives and loves of images*. University of Chicago Press.
- Mitchelstein, Eugenia and Boczkowski, Pablo J (2010). "Online news consumption research: An assessment of past work and an agenda for the future". In: *New media & society* 12.7, pp. 1085–1102.
- Mohammad, Saif (2018). "Obtaining reliable human ratings of valence, arousal, and dominance for 20,000 english words". In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 174–184.

- Mohammad, Saif and Turney, Peter (2010). "Emotions evoked by common words and phrases: Using mechanical turk to create an emotion lexicon". In: *Proceedings of the NAACL HLT 2010 workshop on computational approaches to analysis and generation of emotion in text*, pp. 26–34.
- Mohammad, Saif M (2017). "Word affect intensities". In: *arXiv preprint arXiv:1704.08798*.
- (2018). "Obtaining Reliable Human Ratings of Valence, Arousal, and Dominance for 20,000 English Words". In: *Proceedings of The Annual Conference of the Association for Computational Linguistics (ACL)*. Melbourne, Australia.
- Molyneux, Logan (2015). "What journalists retweet: Opinion, humor, and brand development on Twitter". In: *Journalism* 16.7, pp. 920–935.
- Moore, Fraser (2018). *A former BBC war correspondent explains what news organisations get wrong about reporting conflicts*. Business Insider. URL: <https://www.businessinsider.com/martin-bell-the-biggest-mistakes-made-when-reporting-war-bosnia-srebrenica-2018-2>.
- Moore, Sarah G and McFerran, Brent (2017). "She said, she said: Differential interpersonal similarities predict unique linguistic mimicry in online word of mouth". In: *Journal of the Association for Consumer Research* 2.2, pp. 229–245.
- Mortensen, Mette, Allan, Stuart, and Peters, Chris (2017). "The iconic image in a digital age". In: *Nordicom Review* 38.s2, pp. 71–86.
- Mujica, Constanza and Bachmann, Ingrid (2018). "The impact of melodramatic news coverage on information recall and comprehension". In: *Journalism studies* 19.3, pp. 334–352.
- Munger, Kevin (2020). "All the news that's fit to click: The economics of clickbait media". In: *Political Communication* 37.3, pp. 376–397.
- Muñoz-Torres, Juan Ramón (2012). "Truth and objectivity in journalism: Anatomy of an endless misunderstanding". In: *Journalism studies* 13.4, pp. 566–582.

- Napoli, Philip (2015). "Social media and the public interest: Governance of news platforms in the realm of individual and algorithmic gatekeepers". In: *Telecommunications Policy* 39.9, pp. 751–760.
- Napoli, Philip M (1999). "Deconstructing the diversity principle". In: *Journal of communication* 49.4, pp. 7–34.
- (2011). "Exposure diversity reconsidered". In: *Journal of information policy* 1, pp. 246–259.
- Nechushtai, Efrat and Lewis, Seth C (2019). "What kind of news gatekeepers do we want machines to be? Filter bubbles, fragmentation, and the normative dimensions of algorithmic recommendations". In: *Computers in Human Behavior* 90, pp. 298–307.
- Nelson, Jacob L (2020). "The elusive engagement metric". In: *Measurable Journalism*. Routledge, pp. 140–156.
- (2021). "The next media regime: The pursuit of 'audience engagement' in journalism". In: *Journalism* 22.9, pp. 2350–2367.
- Newman, Nic, Fletcher, Richard, Kalogeropoulos, Antonis, and Nielsen, Rasmus (2019). *Reuters institute digital news report 2019*. Vol. 2019. Reuters Institute for the Study of Journalism.
- Nikunen, Kaarina (2018). *Media solidarities: Emotions, power and justice in the digital age*. Sage.
- Nissim, Malvina and Patti, Viviana (2017). "Semantic aspects in sentiment analysis". In: *Sentiment analysis in social networks*. Elsevier, pp. 31–48.
- Noguera-Vivo, José Manuel (2018). "You get what you give: Sharing as a new radical challenge for journalism". In: *Communication & society* 31.4, pp. 147–158.
- Nöth, Winfried (2007). "Self-reference in the media: The semiotic framework". In: *Self-reference in the Media*, pp. 3–30.

- Novaković, Jasmina Dj, Veljović, Alempije, Ilić, Siniša S, Papić, Željko, and Tomović, Milica (2017). "Evaluation of classification models in machine learning". In: *Theory and Applications of Mathematics & Computer Science* 7.1, p. 39.
- Nummenmaa, Lauri, Glerean, Enrico, Hari, Riitta, and Hietanen, Jari K (2014). "Bodily maps of emotions". In: *Proceedings of the National Academy of Sciences* 111.2, pp. 646–651.
- Nygren, Gunnar (2012). "Autonomy—A Crucial Element of Professionalization". In: *JOURNALISTIKSTUDIEN VID SÖDERTÖRNS HÖGSKOLA* 4, p. 73.
- O’neill, Deirdre and Harcup, Tony (2009). "News values and selectivity". In: *The handbook of journalism studies*. Routledge, pp. 181–194.
- Oakes, Oakes Michael (2019). *Statistics for corpus linguistics*. Edinburgh University Press.
- Oeldorf-Hirsch, Anne and Sundar, S Shyam (2015). "Posting, commenting, and tagging: Effects of sharing news stories on Facebook". In: *Computers in human behavior* 44, pp. 240–249.
- Okrent, Arika (2014). "The listicle as literary form". In: *University of Chicago Magazine* 106.3, pp. 52–53.
- Olausson, Ulrika (2017). "The reinvented journalist: The discursive construction of professional identity on Twitter". In: *Digital Journalism* 5.1, pp. 61–81.
- Olivieri, Alex, Shabani, Shaban, Sokhn, Maria, and Cudré-Mauroux, Philippe (2019). "Creating task-generic features for fake news detection". In: *Proceedings of the 52nd Hawaii International Conference on System Sciences*.
- Orellana-Rodriguez, Claudia and Keane, Mark T (2018). "Attention to news and its dissemination on Twitter: A survey". In: *Computer Science Review* 29, pp. 74–94.
- Orgeret, Kristin Skare (2020). "Discussing emotions in digital journalism". In: *Digital Journalism* 8.2, pp. 292–297.

- Ørmen, Jacob (2015). "A public conversation in private settings". PhD thesis. PhD Dissertation, University of Copenhagen, Copenhagen.
- (2019). "From consumer demand to user engagement: Comparing the popularity and virality of election coverage on the Internet". In: *The International Journal of Press/Politics* 24.1, pp. 49–68.
- Oschatz, Corinna, Emde-Lachmund, Katharina, and Klimmt, Christoph (2021). "The persuasive effect of journalistic storytelling: Experiments on the portrayal of exemplars in the news". In: *Journalism & Mass Communication Quarterly* 98.2, pp. 407–427.
- Östgaard, Einar (1965). "Factors influencing the flow of news". In: *Journal of peace Research* 2.1, pp. 39–63.
- Paivio, Allan (1991). *Images in mind: the evolution of a theory*. Harvester Wheatsheaf.
- Palmer, Daniel (2010). "Emotional archives: Online photo sharing and the cultivation of the self". In: *Photographies* 3.2, pp. 155–171.
- Palomo, Bella, Teruel, Laura, and Blanco-Castilla, Elena (2019). "Data journalism projects based on user-generated content. How La Nacion data transforms active audience into staff". In: *Digital Journalism* 7.9, pp. 1270–1288.
- Pantti, Mervi (2010). "The value of emotion: An examination of television journalists' notions on emotionality". In: *European Journal of Communication* 25.2, pp. 168–181.
- Pantti, Mervi and Sumiala, Johanna (2009). "Till death do us join: Media, mourning rituals and the sacred centre of the society". In: *Media, culture & society* 31.1, pp. 119–135.
- Pantti, Mervi and Wahl-Jorgensen, Karin (2021). *Journalism and emotional work*.
- Papacharissi, Zizi (2015). *Affective publics: Sentiment, technology, and politics*. Oxford University Press.

- Papacharissi, Zizi and Fatima Oliveira, Maria de (2012). "Affective news and networked publics: The rhythms of news storytelling on hashtag Egypt". In: *Journal of communication* 62.2, pp. 266–282.
- Papathanassopoulos, Stylianos (2002). "European television in the digital age: Issues, dynamics and realities". In: *(No Title)*.
- Papathanassopoulos, Stylianos, Karadimitriou, Achilleas, Kostopoulos, Christos, and Archontaki, Ioanna (2021). "Media concentration and independent journalism between austerity and digital disruption". In: *THE MEDIA FOR DEMOCRACY MONITOR 2021*, p. 177.
- Papathanassopoulos, Stylianos and Miconi, Andrea (2023). *The Media Systems in Europe: Continuities and Discontinuities*.
- Parasie, Sylvain and Dagiral, Eric (2013). "Data-driven journalism and the public good: "Computer-assisted-reporters" and "programmer-journalists" in Chicago". In: *New media & society* 15.6, pp. 853–871.
- Pariser, Eli (2011). *The filter bubble: What the Internet is hiding from you*. Penguin UK.
- Park, Deokgun, Sachar, Simranjit, Diakopoulos, Nicholas, and Elmqvist, Niklas (2016). "Supporting comment moderators in identifying high quality online news comments". In: *Conference on Human Factors in Computing Systems - Proceedings*, pp. 1114–1125. DOI: 10.1145/2858036.2858389.
- Park, S, Fisher, C, Lee, JY, McGuinness, K, Sang, Y, O'Neil, M, Jensen, M, McCallum, K, and Fuller, G (2020). *Digital news report: Australia 2020*. Canberra: News & Media Research Centre, University of Canberra.
- Park, Sora, Sang, Yoonmo, Jung, Jaemin, and Stroud, Natalie Jomini (2021). *News Engagement: The Roles of Technological Affordance, Emotion, and Social Endorsement*.
- Paul, Sharoda A, Hong, Lichan, and Chi, Ed H (2012). "Who is authoritative? understanding reputation mechanisms in quora". In: *arXiv preprint arXiv:1204.3724*.

