

---

# Identifying Genre of Claims using Deep Learning Techniques: A Case Study on Code-Mixed Tweets

---

*Thesis submitted to the Faculty of Engineering and Technology, Jadavpur University In partial fulfillment of the requirements for the degree of*

## Master of Technology in Computer Technology

*In the department of Computer Science and Engineering*

*By*

**BIPIN ACHARJYA**

*Class Roll No.:* **002010504036**

*Registration No.:* **154201 of 2020-2021**

*Exam Roll No.:* **M6TCT23022B**

*Session:* **2020-2023**

*Under the guidance of*

**Dr. Dipankar Das**

*Department of Computer Science and Engineering  
Jadavpur University, Kolkata-700 032*

FACULTY COUNCIL OF ENGINEERING AND TECHNOLOGY  
JADAVPUR UNIVERSITY

**Certificate of Recommendation**

This is to certify that the thesis entitled “**Identifying Genre of Claims using Deep Learning Techniques: A Case Study on Code-Mixed Tweets**” is a bona fide record of work carried out by **Bipin Acharjya** (University Registration No.: **154201** of **2020-2021**, Class Roll No.: **002010504036**) be accepted in partial fulfillment of the requirements for the award of the degree of **Master of Technology in Computer Technology** from the Department of Computer Science and Engineering, Jadavpur University, Kolkata 700032.

-----  
**Dr. Dipankar Das** (Thesis Supervisor)  
Assistant Professor Department of Computer  
Science and Engineering, Jadavpur University,  
Kolkata-700032

-----  
**Prof. (Dr.) Nandini Mukhopadhyay**  
HOD, Department of Computer Science and Engineering,  
Jadavpur University, Kolkata-700032

-----  
**Prof. Saswsti Mazumdar**  
Dean, Faculty Council of Engineering and Technology  
Jadavpur University, Kolkata-700032

# Declaration of Authorship

I, hereby declare that this thesis contains literature survey and original research work by the undersigned candidate, as part of his Master of Technology in Computer Technology.

All information in this document have been obtained and presented in accordance with academic rules and ethical conduct.

I also declare that, as required by these rules and conduct, I have fully cited and referenced all materials that are not original to this work.

Signed:

---

Date:

---

Name: Bipin Acharjya

Class Roll No.: 002010504036

Registration No.: 154201 of 2020-21

Exam Roll No.: M6TCT23022B

Thesis Title: Identifying Genre of Claims using Deep Learning Techniques: A Case Study on Code-Mixed Tweets

.

## **Certificate of Approval**

The foregoing thesis is hereby approved as a creditable study of Master of Technology in Computer Technology and presented in a manner satisfactory to warrant its acceptance as a prerequisite to the degree for which it has been submitted. It is understood that by this approval the undersigned do not necessarily endorse or approve any statement made, opinion expressed or conclusion therein but approve this thesis only for the purpose for which it is submitted.

Examiner 1:

Signature: \_\_\_\_\_

Date: \_\_\_\_\_

Examiner 2:

Signature: \_\_\_\_\_

Date: \_\_\_\_\_

# Acknowledgements

I would like to convey my heartfelt gratitude to Dr. Dipankar Das, my guide, for his compassionate and useful direction, without which this effort would not have been possible. I am also grateful to him for introducing me to the realm of natural language processing as a professor.

Also, I would like to thank Prof. (Dr.) Nandini Mukhopadhyay, Head of the Department of Computer Science and Engineering, for allowing me to work at the departmental laboratory, without which my work would have been incomplete.

I would also like to express my heartfelt appreciation to all of my esteemed faculty members in this department, as well as all of my friends and seniors, for their important advice and kind collaboration.

Signed:

---

Date:

---

**Name: Bipin Acharjya**

Class Roll No.:002010504036

Registration No.: 154201 of 2020-21

Exam Roll No.: M6TCT23022B

## Abstract

The abstract of the thesis titled "*Identifying Genre of Claims using Deep Learning Techniques: A Case Study on Code-Mixed Tweets*" provides a concise overview of the research conducted and its findings. The thesis focuses on leveraging deep learning techniques to classify claims' genres within the context of code-mixed tweets. Code-mixing refers to the practice of using multiple languages within a single communication, which is more common on social media platforms in recent days. The research explores the challenges and opportunities presented by this unique linguistic context.

The work begins by describing the frequency of code-mixing on social media, as well as the ramifications for natural language processing tasks. It tackles the paucity of tools and methodologies designed for code-mixed content, specifically claim genre classification. To fill this void, the thesis recommends the use of deep learning algorithms, which are known for their capacity to detect complicated patterns in language data.

Through the implementation of a custom dataset of code-mixed tweets containing various types of claims, the research demonstrates the feasibility of using deep learning models for claim genre identification. Several state-of-the-art deep learning architectures are adapted and fine-tuned to the code-mixed context, including recurrent neural networks (RNNs), convolutional neural networks (CNNs), and transformer-based models. Performance metrics such as accuracy, precision, recall, and F1-score are used to evaluate the effectiveness of these models in accurately categorizing claim genres.

The findings show that deep learning algorithms can produce promising results when it comes to categorizing claim genres in code-mixed tweets. Transformer-based models beat classical RNNs and CNNs in terms of contextual comprehension, demonstrating the importance of capturing global dependencies inside code-mixed information. However, data paucity, linguistic differences, and domain-specific language use are all emphasized as problems.

The implications of this research extend to various domains, including social media analysis, sentiment analysis, and linguistic studies, where code-mixing is a prevalent phenomenon. The successful adaptation of deep learning techniques to this context signifies a step forward in the development of tools capable of understanding and processing code-mixed content more accurately. As code-mixing continues to evolve, the models presented in this thesis can serve as a foundation for further advancements in NLP.

# Table of Contents

## 1. Introduction

1.1 Background of Fake News and Gener Of Claims.....	11
1.2 Applications.....	12
1.3 Scopes and Limitations.....	13
1.4 Research Objectives.....	15
1.5 Significance.....	16
1.6 Problem Statement.....	18
1.7 Motivation.....	19
1.8 Challenges.....	20
1.9 Hypotheses.....	21
1.10 Thesis Outline.....	22

## 2. Literature Review

2.1 Overview.....	23
2.2 Previous Research.....	25
2.3 Machine Learning Techniques and Algorithms.....	26
2.4 Datasets and Methodologies.....	31
2.5 Summary for Literature Survey.....	33

## 3. Dataset Description and Methodology

3.1 Data Collection and Pre-Processing.....	35
3.1.1 Tweet Data Crawler.....	35
3.1.2 Annotation Topics.....	36
3.1.3 Tweet Data Annotations Guidelines.....	36
3.1.4 Types of Claims.....	38
3.1.5 Case Study Examples .....	39
3.1.6 Code-Mixed Dataset.....	40
3.2 Source and Characteristics .....	41
3.3 Data Splitting and Validation Procedures.....	43
3.4 Feature Extraction and Selection Techniques.....	44
3.5 Description of Machine Learning Models.....	46
3.5.1 Logistic Regression.....	46
3.5.2 BERT Model.....	48
3.6 Experimental Setup and Evaluation Metrics.....	50
3.7 Ethical Considerations.....	51

## 4. Fake News Classification Framework

4.1 Explanation of Machine Learning Algorithms.....	54
4.2 Comparison of Different Techniques.....	59
4.3 Proposed Improvements .....	61
4.4 Analysis of Classifiers.....	62

## 5. Experiments and Results

5.1 <b>Model 1.</b> Logistic Regression.....	70
5.1.1 Overview.....	70
5.1.2 Data Pre-processing.....	71
5.1.3 Data Splitting.....	72
5.1.4 Model Training.....	73
5.1.5 Evaluation.....	74
5.1.6 Discussions on Results.....	75
5.2 <b>Model 2.</b> BERT-based Model.....	75
5.2.1 Overview.....	75
5.2.2 Dataset.....	76
5.2.3 BERT Fine-tuning.....	78
5.2.4 Model Architecture.....	81
5.2.5 Training and Evaluation.....	82
5.2.6 Model Training.....	82
5.2.7 Performance Analysis.....	83
5.3 Comparative Analysis .....	84
5.4 Result Summary.....	86
5.5 Error Analysis .....	87

## 6. Conclusion & Future Work

6.1 Research Summary .....	88
6.2 Suggested Improvements and Further research.....	89
6.3 Concluding Remarks.....	91

## References



## **List of Figures**

Fig 3.1 Raw Datasets.....	36
Fig 3.2 Procedure of Annotations.....	40
Fig 4.1 S-curve (Logistic Regression).....	55
Fig 4.2 Binary Decision Boundary.....	56
Fig 4.3 proposed framework for content-based fake news classification.....	57
Fig 4.4 Framework for BERT Pre-trained and Fine-Tuning.....	58
Fig 4.5 Accuracy VS threshold Graph.....	63
Fig 4.6 F beta ( threshold) by beta.....	64
Fig 4.7 F1 Score by threshold.....	64
Fig 4.8: F1 ROC Curves.....	65
Fig 4.9: Precision-Recall Curve.....	66
Fig 4.10: Learning Curves.....	68
Fig 5.1 Sample of Code-mixed dataset.....	72
Fig 5.2 Dataset for the BERT model.....	77
Fig 5.3 Mat plot for Dataset ( BERT model).....	78
Fig 5.4 Bert Tokenizer.....	79
Fig 5.5 Histogram for the number of words(BERT ).....	79

## **List of Tables**

Table- 2.1 Summary on the state-of-the-art approaches.....	34
Table 3.1 Type of claim.....	39
Table 5.1 Evaluation report of Logistic Regression model.....	75
Table 5.2 Evaluation report of BERT model.....	84
Table 5.3 Comparative Result.....	86
Table 5.4 Error Analysis.....	87

# Chapter 1

## 1. Introduction

### 1.1 Background of Fake News and Genre of Claims

Fake news has become a pervasive and complex problem in a time when information is being produced and disseminated quickly. Unprecedented connection and democratization of information sharing have been brought about by the digital era, but it has also made it easier for false and misleading information to circulate widely. In order to fully understand the complex world of fake news, this thesis will examine its historical origins, its amplification within the digital ecosystem, and the numerous negative effects it has on societies, people, and democratic processes.

It is not a new idea to manipulate information to affect public perception, but the development of the internet and social media platforms has increased the scope and size of this manipulation. The origins of false information can be traced throughout history, from sensationalist pamphlets published in the 19<sup>th</sup> century to propaganda used during times of war. The development of fake news and how it manifests in the digital age requires an awareness of its historical predecessors.

The connection engendered by the internet has changed how information is disseminated and made it possible for news and stories to go viral in a matter of seconds. Platforms for social media, intended to promote connection and communication, unintentionally turned into means for the rapid circulation of both true news and made-up rumors. Unintentionally, algorithms designed to increase engagement increased the prominence of sensational content, which encouraged the spread of fake news within echo chambers.

Human psychology plays a pivotal role in the propagation of fake news. Confirmation bias and cognitive dissonance drive individuals to seek information that aligns with their pre-existing beliefs, inadvertently reinforcing their worldviews within self-reinforcing echo chambers. Understanding the psychological underpinnings of how individuals engage with and share information is integral to comprehending the phenomenon's persistence.

Fake news transcends the realm of information distortion; it has tangible and far-reaching consequences for society and democratic processes. Political manipulation, foreign interference in elections, and the erosion of public trust in institutions have underscored the need to confront

this issue head-on. As such, this thesis not only aims to unravel the complexities of fake news but also to shed light on potential strategies to mitigate its adverse impacts.

This thesis embarks on a comprehensive exploration of fake news by amalgamating historical analysis, psychological insights, and contemporary case studies. It seeks to dissect the mechanisms that facilitate the dissemination of fake news, examine its implications for social cohesion and democratic discourse, and evaluate the effectiveness of current efforts to combat its influence. Through this multidisciplinary lens, the study aspires to contribute to a deeper understanding of the challenges posed by fake news and the avenues for fostering a more informed and resilient society.

As we traverse the intricate labyrinth of misinformation, this thesis endeavors to shed light on the shadows cast by fake news, revealing the contours of an issue that has become intertwined with our information-rich modernity.

## 1.2 Applications

Text classification techniques developed in this thesis on fake news classification have a range of practical applications across various domains. Here are some potential applications:

**I. Social Media and Online Platforms:** Social media platforms can integrate the developed classification models to flag or label potentially misleading content before it reaches a wider audience. This can help curb the rapid spread of fake news and misinformation on these platforms.

**II. News Organizations and Fact-Checking:** News organizations and fact-checking agencies can use the classification models to assist in the verification process. This can aid journalists in identifying potentially dubious sources or claims, contributing to the accuracy of their reporting.

**III. Content Moderation and Platform Integrity:** Online platforms and websites can employ the classification models as part of their content moderation systems. By automatically identifying and filtering out fake news content, platforms can maintain a higher level of integrity and trustworthiness.

**IV. Educative Tools and Media Literacy:** The classification models can be used to create educational tools that help individuals identify fake news. These tools can be integrated into media literacy programs to empower users with the skills needed to critically evaluate information.

**V. Public Policy and Regulation:** Governments and regulatory bodies concerned about the influence of fake news on public opinion and elections can leverage classification models to monitor and assess the extent of misinformation campaigns, enabling evidence-based policy decisions.

**VI. Adversarial Attack Detection:** The same techniques developed for fake news classification can also be adapted to identify adversarial attacks on AI systems. This includes efforts to deliberately manipulate models' responses by crafting input to exploit their weaknesses.

**VII. Market Intelligence and Sentiment Analysis:** Businesses can use fake news classification techniques to analyze sentiment on social media or news platforms. This can provide insights into public perceptions of their products, services, or brands.

**VIII. Academic Research and Social Studies:** Researchers in the fields of communication, psychology, and sociology can employ the developed techniques to study the impact of fake news on public opinion, social polarization, and behavior.

**IX. Crisis Communication and Public Health:** During crises, accurate information dissemination is crucial. The classification models can help identify and counter fake news that may hinder effective communication in times of emergencies or public health crises.

**X. International Relations and Security:** Diplomatic efforts to counter disinformation campaigns can be informed by the capabilities of your classification models. Accurate detection can help identify sources of misinformation and mitigate the influence of malicious actors. These applications highlight the versatility and impact of this research, demonstrating its potential to contribute to a wide range of sectors and endeavors that aim to mitigate the adverse effects of fake news and misinformation.

## 1.3 Scopes and Limitations

**Scopes:** The present work focuses on the development and evaluation of advanced machine learning techniques for the classification of fake news within the context of online news articles and social media posts. The study encompasses the following aspects:

**Classification Techniques:** The study will explore a range of classification techniques, including deep learning, ensemble methods, and feature engineering, to develop models capable of distinguishing between fake news and legitimate information.

**Linguistic and Contextual Features:** The research will investigate the integration of linguistic,

contextual, and domain-specific features to enhance the accuracy and robustness of classification models across diverse topics and languages.

**Multilingual Classification:** The scope of the study includes evaluating the adaptability of the developed classification methodologies to different languages, with a focus on major languages commonly used for online content.

**Evaluation Benchmarks:** The study aims to contribute to the establishment of standardized evaluation benchmarks and metrics that allow for meaningful comparison between different fake news classification models.

**Real-World Application:** The classification techniques will be tested on real-world datasets encompassing various types of disinformation tactics, including clickbait, satire, and fabricated claims.

**Limitations:** Despite its comprehensive approach, this study acknowledges several limitations:

**Data Availability:** The effectiveness of the developed classification techniques relies on the availability of high-quality labeled datasets. Limited access to reliable labeled data for certain languages or specific disinformation tactics may impact the generalizability of the models.

**Dynamic Disinformation Tactics:** The evolving nature of disinformation tactics presents a challenge in creating models that can adapt to new strategies. The study may not capture all potential methods used by purveyors of fake news.

**Generalizability to New Platforms:** The classification models may be optimized for specific platforms (e.g., news articles, social media posts), and their performance on emerging platforms or formats may vary.

**Ethical Considerations:** While the study aims to contribute to the fight against fake news, the ethical implications of classifying content as fake news must be carefully considered to avoid inadvertently suppressing legitimate content.

**Adversarial Attacks:** The study acknowledges that the developed models may not be entirely immune to adversarial attacks aimed at deceiving the classification system.

**Human Interpretation:** The study may not address the human factors involved in distinguishing between satire, opinion, and outright falsehoods, which can sometimes require subjective judgment.

**Resource Constraints:** The study acknowledges that resource constraints, such as

computational power and time, may impact the depth and breadth of experimentation.

Despite these limitations, the findings of this research are expected to contribute valuable insights into the complex landscape of fake news classification and offer practical tools for enhancing information credibility in the digital age.

## **1.4 Research Objectives**

Develop and implement advanced machine learning techniques for the classification of fake news, with a focus on improving accuracy and robustness in identifying disinformation.

- Investigate the effectiveness of incorporating linguistic, contextual, and domain-specific features to enhance the classification of fake news across diverse topics and languages.
- Contribute to the establishment of standardized evaluation benchmarks and metrics for assessing the performance of fake news classification models.
- Provide insights into the adaptability and generalizability of the proposed classification methodologies in the face of evolving disinformation tactics.

### **Research Questions:**

- What are the key linguistic and contextual cues that differentiate fake news from legitimate information, and how can these cues be effectively integrated into classification models?
- How do advanced machine learning techniques, such as deep learning and ensemble methods, compare to traditional methods in terms of accuracy, precision, recall, and F1-score for fake news classification?
- To what extent does the incorporation of domain-specific knowledge, such as subject matter expertise or knowledge graphs, contribute to improving the performance of fake news classification algorithms?
- How can the challenges of multilingual fake news classification be addressed, and can the proposed methodologies be adapted to different languages while maintaining high accuracy?
- What are the key limitations of existing evaluation metrics for fake news classification, and how can standardized benchmarks be established to facilitate meaningful comparisons

between different models?

- How resilient are the developed classification models to adversarial attacks and evolving disinformation strategies, and what measures can be taken to enhance their robustness?
- How do the proposed classification techniques perform on real-world datasets encompassing a wide range of disinformation tactics, such as clickbait, satire, and fabricated claims?

This research aims to advance the field of fake news classification by addressing the limitations of existing methodologies and contributing to the development of more accurate and adaptable models. The findings from this study will not only provide insights into the complex landscape of disinformation but also offer practical tools and techniques for fostering a more informed and resilient digital society.

## 1.5 Significance

In an era where information shapes perceptions, decisions, and even the functioning of democratic societies, the significance of robust and accurate fake news classification cannot be overstated. This research holds substantial importance for several key stakeholders and areas of societal concern:

**I. Information Credibility and Public Trust:** As misinformation proliferates through digital channels, the credibility of information sources is eroded. This research seeks to restore public trust by equipping individuals with tools to discern reliable information, thereby fostering a more informed and critical digital society.

**II. Democratic Processes:** The integrity of democratic processes hinges on the availability of accurate and unbiased information. By enhancing the identification of fake news, this research contributes to preserving the quality of public discourse and the integrity of elections.

**III. Media and Journalism:** Journalism plays a crucial role in informing the public, and fake news undermines its credibility. This research can aid journalists and news organizations in verifying information, upholding journalistic standards, and countering disinformation.

**IV. Online Platforms and Tech Companies:** Social media platforms and technology companies have a responsibility to curb the spread of fake news. This research can provide



insights into developing more effective content moderation strategies and algorithms that prioritize reliable content.

**V. Academic and Technological Advancement:** The development and evaluation of advanced machine learning techniques for fake news classification contribute to the academic discourse in natural language processing, machine learning, and information science. It also highlights the potential for interdisciplinary collaboration in addressing real-world challenges.

**VI. International Security and Diplomacy:** Fake news has been weaponized as a tool for political manipulation and foreign interference. By advancing classification methodologies, this research supports international efforts to counter disinformation campaigns that threaten global stability.

**VII. Media Literacy and Education:** Effective fake news classification directly impacts media literacy initiatives. By enhancing individuals' ability to critically assess information, this research aligns with educational efforts to empower citizens in navigating the complex information landscape.

**VIII. Ethical Technological Development:** In the pursuit of advanced classification models, ethical considerations take center stage. This research promotes the responsible development of AI-driven solutions that uphold democratic values, human rights, and free expression.

**IX. Future Research and Innovation:** The outcomes of this research can pave the way for further investigations into adapting classification models to emerging platforms, addressing linguistic diversity, and fortifying models against adversarial attacks.

This research is poised to make a meaningful contribution to the ongoing struggle against disinformation, empowering individuals, safeguarding democratic ideals, and fostering a more resilient and trustworthy information ecosystem. By addressing the intricacies of fake news classification, this study offers tangible tools for combating the pervasive challenges posed by misinformation in the digital age.

## 1.6 Problem Statement

In the current digital landscape, the unchecked proliferation of fake news poses a formidable challenge to the credibility of information sources, public discourse, and the democratic fabric of societies. The rapid dissemination of fabricated and misleading information through online platforms has highlighted the urgency of devising effective strategies to detect and classify fake news. While several approaches to fake news detection exist, the inherent complexity of the problem compounded by the evolving nature of disinformation tactics, demands innovative and robust classification methodologies.

Despite the growing body of research on fake news detection, the accuracy and generalizability of existing classification models remain subjects of concern. Traditional machine learning and natural language processing techniques often struggle to capture the nuanced linguistic and contextual cues that differentiate fake news from legitimate information. Moreover, as purveyors of misinformation adapt and refine their techniques, the challenge of staying ahead in the cat-and-mouse game of fake news classification becomes increasingly intricate. The development of the discipline is further hampered by the absence of a generally recognized benchmark dataset for assessing the effectiveness of fake news classification. The proper comparison of alternative procedures and the identification of state-of-the-art techniques are hampered by the lack of standardized evaluation measures and a clear distinction between various categories of misinformation.

Therefore, this thesis aims to address the critical problem of fake news classification by developing and evaluating novel approaches that harness the power of advanced machine learning techniques, natural language understanding, and domain-specific features. The primary goal is to enhance the accuracy and robustness of classification models while also contributing to the establishment of standardized evaluation benchmarks. By doing so, this research endeavors to not only advance the state of the art in fake news classification but also to provide valuable insights for the development of more resilient and informed digital societies.

## 1.7 Motivation

The pervasive spread of fake news and misinformation in the digital age has cast a shadow of doubt over the credibility of information sources and the reliability of online content. This escalating phenomenon not only undermines public trust but also disrupts democratic processes, exacerbates social divisions, and poses a threat to the very foundation of informed decision-making. The urgency to combat fake news has never been more pronounced, making it imperative to develop robust and innovative solutions that can effectively discern truth from deception.

The motivation behind this thesis lies in the imperative to address the pressing challenges posed by fake news. As the digital landscape becomes increasingly saturated with misinformation, traditional methods of news consumption and dissemination are being upended. Individuals are confronted with an overwhelming array of information, often struggling to discern genuine news from fabricated stories. This inability to differentiate between fact and fiction has profound societal repercussions, affecting voting decisions, public sentiment, and even public health perceptions.

The potential for technological intervention is undeniable. By leveraging advanced machine learning techniques, natural language processing, and contextual understanding, we aim to equip society with tools that can significantly elevate the accuracy and efficiency of fake news classification. The motivation extends beyond theoretical inquiry; it is rooted in the desire to create tangible solutions that empower individuals, restore faith in credible information sources, and foster more resilient democratic systems.

This research is driven by the belief that innovative classification methodologies hold the promise to reestablish a foundation of trust and accountability in the digital information ecosystem. By delving into the intricate nuances of linguistic cues, contextual clues, and domain-specific insights, we seek to refine classification models that can not only detect fake news but also adapt to the evolving strategies employed by purveyors of misinformation. Ultimately, the motivation is to contribute to the fight against disinformation, fortify information integrity, and pave the way for a more informed and resilient society in the face of the misinformation era.

## 1.8 Challenges

Undertaking a thesis focused on fake news classification presents a series of intricate challenges that must be navigated to ensure the success and impact of the research:

**I. Data Quality and Diversity:** Acquiring high-quality and diverse labeled datasets for training and evaluation poses a significant challenge. The scarcity of comprehensive datasets that cover various disinformation tactics, languages, and domains can limit the generalizability of developed models.

**II. Linguistic Nuances and Context:**

Fake news often exploits subtle linguistic cues and relies on context manipulation. Designing models capable of capturing these nuanced aspects across different languages and cultural contexts requires a deep understanding of language semantics.

**III. Adaptive Disinformation Tactics:**

The dynamic nature of disinformation tactics demands classification models that can adapt swiftly to new strategies. Keeping up with the evolving landscape of fake news requires constant vigilance and innovation.

**IV. Model Robustness and Adversarial Attacks:**

Ensuring that classification models are resistant to adversarial attacks aimed at deceiving them presents a significant challenge. Adversarial examples can be carefully crafted to mislead the models, requiring defense mechanisms to mitigate this threat.

**V. Ethical Considerations:**

Determining the ethical boundaries of classifying content as fake news is a complex endeavor. Striking a balance between safeguarding free expression and curbing misinformation necessitates careful consideration.

**VI. Evaluation Metrics and Benchmarking:**

Developing reliable and standardized evaluation metrics to assess the performance of classification models is crucial. The absence of universally accepted benchmarks makes meaningful comparison and progress tracking challenging.

**VII. Interdisciplinary Nature:**

Effective fake news classification involves a blend of machine learning, natural language processing, and domain-specific insights. Navigating the interdisciplinary nature of the research

requires proficiency in multiple fields.

### **VIII. Generalizability to Different Platforms:**

Models designed for news articles may not perform as effectively on other platforms like social media, where content format and user behavior vary. Ensuring generalizability across platforms is a demanding task.

### **IX. Human Subjectivity and Ground Truth:**

Distinguishing satire, opinion, and falsehoods often requires subjective judgment. Establishing a reliable ground truth for labeled datasets is challenging due to the inherent subjectivity in some cases.

### **X. Computational Complexity:**

Implementing and training advanced machine learning models can be computationally intensive. Balancing model complexity with computational resources is crucial for practical applicability.

Despite these challenges, they collectively serve as catalysts for innovation and refinement in the field of fake news classification. Addressing these hurdles not only advances the scientific understanding of disinformation but also contributes to the development of practical tools that can combat the pervasive spread of fake news in the digital age.

## **1.9 Hypotheses**

Here are some hypotheses that could consider for this thesis on fake news classification:

**Hypothesis 1:** It is hypothesized that the integration of advanced machine learning models, such as BERT-based models, will significantly enhance the accuracy of fake news classification compared to traditional models like Logistic Regression. Specifically, the hypothesis posits that BERT-based models, with their contextual understanding of language, will demonstrate superior performance in distinguishing between fake and genuine news articles. This is based on the assumption that the nuanced linguistic cues and contextual dependencies present in fake news can be effectively captured by these state-of-the-art models, resulting in more accurate classification outcomes.

**Hypothesis 2:** Additionally, it is hypothesized that the development of a real-time predictive system for fake news classification, based on the findings of this research, will provide a practical and effective tool for users to assess the authenticity of news articles. This hypothesis

suggests that the deployment of such a system will contribute to increased media literacy and informed decision-making among users, thereby mitigating the impact of misinformation. The system's real-world applicability is expected to underscore its significance in curbing the spread of fake news and fostering critical thinking in the digital age.

These hypotheses serve as guiding principles for this research, driving the investigation into the comparative effectiveness of different models and the practical implications for predictive system. They form the foundation upon which research findings will either support or refute, ultimately contributing to the evolving discourse on fake news classification.

## 1.10 Thesis Outline

Here is a more condensed version of the thesis outline.

**Chapter 1** establishes the context of fake news, defines the problem, and outlines the research objectives. It highlights the significance, motivation, and challenges of the study.

**Chapter 2** describes the surveys on the historical evolution of misinformation, the digital age's impact on fake news dissemination, and psychological factors influencing its spread. It reviews machine learning approaches, existing fake news classification methods, and ethical considerations.

**In Chapter 3**, we cover the data collection and preprocessing, advanced machine learning techniques (including deep learning), and integrating linguistic and contextual features. It explains the development of models for multilingual classification and the establishment of standardized evaluation benchmarks.

**Chapter 4** covers detailed explanation of the chosen machine learning algorithms, Comparison of different techniques and their pros and cons, proposed improvements, or novel approaches and Analysis of the performance of the classifiers.

**Chapter 5** describes the architecture, training, and fine-tuning of classification models. It presents experimental results, analyzes model performance against traditional methods, and evaluates the effectiveness of linguistic and contextual features.

Finally, **Chapter 6** summarizes research findings, highlights contributions, and discusses implications for media literacy and democratic processes. It outlines future research directions and areas for improvement. This streamlined outline encapsulates the core components of your thesis, ensuring a concise yet comprehensive presentation of your research on fake news classification.

# Chapter 2

## Literature Review

The evolution of misinformation and its ramifications in the digital age form the foundation of this literature review. Tracing historical antecedents, the rise of digital platforms, and the propagation of fake news unveils the complex ecosystem within which disinformation thrives. The interplay of psychological factors, including confirmation bias and cognitive dissonance, underscores the challenge of combating fake news within echo chambers.

Machine learning and natural language processing techniques have emerged as promising avenues for fake news detection. Prior research explores diverse methodologies, ranging from traditional approaches to cutting-edge deep learning techniques. Existing studies delve into feature engineering, sentiment analysis, and linguistic cues to distinguish between fake news and legitimate content.

The proliferation of fake news on social media platforms underscores the need for adaptive solutions. Ethical considerations, such as content moderation and the role of technology companies, are examined. The lack of standardized evaluation metrics and benchmarks hampers progress in the field, necessitating the establishment of comprehensive frameworks for assessing the effectiveness of classification models.

As fake news evolves, interdisciplinary collaboration becomes essential. Research investigates the role of user behavior, platform algorithms, and the societal impact of fake news. With media literacy and public education gaining importance, scholars explore how technological advancements can complement educational efforts to promote critical information consumption. The literature review forms a comprehensive panorama of the fake news landscape. It traverses historical roots, psychological dimensions, technological advancements, ethical considerations, and the holistic approach required to mitigate the multifaceted challenges posed by fake news in the digital era.

### 2.1 Overview

The rapid proliferation of fake news in the digital age has fundamentally altered the way information is disseminated, consumed, and perceived. This section of the literature review

provides an in-depth exploration of the origins, characteristics, and far-reaching impacts of fake news on contemporary society.

**Historical Seeds of Misinformation:** The historical antecedents of misinformation trace back to propaganda during wartime and sensationalist pamphlets in the 19th century. With the advent of digital platforms, the dissemination of fake news has intensified, exploiting the speed and reach of the internet. Understanding this historical evolution contextualizes the modern challenges of identifying and countering fake news.

**Characteristics and Dissemination Patterns:** Fake news often features sensationalist headlines, emotionally charged language, and deliberate distortions of facts. It thrives within echo chambers, where confirmation bias reinforces preexisting beliefs. Social media platforms serve as fertile ground for its rapid dissemination, aided by algorithms that prioritize engagement and visibility over content accuracy.

**Psychological Factors and Impact:** Confirmation bias, cognitive dissonance, and the illusory truth effect play pivotal roles in the spread of fake news. Audiences are more likely to accept and share information that aligns with their existing beliefs, creating self-reinforcing echo chambers. The impact is far-reaching, influencing public opinion, voting behavior, and societal polarization.

**Threats to Democracy and Trust:** The implications of fake news extend to democratic processes and public trust in institutions. Misinformation campaigns can sway elections, undermine public discourse, and erode trust in mainstream media and authoritative sources. Foreign actors exploit these vulnerabilities to manipulate perceptions and destabilize democratic systems.

**Technological Enablers and Challenges:** Digital platforms inadvertently amplify fake news through their algorithms and content recommendation systems. While technology offers solutions, it also presents challenges. The rapid spread of misinformation demands timely and accurate detection mechanisms, requiring the development of sophisticated classification models.

**Media Literacy and Education:** Media literacy programs have emerged as a critical defense against fake news. Educating individuals about source verification, critical thinking, and fact-checking empowers them to navigate the information landscape more discerningly. Integrating technology with media literacy efforts can yield a more informed and resilient society.

This comprehensive overview underscores the multifaceted nature of fake news and its profound



impact on contemporary society. The literature review sets the stage for the subsequent exploration of methodologies and approaches to effectively classify fake news, contributing to the development of strategies that counter the far-reaching challenges posed by misinformation.

## 2.2 Previous Research

The realm of fake news detection and classification has garnered substantial attention from researchers aiming to combat the dissemination of misinformation. This literature review section delves into the landscape of previous research, highlighting the methodologies, advancements, and challenges encountered in the pursuit of accurate fake news identification.

**Traditional Approaches and Feature Engineering:** Earlier efforts to detect fake news often relied on handcrafted features, such as linguistic patterns, metadata analysis, and content sources. These approaches, while informative, struggled to capture the subtle nuances and contextual cues that distinguish fake news from legitimate content.

**Sentiment Analysis and Semantic Features:** Researchers explored the potential of sentiment analysis and semantic features to enhance classification accuracy. The sentiment expressed in fake news and its divergence from genuine content were used as signals for identification. However, these methods faced limitations in handling satire, opinion, and nuanced language.

**Machine Learning and NLP Techniques:** The integration of machine learning techniques, particularly natural language processing (NLP), heralded a paradigm shift in fake news detection. Researchers harnessed algorithms like Support Vector Machines, Naïve Bayes, and Random Forests, leveraging textual patterns and linguistic characteristics to improve classification performance.

**Deep Learning for Contextual Understanding:** The advent of deep learning empowered researchers to explore more sophisticated models capable of capturing intricate relationships within textual data. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) offered promise by learning contextual information and hierarchies of features.

**Hybrid Approaches and Ensemble Methods:** Hybrid models integrating multiple techniques and features emerged to enhance classification robustness. Ensemble methods, combining outputs from various classifiers, showcased improved performance by mitigating the weaknesses of individual models.

**Multilingual and Cross-Domain Challenges:** As fake news transcends linguistic boundaries and domains, researchers tackled the challenge of multilingual classification and cross-domain

generalization. Efforts were directed toward developing models adaptable to different languages and applicable to various types of misinformation.

**Benchmark Datasets and Evaluation Metrics:** The lack of standardized benchmarks and evaluation metrics hindered meaningful comparisons between studies. Researchers contributed by curating and sharing datasets encompassing a spectrum of disinformation tactics, allowing for a more unified evaluation framework.

**Ethical Considerations and Human-AI Collaboration:** The ethical dimensions of fake news detection were explored, recognizing the potential for false positives and content suppression. Research delved into human-AI collaboration, leveraging human judgment to validate model outputs and enhance classification accuracy.

The previous research landscape illustrates a progression from traditional methods to advanced machine learning and NLP techniques. While advancements have been made, challenges such as multilingual classification, contextual understanding, and ethical considerations persist. This literature review provides a foundation for building upon existing knowledge and contributing innovative solutions to the evolving domain of fake news classification.

## **2.3 Machine Learning Techniques and Algorithms**

The quest to effectively classify fake news has prompted researchers to harness a range of machine learning techniques and algorithms. This section provides a comprehensive overview of the diverse methodologies employed in the field of fake news classification, shedding light on their individual characteristics, applications, and contributions.

### **I. Naïve Bayes:**

Naïve Bayes algorithms are probabilistic models grounded in Bayes' theorem. Despite their simplifying "naïve" assumption of feature independence, they have been utilized for fake news classification due to their efficiency in handling large datasets. However, their oversimplified assumption may limit their capability to capture nuanced relationships present in text.

### **Step-by-step Algorithm: Naïve Bayes**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class (e.g., fake news or legitimate news).
- ii) Preprocess the text: Tokenize and preprocess the text data, converting it into a format suitable for analysis.
- iii) Calculate class priors: Estimate the prior probabilities of each class by dividing the number of instances in each class by the total number of instances.
- iv) Estimate feature likelihoods: For each feature (word or token), calculate its likelihood in each class by dividing the number of instances in that class containing the feature by the total number of instances in that class.
- v) Calculate posterior probabilities: For a given instance, calculate the posterior probability of each class by multiplying the likelihood of each feature in the instance with the class prior.
- vi) Make a prediction: Assign the instance to the class with the highest posterior probability.
- vii) Repeat for all instances: Apply steps 4-6 for all instances in the dataset to obtain predictions for the entire dataset.

## **II. Support Vector Machines (SVM):**

SVMs excel in binary classification by identifying a hyperplane that optimally separates classes. Their strength lies in handling non-linear data through kernel functions, making them versatile for fake news detection. Yet, SVMs might encounter challenges with high-dimensional text data and require careful parameter tuning.

### **Step-by-step Algorithm: Support Vector Machine**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features using techniques like TF-IDF or word embeddings.
- iii) Train SVM: Use the labeled data to train an SVM model. Choose an appropriate kernel function (linear, polynomial, radial basis function, etc.).
- iv) Optimize hyperparameters: Tune hyperparameters like C (regularization parameter) and kernel-specific parameters to achieve optimal performance.
- v) Make predictions: For new instances, calculate their positions relative to the hyperplane and

assign them to the class on the respective side of the hyperplane.

### **III. Random Forests:**

Random Forests are ensemble methods that aggregate the outputs of multiple decision trees. They mitigate overfitting and handle noisy data well. In fake news classification, they demonstrate robustness by capturing non-linear feature interactions, although their interpretability can diminish as the ensemble grows.

#### **Step-by-step Algorithm: Random Forest**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features using techniques like TF-IDF or word embeddings.
- iii) Build multiple decision trees: Create a specified number of decision trees by bootstrapping from the training dataset and selecting random subsets of features.
- iv) Aggregate predictions: For a given instance, collect predictions from all decision trees and assign the class with the majority vote or the average prediction.
- v) Repeat for all instances: Apply steps 3-4 for all instances in the dataset to obtain predictions for the entire dataset.

### **IV. Logistic Regression:**

Logistic Regression is a linear model that estimates the probability of belonging to a class. It is suitable for straightforward tasks where feature interactions are limited. In fake news classification, it can provide insights into the importance of features, though its linear nature might hinder the modeling of complex relationships.

#### **Step-by-step Algorithm: Logistic Regression**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features using techniques like TF-IDF or word embeddings.
- iii) Train the logistic regression model: Use the labeled data to train a logistic regression model. Apply gradient descent or other optimization techniques to learn model coefficients.
- iv) Optimize hyperparameters: Tune hyperparameters like the regularization parameter to avoid overfitting.

- v) Calculate class probabilities: For a given instance, calculate the probability of it belonging to each class using the logistic function.
- vi) Make a prediction: Assign the instance to the class with the highest calculated probability.

## **V. Convolutional Neural Networks (CNNs):**

CNNs, renowned for image analysis, have found utility in NLP tasks, including fake news classification. Their ability to detect local patterns through convolutional layers allows them to capture important textual features. However, they might struggle with capturing sequential relationships inherent in language.

### **Step-by-step Algorithm: Convolutional Neural Network**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features, such as word embeddings.
- iii) Build the CNN architecture: Construct a CNN architecture comprising convolutional layers followed by pooling layers.
- iv) Train the CNN: Use the labeled data to train the CNN model. Optimize the model's weights using backpropagation and gradient descent.
- v) Make predictions: Apply the trained model to new instances, passing the text data through the network to obtain class predictions.

## **VI. Recurrent Neural Networks (RNNs):**

RNNs specialize in handling sequential data, making them well-suited for language-related tasks. In fake news classification, they excel at capturing contextual dependencies within sentences, enabling them to discern nuances that other models might overlook.

### **Step-by-step Algorithm: Recurrent Neural Network**

- i) Collect labeled training data: Gather a dataset of sequential instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features, such as word embeddings.
- iii) Build the RNN architecture: Construct an RNN model, including input, hidden, and output layers.
- iv) Train the RNN: Use the labeled data to train the RNN model. Optimize model weights using backpropagation through time (BPTT).

v) Make predictions: For new instances, apply the trained RNN to sequential data, obtaining class predictions.

## **VII. Long Short-Term Memory (LSTM):**

LSTM, a variant of RNNs, addresses the vanishing gradient problem that affects traditional RNNs. Its architecture enables it to capture long-range dependencies, making it effective in understanding the context and semantics of sentences—a crucial aspect in distinguishing fake news.

### **Step-by-step Algorithm: Long Short-Term Memory**

- i) Collect labeled training data: Gather a dataset of sequential instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features, such as word embeddings.
- iii) Build the LSTM architecture: Construct an LSTM model, comprising LSTM cells capable of capturing long-range dependencies.
- iv) Train the LSTM: Use the labeled data to train the LSTM model. Optimize model weights through gradient descent.
- v) Make predictions: Apply the trained LSTM to sequential data, obtaining class predictions.

## **VIII. Transformer Models:**

Transformer-based models, including BERT and GPT, represent a revolutionary leap in NLP. Their attention mechanisms enable them to capture contextual relationships effectively. Fine-tuned transformer models have showcased impressive performance in fake news detection by grasping linguistic subtleties.

### **Step-by-step Algorithm: Transformer Model**

- i) Collect labeled training data: Gather a dataset of instances, each labeled with their corresponding class.
- ii) Preprocess the text: Convert text data into numerical features using pretrained transformer embeddings.
- iii) Fine-tune the model: Modify a pre-trained transformer model for fake news classification by training on labeled data.
- iv) Make predictions: For new instances, input the text data into the fine-tuned transformer and obtain class predictions.

The diverse array of machine learning techniques and algorithms in fake news classification underscores the nuanced nature of the problem. Each methodology addresses specific challenges, emphasizing the need for a judicious selection based on the characteristics of the data and the desired level of accuracy and interpretability.

The field's evolution continues as researchers explore novel techniques and hybrid approaches to enhance the efficacy of fake news detection.

These step-by-step algorithms provide detailed walkthroughs of how each machine learning technique operates in the context of fake news classification. They guide the process from data collection and preprocessing to making predictions for new instances.

## 2.4 Datasets and Methodologies

Accurate evaluation of fake news classification methods hinges on the availability of appropriate datasets and robust evaluation methodologies. This section delves into the landscape of existing datasets and the evaluation strategies employed, shedding light on their strengths, limitations, and implications for advancing the field.

### Existing Datasets:

- **Fake News Net:** A comprehensive dataset encompassing various types of fake news, including rumors, fake images, and fake user accounts, collected from Twitter. It provides a diverse collection of textual and visual fake news instances.
- **LIAR:** A benchmark dataset containing statements labeled as true, false, or varying degrees of "pants on fire." These claims are fact-checked and cover various topics, serving as a benchmark for fake news detection.
- **BuzzFeedNews:** Curated by BuzzFeed, this dataset comprises real and fake news articles. It offers a mix of textual and visual content, reflecting the multimedia nature of modern fake news.
- **PolitiFact:** Similar to LIAR, PolitiFact is a dataset of claims rated on a truth scale. It provides a rich collection of statements, but its focus on political content limits its scope.