- Pavlik, John V (2013). "Trends in new media research: A critical review of recent scholarship". In: *Sociology Compass* 7.1, pp. 1–12.
- Pearce, Warren, Özkula, Suay M, Greene, Amanda K, Teeling, Lauren, Bansard, Jennifer S, Omena, Janna Joceli, and Rabello, Elaine Teixeira (2020). "Visual cross-platform analysis: Digital methods to research social media images". In: *Information, Communication & Society* 23.2, pp. 161–180.
- Pedersen, Marius, Bonnier, Nicolas, Hardeberg, Jon Yngve, and Albrechtsen, Fritz (2010). "Attributes of image quality for color prints". In: *Journal of Electronic Imaging* 19.1, p. 011016.
- Pedregosa, F, Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P, Weiss, R., Dubourg, V, Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., and Duchesnay, E. (2011). "Scikit-learn: Machine Learning in Python". In: *Journal of Machine Learning Research* 12, pp. 2825–2830.
- Pellegrini, Tassilo (2012). "Semantic metadata in the news production process: Achievements and challenges". In: *Proceeding of the 16th international academic mindtrek conference*, pp. 125–133.
- Pempek, Tiffany A, Yermolayeva, Yevdokiya A, and Calvert, Sandra L (2009). "College students' social networking experiences on Facebook". In: *Journal of applied developmental psychology* 30.3, pp. 227–238.
- Peng, Tai-Quan, Liang, Hai, and Zhu, Jonathan JH (2019). *Introducing computational social science for Asia-Pacific communication research*.
- Pennington, Jeffrey, Socher, Richard, and Manning, Christopher (2014). "Glove: Global Vectors for Word Representation". In: *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pp. 1532–1543.
- Pereira, Francisco, Mitchell, Tom, and Botvinick, Matthew (2009). "Machine learning classifiers and fMRI: a tutorial overview". In: *Neuroimage* 45.1, S199–S209.

- Pérez-Rosas, Verónica, Kleinberg, Bennett, Lefevre, Alexandra, and Mihalcea, Rada (2017). "Automatic detection of fake news". In: *arXiv preprint arXiv:1708.07104*.
- Peters, Chris (2011). "Emotion aside or emotional side? Crafting an 'experience of involvement' in the news". In: *Journalism* 12.3, pp. 297–316.
- Picone, Ike, Kleut, Jelena, Pavličková, Tereza, Romic, Bojana, Møller Hartley, Jannie, and De Ridder, Sander (2019). "Small acts of engagement: Reconnecting productive audience practices with everyday agency". In: *New Media & Society* 21.9, pp. 2010–2028.
- Pieters, Rik and Wedel, Michel (2004). "Attention capture and transfer in advertising: Brand, pictorial, and text-size effects". In: *Journal of marketing* 68.2, pp. 36–50.
- Pimentel, João Felipe, Murta, Leonardo, Braganholo, Vanessa, and Freire, Juliana (2021). "Understanding and improving the quality and reproducibility of Jupyter notebooks". In: *Empirical Software Engineering* 26.4, pp. 1–55.
- Planer, Rosanna, Godulla, Alexander, Seibert, Daniel, and Pietsch, Patrick (2022). "Journalistic Quality Criteria under the Magnifying Glass: A Content Analysis of the Winning Stories of World Press Photo Foundation's Digital Storytelling Contest". In: *Journalism and Media* 3.4, pp. 594–614.
- Plutchik, Robert (1980a). "A general psychoevolutionary theory of emotion". In: *Theories of emotion*. Elsevier, pp. 3–33.
- (1980b). "Emotion". In: *A psychoevolutionary synthesis*.
- Popović, Virginia and Popović, Predrag (2014). "The twenty-first century, the reign of tabloid journalism". In: *Procedia-Social and Behavioral Sciences* 163, pp. 12–18.
- Porlezza, Colin and Splendore, Sergio (2019). "From open journalism to closed data: Data journalism in Italy". In: *Digital journalism* 7.9, pp. 1230–1252.
- Potter, Mary C, Wyble, Brad, Hagmann, Carl Erick, and McCourt, Emily S (2014). "Attention, Perception, & Psychophysics". In:

- Poulakidakos, S, Sotirakou, C, Armenakis, A, Mandenaki, K, A, Karampela, and Mourlas, C (2023). "Deliverable of IQ Project:" in: 2.2.
- Purba, Kristo Radion, Asirvatham, David, and Murugesan, Raja Kumar (2021). "Instagram post popularity trend analysis and prediction using hashtag, image assessment, and user history features." In: *Int. Arab J. Inf. Technol.* 18.1, pp. 85–94.
- Purcell, Kristen, Rainie, Lee, Mitchell, Amy, Rosenstiel, Tom, and Olmstead, Kenny (2010). "Understanding the participatory news consumer". In: *Pew Internet and American Life Project 1*, pp. 19–21.
- Qian, Crystal J, Tang, Jonathan D, Penza, Matthew A, and Ferri, Christopher M (2017). "Instagram Popularity Prediction via Neural Networks and Regression Analysis". In: *IEEE Transactions on Multimedia* 19.
- Quandt, Thorsten (2008). "News on the World Wide Web? A comparative content analysis of online news in Europe and the United States". In: *Journalism Studies* 9.5, pp. 717–738.
- (2018). "Dark participation". In: *Media and communication* 6.4, pp. 36–48.
- Rashkin, Hannah, Choi, Eunsol, Jang, Jin Yea, Volkova, Svitlana, and Choi, Yejin (2017). "Truth of varying shades: Analyzing language in fake news and political fact-checking". In: *Proceedings of the 2017 conference on empirical methods in natural language processing*, pp. 2931–2937.
- Rayson, Steve (2017). *2017's Most Shared Facebook Content: Viral Posts, Videos and Articles*. Buzzsumo. URL: <https://buzzsumo.com/blog/the-most-shared-facebook-content-posts-videos/>.
- Reddit (n.d.). *Reddit*. Accessed: 2020-04-09. URL: <https://www.reddit.com>.
- Reddy, William M (1997). "Against constructionism: the historical ethnography of emotions". In: *Current anthropology* 38.3, pp. 327–351.

- Redies, Christoph, Grebenkina, Maria, Mohseni, Mahdi, Kaduhm, Ali, and Dobel, Christian (2020). "Global image properties predict ratings of affective pictures". In: *Frontiers in psychology* 11, p. 953.
- Reinemann, Carsten, Stanyer, James, Scherr, Sebastian, and Legnante, Guido (2012). "Hard and soft news: A review of concepts, operationalizations and key findings". In: *Journalism* 13.2, pp. 221–239.
- Reis, Julio CS, Correia, André, Murai, Fabrício, Veloso, Adriano, and Benevenuto, Fabrício (2019). "Supervised learning for fake news detection". In: *IEEE Intelligent Systems* 34.2, pp. 76–81.
- Ribeiro, Marco Tulio, Singh, Sameer, and Guestrin, Carlos (2016). "' Why should I trust you?' Explaining the predictions of any classifier". In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144.
- Richards, Barry (2007). *Emotional governance: Politics, media and terror*. Springer.
- (2009). "News and the emotional public sphere". In: *The Routledge companion to news and journalism*. Routledge, pp. 345–355.
- Rick, Edmonds (2019). *The New York Times sells premium ads based on how an article makes you feel*. Poynter. URL: <https://www.poynter.org/business-work/2019/the-new-york-times-sells-premium-ads-based-on-how-an-article-makes-you-feel/>.
- Rivo, Eduardo, Fuente, Javier de la, Rivo, Ángel, García-Fontán, Eva, Cañizares, Miguel-Ángel, and Gil, Pedro (2012). "Cross-Industry Standard Process for data mining is applicable to the lung cancer surgery domain, improving decision making as well as knowledge and quality management". In: *Clinical and Translational Oncology* 14.1, pp. 73–79.
- Robertson, Craig T and Mourão, Rachel R (2020). "Faking alternative journalism? An analysis of self-presentations of "fake news" sites". In: *Digital Journalism* 8.8, pp. 1011–1029.
- Robinson, James G. (n.d.). *Tow Center for Digital Journalism*. Accessed: 2020-01-07.

- Rogers, Simon (2014). "Data journalism is the new punk". In: *British journalism review* 25.2, pp. 31–34.
- Rony, Md Main Uddin, Hassan, Naeemul, and Yousuf, Mohammad (2017). "Diving deep into clickbaits: Who use them to what extents in which topics with what effects?" In: *Proceedings of the 2017 IEEE/ACM international conference on advances in social networks analysis and mining 2017*, pp. 232–239.
- Rosas, Omar V (2018). "Strategic Avoidance and Strategic Use: A Look into Spanish Online Journalists' Attitudes Toward Emotions in Reporting." In: *International Journal of Communication* (19328036) 12.
- Rosenstiel, Tom, Just, Marion, Belt, Todd, Pertilla, Atiba, Dean, Walter, and Chinni, Dante (2007). *We interrupt this newscast: How to improve local news and win ratings, too*. Cambridge University Press.
- Royal, Cindy (2010). "The journalist as programmer: A case study of the New York Times interactive news technology department". In: *International Symposium on Online Journalism*. Vol. 2. 1. Citeseer, pp. 5–24.
- Rubin, Victoria L, Conroy, Niall, Chen, Yimin, and Cornwell, Sarah (2016). "Fake news or truth? using satirical cues to detect potentially misleading news". In: *Proceedings of the second workshop on computational approaches to deception detection*, pp. 7–17.
- Ruggiero, Thomas E (2000). "Uses and gratifications theory in the 21st century". In: *Mass communication & society* 3.1, pp. 3–37.
- Ruotsalainen, Juho (2018). "Scanning the shape of journalism—Emerging trends, changing culture?" In: *Futures* 104, pp. 14–24.
- Russell, James A (1980). "A circumplex model of affect." In: *Journal of personality and social psychology* 39.6, p. 1161.

- (2003). “Core affect and the psychological construction of emotion.” In: *Psychological review* 110.1, p. 145.
- Salganik, Matthew J (2019). *Bit by bit: Social research in the digital age*. Princeton University Press.
- Saltz, Jeff (2022). *CRISP-DM is Still the Most Popular Framework for Executing Data Science Projects*. URL: <https://www.datascience-pm.com/crisp-dm-still-most-popular/>.
- Sánchez, S. (2021). *The Wall Street Journal uses Narrativa’s AI for its news automation*. *Narrativa*. Retrieved 01.05.2022. URL: <https://www.narrativa.com/the-wall-street-journal-uses-narrativas-ai-for-its-news-automation/>.
- Sánchez Laws, Ana Luisa (2020). “Can immersive journalism enhance empathy?” In: *Digital journalism* 8.2, pp. 213–228.
- Sang, Yoonmo, Lee, Jee Young, Park, Sora, Fisher, Caroline, and Fuller, Glen (2020). “Signalling and expressive interaction: Online news users’ different modes of interaction on digital platforms”. In: *Digital journalism* 8.4, pp. 467–485.
- Scacco, Joshua and Muddiman, Ashley (2016). “Investigating the influence of “clickbait” news headlines”. In: *Engaging News Project Report*.
- Schaudt, Sky and Carpenter, Serena (2009). “The News That’s Fit to Click: An Analysis of Online News Values and Preferences Present in the Most-viewed Stories on azcentral.com.” In: *Southwestern Mass Communication Journal* 24.2.
- Schmidt, Jan-Hinrik (2014). “Twitter and the rise of personal publics”. In: *Twitter and society*, pp. 3–14.
- Schmidt, Thomas R (2021). “‘It’s OK to feel’: The emotionality norm and its evolution in US print journalism”. In: *Journalism* 22.5, pp. 1173–1189.
- Schudson, Michael (2001). “The objectivity norm in American journalism”. In: *Journalism* 2.2, pp. 149–170.

- Schudson, Michael (2008). "News and Democratic Society: past, present, and future". In: *Hedgehog Review* 10.2, pp. 7–21.
- Schwartz, Evan I (2002). *Digital Darwinism: 7 breakthrough business strategies for surviving in the cutthroat Web economy*. Currency.
- Serafeim, Katerina (2012). "The impact of social media on press freedom in Greece: Benefits, challenges and limitations". In: *ESSACHESS Journal for Communication Studies* 5.1, p. 9.
- Shah, Cappella, Joseph N, and Neuman, W Russell (2015). "Big data, digital media, and computational social science: Possibilities and perils". In: *The ANNALS of the American Academy of Political and Social Science* 659.1, pp. 6–13.
- Shah, Hetan (2018). "Algorithmic accountability". In: *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 376.2128, p. 20170362.
- Shao, Chengcheng, Ciampaglia, Giovanni Luca, Varol, Onur, Flammini, Alessandro, and Menczer, Filippo (2017). "The spread of fake news by social bots". In: *arXiv preprint arXiv:1707.07592* 96, p. 104.
- She, Jie, Zhang, Tao, Chen, Qun, Zhang, Jianzhang, Fan, Weiguo, Wang, Hongwei, and Chang, Qingqing (2021). "Which social media posts generate the most buzz? Evidence from WeChat". In: *Internet Research*.
- Shearer, Colin (2000). "The CRISP-DM model: the new blueprint for data mining". In: *Journal of data warehousing* 5.4, pp. 13–22.
- Shehata, Adam and Strömbäck, Jesper (2021). "Learning political news from social media: Network media logic and current affairs news learning in a high-choice media environment". In: *Communication Research* 48.1, pp. 125–147.
- Shin, Jieun and Thorson, Kjerstin (2017). "Partisan selective sharing: The biased diffusion of fact-checking messages on social media". In: *Journal of Communication* 67.2, pp. 233–255.

- Shoemaker, Pamela J and Cohen, Akiba A (2012). *News around the world: Content, practitioners, and the public*. Routledge.
- Shoemaker, Pamela J and Reese, Stephen D (1996). *Mediating the message*. White Plains, NY: Longman.
- Shoemaker, Pamela J and Vos, Timothy (2009). *Gatekeeping theory*. Routledge.
- Shouse, Eric (2005). "Feeling, emotion, affect". In: *M/c journal* 8.6.
- Shu, Kai, Mahudeswaran, Deepak, Wang, Suhang, Lee, Dongwon, and Liu, Huan (2020). "Fakenewsnet: A data repository with news content, social context, and spatiotemporal information for studying fake news on social media". In: *Big Data* 8.3, pp. 171–188.
- Shu, Kai, Sliva, Amy, Wang, Suhang, Tang, Jiliang, and Liu, Huan (2017). "Fake news detection on social media: A data mining perspective. CoRR abs/1708.01967 (2017)". In: *arXiv preprint arXiv:1708.01967*.
- Singer, Jane (2014). "User-generated visibility: Secondary gatekeeping in a shared media space". In: *New media & society* 16.1, pp. 55–73.
- Singer, Jane B (2005). "The political j-blogger: 'Normalizing' a new media form to fit old norms and practices". In: *Journalism* 6.2, pp. 173–198.
- Singer, Jane B, Domingo, David, Heinonen, Ari, Hermida, Alfred, Paulussen, Steve, Quandt, Thorsten, Reich, Zvi, and Vujnovic, Marina (2011). *Participatory journalism: Guarding open gates at online newspapers*. John Wiley & Sons.
- Slattery, Karen L and Hakanen, Ernest A (1994). "Trend: Sensationalism versus public affairs content of local TV news: Pennsylvania revisited". In: *Journal of Broadcasting & Electronic Media* 38.2, pp. 205–216.
- Sloam, James and Henn, Matt (2018). "Youthquake 2017: how the rise of young cosmopolitans in Britain could transform politics". In: *Democratic Audit Blog*.

- Snodgrass, Joan Gay, Wasser, Barry, Finkelstein, Marjorie, and Goldberg, Linda Brainin (1974). "On the fate of visual and verbal memory codes for pictures and words: Evidence for a dual coding mechanism in recognition memory". In: *Journal of verbal learning and verbal behavior* 13.1, pp. 27–37.
- Song, Hyunjin, Tolochko, Petro, Eberl, Jakob-Moritz, Eisele, Olga, Greussing, Esther, Heidenreich, Tobias, Lind, Fabienne, Galyga, Sebastian, and Boomgaarden, Hajo G (2020). "In validations we trust? The impact of imperfect human annotations as a gold standard on the quality of validation of automated content analysis". In: *Political Communication* 37.4, pp. 550–572.
- Soroka, Stuart and McAdams, Stephen (2015). "News, politics, and negativity". In: *Political Communication* 32.1, pp. 1–22.
- Sotirakou, C, Germanakos, P, Holzinger, A, and Mourlas, C (2018). "Feedback Matters! Predicting the Appreciation of Online Articles A Data-Driven Approach". In: *International Cross-Domain Conference for Machine Learning and Knowledge Extraction*. Springer, pp. 147–159.
- Sotirakou, C, Karampela, A, and Mourlas, C (2021). "Evaluating the Role of News Content and Social Media Interactions for Fake News Detection". In: *Disinformation in Open Online Media: Third Multidisciplinary International Symposium, MISDOOM 2021, Virtual Event, September 21–22, 2021, Proceedings* 3. Springer, pp. 128–141.
- Sotirakou, C, Koutromanou, E, and Mourlas, C (2023). "Exploring the Impact of Featured Images in News Stories using Machine Learning". In: *Future Technologies Conference (FTC) 2023, Lecture Notes in Networks and Systems*.
- Sotirakou, C, Trilling, D, Germanakos, P, and Mourlas, C (2019). "Opening the black box of perceived quality: Predicting endorsement on a blog site". In: *2019 IEEE/WIC/ACM International Conference on Web Intelligence (WI)*. IEEE, pp. 388–392.

- Spangher, Alexander (2018). *How Does This Article Make You Feel? Using data science to predict the emotional resonance of New York Times articles for better ad placement*. <https://open.nytimes.com/how-does-this-article-make-you-feel-4684e5e9c47>, Retrieved 20-7-2020.
- Sparks, Colin (1998). "Introduction: Tabloidization and the media". In: *Javnost–The Public* 5.3, pp. 5–10.
- Sparrow, Bartholomew H (2006). "A research agenda for an institutional media". In: *Political communication* 23.2, pp. 145–157.
- Splendore, Sergio, Di Salvo, Philip, Eberwein, Tobias, Groenhart, Harmen, Kus, Michal, and Porlezza, Colin (2016). "Educational strategies in data journalism: A comparative study of six European countries". In: *Journalism* 17.1, pp. 138–152.
- Staab, Joachim Friedrich (1990). "The role of news factors in news selection: A theoretical reconsideration". In: *European Journal of communication* 5.4, pp. 423–443.
- Steensen, Steen (2016). "The intimization of journalism". In: *Handbook of Digital Journalism Studies*, pp. 113–127.
- (2017). "Subjectivity as a journalistic ideal". In:
- Steensen, Steen, Ferrer-Conill, Raul, and Peters, Chris (2020). "(Against a) theory of audience engagement with news". In: *Journalism Studies* 21.12, pp. 1662–1680.
- Stenvall, Maija (2008). "On emotions and the journalistic ideals of factuality and objectivity—Tools for analysis". In: *Journal of Pragmatics* 40.9, pp. 1569–1586.
- Stepaniuk, Krzysztof (2015). "The relation between destination image and social media user engagement—theoretical approach". In: *Procedia-Social and Behavioral Sciences* 213, pp. 616–621.
- Stoneman, Jonathan (2015). "Does open data need journalism?" In:

- Strapparava, Carlo and Mihalcea, Rada (2007). "Semeval-2007 task 14: Affective text". In: *Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)*, pp. 70–74.
- Stray, Jonathan (2010a). *A full-text visualization of the Iraq War Logs.*, Retrieved 11.11.2022. URL: <http://jonathanstray.com/a-full-text-visualization-of-the-iraq-war-logs>.
- (2010b). *Paradise Papers*. URL: <http://jonathanstray.com/a-full-text-visualization-of-the-iraq-war-logs>.
- (2019). "Making artificial intelligence work for investigative journalism". In: *Digital Journalism* 7.8, pp. 1076–1097.
- Stroud, Natalie Jomini (2008). "Media use and political predispositions: Revisiting the concept of selective exposure". In: *Political Behavior* 30.3, pp. 341–366.
- (2011). *Niche news: The politics of news choice*. Oxford University Press on Demand.
- Sturm Wilkerson, Heloisa, Riedl, Martin J, and Whipple, Kelsey N (2021). "Affective affordances: Exploring facebook reactions as emotional responses to hyperpartisan political news". In: *Digital Journalism* 9.8, pp. 1040–1061.
- Szabo, Gabor and Huberman, Bernardo A (2010). "Predicting the popularity of online content". In: *Communications of the ACM* 53.8, pp. 80–88.
- Tacchini, Eugenio, Ballarin, Gabriele, Della Vedova, Marco L, Moret, Stefano, and Alfaro, Luca de (2017). "Some like it hoax: Automated fake news detection in social networks". In: *arXiv preprint arXiv:1704.07506*.
- Tafesse, Wondwesen and Wien, Anders (2017). "A framework for categorizing social media posts". In: *Cogent Business & Management* 4.1, p. 1284390.

- Tandoc, Edson and Thomas (2015). "The Ethics of Web Analytics: Implications of using audience metrics in news construction". In: *Digital Journalism* 3.2, pp. 243–258. ISSN: 2167082X. DOI: 10.1080/21670811.2014.909122.
- Tandoc, Edson and Vos (2016). "The journalist is marketing the news: Social media in the gatekeeping process". In: *Journalism Practice* 10.8, pp. 950–966.
- Tandoc Jr, Edson C (2014). "Journalism is twerking? How web analytics is changing the process of gatekeeping". In: *New media & society* 16.4, pp. 559–575.
- (2019). *Analyzing Analytics: Disrupting Journalism One Click at a Time*. Routledge.
- Tandoc Jr, Edson C and Ferrucci, Patrick R (2017). "Giving in or giving up: What makes journalists use audience feedback in their news work?" In: *Computers in Human Behavior* 68, pp. 149–156.
- Tandoc Jr, Edson C, Lim, Zheng Wei, and Ling, Richard (2018). "Defining "fake news" A typology of scholarly definitions". In: *Digital journalism* 6.2, pp. 137–153.
- Tandoc Jr, Edson C and Oh, Soo-Kwang (2017). "Small departures, big continuities? Norms, values, and routines in The Guardian's big data journalism". In: *Journalism studies* 18.8, pp. 997–1015.
- Tang, Rong, Ng, Kwong Bor, Strzalkowski, Tomek, and Kantor, Paul B (2003). "Automatically predicting information quality in news documents". In: *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology: companion volume of the Proceedings of HLT-NAACL 2003–short papers-Volume 2*. Association for Computational Linguistics, pp. 97–99.
- Tao, Xiaohui, Velasquez-Silva, Juan Domingo, Liu, Jiming, and Zhong, Ning (2020). "Computational Social Science as the ultimate Web Intelligence". In: *World Wide Web* 23.3, pp. 1743–1745.