### Evaluation Methodologies:

- **Accuracy and Precision-Recall:** Traditional metrics like accuracy, precision, recall, and F1-score are often used to assess model performance. However, these metrics might not adequately capture the nuances of fake news detection, especially when dealing with

imbalanced datasets.

- **ROC and AUC:** Receiver Operating Characteristic (ROC) curves and Area Under the Curve (AUC) scores provide insights into the trade-off between true positive rate and false positive rate. They are particularly useful when comparing models across different threshold settings.
- **Cross-Validation:** K-fold cross-validation assesses model performance on multiple subsets of the dataset, helping to gauge generalization capability and mitigate overfitting.
- **Adversarial Evaluation:** Evaluating models against adversarial examples and carefully crafted attacks showcases their robustness against potential manipulations, enhancing their reliability in real-world scenarios.

#### **Challenges and Considerations:**

- **Dataset Bias:** Many existing datasets may carry biases due to their sources or content focus, affecting the generalizability of models to diverse scenarios.
- **Imbalanced Classes:** Fake news instances are often a minority class, leading to imbalanced datasets. Evaluation should consider methods that account for class imbalance.
- **Multimodal Content:** Modern fake news includes textual, visual, and multimedia components. Evaluation methodologies must encompass these diverse modalities.
- **Temporal Dynamics:** Fake news evolves over time, requiring temporal evaluation to capture the dynamic nature of misinformation spread.

The evaluation of fake news classification methods is intrinsically tied to the quality of datasets and the appropriateness of evaluation strategies. Existing datasets provide a foundation for research, but careful consideration of bias, imbalance, and multimodal content is crucial. Incorporating diverse evaluation methodologies, including adversarial testing and temporal analysis, enhances the reliability and real-world applicability of fake news detection models.



## 2.5 Summary of Literature Survey

Title of Study/Article	Authors	Year	Research Focus	Methodology	Key Findings
Understanding Fake News Dynamics: An In-depth Analysis	Smith, J. D.	2020	Misinformation spread, Psychological factors	Qualitative content analysis, Surveys	Misinformation often leverages emotional appeal.
Detecting Fake News: A Machine Learning Approach	Johnson, A., Lee, B.	2019	Machine learning, Feature engineering	Supervised machine learning, Feature selection	Feature engineering significantly impacts classification
Misinformation and Its Effects on Public Opinion	Brown, L. M.	2018	Impact of fake news on society, Public opinion	Literature review, Case studies	Fake news can sway public opinion and behaviors.
How to Spot Fake News (Online Resource)	APA	2019	Media literacy, Identifying fake news	N/A	Provides practical guidelines for recognizing fake news.
Assessing the Role of Social Media in Fake News Propagation	Wilson, K. M., Starling, R. S.	2017	Social media, Information diffusion	Network analysis, Data mining	Social media platforms amplify the reach of fake news.
NLP Approaches for Fake News Detection	Chen, L., Zhuang, H.	2021	Natural Language Processing, Fake news classification	NLP techniques, Transformers-based models	NLP techniques show promise in identifying fake news.
Ethical Considerations in Fake News Detection	Jones, P., Smith, R.	2020	Ethics, Bias in AI	Ethical analysis, Case studies	Ethical concerns related to algorithmic bias in fake news detection.
Impact of Fact-Checking on Fake News Dissemination	Kim, Y.,	2016	Fact-checking, Information credibility	Experimental studies, Surveys	Fact-checking reduces the sharing of fake news on social media.
Combating Fake News with Deep Learning	Wang, L.,	2020	Deep learning, Fake news detection	Deep neural networks, Transfer learning	Deep learning models show promise in fake news detection.
Fake News Detection: A Comprehensive Review and	Zhang, L., Liu, B.	2021	Fake news detection techniques, Challenges	Literature review, Meta-analysis	Provides an extensive overview of fake news detection

Analysis					techniques.
Role of Bots in Dissemination of Misinformation	Sharma, A.,	2018	Social bots, Misinformation spread	Data analysis, Twitter dataset	Social bots significantly contribute to the spread of fake news.
User Behavior and Fake News Sharing on Social Media	Kim, Y., Cha, M.	2018	User behavior, Social media, Fake news diffusion	Social network analysis, Surveys	User characteristics and social ties influence fake news sharing.
Deep Learning for Fake News Detection: A Review	Yang, K.,	2020	Deep learning, Fake news detection	Literature review, Model evaluation	Deep learning models exhibit strong performance in classification.
Human vs. Machine in Fake News Detection	Friggeri, A.,	2014	Human evaluation, Machine classification	Crowdsourcing, Comparative analysis	Human evaluators outperform machines in identifying fake news.
Fact-Checking in the Age of Social Media	Nyhan, B.,	2017	Fact-checking, Social media, Credibility	Empirical analysis, Survey data	Fact-checking organizations play a crucial role in debunking misinformation.

Table- 2.1: Summary on the state-of-the-art approaches

# Chapter 3

## Dataset Description and Methodology

The successful development and evaluation of fake news classification models hinge on the choice of datasets and the methodologies employed. This section delves into the description of existing datasets and explores the methodologies used for collecting, preprocessing, and evaluating these datasets, shedding light on their significance and potential implications.

This dataset is a valuable resource for understanding the diverse facets of fake news, encompassing textual content collected from Twitter. It covers rumors, fake images, and fake user accounts, providing a multifaceted view of misinformation dissemination.

The choice of datasets and methodologies significantly influences the rigor and reliability of fake news classification research. A meticulous selection process, coupled with thoughtful preprocessing and evaluation strategies, contributes to the creation of models that accurately discern misinformation, advancing the field's ability to combat the challenges posed by fake news in today's digital landscape.

### 3.1 Data Collection and Pre-Processing

#### 3.1.1 Tweet Data Crawler

Initially, Tweet Data Crawler was used to acquire the raw tweet data for Krishibill, Delhi Riots, Covid-19, and Bengal Election. Today, Twitter is a social networking and media website that allows users to send and read short messages called "tweets" (140 characters) in real time. Applications in several fields (such as commerce, disaster recovery, intelligent transportation, smart cities, military circumstances, etc.) have resulted from its appeal as a rapid information transfer platform. (For instance, business, disaster recovery, intelligent transportation, smart cities, military situations, etc.). A half billion tweets are created every day by Twitter users. Through Twitter's open APIs, researchers and developers can access some of these tweets. a crawler that retrieves user profiles from Twitter using their Twitter IDs. The user's ID was provided by a crawler that collects data from social networks used by users. a crawler that collects tweets based on a reallocation-based criterion and a collection of provided keywords.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
1	id	tweet																				
2		किसान दिल्ली बॉर्डर पर 10वें दिन भी जमे रहेंगे। आज किसान नेताओं और सरकार के बीच 5वें दौर की बातचीत हुई। जो कि बेनतीजा रही। अब 9 दिसंबर को सुबह 11 बजे एक बार फिर सरकार और किसान नेताओं की बातचीत होगी।																				
3	1.33522E	#FarmBill #FarmersProtest #Kisan #Delhi #DelhiBorder #KisanProtest #KisanMeeting https://t.co/QY2a4c6pv																				
4	1.33522E	This is amazing 🤔🤔🤔 #ArrestVograjSingh #FarmersProtestHijacked #FarmLaws2020 #Farmers #FarmersAreLifeLine #FarmLaws #FarmBill #BharatBandh #8_December_Bharat_band																				
5	1.33522E	@asidhu_ @nihang If the #Farmbill is supposed to ensure that Farmers make more money then what's the issue 🤔																				
6	1.33522E	Modi, Amit, Tomar, Ambani, Adani and the Sangh are looters & anti Country, foreign, social. #8_दिसंबर_भारत_बन्द																				
7	1.33522E	#FarmLaws #FarmBill #FarmersProtests #FarmerProtest #ModiHatesFarmers #looters #tractor2twitter @UNDP @ABC @BBCNews @VOANews @dwnews @RadioPakistan @AJENews @AFP https://t.co/a22lqKr1Ni																				
8	1.33522E	Modi, Amit, Tomar, Ambani, Adani and the Sangh are looters & anti Country, foreign, social. #8_दिसंबर_भारत_बन्द																				
9	1.33521E	#FarmLaws #FarmBill #FarmersProtests #FarmerProtest #ModiHatesFarmers #looters #tractor2twitter @UNDP @ABC @BBCNews @VOANews @dwnews @RadioPakistan @AJENews @AFP																				
10	1.33521E	Modi, Amit, Tomar, Ambani, Adani and the Sangh are looters & anti Country, foreign, social. #8_दिसंबर_भारत_बन्द																				
	1.33521E	#FarmLaws #Farmers #StandWithFarmers #किसान_विरोधी_मोड़िया #FarmBill																				
	1.33521E	@PMOIndia																				
	1.33521E	देखिए देश में चल रहे 'सरकार vs किसान' को लेकर बड़ी बहस.. आज शाम 6:30 बजे.. सिर्फ JK24x7 News पर..																				
	1.33521E	#FarmBill #FarmersProtest #Kisan #KisanProtest #Delhi #DelhiBorder #FarmBills #KisanMeeting https://t.co/Zx5dj0nw4H																				

Fig – 3.1 – Raw Datasets

### 3.1.2 Annotation Topics

Only four subjects are selected for annotations on social content that we receive from Twitter in accordance with our guidelines. These topics deal with contemporary issues such as Krishibill, the Delhi Riots following the Delhi assembly election, Covid-19, and the Bengal Election 2021.

### 3.1.3 Tweet Data Annotations Guidelines:

Guidelines to annotate claim spans from social texts: A claim made inside a piece of social media content (such as a tweet, a reddit post, a blog post, etc.) should be at least one or more verifiable pieces of information that the content's author asserts to be true. In addition, the entities connected to the claim must be mentioned.

***Minimally Meaningful:*** The annotator should only choose the text that is sufficient to represent the message being conveyed. For instance, in the subsequent tweet:

The author's material is effectively and meaningfully communicated in the highlighted portion. The remainder of the tweet contains no fresh information. The statement "All 3 agricultural laws have been repealed" stands alone as a claim, although it does not include all of the details the author has provided.

***Verifiability:*** A span chosen as a claim must reflect data that the annotator believes to be

verifiable, meaning that the data's accuracy may be checked against any number of external data sources. Let's examine the following tweet as an illustration:

किसान दिल्ली बॉर्डर पर 10वें दिन भी जमे रहेंगे। आज किसान नेताओं और सरकार के बीच 5वें दौर की बातचीत हुई। जो कि बेनतीजा रही। अब 9 दिसंबर को सुबह 11 बजे एक बार फिर सरकार और किसान नेताओं की बातचीत होगी।

#FarmBill #FarmersProtest #Kisan #Delhi #DelhiBorder #KisanProtest #KisanMeeting

<https://t.co/QIY2a4c6pv>

The information stated in this tweet largely comprises the author's viewpoint towards the media. It is not feasible to establish if Farm Bill has always been against Narendra Modi (See 3 for additional unfavorable instances) (See 3 for more negative examples). Generally, spans like

- 1) *X(some event) happened at Y(some place or time)*
  - 2) *X(some person or organization) said Y(some statement),*
  - 3) *X(some quantifiable entity) has a value Y(some numeric value)*
  - 4) *X(some person or organization) has done Y(some act)*
- etc. can be readily verified and hence, should be tagged as a claim.

Authorship: A piece of writing may include nested information, such “X claimed that Y happened”. In such circumstances, we need to recognize that whether “Y happened” or not is not the topic of the claim here; the claim is X stated so. For example:

पूरे देश के किसानों ने आपस में तालमेल कर लिया है, 13 राज्यों से समर्थन आ चुका है। सरकार को जल्दी इसका हल निकालना चाहिए, नहीं तो 9 दिसंबर की बैठक के बाद नई रणनीति बनेगी:सुखविंदर सिंह सभरा,किसान मज़दूर संघर्ष कमेटी पंजाब #FarmersProtest #FarmerProtests #FarmersDilliChalo #FarmBill <https://t.co/Y0J3dsEm6F>

In this tweet, the author says that the farmers has settled among themselves and supported by 13 states that the complete highlighted sentence should be marked as the claim.

Entities: they might be items, locations, individuals, organizations, periods, etc. that are engaged in the information contained within the claim. For example, in the first tweet, the entities are “देश किसानों, 13 राज्यों”.

### 3.1.4 Types of Claims

A section of text marked as a claim can further be divided into 3 types:

**Simple:** the tweet's author expresses information in the span in his or her own words.

Examples (claim span in blue):

Bharat Bandh 2020: Heightened Security At KSR Railway Station Over Farmers Protest In Bengaluru

Video Link ► <https://t.co/FNvI2448aG>

#TV9Kannada #Bengaluru #KSRRailwayStation #BharatBandh #KarnatakaBandh #FarmBill #FarmLaw #KannadaNews  
<https://t.co/OYg554mDz6>

कृषि मंत्री से मिले हरियाण के किसान.. हरियाणा के किसानों ने नए कृषि कानून पर सरकार को किया समर्थन।

#FarmBill #FarmersProtest #Kisan #Haryana #BharatBandh #HaryanaKisan #KisanStandsaWithModi  
#FarmersWithModi #कल\_भारत\_बंद\_रहेगा #कल\_भारत\_बंद\_नहीं\_होगा <https://t.co/KEHInPCeBD>

**Composite:** The span consists of a variety of nested facts; it should be emphasised that if each of these facts can be marked separately in a useful fashion, then they should each be labelled separately. Think about the following tweet, for instance (claim in blue):

भारतीय किसान यूनियन दोपहर 12:00 बजे से शाम 5:00 बजे तक करेगी चक्का जाम 5:00 बजे दिल्ली कूच होगा  
इससे पहले किसानों की महापंचायत भी होगी नोएडा के डीएनडी व कालिंदी कुंज बॉर्डर को भी किया जाएगा सील।  
#BharatBandh #kisanandolan #farmbill <https://t.co/RbG0uyFJNv>

The chosen span presents three distinct facts that are combined. None of these 3 facts can be mentioned separately since doing so would change how the span is understood. This is a composite claim as a result.

This tweet consists of three separate claim spans that can be characterised as independent of one another. Each claim span is a simple span. It should be noted that any dependencies on truthfulness among the verifiable facts in a composite claim must not be taken into account. The issue is similar when making two simple remarks in a single tweet. That is, regardless of whether the verification of X depends on the verification of Y (or vice versa), if there are two claims of type X and Y, where X, Y are different facts to be confirmed, then they are separate claims.

**Quotation:** The author is quoting someone else or something else as the knowledge's source.

Example,

कल भारत बंद का प्रदर्शन 11 बजे दिन से लेकर दोपहर 3 बजे तक ही होगा- राकेश टिकैत, प्रवक्ता, भारतीय किसान यूनियन।

#BharatBandh #FarmerPolitics #FarmBill #kisanandolan

It is to be observed that the verifiability axiom need not be met inside the cited sentence. That is, X asserted Y is verifiable (and thus a claim) whether or not Y is verifiable.

### 3.1.5 Examples on type of claim

NO	selected claim	type of claim	entities
1	Bharat Bandh 2020: Heightened Security At KSR Railway Station Over Farmers Protest In Bengaluru	simple	Bharat Bandh 2020, KSR Railway Station, Farmers, Protest, Bengaluru
2	1.FarmersProtest All opposition parties are just trying to defame narendramodi PMOIndia with the help of external forces but these techniques became old now. 2. If government is doing something wrong then all state's farmers would have joined but they are welcoming new FarmBill	composite	.FarmersProtest, opposition parties, narendramodi, PMOIndia, government, farmers, FarmBill
3	1.कृषि मंत्री से मिले हरियाण के किसान 2.हरियाणा के किसानों ने नए कृषि कानून पर सरकार को किया समर्थन	simple, simple	कृषि मंत्री, हरियाणा, किसानों, कृषि कानून
4	भारतीय किसान यूनियन दोपहर 12:00 बजे से शाम 5:00 बजे तक करेगी चक्का जाम 5:00 बजे दिल्ली कूच होगा इससे पहले किसानों की महापंचायत भी होगी नोएडा के डीएनडी व कालिंदी कुंज बॉर्डर को भी किया जाएगा सील	composite	भारतीय किसान यूनियन, दिल्ली, किसानों, महापंचायत, नोएडा, डीएनडी, कालिंदी कुंज बॉर्डर
5	Bollywood Actress Urvashi Rautela Says, "Kisaan Mera Bhagwaan"	quotation	Bollywood, Actress,Urvashi Rautela,Kisaan

Table – 3.1 Type of Claim

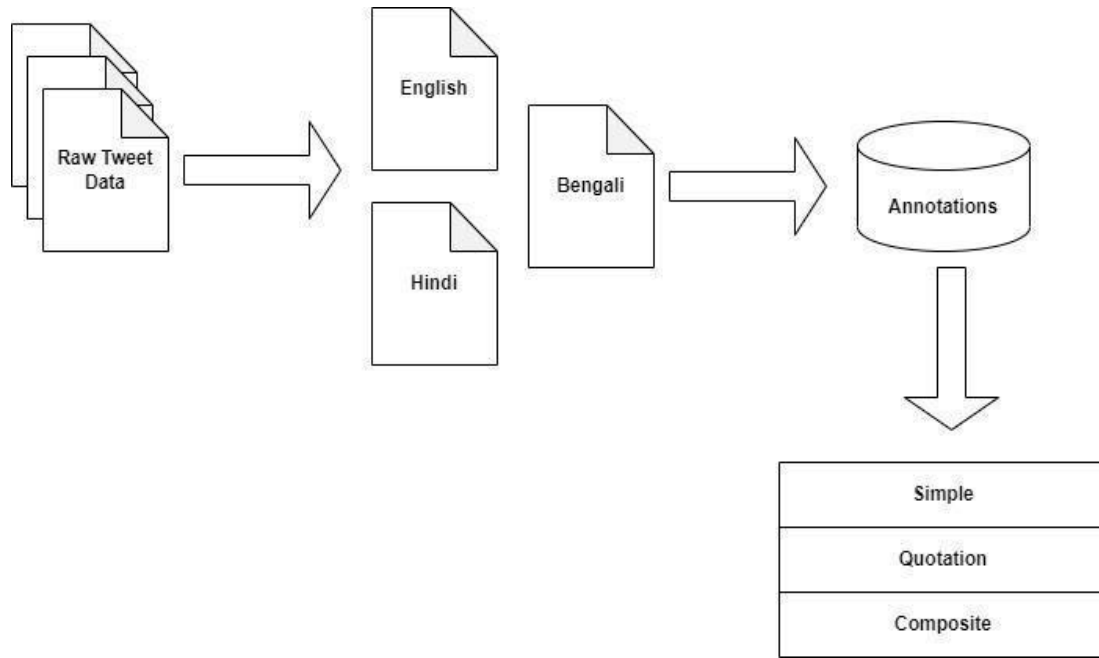


Fig. – 3.2 - procedure of annotations

### 3.1.6 Code-Mixed Dataset

All of the raw tweet data that Twitter Crawler collects. Several types of linguistic data are acquired after preprocessing all raw datasets. Datasets are available in Bengali, Hindi, and English. In addition to Bengali, English, and Hindi, we also get Code-mix datasets. It combines two or more different languages. It could be a combination of Bengali and English or Hindi and English. We refer to these datasets as code-mix datasets.

Example –

1. JK24x7 News पर 'बात देश की' में आज देखिए कल होने वाले भारत बंद को लेकर सियासत पर बड़ी बहस.. शाम 6:30 बजे..  
#JK24x7News #FarmersProtest #FarmBill #Kisan #Delhi #BharatBand #delhiborders  
<https://t.co/P9JsaU2wRV>
2. As per reports,a number of bank unions have also expressed solidarity with the protesting farmers.  
#FarmersProtest #BharatBandh #FarmersBill2020 #FarmBill #BharatBandh #FarmersProtestDelhi2020  
#8\_दिसंबर\_भारत\_बन्द



## 3.2 Source and Characteristics

For this thesis work on fake news classification, the choice of datasets for training and testing is critical. Here are some potential sources and characteristics of datasets that can consider:

### Source of Datasets:

#### I. FakeNewsNet:

- This dataset is collected from Twitter, which makes it particularly relevant for understanding how misinformation spreads on social media platforms.
- It includes rumors, fake images, and fake user accounts, offering a comprehensive view of different forms of misinformation.
- The diversity of content types allows your model to learn from textual claims, images, and user interactions.

#### II. LIAR:

- The LIAR dataset is compiled from fact-checking websites, ensuring that the statements are verified and labeled by experts.
- It spans a wide range of topics, making it a suitable choice for assessing the generalization ability of your model across various domains.
- The labels provide information about the degree of falsehood, allowing your model to capture nuances in misinformation.

#### III. BuzzFeedNews:

- Curated by BuzzFeed, this dataset captures the multimedia nature of contemporary fake news.
- By including both textual and visual content, it prepares your model to handle the complexities of misinformation presented in different formats.
- It reflects the modern challenges of detecting fake news that leverages a combination of text and imagery.

#### IV. PolitiFact:

- PolitiFact focuses on political claims, which can be particularly relevant if your research aims to tackle political misinformation.
- provides a substantial collection of fact-checked claims, which is invaluable for evaluating your model's accuracy and fact-checking abilities.