- Tatar, Alexandru, De Amorim, Marcelo Dias, Fdida, Serge, and Antoniadis, Panayotis (2014). “A survey on predicting the popularity of web content”. In: *Journal of Internet Services and Applications* 5.1, pp. 1–20.
- Tatar, Alexandru, Leguay, Jérémie, Antoniadis, Panayotis, Limbourg, Arnaud, Amorim, Marcelo Dias de, and Fdida, Serge (2011). “Predicting the popularity of online articles based on user comments”. In: *Proceedings of the International Conference on Web Intelligence, Mining and Semantics*, pp. 1–8.
- Tausczik, Yla R and Pennebaker, James W (2010). “The psychological meaning of words: LIWC and computerized text analysis methods”. In: *Journal of language and social psychology* 29.1, pp. 24–54.
- Tejedor, Santiago and Vila, Pere (2021). “Exo Journalism: A Conceptual Approach to a Hybrid Formula between Journalism and Artificial Intelligence”. In: *Journalism and Media* 2.4, pp. 830–840.
- Tenenboim, Ori and Cohen, Akiba A (2015). “What prompts users to click and comment: A longitudinal study of online news”. In: *Journalism* 16.2, pp. 198–217.
- Tenk, Yasmin (2021). “Five golden rules for using images to engage readers with your story”. In: *journalism.co.uk*. URL: <https://www.journalism.co.uk/news/five-rules-to-make-your-article-images-more-engaging/s2/a789231/>.
- Terence Parr, Jeremy Howard (2018). *The Mechanics of Machine Learning*. URL: <https://mlbook.explained.ai/> (visited on 09/01/2022).
- Thompson, Alex (2016). *Journalists and Trump voters live in separate online bubbles, MIT analysis shows*. Vice News. URL: <https://www.vice.com/en/article/d3xamx/journalists-and-trump-voters-live-in-separate-online-bubbles-mit-analysis-shows>.

- Thomson, TJ and Greenwood, Keith (2017). "I "like" that: Exploring the characteristics that promote social media engagement with news photographs". In: *Visual Communication Quarterly* 24.4, pp. 203–218.
- Thorpe, Simon, Fize, Denis, and Marlot, Catherine (1996). "Speed of processing in the human visual system". In: *nature* 381.6582, pp. 520–522.
- Thorson, Kjerstin and Wells, Chris (2016a). "Curated flows: A framework for mapping media exposure in the digital age". In: *Communication Theory* 26.3, pp. 309–328. ISSN: 10503293. DOI: 10.1111/comt.12087.
- (2016b). "Curated flows: A framework for mapping media exposure in the digital age". In: *Communication Theory* 26.3, pp. 309–328.
- Tierney, John (2013). "Good news beats bad on social networks". In: *The New York Times* 18.
- Tifentale, Alise and Manovich, Lev (2015). "Selfiecity: Exploring photography and self-fashioning in social media". In: *Postdigital aesthetics*. Springer, pp. 109–122.
- Tolochko, Petro and Boomgaarden, Hajo G (2018). "Analysis of linguistic complexity in professional and citizen media". In: *Journalism Studies* 19.12, pp. 1786–1803.
- Trappel, Josef and Tomaz, Tales (2021). *The Media for Democracy Monitor 2021: How leading news media survive digital transformation (Vol. 1)*. Nordicom, University of Gothenburg.
- Treder, Matthias Sebastian (2010). "Behind the looking-glass: A review on human symmetry perception". In: *Symmetry* 2.3, pp. 1510–1543.
- Trilling, Damian, Tolochko, Petro, and Burscher, Björn (2017). "From newsworthiness to shareworthiness: How to predict news sharing based on article characteristics". In: *Journalism & Mass Communication Quarterly* 94.1, pp. 38–60.
- Tromble, Rebekah (2019). "The (Mis) Informed Citizen: Indicators for Examining the Quality of Online News". In: *Available at SSRN 3374237*.

- Tsagkias, Manos, Weerkamp, Wouter, and De Rijke, Maarten (2009). "Predicting the volume of comments on online news stories". In: *Proceedings of the 18th ACM conference on Information and knowledge management*, pp. 1765–1768.
- Tuchman, Gaye (1972). "Objectivity as strategic ritual: An examination of newsmen's notions of objectivity". In: *American Journal of sociology* 77.4, pp. 660–679.
- Turner, Graeme (2013). *Understanding celebrity*. Sage.
- Urban, Juliane and Schweiger, Wolfgang (2014). "News quality from the recipients' perspective: Investigating recipients' ability to judge the normative quality of news". In: *Journalism Studies* 15.6, pp. 821–840.
- Usher, Nikki (2016). *Interactive journalism: Hackers, data, and code*. University of Illinois Press.
- Uskali, Turo I and Kuutti, Heikki (2015). "Models and streams of data journalism". In: *The journal of media innovations* 2.1, pp. 77–88.
- Valenzuela, Sebastián, Piña, Martina, and Ramírez, Josefina (2017a). "Behavioral effects of framing on social media users: How conflict, economic, human interest, and morality frames drive news sharing". In: *Journal of communication* 67.5, pp. 803–826.
- (2017b). "Behavioral effects of framing on social media users: How conflict, economic, human interest, and morality frames drive news sharing". In: *Journal of Communication* 67.5, pp. 803–826. ISSN: 00219916. DOI: 10.1111/jcom.12325.
- Van der Haak, Bregtje, Parks, Michael, and Castells, Manuel (2012). "The future of journalism: Networked journalism". In: *International journal of communication* 6, p. 16.
- Van Der Wurff, Richard and Schoenbach, Klaus (2014). "Civic and citizen demands of news media and journalists: What does the audience expect from good journalism?" In: *Journalism & mass communication quarterly* 91.3, pp. 433–451.

- Van Dijck, José (2008). "Digital photography: communication, identity, memory". In: *Visual communication* 7.1, pp. 57–76.
- Vaswani, Ashish, Shazeer, Noam, Parmar, Niki, Uszkoreit, Jakob, Jones, Llion, Gomez, Aidan N, Kaiser, Łukasz, and Polosukhin, Illia (2017). "Attention is all you need". In: *Advances in neural information processing systems*, pp. 6000–6010.
- Vis, Farida and Goriunova, Olga (2015). "The iconic image on social media: A rapid research response to the death of Aylan Kurdi". In: *Visual social media lab*.
- Vos and Thomas (2018). "The discursive construction of journalistic authority in a post-truth age". In: *Journalism Studies* 19.13, pp. 2001–2010.
- Vosoughi, Soroush, Roy, Deb, and Aral, Sinan (2018). "The spread of true and false news online". In: *Science* 359.6380, pp. 1146–1151.
- Vu, Hong Tien (2014). "The online audience as gatekeeper: The influence of reader metrics on news editorial selection". In: *Journalism* 15.8, pp. 1094–1110.
- Wahl-Jorgensen, Karin (2013). "The strategic ritual of emotionality: A case study of Pulitzer Prize-winning articles". In: *Journalism* 14.1, pp. 129–145.
- (2018). "News Media and the Emotional Public Sphere| Toward a Typology of Mediated Anger: Routine Coverage of Protest and Political Emotion". In: *International Journal of Communication* 12, p. 17.
- (2019). "Karin Wahl-Jorgensen Emotions, Media and Politics." In: *Cambio* 9.17, pp. 139–143.
- (2020). "An emotional turn in journalism studies?" In: *Digital journalism* 8.2, pp. 175–194.
- Wahl-Jorgensen, Karin and Pantti, Mervi (2021). *Introduction: The emotional turn in journalism*.

- Wallach, Hanna (2018). "Computational social science computer science+ social data". In: *Communications of the ACM* 61.3, pp. 42–44.
- Wang, Qun (2018). "Dimensional Field Theory: The adoption of audience metrics in the journalistic field and cross-field influences". In: *Digital journalism* 6.4, pp. 472–491.
- Wang, William Yang (2017). "' liar, liar pants on fire': A new benchmark dataset for fake news detection". In: *arXiv preprint arXiv:1705.00648*.
- Warriner, Amy Beth, Kuperman, Victor, and Brysbaert, Marc (2013). "Norms of valence, arousal, and dominance for 13,915 English lemmas". In: *Behavior research methods* 45.4, pp. 1191–1207.
- Waterson, J. (2018). *Financial Times tool warns if articles quote too many men*. <https://www.theguardian.com/media/2018/nov/14/financial-times-tool-warns-if-articles-quote-too-many-men>, Retrieved 9-04-2022.
- Weaver, David Hugh and Wu, Wei (1998). *The global journalist: News people around the world*. Hampton Press (NJ).
- Weber, Patrick (2014). "Discussions in the comments section: Factors influencing participation and interactivity in online newspapers' reader comments". In: *New media & society* 16.6, pp. 941–957.
- Weber, Wibke, Engebretsen, Martin, and Kennedy, Helen (2018). "Data stories: Rethinking journalistic storytelling in the context of data journalism". In: *Studies in communication sciences* 2018.1, pp. 191–206.
- Weeks, Brian E and Holbert, R Lance (2013). "Predicting dissemination of news content in social media: A focus on reception, friending, and partisanship". In: *Journalism & mass communication quarterly* 90.2, pp. 212–232.
- Welbers, Kasper, Atteveldt, Wouter van, Kleinnijenhuis, Jan, Ruigrok, Nel, and Schaper, Joep (2016). "News selection criteria in the digital age: Professional norms versus online audi-

- ence metrics”. In: *Journalism: Theory, Practice & Criticism* 17.8, pp. 1037–1053. ISSN: 1464-8849. DOI: 10.1177/1464884915595474.
- Wendelin, Manuel, Engelmann, Ines, and Neubarth, Julia (2017). “User rankings and journalistic news selection: Comparing news values and topics”. In: *Journalism studies* 18.2, pp. 135–153.
- Westerståhl, Jörgen (1983). “Objective news reporting: General premises”. In: *Communication research* 10.3, pp. 403–424.
- Whittaker, Zack (2022). *Web scraping is legal, US appeals court reaffirms*. Techcrunch. URL: <https://techcrunch.com/2022/04/18/web-scraping-legal-court/> (visited on 10/01/2022).
- Willens, M. (2019). *Forbes is building more AI tools for its reporters*. <https://digiday.com/media/forbes-built-a-robot-to-pre-write-articles-for-its-contributors>, Retrieved 9-04-2022.
- Wilson, Theresa, Wiebe, Janyce, and Hoffmann, Paul (2005). “Recognizing contextual polarity in phrase-level sentiment analysis”. In: *Proceedings of human language technology conference and conference on empirical methods in natural language processing*, pp. 347–354.
- Wojcieszak, Magdalena, Bimber, Bruce, Feldman, Lauren, and Stroud, Natalie Jomini (2016). “Partisan news and political participation: Exploring mediated relationships”. In: *Political Communication* 33.2, pp. 241–260.
- Wojdyski, Bartosz W (2015). “Interactive data graphics and information processing: The moderating role of involvement.” In: *Journal of Media Psychology: Theories, Methods, and Applications* 27.1, p. 11.
- Wollebæk, Dag, Karlsen, Rune, Steen-Johnsen, Kari, and Enjolras, Bernard (2019). “Anger, fear, and echo chambers: The emotional basis for online behavior”. In: *Social Media+ Society* 5.2, p. 2056305119829859.

- Xiao, Li and Ding, Min (2014). "Just the faces: Exploring the effects of facial features in print advertising". In: *Marketing Science* 33.3, pp. 338–352.
- Xu, Weiai Wayne, Sang, Yoonmo, and Kim, Christopher (2020). "What drives hyper-partisan news sharing: Exploring the role of source, style, and content". In: *Digital Journalism* 8.4, pp. 486–505.
- Yoo, Sung Woo, Kim, Ji Won, and Gil de Zúñiga, Homero (2017). "Cognitive benefits for senders: Antecedents and effects of political expression on social media". In: *Journalism & mass communication quarterly* 94.1, pp. 17–37.
- Zailskaitė-Jakštė, Ligita, Ostreika, Armantas, Jakštas, Adomas, Stanevičienė, Evelina, and Damaševičius, Robertas (2017). "Brand communication in social media: The use of image colours in popular posts". In: *2017 40th International Convention on Information and Communication Technology, Electronics and Microelectronics (MIPRO)*. IEEE, pp. 1373–1378.
- Zaki, Jamil and Williams, W Craig (2013). "Interpersonal emotion regulation." In: *Emotion* 13.5, p. 803.
- Zamith, Rodrigo (2019). "Transparency, interactivity, diversity, and information provenance in everyday data journalism". In: *Digital journalism* 7.4, pp. 470–489.
- Zanzotto, Fabio Massimo (2019). "Human-in-the-loop artificial intelligence". In: *Journal of Artificial Intelligence Research* 64, pp. 243–252.
- Zelizer, Barbie (1998). *Remembering to forget: Holocaust memory through the camera's eye*. University of Chicago Press.
- Zelizer, Barbie and Allan, Stuart (2011). *Journalism after september 11*. Taylor & Francis.
- Zhang, Jun, Wang, Wei, Xia, Feng, Lin, Yu-Ru, and Tong, Hanghang (2020). "Data-driven computational social science: A survey". In: *Big Data Research* 21, p. 100145.

- Zhang, Shuling and Feng, Jieyun (2019). "A Step Forward? Exploring the diffusion of data journalism as journalistic innovations in China". In: *Journalism studies* 20.9, pp. 1281–1300.
- Zhang, Shunyuan, Lee, Dokyun, Singh, Param Vir, and Srinivasan, Kannan (2017). "How much is an image worth? Airbnb property demand estimation leveraging large scale image analytics". In: *Airbnb Property Demand Estimation Leveraging Large Scale Image Analytics* (May 25, p. 2017).
- Zhang, Zhongping, Chen, Tianlang, Zhou, Zheng, Li, Jiabin, and Luo, Jiebo (2018). "How to become instagram famous: Post popularity prediction with dual-attention". In: *2018 IEEE international conference on big data (big data)*. IEEE, pp. 2383–2392.
- Zhao, Shanyang, Grasmuck, Sherri, and Martin, Jason (2008). "Identity construction on Facebook: Digital empowerment in anchored relationships". In: *Computers in human behavior* 24.5, pp. 1816–1836.
- Zhou, Wang, Wei-Wei, Li, Yan, Jin, Kai-Rui, Wang, Xuan-Yi, Wang, Zi-Wei, Chen, Yi-Shan, Wang, Shao-Jia, Hu, Jing, Zhang, Hui-Na, et al. (2019). "In-depth mining of clinical data: the construction of clinical prediction model with R". In: *Annals of translational medicine* 7.23.
- Zhou, Y. (2017). *The media's language about killers in mass shootings analyzed*. Quartz., Retrieved 01.10.2022. URL: <https://qz.com/1099083/analysis-of-141-hours-of-cable-news-reveals-how-mass-killers-are-really-portrayed/>.
- Ziegele, Marc, Breiner, Timo, and Quiring, Oliver (2014). "What creates interactivity in online news discussions? An exploratory analysis of discussion factors in user comments on news items". In: *Journal of Communication* 64.6, pp. 1111–1138.
- Zohourian, Alireza, Sajedi, Hedieh, and Yavary, Arefeh (2018). "Popularity prediction of images and videos on Instagram". In: *2018 4th International Conference on Web Research (ICWR)*. IEEE, pp. 111–117.

Appendix A

Biographical Sketch

Catherine Sotirakou is a computational journalist with over a decade of experience in the media industry in Greece. She holds a Master's degree in Digital Media and Interactive Environments from the National and Kapodistrian University of Athens (NKUA) and a Bachelor's in Journalism from the Aristotle University of Thessaloniki (AUTH). From 2012 she worked as a broadcast journalist at Alpha TV in Greece and later pursued a Ph.D. in Communication and Media Studies at NKUA. In 2017, she received a Stavros Niarchos scholarship to attend the Lede Program at Columbia University, a intensive certification program in data journalism that combines journalism and computer science courses. After completing the program, she returned to Greece and was promoted to Tech and Innovation Consultant at Alpha TV, where she led the digital transformation of the television channel. In addition, Catherine has attended several summer schools on artificial intelligence and machine learning, including the "Journalism in the era of algorithms and artificial intelligence" program at the European University Institute, the "Summer School on Methods for Computational Social Science: Methods for analyzing and modeling textual data" at the University of Southern California, and the "Lisbon Machine Learning School" at the Instituto Superior Técnico.

In 2019 October, she became an Erasmus student at the University of Vienna's Department of Communication. She worked closely on her dissertation with Professor Hajo Boomgaarden and also participated in research projects at the Computational Communication Sci-

ence Lab. During her time at the University of Vienna, she took a course on “Introduction to Machine Learning” taught by Professor David Steyrl. In late 2020, she returned to Greece and worked on various national, international, and European Commission projects with Professor Mourlas. These projects included “Fact-Checking” and “CALYPSO”, which focused on disinformation, the localization of the “Workbench” platform for data journalism, and “IQ Journalism”, which focused on artificial intelligence and journalism. In the spring of 2022, she received a Greek-British short-term scholarship and became a Research Fellow at the London School of Economics and Political Science’s Department of Media and Communications. While there, she worked with Professor Charlie Beckett on the JournalismAI project, a global initiative that aims to help news organizations use artificial intelligence responsibly. Additionally, since 2014 she has worked as a Teaching Assistant and Research Associate at the New Technologies Laboratory in Communication, Education, and the Mass Media (NKUA).

Appendix B

Code

B.1 Study: 1

In this study in order to create the two buckets for the high and low claps, it was necessary to first remove the outliers, specifically the articles with hundred thousands of claps. This was done by creating a boxplot of the target variable (claps) and adding a horizontal line at the 99.8th percentile of the values in claps. The boxplot visualized the distribution of the values, and highlighted the values that were significantly higher than the majority of the data. Then those outliers were removed from the data. The code block below presents the process.

Listing B.1: Removing the outliers from the target variable.

```
%Import the libraries
from pdpbox import pdp
from plotnine import *
import matplotlib.pyplot as plt
%matplotlib inline

%Create the boxplot
claps = pd.Series(df['claps'])
plt.boxplot(claps.values,0,'rs')
plt.hlines(y = claps.quantile(0.998), label = '99.8 percentile', xmin =
    0,xmax = 100, linestyle='dashed' )

%Remove the outliers over the 99.8th percentile.
data = pd.DataFrame()
quant_998 = claps.quantile([0.998])
print(quant_998)
data = df[ (df["claps"] < int(quant_998)) ].copy()
```

```
data.reset_index(drop=True, inplace = True)
data.shape
```

Furthermore, the code related to the best Classifier, the XGBoost is provided here. At first several variables are defined like the number of iterations to run the model for, the learning rate, and the number of rounds for early stopping. Next, the XGBClassifier object is build based on certain parameters. Lastly the model is trained and then able to make predictions on the validation set using the predict() method. More specifically, the hyperparameters included in the model initialization are:

1. `n_estimators`: the maximum number of rounds (iterations) that the model will run for.
2. `max_depth`: the maximum depth of each tree in the model.
3. `objective`: the objective function to be minimized by the model. In this case, it is set to “binary:logistic” for a binary classification task.
4. `learning_rate`: a technique to slow down the learning in the gradient boosting model, so the new trees do not overfit the training dataset.
5. `subsample`: the fraction of the training data to be used in each round.
6. `min_child_weight`: the minimum sum of instance weights needed in a child.
7. `colsample_bytree`: the fraction of columns to be randomly sampled for each tree.
8. `scale_pos_weight`: a technique to control the balance of positive and negative weights, useful for unbalanced classes.
9. `gamma`: the minimum loss reduction required to make a split.
10. `reg_alpha`: the L1 regularization term.
11. `reg_lambda`: the L2 regularization term. When those terms on weights are increased it makes the model more conservative.

Listing B.2: Python code for building the XGBoost Classifier.

```
from xgboost import XGBClassifier
```

```

from sklearn.model_selection import train_test_split

MAX_ROUNDS = 180
LEARNING_RATE = 0.1

%Set up classifier
model = XGBClassifier(
    n_estimators=MAX_ROUNDS,
    max_depth=8,
    objective="binary:logistic",
    learning_rate=LEARNING_RATE,
    subsample=.8,
    min_child_weight=6,
    colsample_bytree=.8,
    scale_pos_weight=1.6,
    gamma=10,
    reg_alpha=8,
    reg_lambda=1.3,
)

fit_model_xgb = model.fit(X_train,y_train)
val_pred_xgb = fit_model_xgb.predict(X_val)

```

For the evaluation and interpretation of the model, many different libraries were used, such as ELI5¹, TreeInterpreter² and SHAP³. Specifically, the treeinterpreter is a tool for interpreting the predictions of tree-based machine learning models. It can be used to decompose the predictions of a model into the contributions of each feature to the prediction. The code for the explanation of the Random Forest Classifier, and specifically for the features that are more significant for the high engagement claps class is provided below.