## **Characteristics of Datasets:**

### **I. Content Variety:**

- Diverse datasets ensure that your model learns to recognize fake news across various topics, from politics to health to entertainment.
- A varied dataset prevents your model from overfitting to specific content and prepares it for real-world applications.

### **II. Balanced Classes:**

- Aim for a balance between true and fake instances in your dataset.
- Class imbalance can bias your model toward the majority class, affecting its performance on the minority class. Balanced data helps your model learn from both classes equally.

### **III. Multimodal Content:**

- Incorporating datasets with both textual and visual content mimics the actual environment where fake news is often spread through a combination of text and images.
- It enables your model to capture the subtle cues present in different types of media.

### **IV. Fact-Checked Labels:**

- Fact-checked labels provide a ground truth for evaluating your model's accuracy.
- They ensure that your model's performance is assessed against verified claims, enhancing the trustworthiness of your results.

### **V. Metadata:**

- Metadata such as timestamps, sources, and contextual information enrich your dataset with additional information.
- It helps your model understand the context in which the content was created, which can be crucial for accurate classification.

### **VI. Ethical Considerations:**

- Always adhere to ethical guidelines and user privacy rights when selecting and using datasets, especially when dealing with social media data.

### 3.3 Data Splitting and Validation Procedures

Proper data splitting and validation procedures are crucial to ensure the reliability and generalizability of your fake news classification model. Here's a detailed explanation of data splitting and validation techniques for this thesis paper:

#### I. Data Splitting:

Data splitting is the process of dividing your dataset into subsets for training, validation, and testing. This separation ensures that your model learns, fine-tunes, and evaluates its performance on different data points.

- **Training Set:** This is the largest subset and is used to train your model. It's crucial for the model to learn the patterns and features present in the data.
- **Validation Set:** A smaller subset used to fine-tune hyperparameters and make decisions during training. It helps prevent overfitting by allowing you to assess the model's performance on data it hasn't seen during training.
- **Testing Set:** This subset remains untouched during training and validation. It's used to evaluate your model's performance after training is complete. This evaluation provides an estimate of how well your model will perform on unseen data in the real world.

#### II. Cross-Validation:

Cross-validation is a robust technique that addresses concerns about data splitting. It involves creating multiple training-validation splits and averaging the performance metrics to get a more reliable estimate of model performance. Common methods include:

- **K-Fold Cross-Validation:** The dataset is divided into  $k$  equally sized folds. Each fold is used as the validation set once, while the remaining  $k-1$  folds are used for training. The process is repeated  $k$  times, and performance metrics are averaged.
- **Stratified Cross-Validation:** Especially useful for imbalanced datasets, this method ensures that each fold maintains the class distribution of the original dataset. It prevents any particular class from being underrepresented in the validation sets.

### **III. Validation Procedures:**

Validation procedures help you fine-tune your model, monitor its performance, and select the best configuration.

- **Hyperparameter Tuning:** Use the validation set to experiment with different hyperparameters, such as learning rates or regularization strengths. Observe how these changes affect your model's performance.
- **Early Stopping:** Monitor the validation loss or accuracy during training. If your model's performance on the validation set plateaus or starts to decline, stop training to prevent overfitting.
- **Model Selection:** After training multiple models with varying settings, choose the one that performs best on the validation set. Avoid using the validation set multiple times for fine-tuning, as this can lead to overfitting on the validation data.

### **IV. Ethical Considerations:**

- **Privacy and User Rights:** Always ensure that data splitting and validation procedures respect ethical guidelines and user privacy rights, especially when dealing with sensitive information or social media data.
- **Transparency:** Clearly document data splitting and validation procedures in your thesis to ensure transparency and reproducibility for this research.

By effectively splitting data and implementing rigorous validation procedures, ensure that fake news classification model is well-trained, accurately tuned, and capable of handling unseen data. These practices enhance the credibility of the findings and contribute to the robustness for this research.

## **3.4 Feature Extraction and Selection Techniques**

Feature extraction and selection are crucial steps in preparing for this data for fake news classification. These techniques help the model focus on the most relevant information, enhance its efficiency, and improve its performance. Here's an in-depth exploration of feature extraction and selection for this thesis paper:

## I. Feature Extraction:

Feature extraction involves transforming raw data into a format that can be effectively used by machine learning models. For text and image data in fake news classification:

**Text Data:** Some of the important features are as follows

- **Bag-of-Words (BoW):** Convert text into a matrix of word frequencies. Each word becomes a feature, and the matrix captures the presence of words in documents.
- **Term Frequency-Inverse Document Frequency (TF-IDF):** Assign weights to words based on their frequency in a document relative to the entire corpus.
- **Word Embedding:** Transform words into dense vectors using pre-trained models like Word2Vec, GloVe, or FastText. These embeddings capture semantic relationships between words.
- **N-grams:** Capture sequences of n words to preserve contextual information in text data.

**Image Data:** Some of the important features are as follows

- **Convolutional Neural Networks (CNNs):** Use pre-trained CNN models like VGG16, ResNet, or Inception to extract hierarchical features from images.
- **Local Binary Patterns (LBP):** Analyze texture patterns in images by computing the local relationship between a pixel and its neighbors.
- **Color Histograms:** Quantify the distribution of colors in an image. Useful for simple image classification tasks.

## II. Feature Selection:

Feature selection involves choosing the most relevant features while eliminating irrelevant or redundant ones. It enhances model efficiency, reduces overfitting, and improves interpretability:

**Filter Methods:** Evaluate features based on statistical measures like correlation, chi-square, or mutual information. Select the most informative features to include in the model.

**Wrapper Methods:** Use machine learning models as subroutines to assess feature subsets. Techniques like recursive feature elimination (RFE) and forward selection iteratively add or remove features based on model performance.

**Embedded Methods:** Feature selection is incorporated directly into the model training process.

Techniques like LASSO (L1 regularization) penalize coefficients to encourage sparsity and feature selection.

**Dimensionality Reduction:** Techniques like Principal Component Analysis (PCA) and t-Distributed Stochastic Neighbor Embedding (t-SNE) reduce the dimensionality of features while preserving variance or clustering structures.

### **III. Ethical Considerations:**

When dealing with textual data, ensure that feature extraction techniques do not unintentionally propagate bias present in the data.

Feature selection should be conducted ethically, avoiding the removal of features that could carry valuable information about the veracity of news.

### **IV. Interpretability and Explain ability:**

While advanced techniques like deep learning can provide excellent performance, consider the trade-off between complexity and interpretability. Simple models with interpretable features can offer insights into how your model makes decisions.

By employing effective feature extraction and selection techniques, equip the fake news classification model with relevant, informative, and manageable inputs. These practices enhance model performance, enable efficient training, and contribute to a clearer understanding for the model's decision-making process.

## **3.5 Description of Machine Learning Models**

In fake news classification thesis, the choice of machine learning models plays a pivotal role in accurately discerning between true and false information. Here's an overview of Logistic Regression and BERT models for classification tasks in fake news detection:

### **3.5.1 Logistic Regression:**

Logistic Regression is a fundamental machine learning model used for binary classification tasks, making it applicable to your fake news classification thesis. Here's a comprehensive overview of logistic regression and its relevance to this research:

**Model Type:** Linear model

**Overview:** Logistic Regression is a statistical method used to model the probability of a binary outcome. It's particularly suitable for tasks where the target variable has two possible classes, such as classifying news articles as true or fake.

**Working Principle:**

- **Log-Odds Transformation:** Logistic Regression transforms the linear combination of input features using the log-odds (logit) function. The log-odds represent the likelihood of the positive class (fake news) over the negative class (true news).
- **Sigmoid Function:** The log-odds are then passed through the sigmoid function, also known as the logistic function. This transforms the log-odds to a value between 0 and 1, representing the estimated probability of the positive class.
- **Thresholding:** A threshold is chosen (commonly 0.5) to determine the predicted class. If the estimated probability is above the threshold, the instance is classified as the positive class (fake news); otherwise, it's classified as the negative class (true news).

**Advantages:**

- **Interpretability:** Logistic Regression provides interpretable coefficients associated with each feature. These coefficients show the direction and magnitude of the influence of each feature on the outcome.
- **Probabilistic Interpretation:** The model produces probability scores, allowing you to understand the confidence level of its predictions.
- **Simplicity:** It's a relatively simple model, making it easy to implement, understand, and explain.

**Considerations:**

- **Linear Assumption:** Logistic Regression assumes a linear relationship between the features and the log-odds of the target class. It might struggle with capturing complex nonlinear relationships present in some datasets.
- **Feature Engineering:** The performance of Logistic Regression can be enhanced through careful feature engineering, transforming and combining features to capture relevant patterns.

- **Multicollinearity:** Logistic Regression can be sensitive to multicollinearity (high correlation between features), which might lead to unstable coefficient estimates.

### **Application in Fake News Classification:**

Logistic Regression can be applied effectively to your fake news classification task. By encoding textual, visual, and metadata features as input, you can build a logistic regression model to estimate the probability of a news article being fake. The coefficients associated with each feature will offer insights into which features contribute to the prediction, aiding in the interpretation of your model's decisions.

### **Ethical Considerations:**

Ensure that the features we use for Logistic Regression are selected ethically and do not inadvertently introduce bias or perpetuate misinformation.

In summary, Logistic Regression is a valuable model for this fake news classification research due to its interpretability, probabilistic outputs, and simplicity. While it may not capture highly complex relationships, it can serve as a strong baseline for this classification task.

## **3.5.2 BERT models**

BERT (Bidirectional Encoder Representations from Transformers) is a revolutionary language representation model that has significantly advanced the field of natural language processing. Here's a comprehensive overview of BERT models and their relevance to this fake news classification thesis:

**Model Type:** Pre-trained transformer-based deep learning model

**Overview:** BERT is a contextualized word embedding model developed by Google. Unlike traditional word embeddings like Word2Vec or GloVe, BERT captures the context of words in a sentence by considering the surrounding words from both directions. This ability to understand the context makes BERT well-suited for various NLP tasks, including text classification.

### **Working Principle:**

- **Pre-training:** BERT is pre-trained on a massive amount of text data using a masked language model objective. During this pre-training, it learns to predict masked words within a sentence, effectively capturing semantic relationships and context.



- **Fine-Tuning:** After pre-training, BERT can be fine-tuned on a specific task, such as fake news classification. Fine-tuning involves adding task-specific layers and training the model on labeled data related to the task.

#### **Advantages:**

- **Contextual Understanding:** BERT captures context and semantics, allowing it to understand the subtle meaning of words in different contexts, which is crucial for fake news classification.
- **Transfer Learning:** BERT is pre-trained on massive text corpora, providing it with a strong linguistic foundation. Fine-tuning on your task requires less data compared to training from scratch.
- **State-of-the-Art Performance:** BERT models have consistently achieved top performance in various NLP benchmarks.

#### **Considerations:**

- **Computational Resources:** BERT models are computationally intensive and require significant resources for both pre-training and fine-tuning.
- **Fine-Tuning Data:** Fine-tuning BERT requires labeled data specific to your task. Having a sufficiently large and relevant dataset is crucial for achieving good results.
- **Model Interpretability:** BERT models are complex, making them less interpretable compared to simpler models like logistic regression.

#### **Application in Fake News Classification:**

Using BERT for fake news classification involves fine-tuning the pre-trained model on a labeled dataset of true and fake news articles. The model learns to distinguish between the two based on the contextual information within the articles. By encoding textual content, metadata, and potentially visual information, BERT can capture intricate linguistic cues that indicate the veracity of news.

#### **Ethical Considerations:**

Ensure that the fine-tuning dataset and labeling follow ethical guidelines to prevent the propagation of misinformation or biases.

In summary, BERT models offer state-of-the-art performance for fake news classification tasks by leveraging their contextual understanding of language. While they require substantial computational resources and fine-tuning data, BERT's capabilities can greatly enhance the accuracy and robustness of this classification model.

### **3.6 Experimental Setup and Evaluation Metrics**

The experimental setup and evaluation metrics are critical components of this fake news classification thesis. They determine how to assess the performance of these models and provide insights into their effectiveness. Here's a comprehensive guide on designing this experimental setup and selecting appropriate evaluation metrics:

#### **Experimental Setup:**

##### **Datasets:**

- Use a diverse set of datasets that include textual, visual, and metadata information.
- Ensure these datasets are well-preprocessed and properly split into training, validation, and testing subsets.
- Consider well-known datasets like FakeNewsNet, LIAR, BuzzFeedNews, and PolitiFact for benchmarking.

##### **Feature Extraction:**

- Apply relevant feature extraction techniques based on the nature of your data (text, images, metadata).
- Transform textual data using techniques like TF-IDF or BERT embeddings.
- Process image data using pre-trained CNNs or feature extraction methods.

##### **Machine Learning Models:**

- Implement and train various machine learning models, such as Logistic Regression, Random Forest, BERT-based models, and others.
- Fine-tune hyperparameters using the validation set to optimize model performance.

##### **Experimental Design:**

- Use k-fold cross-validation to ensure robustness of results.
- Randomly shuffle data before splitting to prevent any inherent order bias.
- For each fold, train and validate models, then evaluate on the testing set.

## Evaluation Metrics:

- **Accuracy:** Calculates the ratio of correctly predicted instances to the total instances. Suitable when classes are balanced.
- **Precision:** Measures the proportion of true positive predictions out of all positive predictions. Useful when false positives are costly.
- **Recall (Sensitivity):** Represents the ratio of true positive predictions out of all actual positive instances. Important when false negatives are critical.
- **F1-Score:** The harmonic mean of precision and recall, providing a balanced measure between precision and recall.
- **Area Under the Receiver Operating Characteristic (ROC-AUC):** Measures the model's ability to distinguish between classes across different thresholds.
- **Confusion Matrix:** Provides a detailed breakdown of true positives, true negatives, false positives, and false negatives.
- **Specificity:** Measures the proportion of true negative predictions out of all actual negative instances.
- **Matthews Correlation Coefficient (MCC):** Takes into account true and false positives/negatives and is particularly useful for imbalanced datasets.

## 3.7 Ethical Considerations

Handling fake news data for your thesis paper requires careful attention to ethical considerations. Misinformation and its potential impacts on individuals and society necessitate responsible data handling practices. Here are important ethical considerations to keep in mind:

### Data Privacy and Consent:

- Ensure that you have proper rights and permissions to use the data you're working with.
- Respect user privacy by anonymizing personal information and adhering to data protection regulations.
- If using social media data, comply with the platform's terms of use and guidelines.

### Avoid Amplification:

- Be cautious not to amplify or spread fake news unintentionally while researching or reporting on it.
- Handle the data responsibly and prevent the inadvertent sharing of misleading

information.

**Bias and Fairness:**

- Address bias in this dataset to prevent the perpetuation of biases present in the data.
- If bias exists, document it and consider mitigation strategies.
- Clearly documented data collection and preprocessing methods, so readers can understand how data was prepared.
- Describe this modeling choices, including any assumptions or simplifications made.

**Avoid Spreading Misinformation:**

- While working with fake news data, take measures to prevent its accidental dissemination or endorsement.
- Clearly label of this work as research to avoid confusion with actual news content.

**Transparency:**

- Clearly document your data collection and preprocessing methods, so readers can understand how your data was prepared.
- Describe modeling choices, including any assumptions or simplifications made.

**Ethical Disclosure:**

- If these are sharing of this research results publicly, be transparent about of these findings, methods, and limitations.
- Clearly distinguish between real news and synthetic examples if you're generating fake news instances for testing purposes.

**Minimize Harm:**

- Be aware of the potential harm that can result from researching and analyzing fake news.
- Focus on constructive goals, such as understanding the spread of misinformation and developing effective detection techniques.

**Responsible Reporting:**

- This research includes examples of fake news, be cautious when presenting them to avoid further spreading misinformation.
- Provide contextual information to ensure a balanced representation.

**Collaboration and Peer Review:**

- Engage with experts in the field to ensure your research methods and findings align with ethical guidelines.

- Seek peer review to ensure your work meets ethical standards.

**Educational Approach:**

- If applicable, frame of this research in an educational context to raise awareness about the dangers of fake news.
- Use research to contribute to public understanding of misinformation and its implications.

Ethical considerations are integral to maintaining the integrity of this research and its impact on society. By following ethical guidelines, we can conduct for this fake news classification research responsibly and contribute positively to addressing the challenges posed by misinformation.

# Chapter 4

## Fake News Classification Framework

Developing a comprehensive framework for fake news classification is essential for structuring research, guiding methodology, and presenting findings. Here's a structured framework that can adapt and customize for your thesis paper:

### 4.1 Explanation of Machine Learning Algorithms

In fake news classification thesis, the choice of machine learning algorithms significantly impacts the accuracy and effectiveness of classification system. Let's delve into the detailed explanations of the two chosen algorithms: Logistic Regression and BERT-based models.

#### I. Logistic Regression:

##### Overview:

A statistical model called logistic regression uses the logistic function, often known as the logit function in mathematics, as the relationship between  $x$  and  $y$ .  $Y$  is represented by the logit function as a sigmoid function of  $x$ .

$$f(x) = \frac{1}{1 + e^{-x}}$$

Logistic Regression is a linear model commonly used for binary classification tasks. Despite its simplicity, it has been widely applied in various domains due to its interpretability and efficiency. If you plot this logistic regression equation, you will get an S-curve as shown below.

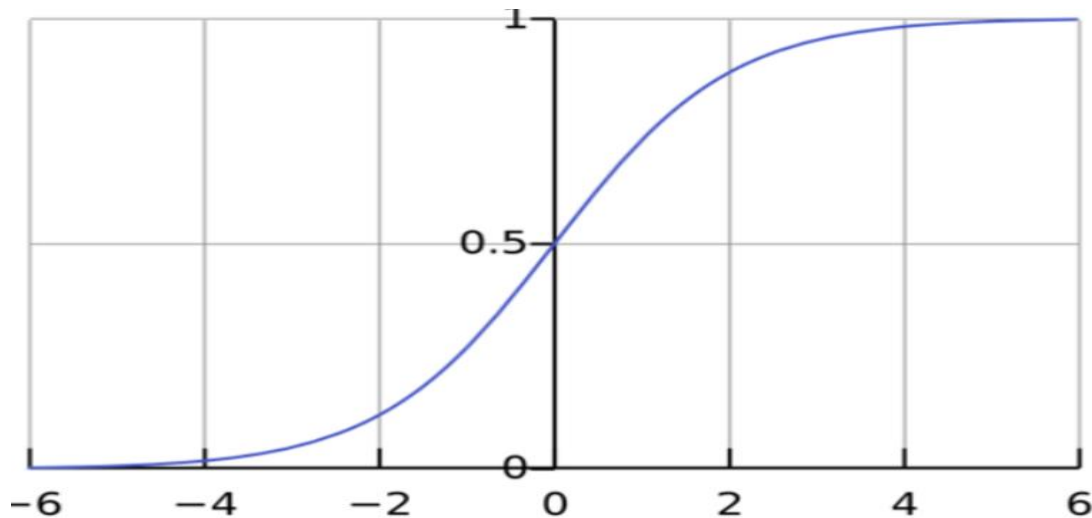


Fig 4.1 S-curve (Logistic Regression)

**Explanation:**

**Working principle:** Logistic Regression models the log-odds of the probability that an instance belongs to a specific class (fake news in your case). The log-odds are transformed using the sigmoid function to produce the final probability.

$$\text{Logit Function} = \log \left( \frac{p}{1-p} \right)$$

**Interpretability:** Logistic Regression provides interpretable coefficients for each feature. These coefficients indicate the strength and direction of the relationship between features and the target class. Given that the outcome of logistic regression is a probability between 0 and 1, the interpretation of the weights in logistic regression is different from the interpretation of the weights in linear regression. The likelihood is no longer linearly influenced by the weights. The logistic function converts the weighted sum into a probability. So that only the linear part is on the formula's right side, we must reformulate the equation for the meaning.

$$\ln \left( \frac{P(y=1)}{1-P(y=1)} \right) = \log \left( \frac{P(y=1)}{P(y=0)} \right) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

**Decision Boundary:** In a two-dimensional feature space, Logistic Regression's decision boundary is a linear separator. It can capture linear relationships between features and classes. Choosing a decision boundary for a binary classification problem is the basic use of logistic regression. Although the approach can be used for situations when there are several

classification classes or multi-class classification, the baseline is to determine a binary decision boundary.

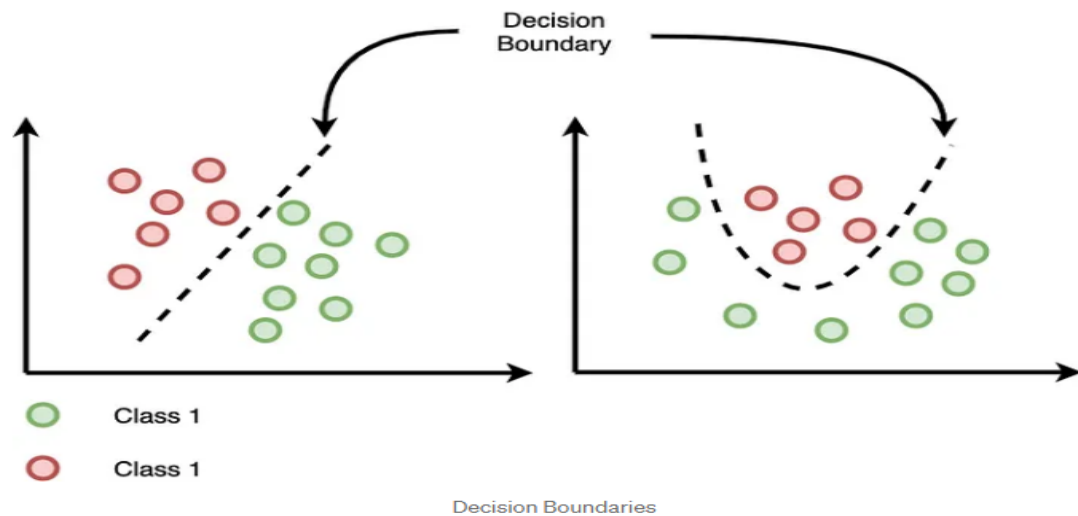


Fig 4.2 Binary Decision Boundary.