Listing B.3: Interpretation of the model with TreeInterpreter.

```

from treeinterpreter import treeinterpreter as ti

%The predict function takes the Random Forest model and the features
and returns the predictions, bias, and feature contributions as
output.

prediction, bias, contributions = ti.predict(rf, X_val[features])
prediction1, bias1, contributions1 = ti.predict(rf, X_val[features])

L = []

for pred, bias, contr in zip(prediction1, bias1, contributions1):
    d = {"pred0" : round(pred[0],2)}

```

¹<https://eli5.readthedocs.io/en/latest/>

²<https://pypi.org/project/treeinterpreter/>

³<https://shap.readthedocs.io/en/latest/>

```

    for c, feature in zip(contr, features):
        d.update({feature: round(c[0], 2)})
#     print(feature, round(c[0], 2))
    L.append(d)
ContributionsClass_df = pd.DataFrame(L)

%Select the top 10 features in terms of absolute mean contribution for
samples with a positive prediction, plot and save the figure.

plot = ContributionsClassOne_df[ContributionsClassOne_df['pred1']>0.5][
features].mean().reindex(ContributionsClassOne_df[
ContributionsClassOne_df['pred1']>0.5][features].mean().
    abs().sort_values(ascending=True).index).tail(10).plot(kind=
'barh', color= 'grey')
fig = plot.get_figure()

fig.savefig('News.eps', format='eps', dpi=1000, bbox_inches='tight')

```

The plot is shown in Figure B.1.

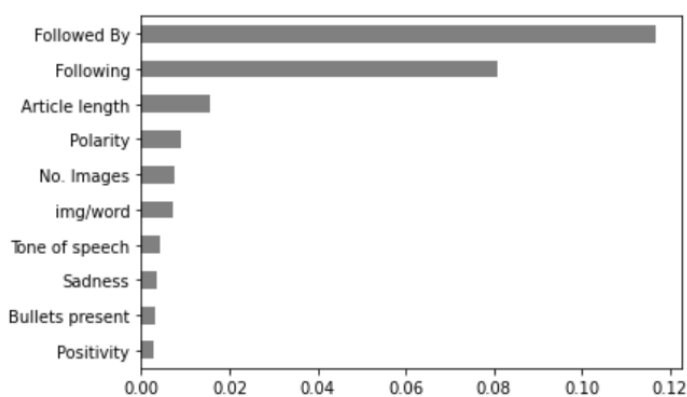


Figure B.1: Feature importance for the high-claps class.

Similar to Treeinterpreter, SHAP is also a tool to explaining the predictions of algorithms. In the following code block, the TreeExplainer class from SHAP is used to explain the Random Forest model. The “shap values” represent the contribution of each feature to the prediction made by the model. A feature that pushes the prediction higher is shown in red, while a feature that pushes the prediction lower is shown in blue.

Listing B.4: Interpretation of the model with SHAP.

```

import shap

% Load JS visualization code to notebook
shap.initjs()

explainer = shap.TreeExplainer(rf)

```

```
shap_values = explainer.shap_values(X_valid_sample.values.astype(int))
shap.summary_plot(shap_values, X_valid_sample)
```

As it can be observed in Figure B.2, a lot of SHAP explanations are stacked horizontally and reveal the explanations for an entire dataset. That way the most important features can be spotted, along with their distribution of the impacts each feature has on the model output. The color of each bar represents the value of the feature (red for high values and blue for low values). In the figure, the feature “Followed by” is the most impactful and the higher it gets it pushes the prediction towards the high claps class. In other words, the more followers an author has the higher the claps they are going to receive on *Medium.com*. The same goes for the article length. The other features are not so significant and their magnitudes on the output are not very clear. More details about the visualization of SHAP can be found on the Github Page⁴.

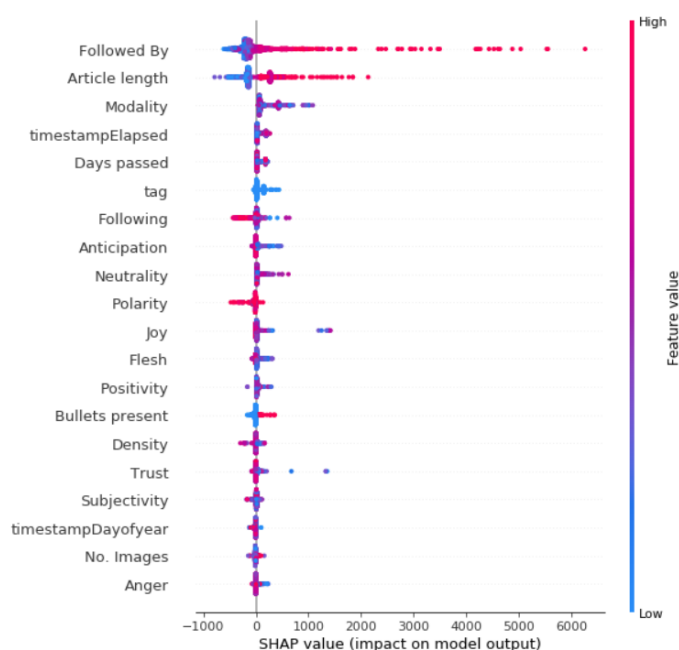


Figure B.2: The contribution of each feature to the prediction made by the model.

B.2 Study: 2

Code, tables and figures for the study 6.2.

⁴<https://github.com/slundberg/shap>

In the beginning, the correlation between the target variable and the features was calculated, and showed that the features were moderately correlated with the “quality” column (Table B.1). Also, the target variable was analyzed to check if the dataset was balanced.

Listing B.5: Correlation with target variable.

```
print(dataset.corr()["quality"].abs().sort_values(ascending=False))
```

Listing B.6: Analysis of the target variable.

```
y = dataset["quality"]
sns.countplot(y)
target_temp = dataset.quality.value_counts()
print(target_temp)
```

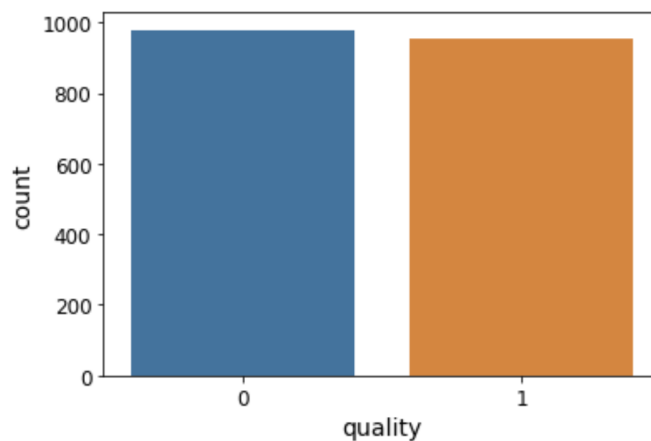


Figure B.3: The two buckets of the Quality variable.

Table B.1: Variable Correlation

Variable	Correlation
Quality	1
Difficulty	0.462214
Diversity	0.400494
Adjectives	0.371633
Trust	0.362439
Length	0.337049
Self Disclosure	0.329908
Numbers	0.269473
No Celebs	0.231839
No Crime	0.20084
Joy intensity	0.18572
#title_words	0.146766
Anticipation	0.130398
Annoyed	0.120989
Inspired	0.115509
Strong subjectivity	0.089355
Fear intensity	0.074113
Sadness intensity	0.065716
Title Polarity	0.065697
img/word	0.064588
No Sensual	0.05405
Anger intensity	0.038973
Mistakes	0.021224
No Animals	0.021158
Disgust	0.00131

The text vectorizers were used as a baseline for the machine learning models. An example of the code can be seen below, where after a basic preprocessing of the text the a TF-IDF (Term Frequency-Inverse Document Frequency) object is created and fitted to the text data, transforming them into numerical representations. Then, the TruncatedSVD is used for reducing the dimensions of the data, and various functions from sklearn for splitting the dataset into training and test sets. Afterwards the Decision Tree Classifier is created, and evaluated.

Listing B.7: Python code for the creation of the TF-IDF Vectorizer.

```
% Import necessary libraries
import pandas as pd
import numpy as np
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.decomposition import TruncatedSVD
from sklearn.tree import DecisionTreeClassifier
```

```

from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report

% Insert the dataset
df = pd.read_csv("newsstories.csv")

% Preprocess the text by lowercasing, removing punctuation, and
  stemming

def clean_text(article):
    (...more code)

df['text'] = df['text'].apply(clean_text)

% Create a TfidfVectorizer object and fit it to the dataset
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(df['text'])

% Create a TruncatedSVD object and fit it to the dataset
svd = TruncatedSVD(n_components=50)
X = svd.fit_transform(X)

% Split the dataset into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, df['quality'],
    test_size=0.2, random_state=42)

% Create a DecisionTreeClassifier object and fit it to the training
  data
clf = DecisionTreeClassifier()
clf.fit(X_train, y_train)

% Predict the labels for the test set
y_pred = clf.predict(X_test)

% Evaluate the classifier
print(classification_report(y_test, y_pred))

```

For the models that consider the dimensions of the theoretical framework before inserting the data to the algorithms, a dendrogram was created to remove the redundant features. An example of the code is provided below:

Listing B.8: Removing redundant features.

```

from scipy.cluster import hierarchy as hc
from scipy import linalg, optimize
import numpy as np
import matplotlib as mpl
import matplotlib.pyplot as plt
import scipy

%The Pearson correlation between the variables in the dataframe df_keep
  is calculated using the spearmanr function.

```

```

corr = np.round(scipy.stats.spearmanr(df_keep).correlation, 4)
corr_condensed = hc.distance.squareform(1-corr)
z = hc.linkage(corr_condensed, method='average')
fig = plt.figure(figsize=(20,15))
dendrogram = hc.dendrogram(z, labels=df_keep.columns, orientation='left',
    leaf_font_size=19)
plt.show()

```

The result of the code is shown in Figure B.4

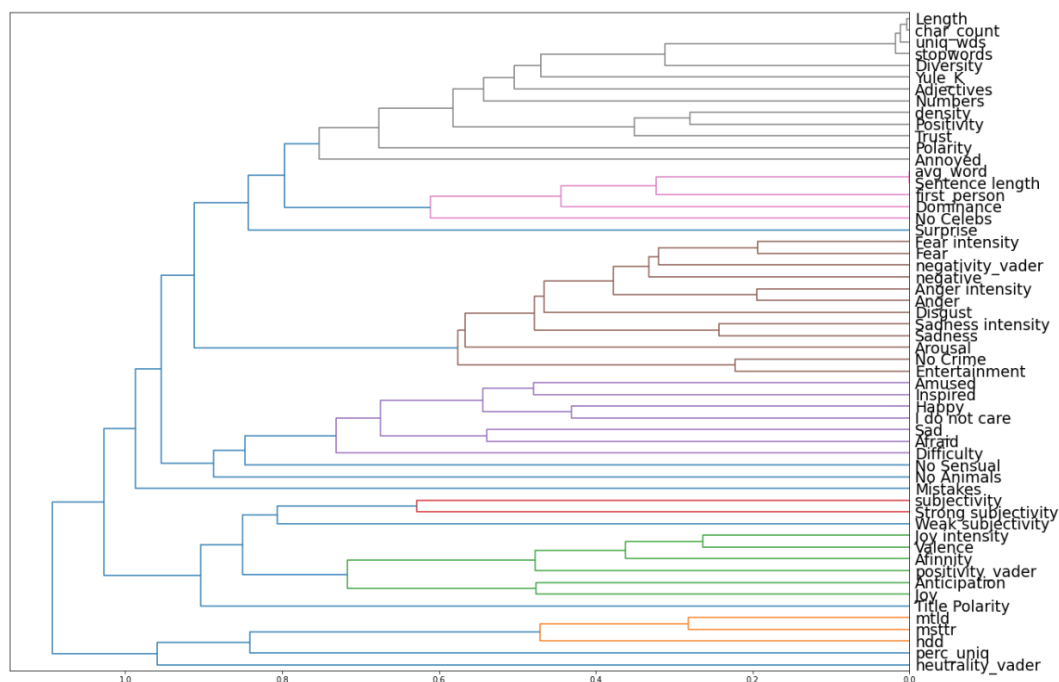


Figure B.4: Dendrogram

Then, some partial dependence plots were created to understand the relationship between a single feature and the target variable (quality). With this technique is visible how the model changes as the value of a single feature changes, while holding all other features constant. Along with feature importance it is possible to find the features that are most important for predicting the target variable and how they contribute to the overall prediction. After removing some of the related features to make the algorithm more simple but with no significant impact to the accuracy, different classification models were build. An example of the creation of the Logistic Regression model is shown below.