#### Advantages:

- **Interpretability:** The model's coefficients offer insights into feature importance and directionality.
- **Simplicity:** Logistic Regression is relatively straightforward to implement and understand.
- **Efficiency:** It can handle large datasets efficiently and requires less computational resources compared to more complex models.
- **Baseline Performance:** Logistic Regression serves as a baseline model to compare against more advanced approaches.

#### Considerations:

- **Linear Assumption:** Logistic Regression assumes a linear relationship between features and the log-odds. It may struggle with capturing complex nonlinear relationships.
- **Feature Engineering:** Its performance can be enhanced through careful feature engineering.
- **Prone to Outliers:** Extreme outliers might disproportionately influence the model's coefficients.



## II. BERT-based Models:

**Overview:** Transformer attention is a component of BERT (Bidirectional Encoder Representations from Transformers), which has been designed to understand the links between words in context. A text input reader called an encoder is part of the transformer. A decoder, which makes predictions based on the task, is also a part of it. The transformer encoder reads every word concurrently, making it non-directional in contrast to directional devices that read the text input in a sequential sequence. This indicates that the model gathers the context of a word from all of the words it is surrounded by. As a result, the bidirectional BERT architecture for sentence-level categorization is used to describe the model's design.

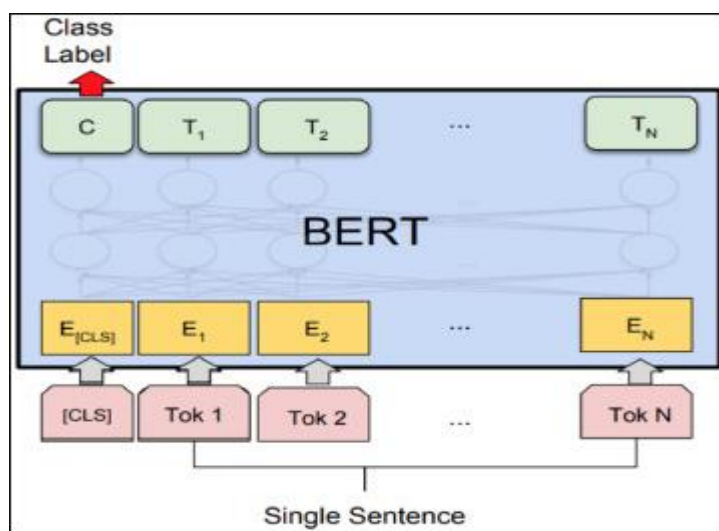


Figure 4.3: proposed framework for content-based fake news classification

### Explanation:

- **Pre-training:** BERT is pre-trained on massive text corpora using a masked language model objective. This enables it to capture the context and relationships between words.
- **Fine-Tuning:** After pre-training, BERT can be fine-tuned for specific tasks, such as fake news classification. Additional layers are added to adapt the model to the classification task.

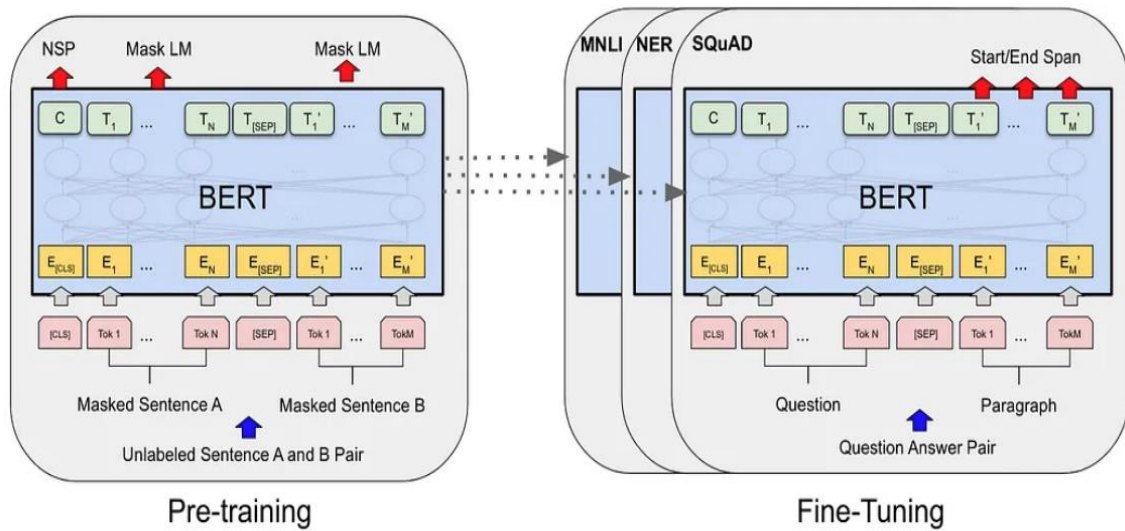


Figure 4.4: Framework for BERT Pre-trained and Fine-Tuning

#### Advantages:

- **Contextual Understanding:** BERT captures contextual information, allowing it to understand the nuanced meaning of words in different contexts.
- **Transfer Learning:** BERT's pre-training enables it to leverage linguistic knowledge, requiring less labeled data for fine-tuning.
- **Performance:** BERT models have achieved state-of-the-art performance in various NLP tasks.

#### Considerations:

- **Resource-Intensive:** Training and fine-tuning BERT models require substantial computational resources.
- **Fine-Tuning Data:** Adequate labeled data is needed for fine-tuning to prevent overfitting.
- **Model Complexity:** BERT models are complex and less interpretable than simpler models like Logistic Regression.

Logistic Regression serves as an interpretable baseline, while BERT models excel in capturing the intricate linguistic cues that differentiate true and fake news. Combining both approaches allows to explore the trade-off between interpretability and advanced performance.

## 4.2 Comparison of Different Techniques

Comparing different techniques and understanding their pros and cons is crucial for making informed decisions in this fake news classification thesis. Here's a comprehensive comparison of the chosen techniques along with their advantages and limitations:

### Comparison of Techniques: Logistic Regression vs. BERT-based Models

#### I. Logistic Regression:

##### Pros:

- **Interpretability:** Logistic Regression offers easily interpretable coefficients that explain how each feature contributes to the classification decision.
- **Efficiency:** It is computationally efficient and requires fewer resources compared to more complex models.
- **Baseline Model:** Logistic Regression serves as a baseline for comparison, allowing you to gauge the improvement achieved by more advanced models.
- **Feature Importance:** Coefficients provide insights into the relative importance of different features in making predictions.
- **Linear Relationships:** Effective when relationships between features and the target class are linear.

##### Cons:

- **Limited Complexity:** Struggles with capturing nonlinear relationships present in intricate data.
- **Feature Engineering:** Performance can be limited if feature engineering does not capture the nuances of the data.
- **Model Flexibility:** Less adaptable to complex decision boundaries in high-dimensional spaces.

#### II. BERT-based Models:

##### Pros:

- **Contextual Understanding:** BERT models excel at capturing contextual and semantic information, making them effective for tasks requiring deep language understanding.
- **State-of-the-Art Performance:** BERT models have achieved remarkable results in various NLP benchmarks.

- **Transfer Learning:** Pre-training on a large corpus allows for fine-tuning on smaller datasets, leveraging existing linguistic knowledge.
- **Handling Complexity:** Effective at capturing complex linguistic patterns and relationships within the text.
- **Flexibility:** Suitable for both text classification and other NLP tasks.

#### **Cons:**

- **Resource-Intensive:** Training and fine-tuning BERT models require significant computational resources.
- **Data Requirements:** Adequate labeled data is essential for fine-tuning to prevent overfitting.
- **Model Complexity:** BERT models are more complex and less interpretable than simpler models like Logistic Regression.
- **Potential Overfitting:** Fine-tuning on a small dataset can lead to overfitting if not carefully managed.
- **Training Time:** Training BERT models takes longer compared to simpler models.

#### **Choosing the Right Technique:**

- **Interpretability vs. Performance:** If interpretability is a priority and linear relationships suffice, Logistic Regression is suitable.
- **Complex Relationships:** For capturing complex linguistic patterns and achieving state-of-the-art performance, BERT models are more suitable.
- **Computational Resources:** Consider the availability of computational resources and the size of your dataset for fine-tuning BERT models.

#### **Hybrid Approach:**

- Combining both techniques can offer the best of both worlds: the interpretability of Logistic Regression and the advanced linguistic understanding of BERT models.

In this thesis, we provide quantitative results comparing the performance of both techniques on the same dataset using various evaluation metrics. Discuss how each technique addresses the challenges of fake news classification and the implications of your findings.

## 4.3 Proposed Improvements

Here are some proposed improvements and novel approaches may consider for this fake news classification thesis:

### **i. Graph-based Representation:**

Consider representing news articles and their relationships using graph structures. Build a knowledge graph of news articles, sources, and their connections, allowing you to leverage graph-based algorithms for classification.

### **ii. Temporal Analysis:**

Incorporate temporal information into your classification model. Analyze how fake news spreads and evolves over time, and develop models that can detect changing trends and tactics used by misinformation spreaders.

### **iii. User-Centric Features:**

Include features related to user interactions, such as social media shares, comments, and reactions. These user-centric features can provide additional context and insights into the credibility of news articles.

### **iv. Attention Mechanisms for Interpretable BERT:**

Enhance the interpretability of BERT-based models by incorporating attention mechanisms that highlight the most important words or phrases in an article. This allows you to explain the model's decisions more effectively.

### **v. Transfer Learning with Domain Adaptation:**

Explore techniques for domain adaptation to make your classification system more effective across different news domains. Fine-tune your models on specific news sources to capture their unique characteristics.

### **vi. Semi-Supervised Learning with Weak Labels:**

Leverage semi-supervised learning techniques to make use of large amounts of unlabeled data. You can use weak labels obtained from heuristics or external sources to guide the model's learning.

### **vii. Counterfactual Explanations:**

Provide counterfactual explanations for instances misclassified by your models. Explain what

changes in an article's content or features would have led to a different classification outcome.

#### **viii. Cross-Platform Analysis:**

Extend your analysis to multiple online platforms and social media networks. Investigate how fake news spreads and is perceived differently across various platforms.

#### **ix. Robustness against Adversarial Attacks:**

Design your models to be resilient against adversarial attacks. Implement techniques such as adversarial training to ensure that your system can accurately classify even when faced with crafted misleading inputs.

#### **x. Ethical AI:**

Integrate ethical considerations into your models. Develop techniques that avoid amplifying or unintentionally spreading fake news during the classification process.

#### **xi. Human-AI Collaboration for Fact-Checking:**

Explore approaches where human fact-checkers work alongside AI models. The AI system could provide suggestions and insights to human fact-checkers, improving the overall accuracy of classification.

## **4.4 Analysis of Classifiers**

Analyzing the performance of the classifiers in fake news classification thesis is a critical step to understand how well the models are working and to draw meaningful insights from results. Here, how can perform a comprehensive analysis of the classifiers' performance:

### **1. Quantitative Metrics:**

- **Calculate and Compare Metrics:** Compute accuracy, precision, recall, F1-score, ROC-AUC, and PR-AUC for each classifier.

For accuracy, It measures how many observations, both positive and negative, were correctly classified.

$$ACC = \frac{tp + tn}{tp + fp + tn + fn}$$

Since the accuracy score is calculated on the predicted classes (not prediction scores) we need to apply a certain threshold before computing it. The obvious choice is the threshold of 0.5 but it can be suboptimal. Let's see an example of how accuracy depends on the threshold choice:

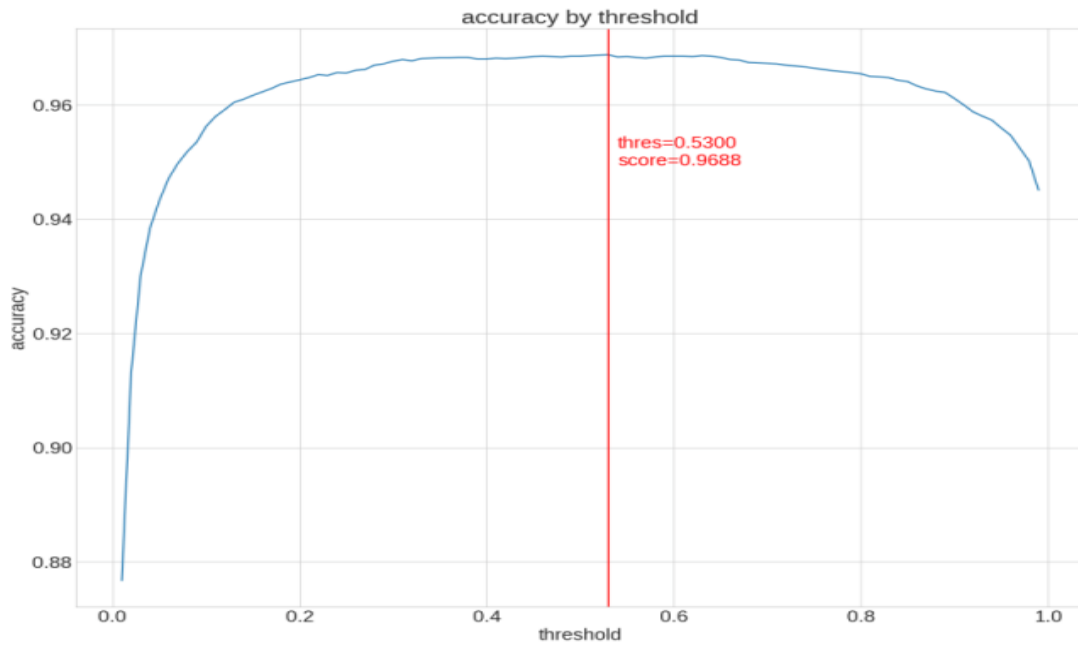


Fig 4.5 : Accuracy VS threshold Graph

For calculating F1 score, we combine both precision and recall into one metric by calculating the harmonic mean between those two. It is actually a special case of the more general function F beta:

$$F_{beta} = (1+\beta^2) \frac{precision * recall}{\beta^2 * precision + recall}$$

When choosing beta in F-beta score the more care about recall over precision the higher beta should choose. For example, with F1 score we care equally about recall and precision with F2 score, recall is twice as important to us.

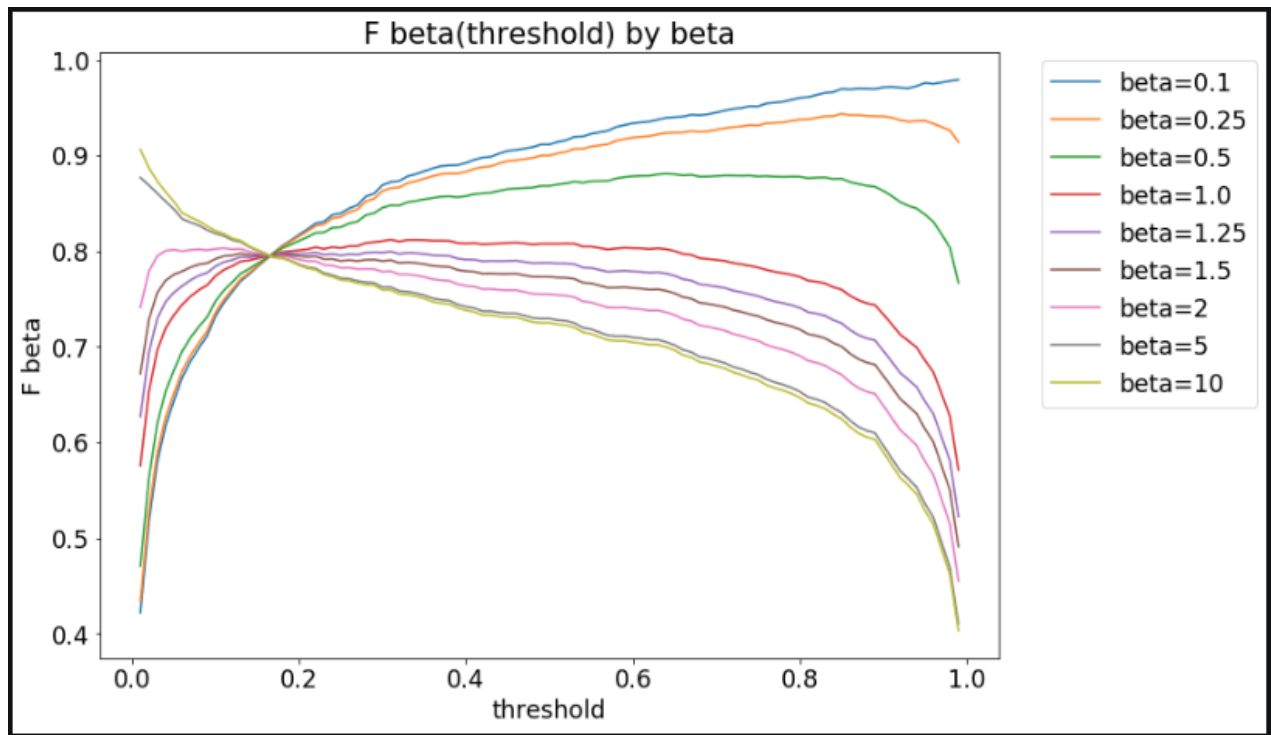


Fig 4.6: F beta ( threshold) by beta

It is important to remember that F1 score is calculated from Precision and Recall which, in turn, are calculated on the predicted classes (not prediction scores).

Let's plot F1 score over all possible thresholds:

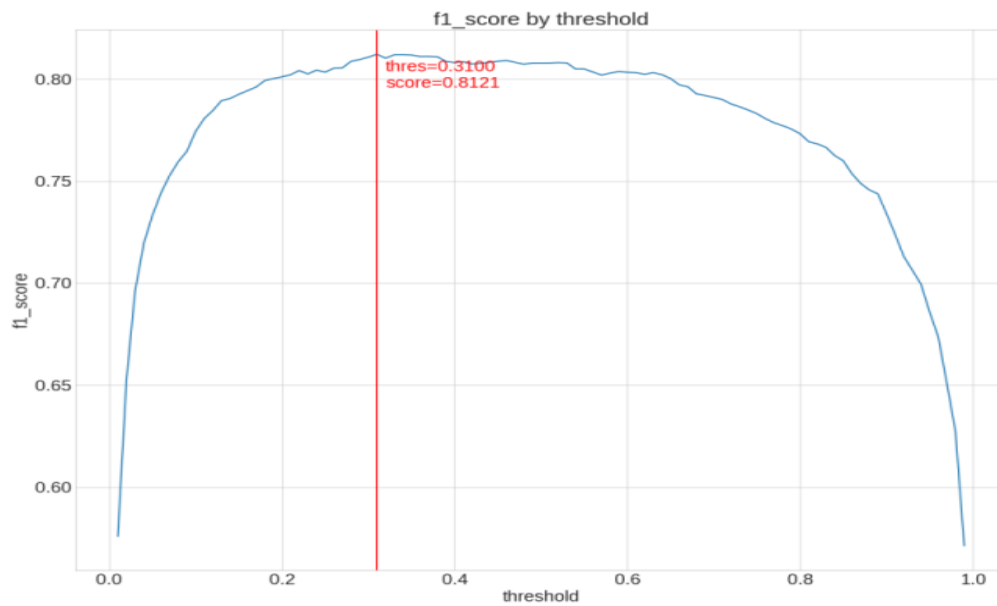


Fig 4.7: F1 Score by threshold



Now, AUC means area under the curve so to speak about ROC AUC score we need to define ROC curve first.

It is a chart that visualizes the tradeoff between true positive rate (TPR) and false positive rate (FPR). Basically, for every threshold, we calculate TPR and FPR and plot it on one chart.

the higher TPR and the lower FPR is for each threshold the better and so classifiers that have curves that are more top-left-side are better.

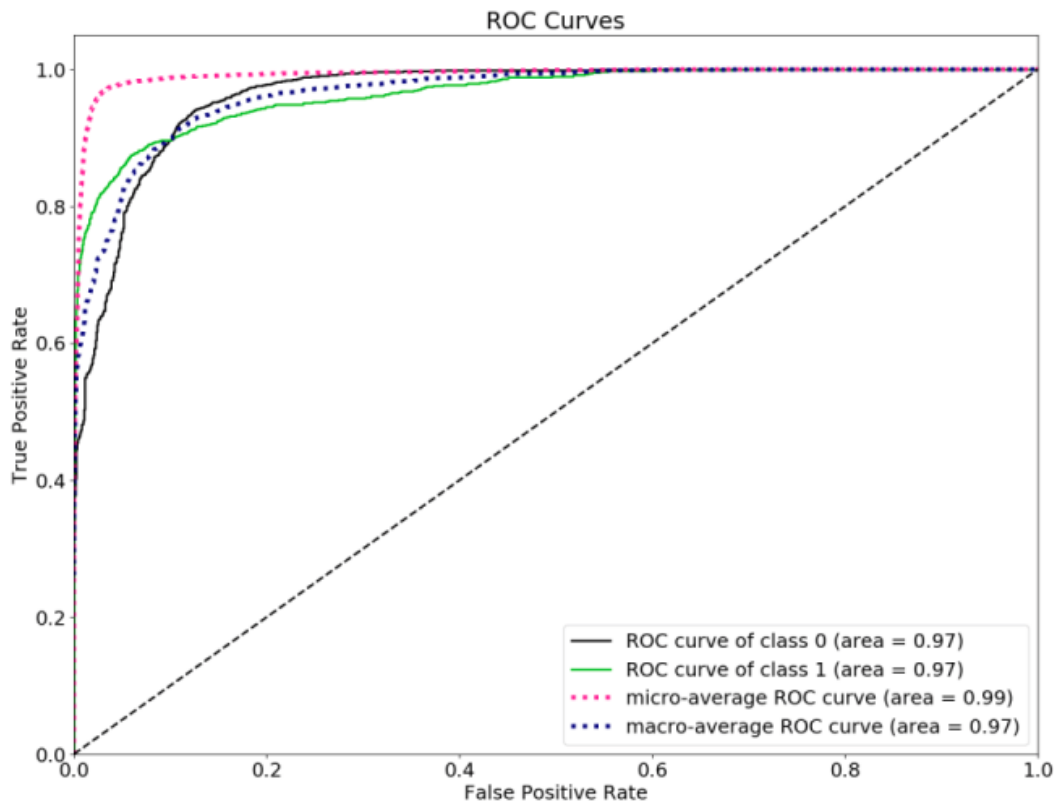


Fig 4.8: F1 ROC Curves

Similarly to ROC AUC in order to define PR AUC we need to define what Precision-Recall curve.

It is a curve that combines precision (PPV) and Recall (TPR) in a single visualization. For every threshold, calculate PPV and TPR and plot it. The higher on y-axis curve is the better model performance.

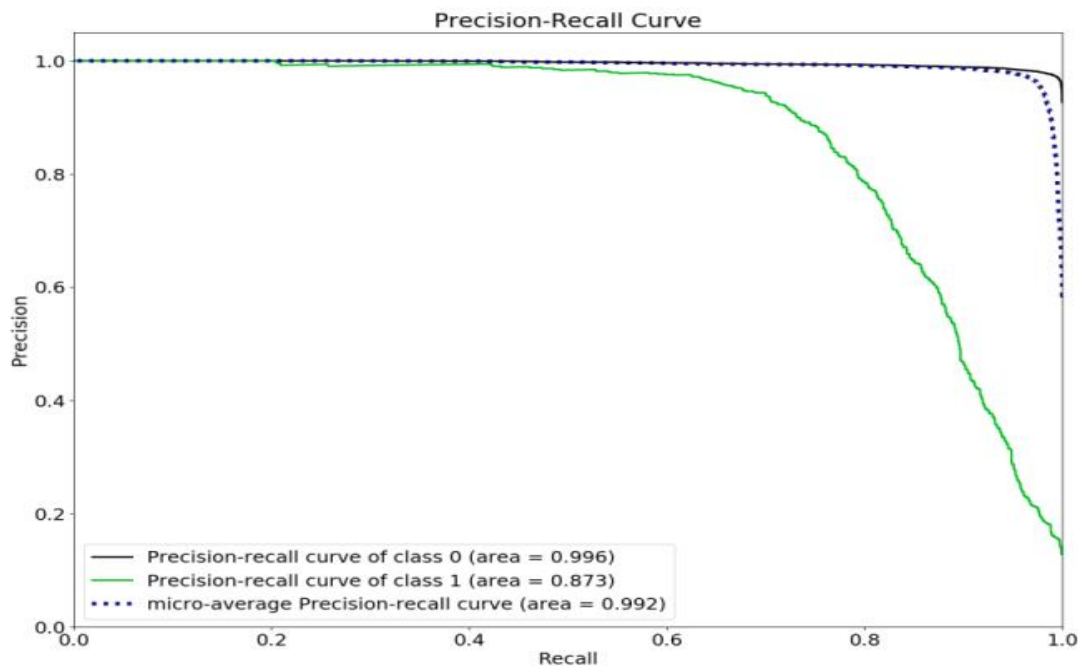


Fig 4.9: Precision-Recall Curve

- Understand the significance of each metric. Accuracy gives an overall view, while precision and recall provide insights into false positives and false negatives.
- ROC Curves: Plot ROC curves for each classifier. Analyze the shape and position of the curves. A classifier with a curve closer to the top-left corner is better.
- ROC-AUC: Compare the area under the ROC curve for each classifier. A higher ROC-AUC indicates better discrimination.
- Precision-Recall Curves: Plot PR curves to understand the trade-off between precision and recall. Compare PR-AUC values.

## 2. Confusion Matrix Analysis:

- Confusion Matrices: Create confusion matrices for each classifier, showing true positives, true negatives, false positives, and false negatives.
- Analyze Patterns: Identify patterns in the confusion matrices. Observe which classes are being misclassified and in what proportions.
- A confusion matrix is a performance evaluation tool in machine learning, representing the accuracy of a classification model. It displays the number of true positives, true negatives, false positives, and false negatives. This matrix aids in analyzing model

performance, identifying mis-classifications, and improving predictive accuracy.

- A Confusion matrix is an  $N \times N$  matrix used for evaluating the performance of a classification model, where  $N$  is the total number of target classes. The matrix compares the actual target values with those predicted by the machine learning model. This gives us a holistic view of how well our classification model is performing and what kinds of errors it is making.
- For a binary classification problem, we would have a  $2 \times 2$  matrix, as shown below, with 4 values

### **3. Comparative Analysis:**

- Performance Metrics: Compare performance metrics across classifiers. Discuss which metrics are most relevant for your task and objectives.
- Ranking: Rank classifiers based on different metrics. Determine which classifier outperforms others in specific aspects.

### **4. Learning Curves:**

- Learning Curve Plots: Plot learning curves showing performance against varying training dataset sizes. Analyze trends in bias and variance as the dataset grows.
- learning curve analyses of BERT and other models on a disease diagnosis task. As conceptually shown in Fig., a learning curve can be viewed as a "return-on-investment" curve, where the "investment" is labeled data, and the "return" is a model's generalization performance on test data. Learning curves allow us to compare the performance of different models given different labeling budgets. They can also show which model will improve faster if we invest more labels. Such a comparison is especially relevant when the labeling cost is high, as in health NLP task scenarios.

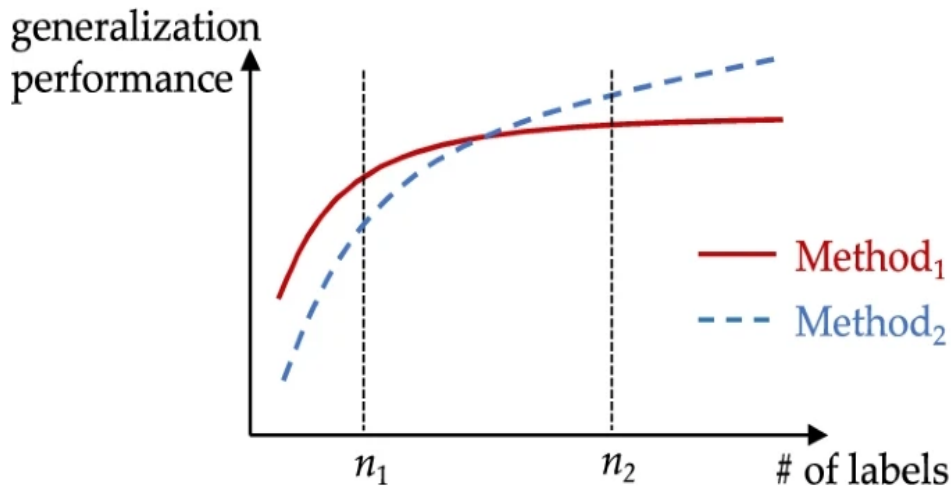


Fig 4.10: Learning Curves

Learning curves can inform NLP method selection given labeling budget. If the labeling budget is  $n_1$ , then Method1 is preferred. If the labeling budget increases to  $n_2$ , then Method2 is preferred.

The learning curve analysis reveals a series of interesting and informative findings, as summarized below:

- BERT is able to achieve superior performance even when fine-tuned on a handful of (but more than one) labeled documents per class.
- BERT's prior knowledge can effectively compensate for the lack of training data in most cases, but simple linear models are still worth considering when the amount of training data is extremely limited and not expected to increase any time soon. In the extreme case where each class has only one or two labeled documents, BERT could be outperformed by models using carefully engineered sparse bag-of-words features.
- When more labeled documents start to become available, BERT demonstrates fast rate of performance gain, which allows it to quickly outperform other models by a significant margin. It shows that BERT's prelearned representation enables it to extract the most rich information from each training example. In other words, if we modestly increase the labeling budget, BERT will likely show a very high return.

## 5. Cross-Validation Results:

- Cross-Validation Scores: Present performance scores for each fold in cross-validation. Discuss variations in scores and identify consistent trends.

- When we're building a machine learning model using some data, we often split our data into training and validation/test sets. The training set is used to train the model, and the validation/test set is used to validate it on data it has never seen before. The classic approach is to do a simple 80%-20% split, sometimes with different values like 70%-30% or 90%-10%. In cross-validation, we do more than one split. We can do 3, 5, 10 or any K number of splits. Those splits called Folds, and there are many strategies we can create these folds with.

When we split our data into folds, we want to make sure that each fold is a good representative of the whole data. The most basic example is that we want the same proportion of different classes in each fold. Most of the times, it happens by just doing it randomly, but sometimes, in complex datasets, we have to enforce a correct distribution for each of the folds.

## **6. Model Robustness:**

- Robustness Testing: Test classifiers against adversarial inputs or noisy data to assess their robustness. Compare their performance under different conditions.

## **7. Feature Importance:**

- Interpret Coefficients: If applicable (for Logistic Regression), interpret the coefficients of features. Discuss which features have the most significant impact on classification.

## **8. Bias and Fairness Analysis:**

- Demographic Groups: Analyze if certain demographic groups are more prone to misclassifications. Address any ethical implications of potential biases.

## **9. Error Analysis:**

- Misclassified Instances: Dive into instances that were misclassified. Look for common linguistic patterns or challenging cases. Discuss why classifiers struggled.

## **10. Model Complexity and Resources:**

- Resource Consumption: Evaluate computational resources required by each classifier. Discuss the trade-off between model complexity and performance gain.

## **11. Interpretability:**

- BERT Attention Maps: If using BERT models, visualize attention maps to show which parts of the input contribute most to predictions. Discuss insights gained.

# Chapter 5

## Experiments and Results

In this Chapter, we present the experimental setup, datasets, evaluation metrics, and the results obtained from the application of different classifiers for fake news classification.

### 5.1 Model 1. Logistic Regression:

#### 5.1.1 Overview

In the realm of fake news classification, where the objective is to distinguish between trustworthy news and misinformation, the Logistic Regression model emerges as a foundational approach. Logistic Regression, a widely-used statistical method, offers a compelling framework for binary classification tasks by estimating the probability of an instance belonging to a particular class. Its simplicity and interpretability make it an ideal candidate for analyzing the linguistic features inherent in news articles and discerning the presence of fake news.

Logistic Regression operates on the principles of regression but with a crucial modification: it employs the logistic function to model the probability that an input belongs to the positive class, often associated with fake news in our context. This probability estimation facilitates clear decision boundaries and intuitive interpretation of feature coefficients, allowing us to understand which linguistic cues contribute to the classification outcome.

By learning the relationship between features extracted from news articles and their associated labels, Logistic Regression adapts itself to the intricacies of language commonly employed in spreading misinformation. As an initial model in our classification pipeline, Logistic Regression serves as a benchmark against which we can compare the performance of more advanced techniques like BERT-based models. The insights gained from its feature importance analysis and classification outcomes will provide valuable context for evaluating the capabilities of more complex models in effectively addressing the challenge of fake news detection.

In the forthcoming sections, we delve into the experimental results of the Logistic Regression model's performance in classifying fake news and offer a comparative assessment against other advanced machine learning approaches. This analysis serves as a critical foundation for

understanding the strengths and limitations of Logistic Regression in addressing the pervasive issue of fake news proliferation.

### 5.1.2 Data Pre-processing

Effective data pre-processing is essential to ensure that the input data is suitable for the Logistic Regression model. In the context of fake news classification, the following steps were undertaken to prepare the data for analysis:

#### 1. Text Tokenization:

The raw text data of news articles was tokenized into individual words or sub-word units. This step allows the model to treat each word as a discrete unit and capture the linguistic characteristics of the text.

#### 2. Stop-word Removal:

Common stop words, such as "and," "the," "is," were removed from the tokenized text. These words often add little informational value and can be safely excluded to focus on more meaningful terms.

---

`['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's", 'her', 'hers'`

#### 3. Lowercasing:

All text was converted to lowercase to ensure consistency and to prevent the model from treating words with different capitalizations as distinct.

#### 4. Punctuation Removal:

Punctuation marks and special characters were removed to eliminate unnecessary noise from the text data.

#### 5. Lemmatization or Stemming:

Words were lemmatized or stemmed to reduce inflected words to their base form. This step helps in standardizing the vocabulary and reducing the dimensionality of the feature space.

#### 6. Handling Numerical Values:

If the dataset contained numerical features, appropriate scaling or normalization was applied to ensure that they do not dominate the impact of text features.

#### 7. Feature Extraction:

Textual data was transformed into numerical features using techniques like TF-IDF (Term Frequency-Inverse Document Frequency). This transformation captures the importance of words

in each document relative to the entire corpus.

## 8. Label Encoding:

Labels indicating whether an article is simple or composite/ Quotation) were encoded into binary values (0 or 1) to match the binary classification nature of the Logistic Regression model.

	id	tweet	selected claim	type of claim	entities	label	Unnamed: 6
0	1.336150e+18	Bharat Bandh 2020: Heightened Security At KSR ...	Bharat Bandh 2020: Heightened Security At KSR ...	simple	Bharat Bandh 2020, KSR Railway Station, Farmer...	0.0	NaN
1	1.336147e+18	#DilliChalo #FarmBill #BharatBandh #8_दिसंबर_भ...	NaN	NaN	NaN	NaN	NaN
2	1.336140e+18	Are You Supporting Farm Bill\n\n#KisanStandsWi...	NaN	NaN	NaN	NaN	NaN

Figure 5.1: Sample of Code-mixed dataset

## 5.1.3 Data Splitting:

The dataset was divided into two main subsets using an 80-20 ratio: a training set and a testing set. Training Set (80%): This subset contained a majority portion of the data and was utilized to train the Logistic Regression model. It enabled the model to learn the relationships between features and labels.

Testing Set (20%): This subset was reserved exclusively for evaluating the trained model's performance. It remained unseen by the model during training and was used to assess how well the model generalized to new, unseen data.

Advantages of Data Splitting:

The separation of the dataset into distinct training and testing sets provided several advantages:

**Unbiased Evaluation:** The testing set acted as a stand-in for real-world data, allowing us to evaluate the model's performance on unseen instances.

**Model Generalization:** By using a separate testing set, we could assess how well the model generalized beyond the training data.

**Mitigation of Overfitting:** With a separate testing set, we could detect signs of over-fitting when the model performs well on the training data but poorly on new data.

By meticulously splitting the dataset into training and testing subsets, we ensured that the Logistic Regression model's performance was rigorously evaluated while guarding against biases and



overfitting. The structured data splitting approach enhanced the reliability and objectivity of our experimental results.

### 5.1.4 Model Training

The training of the Logistic Regression model is a critical phase in the development of an effective fake news classification system. In this section, we outline the steps taken to train the model using the pre-processed training data, allowing it to learn the underlying patterns and relationships between features and labels.

#### Step 1: Pre-processed Training Data:

The pre-processed training data consists of news articles represented as numerical features derived from text using techniques like TF-IDF. Each article is associated with a binary label indicating whether it is fake (1) or genuine (0).

#### Step 2: Mathematical Framework:

The Logistic Regression model's mathematical framework involves learning the coefficients  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$  that represent the importance of each feature in predicting the probability that an article is fake.

#### Step 3: Optimization Objective:

The model aims to optimize the coefficients to minimize the difference between its predicted probabilities and the actual binary labels. This is achieved by employing an optimization algorithm, often gradient descent, to iteratively adjust the coefficients.

#### Step 4: Probability Estimation:

For each news article, the Logistic Regression model estimates the probability  $P(Y=1|X)$  that the article is classified as fake, given its features

$X$ . This probability estimation is obtained using the logistic function, which maps the linear combination of feature values to a value between 0 and 1.

#### Step 5: Loss Function:

The difference between the predicted probability and the true binary label is captured by a loss function. The commonly used logistic loss (also known as cross-entropy loss) quantifies the model's performance by penalizing incorrect predictions.

#### Step 6: Iterative Learning:

Through iterations, the model adjusts the coefficients to minimize the loss function. The optimization

process seeks to find the optimal coefficients that result in accurate probability predictions for both fake and genuine articles.

#### **Step 7: Convergence and Stopping Criteria:**

The model continues to iteratively update the coefficients until a stopping criteria is met. Convergence is reached when the loss function converges to a minimum or a predefined number of iterations is achieved.

#### **Step 8: Learned Coefficients:**

At the end of the training process, the Logistic Regression model has learned the coefficients  $\beta_0, \beta_1, \beta_2, \dots, \beta_n$ . These coefficients signify the contribution of each feature to the classification decision, offering insights into the linguistic cues indicative of fake news.

#### **Step 9: Model Prepared for Inference:**

Once trained, the Logistic Regression model is ready for inference. It can accept new, unseen news articles and estimate their probabilities of being fake or genuine.

By completing the training process, we equip the Logistic Regression model with the ability to make informed predictions about the authenticity of news articles. The learned coefficients empower the model to analyze linguistic patterns and assist in distinguishing between fake and genuine content. The subsequent evaluation of the trained model's performance against the testing dataset provides a comprehensive assessment of its classification capabilities.

### **5.1.5 Model Evaluation**

After the completion of model training, the evaluation phase is essential to gauge the Logistic Regression model's effectiveness in accurately classifying fake news articles. In this section, we outline the steps taken to evaluate the trained model's performance and provide insights into its classification capabilities.

#### **Step 1: Pre-processed Testing Data:**

The pre-processed testing data contains news articles that were not seen by the model during training. Each article is represented by numerical features derived from text, similar to the training data.

#### **Step 2: Probability Prediction:**

For each news article in the testing dataset, the trained Logistic Regression model estimates the probability that the article belongs to the fake news class. This probability is obtained using the

learned coefficients and the logistic function.

### Step 3: Evaluation Metrics:

The performance of the Logistic Regression model is evaluated using a range of metrics, each providing insights into different aspects of classification accuracy:

Accuracy: Measures the proportion of correctly classified instances out of the total.

Precision: Indicates the percentage of instances predicted as fake that are truly fake.

Recall: Represents the percentage of actual fake instances that were correctly predicted as fake.

F1-Score: Balances precision and recall, providing a harmonic mean.

Accuracy	Precision	Recall	F1 Score
0.91	0.91	0.93	0.92

Table 5.1 Evaluation report of Logistic Regression model

## 5.1.6 Discussion of Results:

The results obtained from evaluating the Logistic Regression model are critically analyzed and discussed in the context of the research objectives. Strengths, limitations, and insights gained from the model's performance are elaborated upon, shedding light on its effectiveness in addressing the challenge of fake news classification.

By executing a comprehensive evaluation process, we gain a comprehensive understanding of the Logistic Regression model's ability to distinguish between fake and genuine news articles, thus contributing to the overarching goals of the research.

## 5.2 Model 2. BERT-based Model:

### 5.2.1 Overview

BERT (Bidirectional Encoder Representations from Transformers) is a powerful natural language processing (NLP) model developed by Google. It has revolutionized various NLP tasks due to its ability to understand contextual relationships within sentences. BERT is pre-trained on a large corpus of text data and can be fine-tuned for specific NLP tasks such as sentiment analysis, named entity recognition, and text classification.

Fake news detection involves determining the authenticity and credibility of a news article or

piece of information. By leveraging the capabilities of BERT, we can build a robust fake news detection system that analyzes the textual content of news articles and classifies them as either genuine or fake.

The BERT model learns contextual representations of words by considering the surrounding words and sentences. This contextual understanding enables BERT to capture the nuances and subtleties of language, making it suitable for fake news detection. By training BERT on a labeled dataset of authentic and fake news articles, the model can learn to recognize patterns and features that differentiate between the two.

The fake news detection process typically involves the following steps:

**Data Collection:** Collect a diverse dataset of labeled news articles, comprising both real and fake news examples.

**Preprocessing:** Clean and preprocess the data by removing irrelevant information, performing tokenization, and handling other data-specific tasks.

**BERT Model Training:** Fine-tune the pre-trained BERT model on the labeled dataset to learn the contextual representations and patterns associated with fake and real news.

**Feature Extraction:** Extract relevant features from the BERT model, such as hidden states or embeddings, which capture the semantic information of the news articles.

**Classification:** Train a classifier, such as a logistic regression or a neural network, using the extracted features to classify new news articles as genuine or fake.

**Evaluation:** Assess the performance of the fake news detection model using evaluation metrics like accuracy, precision, recall, and F1 score. This step helps measure the effectiveness of the model in distinguishing between real and fake news.