Listing B.9: The creation of the Logistic Regression Classifier.

```
from sklearn.model_selection import train_test_split
%The predictors (i.e., the feature columns) are extracted from the
  dataset by dropping the "quality" column, while the target variable
  is extracted from the "quality" column of the dataset dataframe.

predictors = dataset.drop("quality",axis=1)
target = dataset["quality"]

%The predictors and target are split into training and test sets using
  the train_test_split function, with a test size of 20%.

X_train,X_test,Y_train,Y_test = train_test_split(predictors,target,
  test_size=0.20,random_state=0)

%The accuracy_score and classification_report functions from sklearn.
  metrics are imported.

from sklearn.metrics import accuracy_score, classification_report
from sklearn.linear_model import LogisticRegression

%The logistic regression model is fitted to the training data using the
  fit method.

lr = LogisticRegression()

lr.fit(X_train,Y_train)

%The model is used to make predictions on the test set using the
  predict method and the predicted labels are stored in Y_pred_lr.

Y_pred_lr = lr.predict(X_test)

%The f1_score function is used to calculate the F1 score for the
  predicted labels and the true labels for the test set, with the
  average parameter set to "weighted".

f1_score(Y_test, Y_pred_lr, average='weighted')
```

Another example for the Random Forest Classifier is shown in the code block that follows.

Listing B.10: The creation of the Random Forest Classifier.

```
from sklearn.ensemble import RandomForestClassifier

%The random forest model is fit to the training data using the fit
  method.

rf = RandomForestClassifier(n_estimators=2000, min_samples_leaf=4,
  max_features=0.5, n_jobs=-1, oob_score=True)
rf.fit(X_train,Y_train)

%The f1_score function is used to calculate the F1 score for the
  predicted labels and the true labels for the test set, with the
  average parameter set to "weighted".
```

```
f1_score(Y_test, Y_pred_rf, average='weighted')

% The feature importance is calculated.

fi = pd.DataFrame(rf.feature_importances_, X_train.columns)
fi.columns = ['Importance']
fi.sort_values(by = 'Importance', ascending=False)[0:20]
```

The results from the Random Forest Classifier on the whole dataset, with an F-1 of 0.80 are presented in Figure B.5.

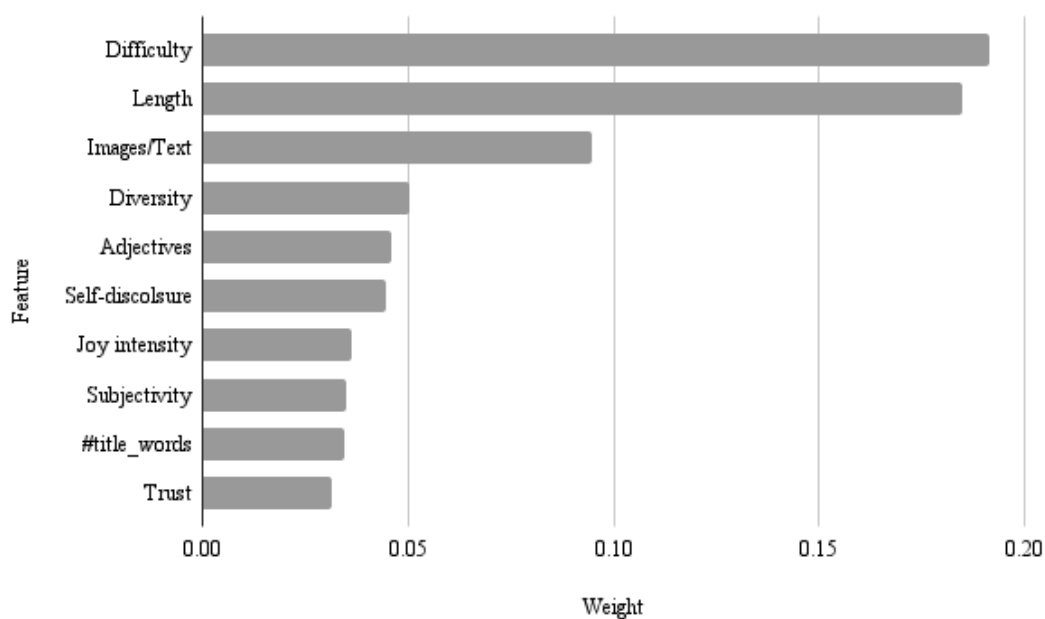


Figure B.5: Feature Importance from the Random Forest Classifier for the whole dataset

Both ELI5 and the dtreeviz libraries were used for the interpretation of the model, but here only code related to the second library is explained in detail, since ELI5 was used in other studies as well and has been presented earlier.

Listing B.11: Interpretation of the model with dtreeviz library.

```
import sys
import os
import sklearn
from sklearn.tree import DecisionTreeClassifier, DecisionTreeRegressor
import xgboost as xgb
from xgboost import plot_importance, plot_tree, plotting
from dtreeviz import trees
import graphviz
```

```

import matplotlib.pyplot as plt
from matplotlib.pylab import rcParams
import pandas as pd
import numpy as np
from dtreeviz.models.shadow_decision_tree import ShadowDecTree
from dtreeviz.models.xgb_decision_tree import ShadowXGBDTree

random_state = 42

%Import dataset.
dataset = pd.read_csv("low_high_features.csv")

%Set the features.
features = ['Difficulty', 'Mistakes', 'Length', 'Title Polarity', '
Diversity', 'Adjectives', 'Joy intensity', 'Fear intensity', '
Sadness intensity', 'Anger intensity', 'Strong subjectivity', '
first_person', 'Disgust', 'Anticipation', 'Trust', 'Surprise', '
Sadness', 'Positivity', 'No Crime', 'No Sensual', 'No Animals', 'No
Celebs', 'img/word', 'Numbers', '#title_words']
target = "quality"

%Train the model the XGBoost model

dtrain = xgb.DMatrix(dataset[features], dataset[target])

params = {"max_depth":4, "eta":0.05, "objective":"binary:logistic", "
subsample":1, "random_state":42}
xgb_model = xgb.train(params=params, dtrain=dtrain, num_boost_round=8)

%Initialize the shadow tree for the interpretation.

d = dataset[features + [target]]
d_matrix = xgb.DMatrix(d)

xgb_shadow = ShadowXGBDTree(xgb_model, 1, d[features], d[target],
features, target, class_names=[0, 1])

%Visualize the leaf samples.
trees.ctreeviz_leaf_samples(xgb_shadow)

%Visualize the decision tree.

trees.dtreeviz(xgb_model, d[features], d[target], features, target,
class_names=[0, 1], tree_index=6)

```

The visualization of the leaf samples where the biggest leaves are shown along with their distribution is shown in Figure B.6, and the whole tree in Figure 6.7 in of the Study 6.2.

B.3 Study: 3

Code, tables and figures for the study 6.3.

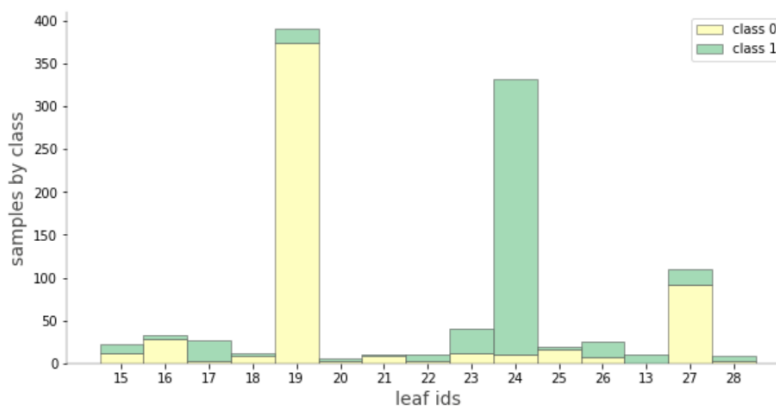


Figure B.6: Visualization of the leaf samples using the dtreeviz library.

In order to eliminate features that are correlated with each other, they were visualized using a correlation matrix and those that had a correlation of 0.70 or higher were removed.

Listing B.12: Correlation Matrix.

```
import math
import matplotlib.pyplot as plt
%matplotlib inline

correlationMatrix = data[features].corr().abs()

plt.subplots(figsize=(15, 15))
sns.heatmap(correlationMatrix, annot=True)

# Mask unimportant features
sns.heatmap(correlationMatrix, mask=correlationMatrix < 1, cbar=False)
plt.show()
```

Additionally, the Local Interpretable Model-agnostic Explanations (LIME)⁵ was used to explain a single prediction by providing an understanding of which characteristics of the news story were important to the decision made. Specifically, LIME assigns a feature importance value to each variable, which can be used to visualize the influence of each feature on the model's decision. However, it is not possible to generalize the results and draw conclusions by visualizing single samples, so this library was not used further. An example can be seen in Figure B.8.