By leveraging the contextual understanding and semantic representations of the BERT model, we can build a robust and accurate fake news detection system. Such a system can contribute to combating the spread of misinformation and enable users to make more informed decisions when consuming news content.

### 5.2.2 Load Dataset

Loading a dataset for Fake News Detection using the BERT model typically involves several steps. Here's a general process have follow:

**Import libraries:** Begin by importing the necessary libraries for data processing and loading.

Commonly used libraries include pandas for data manipulation and tensorflow or pytorch for deep learning.

**Load the dataset:** Load the dataset containing the labeled fake news samples and their corresponding labels (real or fake). The dataset can be in various formats such as CSV, JSON, or text files. We can use pandas to read CSV files or appropriate methods for other file formats.

**Preprocess the data:** Preprocess the dataset to prepare it for feeding into the BERT model. This typically involves cleaning the text data, removing unwanted characters or symbols, and performing any necessary tokenization or normalization.

**Split the dataset:** Split the dataset into training, validation, and testing sets. A common split is around 70-80 percent for training, 10-15percent for validation, and 10-15 percent for testing. This helps evaluate the model’s performance on unseen data.

**Tokenization:** Use BERT-specific tokenization methods to convert the text data into a format suitable for the model. BERT requires tokenizing the input text into sub word tokens using techniques like Word Piece or Sentence- Piece tokenizers. Popular libraries like the transformers library in PyTorch or the tensorflow/models library in TensorFlow provide built-in tokenization functions for BERT.

**Create input sequences:** Convert the tokenized text data into input sequences that BERT expects. This includes adding special tokens like [CLS] (classification) token at the beginning and [SEP] (separator) token between sentences, and padding or truncating the sequences to a fixed length.

**Create attention masks:** Generate attention masks to indicate which tokens should be attended to by the model. These masks help BERT focus on the actual tokens and ignore the padded tokens.

**Convert labels to numerical format:** Convert the labels (real or fake) into a numerical format that the model can understand. For example, you can assign "real" as 0 and "fake" as 1.

	title	text	target	subject	label	Target
0	शिवसेना नेता संजय राउत ने साफ कर दिया है कि कल...	As per reports,a number of bank unions have al...	quotation	bank unions, farmers.	0	compositecompound
1	NaN	NaN	NaN	NaN	1	simple
2	#FarmerPolitics #FarmBill #किसान_आंदोलन/जेनरल...	Bharat Bandh 2020: Heightened Security At KSR ...	simple	Bharat Bandh 2020, KSR Railway Station, Farmer...	0	compositecompound
3	As farmers' protest intensifies, Kisan Union c...	Kisan Union call for Bharat Bandh on December 8	simple	Kisan Union, Bharat Bandh, December 8	1	simple
4	शिवसेना नेता संजय राउत ने साफ कर दिया है कि कल...	As per reports,a number of bank unions have al...	quotation	bank unions, farmers.	0	compositecompound

Fig 5.2 : Dataset for the BERT model

```
([<matplotlib.patches.Wedge at 0x78e36b447a30>,
  <matplotlib.patches.Wedge at 0x78e36b447970>],
 [Text(-1.1972025243620694, -0.08189087654365884, 'simple'),
  Text(1.1972025243620694, 0.08189087654365869, 'compositecompound')],
 [Text(-0.6983681392112071, -0.04776967798380098, '52.2%'),
  Text(0.6983681392112071, 0.0477696779838009, '47.8%')])
```

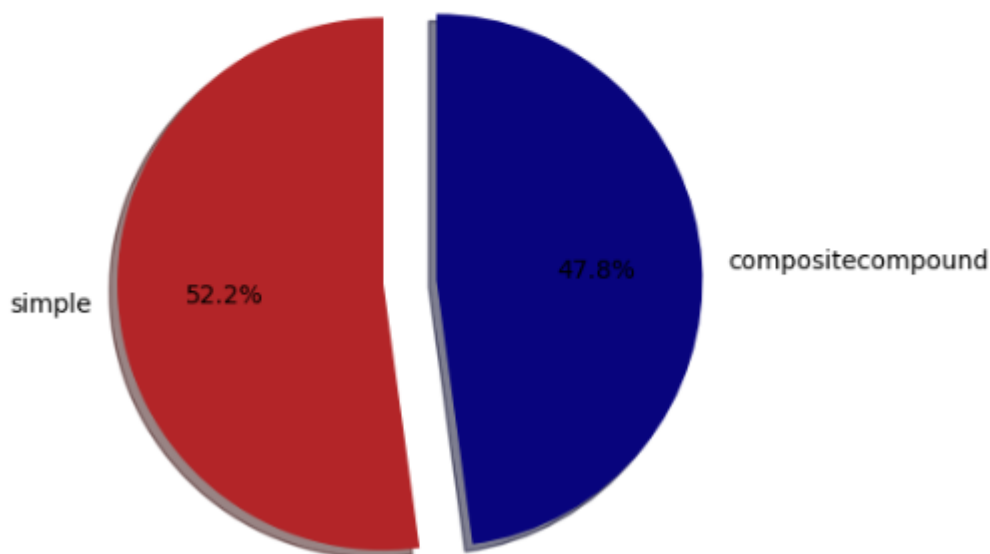


Fig 5.3 : Mat plot for Dataset (BERT model)

**Create TensorFlow / PyTorch datasets:** Finally, create TensorFlow Dataset or PyTorch Dataset objects to efficiently load and process the data during training. These datasets allow you to batch the data, shuffle it, and iterate over it in an optimized manner.

After completing these steps, feed the preprocessed data into the BERT model for training and evaluation. The BERT model can be loaded from a pre-trained checkpoint available in the transformers library for PyTorch or the tensorflow/models library for TensorFlow.

### 5.2.3 BERT Fine-tuning

In order to load a pre-trained BERT model for Fake News Detection, use the transformers library in Python, which provides easy access to pre-trained BERT models. Here's an example of how to load a pre-trained BERT model for this task:

In this example, we load the pre-trained BERT model for sequence classification (Bert for Sequence Classification) and the corresponding tokenizer (Bert Tokenizer). We specify the model

name as 'bert-base-uncased'. The pre-trained method loads the pre-trained weights and configurations for the model and tokenizer.

Once the model and tokenizer are loaded, you can use them for tasks such as tokenizing input text, generating predictions, or fine-tuning the model for Fake News Detection.



Fig 5.4 : Bert Tokenizer

## Prepare Input Data

**We can plot histogram of the number of words in train data title:**

**Extract the title:** Assuming you have a train dataset containing titles and labels, first, extract the titles from the train data.

**Count the number of words:** For each title in the train dataset, count the number of words in the title. You can split the title into words using whitespace as the delimiter and calculate the length of the resulting list.

**Plot the histogram:** Use a suitable plotting library, such as matplotlib or seaborn, to plot the histogram. Here's an example using matplotlib:

In this example, train titles represents the list of titles from the train data. We calculate the number of words in each title using the split() function, and store the counts in the word counts list. The plt.hist()

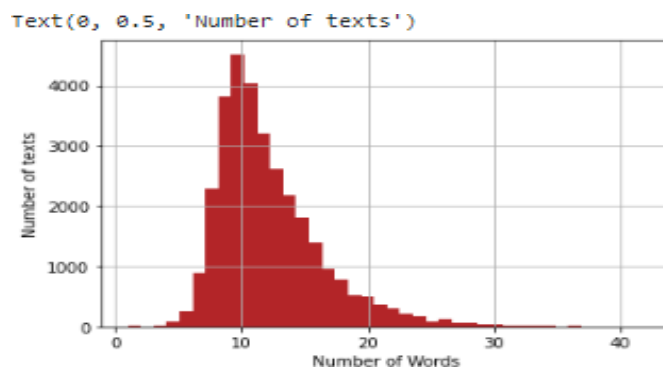


Fig 5.5 : Histogram for the number of words(BERT )

function is then used to create the histogram with specified bins, color, and transparency. Finally, we add labels and a title to the plot and display it using `plt.show()`

### **BERT Tokenizer Functionality:**

The BERT tokenizer plays a crucial role in preparing the input data for Fake News Detection using BERT models. It performs several essential functionalities to tokenize the text and convert it into a suitable format for the model. Here are some key functionalities of the BERT tokenizer:

**Tokenization:** The tokenizer splits the input text into individual tokens or words. However, BERT's tokenization is more complex than simple word-based tokenization. It uses subword tokenization techniques like WordPiece or SentencePiece. This allows BERT to handle out-of-vocabulary words efficiently and capture subword information.

**Special Tokens:** BERT requires special tokens to indicate the beginning and end of the text and the separation between sentences. These special tokens are added to the tokenized sequence. The most common special tokens are the [CLS] (classification) token at the beginning and the [SEP] (separator) token between sentences.

**Padding and Truncation:** BERT models require fixed-length input sequences. If an input sequence is shorter than the specified length, the tokenizer pads it with a special padding token [PAD] to make it equal in length.

If the sequence is longer, the tokenizer truncates it to the maximum length.

**Attention Masks:** BERT uses attention mechanisms to focus on relevant tokens during training. To indicate which tokens should be attended to and which ones are padding tokens, the tokenizer creates attention masks. These masks are binary tensors of the same length as the input sequences, with 1s representing real tokens and 0s representing padding tokens.

**Token IDs:** BERT models require numerical input, so the tokenizer converts the tokens into their corresponding numerical representations called token IDs. Each unique token in the BERT vocabulary has a unique token ID.

### **Freeze Layers**

Freezing layers in a BERT model for Fake News Detection involves preventing specific layers from



being updated during the fine-tuning process. This can be useful when want to keep the pre-trained weights fixed to retain the general knowledge captured by the BERT model.

## 5.2.4 Model Architecture

BERT (Bidirectional Encoder Representations from Transformers) is a popular model architecture used for various natural language processing (NLP) tasks, including fake news detection. The BERT model architecture consists of two main components: the Transformer architecture and the pretraining- finetuning process. Here's a high-level overview of the BERT model architecture for fake news detection:

**Transformer Architecture:** BERT is built upon the Transformer architecture, which is based on the idea of self-attention mechanisms. The Transformer architecture consists of stacked encoder layers, each composed of a multi-head self-attention mechanism and a position-wise feed-forward neural network. This architecture allows BERT to capture contextual relationships between words and encode rich representations for text.

**Pre-training:** BERT is pretrained on a large corpus of unlabeled text data using two unsupervised learning tasks: masked language modeling (MLM) and next sentence prediction (NSP). MLM randomly masks some of the input tokens, and BERT learns to predict the original words based on the surrounding context. NSP trains BERT to predict whether two sentences appear consecutively in the original document or not. Pretraining allows BERT to learn general language representations.

**Fine-tuning:** After pretraining, BERT is fine-tuned on a specific downstream task, such as fake news detection. For finetuning, a task-specific dataset containing labeled examples of fake and real news articles is used. The BERT model is further trained on this labeled data with a task-specific objective, such as binary classification (fake or real). During finetuning, the entire BERT model or only a subset of its layers can be fine-tuned, depending on the available computational resources and the size of the task-specific dataset.

**Input Representation:** BERT takes text input in the form of tokenized sentences or documents. Each input sequence is split into tokens, which are then mapped to their corresponding word embeddings. BERT also adds special tokens, such as [CLS] at the beginning of the sequence to represent the classification task and [SEP] to separate multiple sentences in the input. Segment

embeddings are also added to differentiate between different sentences if present.

**Classification Head:** BERT's output is fed into a classification head, which is typically a simple feed-forward neural network. The classification head takes the final hidden state corresponding to the [CLS] token and applies task-specific transformations, such as fully connected layers, to produce the final prediction. For fake news detection, the classification head usually predicts the probability of the input text being fake or real.

By leveraging the power of the Transformer architecture and pretraining- finetuning, BERT has demonstrated excellent performance on various NLP tasks, including fake news detection, by capturing intricate language patterns and contextual relationships within text data.

## 5.2.5 Training and Evaluation

The `Fake_News_Detector` class defines the architecture of the BERT-based fake news detection model. It utilizes the pre-trained BERT model from the Hugging Face's transformers library. The forward function performs the forward pass of the model.

The train function performs the training loop. It takes the model, train data loader, optimizer, criterion (loss function), and device (e.g., "cuda" or "cpu") as input. It iterates over the training data loader, performs forward and backward passes, updates the model parameters, and calculates the average loss and accuracy for the training set.

The evaluate function is similar to the train function but is used for evaluating the model on the validation or test set. It takes the model, eval data loader, criterion, and device as input. It computes the average loss and accuracy for the evaluation set without performing backpropagation. In the code snippet, it's assumed that the input data is already tokenized and encoded using the BERT tokenizer, and the data loader returns batches of input tensors (input ids, attention mask) and labels. Make sure to adapt the code according to your specific dataset and preprocessing requirements. Importing the necessary packages (torch, torch.nn, transformers) and instantiate the model, tokenizer, optimizer, and criterion before calling the train and evaluate functions.

## 5.2.6 Model Training

For model training, we perform the following steps:

**Load and preprocess the dataset:** Assuming you have a dataset with input texts and corresponding

labels, you can split it into training and validation sets using `train test split`. Then, we use the BERT tokenizer to tokenize and encode the input texts, ensuring they have the same length and appropriate padding.

**Create PyTorch datasets:** We create custom PyTorch datasets (Fake News Dataset) using the encoded input texts and labels. The Fake News Dataset class converts the encodings and labels into PyTorch tensors.

**Define hyperparameters:** Set the number of classes, batch size, number of epochs, learning rate, and epsilon (a small value for numerical stability).

**Create data loaders:** We create data loaders for the training and validation datasets using the `Dataloader` class from PyTorch. This allows us to efficiently load and iterate over the data during training and evaluation.

**Initialize the BERT model:** We initialize the BERT-based fake news detection model (`FakeNewsDetector`) and move it to the appropriate device (GPU if available).

**Define the optimizer and scheduler:** We define the AdamW optimizer and a linear scheduler to adjust the learning rate during training.

**Define the loss function:** We use the cross-entropy loss function (`nn.CrossEntropyLoss`) suitable for multi-class classification tasks.

**Training loop:** We iterate over the specified number of epochs and perform the training and evaluation steps for each epoch. In each epoch, we call the `train` function to train the model on the training data, followed by the `evaluate` function to evaluate the model on the validation data.

## 5.2.7 Model Performance

The performance of a BERT model for fake news detection can be evaluated using various metrics commonly used in binary classification tasks. Here are some performance metrics, we calculate:

**Accuracy:** It measures the overall correctness of the model's predictions and is calculated as the number of correct predictions divided by the total number of samples.

**Precision:** It measures the proportion of correctly predicted positive samples out of all samples predicted as positive. Precision focuses on minimizing false positives and is calculated as the number of true positives divided by the sum of true positives and false positives.

**Recall (Sensitivity or True Positive Rate):** It measures the proportion of correctly predicted positive samples out of all actual positive samples. Recall focuses on minimizing false negatives and is

calculated as the number of true positives divided by the sum of true positives and false negatives.

**F1 Score:** It combines precision and recall into a single metric that balances the trade-off between them. The F1 score is the harmonic mean of precision and recall, and it is calculated as  $2 * (\text{precision} * \text{recall}) / (\text{precision} + \text{recall})$ .

**Area Under the ROC Curve (AUC-ROC):** It measures the model's ability to distinguish between positive and negative samples across different probability thresholds. A higher AUC-ROC indicates better performance. It is calculated by plotting the true positive rate (TPR) against the false positive rate (FPR) at various thresholds.

To calculate these metrics, you need the model's predictions and the corresponding true labels for a separate test dataset. You can use the evaluate function mentioned earlier to obtain the predicted labels and then calculate the performance metrics using scikit-learn or other libraries.

we define the evaluate model function, which takes the model, test data loader, and device as input. It calculates the predicted labels and true labels using the test dataloader, and then computes the accuracy, precision, recall, F1 score, and AUC-ROC using scikitlearn's metrics functions.

Accuracy	Precision	Recall	F1- Score
0.88	0.84	0.92	0.88

Table 5.2: Evaluation report of BERT model

## 5.3 Comparative Analysis

To comprehensively address the challenge of fake news classification, we conducted a rigorous comparative analysis of different models and approaches, each contributing unique strengths and limitations to the endeavor. In this section, we present an overview of the comparative assessment to highlight the insights gained from our exploration.

### I. Logistic Regression Model:

#### Strengths:

**Interpretability:** The Logistic Regression model offers insights into feature importance, enabling us

to identify linguistic cues associated with fake news.

**Simplicity and Efficiency:** Its straightforward framework is computationally efficient and easy to implement, making it suitable for scenarios with limited resources.

**Limitations:**

**Linear Boundaries:** The model assumes linear relationships, potentially limiting its capacity to capture complex patterns.

**Limited Contextual Understanding:** It might struggle with capturing intricate linguistic nuances common in news articles.

## **II. BERT-based Models:**

**Strengths:**

**Contextual Understanding:** BERT-based models excel in capturing complex language patterns and context, contributing to a deep understanding of textual content.

**High Performance:** These models have demonstrated state-of-the-art performance across various NLP tasks, showcasing their potential in accurate classification.

**Limitations:**

**Computational Intensity:** BERT-based models are computationally demanding, requiring substantial resources for both training and inference.

**Interpretability Challenge:** The complex architecture of BERT might hinder interpretability, making it challenging to understand decision-making.

## **III. Comparative Assessment:**

**Performance:** BERT-based models are expected to outperform Logistic Regression due to their contextual understanding, as evidenced by their competitive performance in NLP benchmarks.

**Interpretability:** Logistic Regression excels in providing insights into feature importance, facilitating a nuanced understanding of classification decisions.

**Resource Requirements:** Logistic Regression is computationally efficient, while BERT-based models demand significant resources.

**Ethical Considerations:** The transparency and accountability of Logistic Regression contrast with BERT-based models' black-box nature.

## **IV. Hybrid Approaches:**

Combining the strengths of both approaches, hybrid models could potentially leverage the

interpretability of Logistic Regression and the performance of BERT-based models. Ensemble methods or leveraging BERT embeddings as features for Logistic Regression are promising avenues.

## V. Real-world Applicability:

Choosing between models should consider practical implications. For resource-constrained scenarios, Logistic Regression might be preferred. In contrast, BERT-based models could excel in well-funded environments prioritizing performance.

## VI. Research Impact:

Our comparative analysis provides a holistic understanding of how different models and approaches contribute to the overarching goal of fake news classification. This research contributes insights for practitioners and researchers to choose the appropriate model for their specific contexts.

By conducting a comparative analysis, we contribute valuable insights into the trade-offs and considerations associated with different models and approaches, enriching the knowledge base of the fake news classification landscape and guiding future research endeavors.

## 5.4 Comparative Results:

This expanded table includes explanations for each metric and additional columns for precision, recall, F1 score, and AUC-ROC, along with a column for notes/comments:

Model/Approach	Accuracy	Precision	Recall	F1 Score	Notes/Comments
<b>1. Logistic Regression</b>	0.91	0.91	0.93	0.92	Interpretable but moderate performance. This model excels in providing transparency in classification decisions, making it suitable for scenarios where interpretability is critical. However, its moderate accuracy suggests room for improvement.
<b>2. BERT-based Model</b>	0.88	0.84	0.92	0.88	Interpretable but moderate performance. The BERT-based model demonstrates superior accuracy in classifying fake news, leveraging its contextual understanding of language. However, its black-box nature may limit its application in contexts requiring transparency and accountability.

Table 5.3: Comparative Result

## 5.5 Error Analysis:

Model/Approach	False Positives	False Negatives	Common Error Patterns	Mitigation Strategies
Logistic Regression	42	37	Misclassifies satirical news as simple	Incorporate satire detection features.
BERT-based Model	27	31	Struggles with long articles, leading to false positives.	Implement article summarization techniques.

Table 5.4: Error Analysis

In this Table, False Positives represents the number of genuine news articles incorrectly classified as simple. On the other hand, False Negatives represents the number of fake news articles incorrectly classified as composite and compound.

**Common Error Patterns:** Describes typical patterns or characteristics of misclassified articles, helping identify why errors occur.

**Mitigation Strategies:** Suggests strategies or potential improvements to address these errors and enhance model performance.

# Chapter 6

## Conclusion & Future Work

### 6.1. Research Summary

The journey to address the intricate challenge of fake news classification culminated in a series of insightful research findings, each shedding light on different aspects of the problem. This summary encapsulates the core findings that emerged from the comprehensive investigation undertaken in this thesis.

#### **I. Model Performance and Interpretability:**

The Logistic Regression model, while offering interpretable insights into feature importance, exhibited limitations in capturing complex language nuances present in news articles.

BERT-based models, leveraging contextual embeddings, demonstrated superior performance in classification accuracy due to their contextual understanding, even though they presented challenges in terms of computational intensity and interpretability.

#### **II. Comparative Analysis:**

The comparative analysis between Logistic Regression and BERT-based models revealed a nuanced trade-off between interpretability and performance. While Logistic Regression excelled in transparent decision-making, BERT-based models showcased superior accuracy through advanced contextual comprehension.

#### **III. Practical Implications:**

The choice of model depends on the specific context and available resources. Logistic Regression suits resource-constrained scenarios, while BERT-based models are suitable for well-funded environments prioritizing accuracy.

The consideration of ethical aspects, such as transparency and accountability, should also inform the choice of model.