⁵<https://github.com/marcotcr/lime>

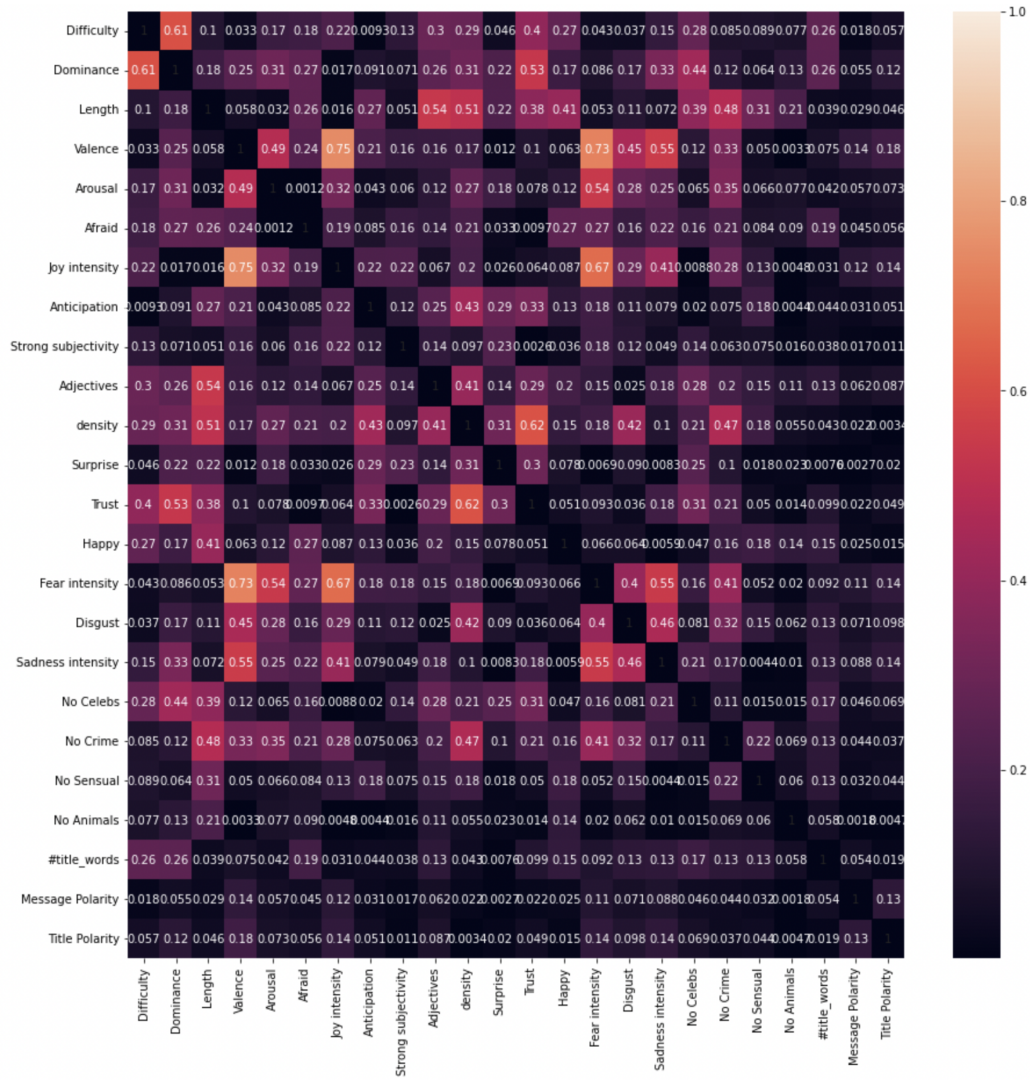


Figure B.7: Visualization of the correlation matrix.

Probability of success: 0.9659925361414831

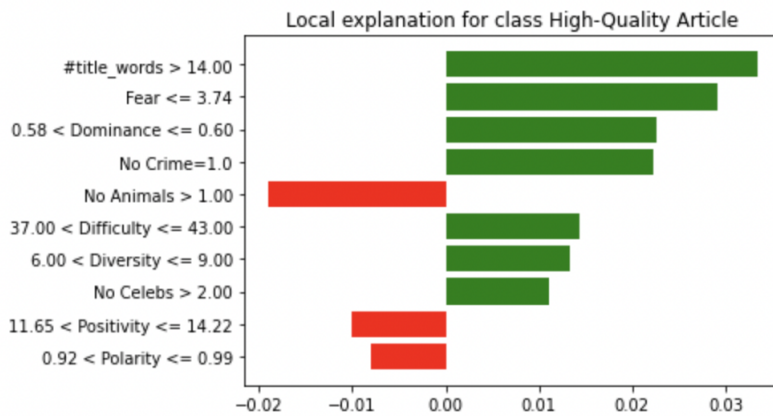


Figure B.8: Explanation of a prediction using LIME.

Listing B.13: Permutation Importance calculation with ELI5.

```

from eli5.sklearn import PermutationImportance
perm = PermutationImportance(rf).fit(X_train,y_train)
eli5.show_weights(perm, top = 10, feature_names = X_train.columns.
    tolist() )

```

Weight	Feature
0.1084 ± 0.0089	Length
0.0859 ± 0.0077	#title_words
0.0252 ± 0.0072	Difficulty
0.0228 ± 0.0011	Dominance
0.0166 ± 0.0031	No Celebs
0.0069 ± 0.0030	Anticipation
0.0055 ± 0.0019	Arousal
0.0044 ± 0.0021	Fear intensity
0.0040 ± 0.0034	Trust
0.0037 ± 0.0011	Adjectives

Figure B.9: Permutation importance using ELI5 library.

Instead, the ELI5 library provides clearer explanations and displays the weights of each feature. Additionally, it allows for the examination of the features responsible for not just one prediction, but for the entire model. Therefore, it was the preferred method for calculating permutation importance in most of the studies. An example can be seen in Figure B.9.

For the studies four and five, the same methods showed here were employed, therefore the code is not included in the appendix to avoid repetitions.

Appendix C

Lexicons for Feature Engineering

Lexicons used for feature creation in Chapter 5.

Subjectivity feature: For this feature in Study 1 the Python library TextBlob (Loria 2018) was used which returns a float within the range [0.0, 1.0] where 0.0 is very objective and 1.0 is very subjective. For the rest of the studies the Subjectivity Lexicon (Wilson et al. 2005) which includes a list of subjectivity clues to calculate weak and strong subjectivity in texts. The score is the total number of subjectivity clues divided by the number of words in the text.

Modality: Grammatical modality is signaled by grammatical moods that express a speaker's general intentions and is implemented grammatically through three moods namely indicative, imperative, and subjunctive. To capture the degree of certainty of an author, we used pattern.en (De Smedt et al. 2012), which measures modality as a value between -1.0 and 1.0, where values > 0.5 represent certainty.

Readability: For this feature in Study 1, the Flesch-Kincaid Grade Level (Kincaid et al. 1975) was used. For evaluating the readability of diverse texts, the "Flesch-Kincaid Grade Level Formula" assigns a score that corresponds to a U.S. grade level. If the calculation yields a number greater than 10, it may also refer to the average number of years of schooling needed to comprehend this book. The formula shown in Figure C.1 was used to determine the grade level:

$$0.39 \left(\frac{\text{total words}}{\text{total sentences}} \right) + 11.8 \left(\frac{\text{total syllables}}{\text{total words}} \right) - 15.59$$

Figure C.1: Formula for “Flesch-Kincaid Grade Level”

For the rest of the studies the “Flesch Reading Ease Score” (Flesch 1948) was used. Higher numbers suggest easier-to-read information, while lower numbers indicate more difficult-to-read passages. The formula is shown in Figure C.2:

$$206.835 - 1.015 \left(\frac{\text{total words}}{\text{total sentences}} \right) - 84.6 \left(\frac{\text{total syllables}}{\text{total words}} \right)$$

Figure C.2: Formula for “Flesch Reading Ease Score”

The scores can be explained further according to Wikipedia in the following Figure C.3:

Score	School level (US)	Notes
100.00–90.00	5th grade	Very easy to read. Easily understood by an average 11-year-old student.
90.0–80.0	6th grade	Easy to read. Conversational English for consumers.
80.0–70.0	7th grade	Fairly easy to read.
70.0–60.0	8th & 9th grade	Plain English. Easily understood by 13- to 15-year-old students.
60.0–50.0	10th to 12th grade	Fairly difficult to read.
50.0–30.0	College	Difficult to read.
30.0–10.0	College graduate	Very difficult to read. Best understood by university graduates.
10.0–0.0	Professional	Extremely difficult to read. Best understood by university graduates.

Figure C.3: Flesch Reading Ease Scores

Sentiment: For the sentiment multiple lexicons that detect polarity (e.g. a positive or negative opinion) were used. The Vader (Hutto and E. Gilbert 2014) package VADER (Valence Aware Dictionary and sEntiment Reasoner) is a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.

Textblob¹ sentiment analysis library was also used for capturing polarity in headlines, since it is very good for short texts. The score is calculated from -1 to 1 respectively to detect negative or positive texts.

¹<https://textblob.readthedocs.io/en/dev/>

Emotions: The NRC EMOLEX lexicon is a list of words that captures the eight basic emotions: joy, sadness, anger, fear, trust, surprise, disgust, and anticipation (Nissim and Patti 2017).

The NRC Affect Intensity Lexicon (Saif M. Mohammad 2018), focuses on English words and their associations with four basic emotions (anger, fear, sadness, joy).

The NRC VAD Lexicon was used which identifies the sentiments of valence, arousal, and dominance (S. Mohammad 2018) and it is based on Russell's theory. The dimensions are:

Valence is the positive–negative or pleasure–displeasure dimension; Arousal is the excited–calm or active–passive dimension; and Dominance is the powerful–weak or “have full control”–“have no control” dimension.

For every word the scores range from 0 (lowest V/A/D) to 1 (highest V/A/D).

Appendix D

Dataset Cleaning

List with keywords excluded from the annotated dataset used in the second study, so it would not be possible for the annotators to recognise the news organisation which published each news story.

Daily Star Online, Daily Mail Australia, Mirror Online, The Sun, Daily Mail, Sunday Mirror, DailyMail.com, Email webnews@mirror.co.uk, BBC, CNN, NBC, Daily Star Online, CBS, Reuters, (Reuters), Fox News, The New York Times, The Times, Sky News, The Mirror, The Daily Star, The Star, the Guardian, ABC, Press Association, MSNBC, Fox News, Daily-.com, The Independent, Express.co.uk, Telegraph, Enquirer, USA TODAY, Daily Telegraph, New York Daily.

Get email updates with the day's biggest stories

The Daily Star's FREE newsletter is spectacular! Sign up today for the best stories straight to your inbox

Get US and UK politics insight with our free daily email briefing straight to your inbox Missing out on the latest scoops? Sign up for – Playbook and get the latest news, every morning — in your inbox. [What you need to know to start the day: Get New York Today in your inbox.]

Today's coverage from Post correspondents around the world Like Washington Post World

on Facebook and stay updated on foreign news

You can WhatsApp us on 07810 791 502.

A daily play-by-play of congressional news in your inbox. By signing up you agree to receive email newsletters or updates from – and you agree to our privacy policy and terms of service. You can unsubscribe at any time and you can contact us here. This sign-up form is protected by reCAPTCHA and the Google Privacy Policy and Terms of Service apply. Sign up to FREE email alerts with news to brighten your day Get email updates with the day's biggest stories —> . DON'T MISS READ MORE: The Daily Star's FREE newsletter is spectacular! Sign up today for the best stories straight to your inbox Missing out on the latest scoops? Sign up for – Playbook and get the latest news, every morning — in your inbox.