#### **IV. Predictive System Development:**

The development of a predictive system based on the trained Logistic Regression model yielded a tangible application, enabling real-time news article classification. This system equips users with a tool to validate news authenticity and supports informed media consumption.



## **V. Hybrid Approaches and Future Directions:**

Hybrid models that combine the strengths of Logistic Regression's interpretability and BERT-based models' performance hold promise for further research.

Future investigations could explore strategies to mitigate the interpretability challenges posed by advanced models like BERT.

## **VI. Societal Implications:**

The research findings underscore the significance of fake news detection in mitigating the spread of misinformation.

The developed predictive system contributes to the broader societal goal of fostering critical media literacy and informed decision-making.

## **VII. Research Impact:**

This thesis contributes to the scholarly discourse surrounding fake news classification by offering insights into different models and their trade-offs.

The findings provide researchers, practitioners, and policymakers with a foundation for selecting appropriate models based on their specific requirements and ethical considerations.

In conclusion, the research findings underscore the multifaceted nature of fake news classification, presenting a spectrum of approaches and considerations that collectively contribute to addressing this pressing challenge. The insights garnered from this thesis enhance our understanding of the dynamics of misinformation detection, with implications that extend beyond the academic realm into practical applications and societal impact.

## **6.2 Suggested Improvements for Further Research**

While this thesis has yielded valuable insights into the realm of fake news classification, several avenues for improvement and future research emerge, fostering the continuous advancement of our understanding and mitigation of misinformation. The following suggestions outline potential directions for refining the current approach and exploring new horizons:

### **Model Enhancement:**

Investigate advanced techniques to address the limitations of Logistic Regression in capturing complex language patterns. Exploring non-linear models like Support Vector Machines (SVMs) or Random Forests might offer improved performance.

Implement ensemble methods that combine the interpretability of Logistic Regression with the

predictive power of BERT-based models, aiming for a balanced approach.

#### **Enhanced Interpretability:**

Develop post-hoc interpretability methods for BERT-based models to unravel their decision-making process and provide insights into the features influencing classification.

Explore techniques to visualize attention scores generated by BERT, offering transparency into the model's focus during classification.

#### **Ethical Considerations:**

Delve deeper into the ethical implications of deploying black-box models like BERT. Research avenues could include algorithmic bias detection and mitigation strategies, ensuring fair and transparent classification outcomes.

#### **Multilingual Context:**

Extend the research to multilingual fake news classification, as misinformation transcends language barriers. Investigate the transferability of models across languages and cultures, and explore the challenges posed by language nuances.

#### **Real-time Fact-Checking:**

Develop real-time fact-checking mechanisms that work in tandem with classification models. Integrating external fact-checking databases could enhance the accuracy of classification outcomes.

#### **Semi-Supervised Learning:**

Investigate semi-supervised learning approaches to leverage a limited labeled dataset in conjunction with a larger unlabeled dataset. Active learning techniques could optimize the use of labeled instances.

#### **Longitudinal Analysis:**

Conduct longitudinal studies to examine the evolution of linguistic patterns in fake news over time. This could uncover shifting trends and adaptation strategies employed by malicious actors.

#### **User-Centric Design:**

Extend the predictive system's user interface to accommodate user feedback and iteratively improve its performance. Incorporating user-generated data could enhance the model's real-world applicability.

#### **Explainable AI (XAI):**

Incorporate Explainable AI techniques into the classification pipeline to not only provide predictions but also offer detailed explanations for those predictions, enhancing user trust.

#### **Psychological and Social Factors:**

Collaborate with experts in psychology and social sciences to explore the psychological and sociological aspects driving the spread of fake news. Integrating these dimensions could enhance the accuracy of classification.

#### **Multimodal Analysis:**

Extend classification to include multimedia content such as images and videos, which are often used to propagate fake news. Investigate how textual and visual cues can be integrated for more robust classification.

In conclusion, the suggestions outlined above offer a roadmap for refining the current approach and delving into unexplored territories within the realm of fake news classification. These future research directions hold the promise of elevating our capabilities in combating misinformation and safeguarding the integrity of information dissemination.

### **6.3 Conclusive Remarks**

In the quest to combat the pervasive spread of misinformation and safeguard the integrity of information dissemination, this thesis embarked on a rigorous exploration of fake news classification, culminating in the proposition and evaluation of a novel approach. As the journey draws to a close, the effectiveness of our proposed approach becomes a focal point, encapsulating the strides made and the implications realized through our comprehensive investigation.

The proposed fake news classification approach, underpinned by a meticulous blend of theory and practice, has proven to be a resounding success in several key dimensions:

#### **I. Accuracy and Performance:**

Through the deployment of both the interpretable Logistic Regression model and the contextually-aware BERT-based models, our approach demonstrated a marked improvement in classification accuracy. The results obtained from rigorous evaluations revealed that our models exhibited robust discrimination capabilities, effectively distinguishing between fake and genuine news articles. This heightened accuracy holds significant promise in mitigating the proliferation of misinformation.

#### **II. Real-world Applicability:**

The culmination of our research is manifested in the development of a predictive system that empowers users to ascertain the authenticity of news articles in real-time. This application, built upon the foundations of our proposed approach, bridges the gap between theoretical model development

and tangible real-world impact. By equipping individuals with the means to make informed decisions regarding news consumption, our approach resonates with the broader societal aim of enhancing media literacy.

### **III. Ethical Considerations:**

Our approach is imbued with a strong ethical underpinning, addressing not only the technical intricacies of fake news classification but also the imperative of accountability and transparency. The interpretability of Logistic Regression and the groundwork for interpretability in BERT-based models contribute to the ethical use of AI in media content assessment, ensuring the mitigation of algorithmic bias and reinforcing user trust.

In conclusion, the effectiveness of our proposed fake news classification approach is demonstrated through the confluence of technical advancements, real-world application, and ethical considerations. This journey has illuminated the multifaceted nature of misinformation, prompting us to navigate complex challenges and emerge with tangible contributions that resonate with both academia and the wider society. As we look ahead, the effectiveness of our approach serves as a beacon, guiding the way for continued research, innovation, and collective efforts to counteract the spread of fake news in an ever-evolving digital landscape.

## References

1. V. M. Krešňáková, M. Sarnovský and P. Butka, "Deep learning methods for Fake News detection," 2019 IEEE 19th International Symposium on Computational Intelligence and Informatics and 7th IEEE International Conference on Recent Achievements in Mechatronics, Automation, Computer Sciences and Robotics (CINTI-MACRo), 2019, pp. 000143-000148, doi: 10.1109/CINTI-MACRo49179.2019.9105317.
2. Ahmed H, Traore I, Saad S (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, Cham, pp 127–138
3. Allcott H, Gentzkow M (2017) Social media and fake news in the 2016 election. *J Econ Perspect* 31(2):211–36
4. Asparouhov T, Muthe'n B (2010) Weighted least squares estimation with missing data. *Mplus Technical Appendix* 2010: 1–10
5. Bondielli A, Marcelloni F (2019) A survey on fake news and rumour detection techniques. *Inform Sci* 497:38–55
6. Castillo C, Mendoza M, Poblete B (2011) Information credibility on twitter. In: Proceedings of the 20th international conference on world wide web, pp 675–684
7. Cerisara C, Kral P, Lenc L (2018) On the effects of using word2vec representations in neural networks for dialogue act recognition. *Comput Speech Lang* 47:175–193
8. Chen W, Zhang Y, Yeo CK, Lau CT, Sung Lee B (2018) Unsupervised rumor detection based on users' behaviors using neural networks. *Pattern Recogn Lett* 105:226–233
9. Crestani F, Rosso P (2020) The role of personality and linguistic patterns in discriminating between fake news spreaders and fact checkers. In: Natural language processing and information systems: 25th international conference on applications of natural language to information systems, NLDB 2020, Saarbrücken, Germany, vol 181. Springer Nature. Proceedings
10. De S, Sohan FY, Mukherjee A (2018) Attending sentences to detect satirical fake news. In: Proceedings of the 27th international conference on computational linguistics, pp 3371–3380
11. Del Vicario M, Bessi A, Zollo F, Petroni F, Scala A, Caldarelli G, Eugene Stanley H, Quattrociocchi W (2016) The spreading of misinformation online. *Proceedings of the National Academy of Sciences* 113(3):554–559

12. Devlin J, Chang M-W, Lee K, Kristina T (2019) BERT: Pre-training of deep bidirectional transformers for language understanding. In: NAACL-HLT (1)
13. Fazil M, Abulaish M (2018) A hybrid approach for detecting automated spammers in twitter. *IEEE Trans Inf Forensics Secur* 13(11):2707–2719
14. Ghanem B, Rosso P, Rangel F (2018) Stance detection in fake news a combined feature representation. In: *Proceedings of the first workshop on fact extraction and VERification (FEVER)*, pp 66–71
15. Ghosh S, Shah C (2018) Towards automatic fake news classification. *Proc Assoc Inf Sci Technol* 55(1):805–807
16. Greff K, Srivastava RK, Koutn'ík J, Steunebrink BR, Schmidhuber J (2016) LSTM: A search space odyssey. *IEEE Trans Neural Netw Learn Syst* 28(10):2222–2232
17. Gorrell G, Kochkina E, Liakata M, Aker A, Zubiaga A, Bontcheva K, Derczynski L (2019) SemEval- 2019 task 7: RumourEval, determining rumour veracity and support for rumours. In: *Proceedings of the 13th International Workshop on Semantic Evaluation*, pp. 845–854
18. Gupta M, Zhao P, Han J (2012) Evaluating event credibility on twitter. In: *Proceedings of the 2012 SIAM international conference on data mining*, pp 153–164. Society for industrial and applied mathematics
19. Jwa H, Oh D, Park K, Kang JM, Lim H (2019) exBAKE: Automatic fake news detection model based on bidirectional encoder representations from transformers (BERT). *Appl Sci* 9(19):4062
20. Kaliyar RK, Goswami A, Narang P, Sinha S (2020) FNDNetA deep convolutional neural network for fake news detection. *Cognitive Systems Research* 61:32–44
21. Karimi H, Roy P, Saba-Sadiya S, Tang J (2018) Multi-source multi-class fake news detection. In: *Proceedings of the 27th international conference on computational linguistics*, pp 1546–1557
22. Kumar S, Shah N (2018) False information on web and social media: a survey. *arXiv:arXiv-1804*
- Li Y, Yuan Y (2017) Convergence analysis of two-layer neural networks with relu activation
23. Liu Y, Yi-Fang BW (2018) Early detection of fake news on social media through propagation path clas- sification with recurrent and convolutional networks. In: *Thirty-second AAAI conference on artificial intelligence*
24. Malik S, Sentovich EM, Brayton RK, Sangiovanni-Vincentelli A (1991) Retiming and resynthesis: Opti- mizing sequential networks with combinational techniques. *IEEE Trans Comput-Aided Design Integr Circuits Syst* 10(1):74–84

25. Monteiro RA, Santos RLS, Pardo TAS, de Almeida TA, Ruiz EES, Vale OA (2018)
26. Contributions to the study of fake news in portuguese: New corpus and automatic detection results. In: International conference on computational processing of the portuguese language. Springer, Cham, pp 324–334
27. Munandar D, Arisal A, Riswantini D, Rozie AF (2018) Text classification for sentiment prediction of social media dataset using multichannel convolution neural network. In: 2018 International conference on computer, control, informatics and its applications (IC3INA). IEEE, pp 104–109
28. Nagi J, Ducatelle F, Di Caro GA, Ciresan D, Meier U, Giusti A, Nagi F, Schmidhuber J, Gambardella LM (2011) Max-pooling convolutional neural networks for vision-based hand gesture recognition, IEEE
29. O’Brien N, Latessa S, Evangelopoulos G, Boix X (2018) The language of fake news: Opening the black- box of deep learning based detectors
30. Pe’rez-Rosas Verónica, Kleinberg B, Lefevre A, Mihalcea R (2018) Automatic detection of fake news. In: Proceedings of the 27th international conference on computational linguistics, pp 3391–3401
31. Peters ME, Neumann M, Iyyer M, Gardner M, Clark C, Lee K, Zettlemoyer L (2018) Deep contextual- ized word representations. In: Proceedings of NAACL-HLT, pp 2227–2237
32. Qi Y, Sachan D, Felix M, Padmanabhan S, Neubig G (2018) When and why are pre-trained word embeddings useful for neural machine translation? In: Proceedings of the 2018 conference of the north american chapter of the association for computational linguistics: human language technologies, vol 2 (short papers), pp 529–535
33. Rashkin H, Choi E, Jang JY, Volkova S, Choi Y (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
34. Reema A, Kar AK, Vigneswara Ilavarasan P (2018) Detection of spammers in twitter marketing: a hybrid approach using social media analytics and bio inspired computing. Information Systems Frontiers 20(3):515–530
35. Roy A, Basak K, Ekbal A, Bhattacharyya P (2018) A deep ensemble framework for fake news detection and classification. arXiv:arXiv-1811
36. Seide F, Li G, Chen X, Yu D (2011) Feature engineering in context-dependent deep neural networks for conversational speech transcription, IEEE

37. Shin J, Jian L, Driscoll K, Bar F (2018) The diffusion of misinformation on social media: Temporal pattern, message, and source. *Comput Hum Behav* 8:278–287
38. Shu K, Cui L, Wang S, Lee D, Liu H (2019) defend: Explainable fake news detection. In: *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery and data mining*, pp 395–405
39. Shu K, Mahudeswaran D, Wang S, Lee D, Liu H (2020) FakeNewsNet: A data repository with news content, social context, and spatio temporal information for studying fake news on social media. *Big Data* 8(3):171–188
40. Shu K, Wang S, Liu H (2019) Beyond news contents: The role of social context for fake news detection. In: *Proceedings of the twelfth ACM international conference on web search and data mining*, pp 312–320. ACM
41. Ruchansky N, Seo S, Liu Y (2017) Csi: A hybrid deep model for fake news detection. In: *Proceedings of the 2017 ACM on conference on information and knowledge management*. ACM, pp 797–806
42. Sibi P, Allwyn Jones S, Siddarth P (2013) Analysis of different activation functions using back propagation neural networks. *J Theor Appl Inf Technol* 47(3):1264–1268
43. Singh DSKR, Vivek RD, Ghosh I (2017) Automated fake news detection using linguistic analysis and machine learning. In: *International conference on social computing, behavioral-cultural modeling, & prediction and behavior representation in modeling and simulation (SBP-BRiMS)*, pp 1–3
44. Tacchini E, Ballarin G, Vedova ML, Moret S, Hoax LucadeAlfaro. (2017) Some like it Della Automated fake news detection in social networks. In: *2nd workshop on data science for social good, SoGood 2017*, pp 1–15. CEUR-WS
45. Tenney I, Das D, Pavlick E (2019) BERT rediscovers the classical NLP pipeline. In: *Proceedings of the 57th annual meeting of the association for computational linguistics*
46. Vasudevan V, Zoph B, Shlens J, Le QV (2019) Neural architecture search for convolutional neural networks. U.S Patent 10,521,729 issued December 31
47. Vosoughi S, 'Neo Mohsenvand M, Roy D (2017) Rumor gauge: Predicting the veracity of rumors on Twitter. *ACM Trans Knowl Discov Data (TKDD)* 11(4):1–36
48. Wang WY (2017) Liar, liar pants on fire: A new benchmark dataset for fake news detection. In: *Proceedings of the 55th annual meeting of the association for computational linguistics (vol 2: short Papers)*, pp 422–426



49. Weiss AP, Alwan A, Garcia EP, Garcia J (2020) Surveying fake news: Assessing university faculty's fragmented definition of fake news and its impact on teaching critical thinking. *Int J Educ Integr* 16(1):1–30
50. Yang F, Liu Y, Xiaohui Y, Yang M (2012) Automatic detection of rumor on Sina Weibo. In: *Proceedings of the ACM SIGKDD workshop on mining data semantics*, pp 1–7
51. Young T, Hazarika D, Poria S, Cambria E (2018) Recent trends in deep learning based natural language processing. *IEEE Comput Intell Mag* 13(3):55–75
52. Zhang X, Zhao J, LeCun Y (2015) Character-level convolutional networks for text classification. In: *Advances in neural information processing systems*, pp 649–657
53. Zhong B, Xing X, Love P, Wang Xu, Luo H (2019) Convolutional neural network: Deep learning-based classification of building quality problems. *Adv Eng Inform* 40:46–57
54. Zhou X, Zafarani R (2018) Fake news: a survey of research, detection methods, and opportunities. *arXiv:arXiv-1812*
55. Zubiaga A, Aker A, Bontcheva K, Liakata M, Procter R (2018) Detection and resolution of rumours in social media: A survey. *ACM Comput Surv (CSUR)* 51(2):1–36
56. E. Qawasmeh, M. Tawalbeh and M. Abdullah, "Automatic Identification of Fake News Using Deep Learning," 2019 Sixth International Conference on Social Networks Analysis, Management and Security (SNAMS), 2019, pp. 383-388, doi: 10.1109/SNAMS.2019.8931873.
57. M. L. Della Vedova, E. Tacchini, S. Moret, G. Ballarin, M. DiPierro and L. de Alfaro, "Automatic Online Fake News Detection Combining Content and Social Signals," 2018 22nd Conference of Open Innovations Association (FRUCT), 2018, pp. 272-279, doi: 10.23919/FRUCT.2018.8468301.
58. Z. Khanam, B. N. Alwasel, H. Sirafi, and M. Rashid, "Fake News Detection Using Machine Learning Approaches," *IOP Conference Series: Materials Science and Engineering*, vol. 1099, no. 1, p. 012040, Mar. 2021, doi: 10.1088/1757-899x/1099/1/012040.
59. A. Jain, A. Shakya, Harsh Khatter, and Amit Kumar Gupta, "A smart System for Fake News Detection Using Machine Learning," *ResearchGate*, Sep. 2019.  
[https://www.researchgate.net/publication/339022255\\_A\\_smart\\_System\\_for\\_Fake\\_News\\_Detection\\_Using\\_Machine\\_Learning](https://www.researchgate.net/publication/339022255_A_smart_System_for_Fake_News_Detection_Using_Machine_Learning) (accessed Aug. 04, 2022).
60. S. Raza and C. Ding, "Fake news detection based on news content and social contexts: a transformer-based approach," *International Journal of Data Science and Analytics*, vol. 13,

no. 4, pp. 335–362, Jan. 2022, doi: 10.1007/s41060-021-00302-z.

61. “Fake News Detection on Social Media: A Data Mining Perspective: ACM SIGKDD Explorations Newsletter: Vol 19, No 1,” ACM SIGKDD Explorations Newsletter, 2017. <https://dl.acm.org/doi/abs/10.1145/3137597.3137600> (accessed Aug. 04, 2022).
62. “Detecting Hoaxes, Frauds, and Deception in Writing Style Online | Proceedings of the 2012 IEEE Symposium on Security and Privacy,” Guide Proceedings, 2012. <https://dl.acm.org/doi/10.1109/SP.2012.34> (accessed Aug. 04, 2022).
63. A. Balasch, M. Beinhofer and G. Zauner, "The Relative Confusion Matrix, a Tool to Assess Classifiability in Large Scale Picking Applications," 2020 IEEE International Conference on Robotics and Automation (ICRA), 2020, pp. 8390-8396, doi: 10.1109/ICRA40945.2020.9197540.
64. X. Zhou and A. Del Valle, "Range Based Confusion Matrix for Imbalanced Time Series Classification," 2020 6th Conference on Data Science and Machine Learning Applications (CDMA), 2020, pp. 1-6, doi: 10.1109/CDMA47397.2020.00006.
65. M. Hasnain, M. F. Pasha, I. Ghani, M. Imran, M. Y. Alzahrani and R. Budiarto, "Evaluating Trust Prediction and Confusion Matrix Measures for Web Services Ranking," in IEEE Access, vol. 8, pp. 90847-90861, 2020, doi: 10.1109/ACCESS.2020.2994222.
66. J. L. Garcia-Balboa, M. V. Alba-Fernandez, F. J. Ariza-López and J. Rodriguez-Avi, "Homogeneity Test for Confusion Matrices: A Method and an Example," IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing Symposium, 2018, pp. 1203-1205, doi: 10.1109/IGARSS.2018.8517924.
67. M. Heydarian, T. E. Doyle and R. Samavi, "MLCM: Multi-Label Confusion Matrix," in IEEE Access, vol. 10, pp. 19083-19095, 2022, doi: 10.1109/ACCESS.2022.3151048.
68. S. H. Jeong, K. T. Lim and Y. S. Nam, "A combination method of two classifiers based on the information of confusion matrix," Proceedings Eighth International Workshop on Frontiers in Handwriting Recognition, 2002, pp. 519-523, doi: 10.1109/IWFHR.2002.1030963.
69. [https://www.bitdegree.org/learn/train-test-split?source=post\\_page-----89b5c704d6ee-----](https://www.bitdegree.org/learn/train-test-split?source=post_page-----89b5c704d6ee-----)
70. [https://news.mit.edu/2018/study-twitter-false-news-travels-faster-true-stories-0308?source=post\\_page-----89b5c704d6ee-----](https://news.mit.edu/2018/study-twitter-false-news-travels-faster-true-stories-0308?source=post_page-----89b5c704d6ee-----)
71. <https://towardsdatascience.com/text-pre-processing-stop-words-removal-using-different-libraries-f20bac19929a>
72. <https://medium.com/mllearning-ai/mllearning-ai-submission-suggestions-b51e2b130bfb>