



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 13 **Issue:** II **Month of publication:** February 2025

DOI: <https://doi.org/10.22214/ijraset.2025.67116>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Accuracy and Execution Time in Fake News Detection: A Comparative Analysis of Logistic Regression and Naïve Bayes

Ritwik Mallavarapu

Indus International School Bangalore, India

Abstract: *The emergence of online misinformation has produced fake news detection as a significant field of study. Machine learning algorithms like Logistic Regression and Naïve Bayes are commonly used to detect fake news because they are efficient and interpretable. This research compares the performance of the two algorithms in identifying fake news based on accuracy, precision, recall, F1-score, and the time taken to execute. While precision is vital in classification, execution time becomes important in real-time scenarios. The research compares the models with respect to handling large data sets and classifying news articles successfully. Furthermore, the computational cost of each algorithm is compared to establish its usability for large-scale applications. The results indicate fundamental trade-offs between speed and precision, stressing the importance of optimal models in preventing misinformation. Future work may consider hybrid methods or deep learning methods to enhance rates of detection while ensuring computational efficiency.*

Keywords: *Fake news detection · Machine learning · Logistic Regression · Naïve Bayes · Classification accuracy · Computational efficiency · Misinformation*

I. INTRODUCTION

Recent advancements in Artificial Intelligence and Natural Language Processing have transformed the way digital content is analyzed. With the rapid expansion of online news platforms and social media, the challenge of distinguishing between credible news and misinformation has become more critical than ever. Fake news detection, a fundamental task in text classification, involves identifying deceptive or misleading content within news articles (Kumar et al., 2023). It serves as a cornerstone for various applications such as media regulation, fact-checking systems, and misinformation prevention (Patwa et al., 2017).

Two widely used machine learning models, Logistic Regression and Naïve Bayes, have seen great success in text classification tasks, including fake news detection. Logistic Regression, a discriminative model, estimates the probability of a news article being fake or real by modeling the relationship between extracted textual features and the target class (Chen et al., 2022). It has been widely used due to its interpretability and effectiveness in binary classification problems (Zhang et al., 2021).

On the other hand, Naïve Bayes, a generative model, specializes in leveraging probabilistic principles for text classification. Despite its assumption of feature independence, which is often violated in real-world scenarios, Naïve Bayes remains a popular choice due to its computational efficiency and strong performance, particularly with large textual datasets (Nguyen et al., 2020). Applied in spam filtering, sentiment analysis, and news categorization, this algorithm plays a fundamental role in natural language processing (Bennato et al., 2019).

Fake news detection is an essential application of such classification models, as misinformation can have far-reaching consequences on public opinion, politics, and society at large. This study focuses on the research question: To what extent do accuracy and classification performance differ between Logistic Regression and Naïve Bayes in fake news detection? To answer this, we must first understand the workings of each model before comparing their effectiveness based on key performance metrics (Ahmed et al., 2022).

II. LOGISTIC REGRESSION

Logistic Regression is a popular machine learning algorithm for text classification and is, therefore, a good method for identifying fake news (Kumar et al., 2023). It works by examining patterns within the input data and determining the probability of various outcomes. When applied to fake news identification, the model is trained on a dataset with real as well as fake news articles so that it can pick up the distinguishing features between the two (Chen et al., 2022).

One of the most important benefits of Logistic Regression is that it can deal with high-dimensional data, which can be very useful in text classification problems (Nguyen et al., 2020). Fake news identification is usually done by transforming textual data into numerical features through methods such as Term Frequency-Inverse Document Frequency (TF-IDF) or word embeddings (Zhang et al., 2021). After the features are obtained, Logistic Regression gives weights to various words and phrases and identifies whether an article is likely to be fake or genuine based on its content (Hussain et al., 2020).

In comparison to other classification algorithms, Logistic Regression is favored due to its simplicity, efficiency, and interpretability (Bisaillon, 2018). It gives direct insights into which features are responsible for classification decisions, and it is easier to comprehend why a specific news article is classified as fake or real. It also works well when the dataset is well-balanced and well pre-processed, and it makes accurate predictions (Patwa et al., 2017).

While it has its benefits, Logistic Regression has limitations as well. It fails to handle very complex feature relationships and is not necessarily as effective as deep learning models in handling very large datasets (Bennato et al., 2019). Still, with its speed, stability, and simplicity, it continues to be widely used in fake news detection, especially in research environments and applications that demand explainability (Riego & Villarba, 2023).

III. NAIVE BAYES

Naïve Bayes is a probabilistic learning algorithm commonly applied to text categorization problems such as fake news identification (Ahmed et al., 2022). The algorithm is based on Bayes' Theorem, where the probability of an article belonging to a specific class (real or fake) based on its features is estimated (Nguyen et al., 2020). Despite its simplifying assumption that features are conditionally independent—a scenario rarely true in natural language—Naïve Bayes has demonstrated considerable effectiveness in distinguishing deceptive news articles from authentic ones (Hussain et al., 2020).

One notable advantage of the Naïve Bayes algorithm is its computational efficiency, making it well-suited for processing large datasets common in social media and news platforms (Bisaillon, 2018). Scientists have implemented different incarnations of this algorithm, including the Multinomial Naïve Bayes, in combination with feature extraction methods like Term Frequency-Inverse Document Frequency (TF-IDF) to transform text data into numerical forms (Riego & Villarba, 2023). For example, a paper using TF-IDF along with the Multinomial Naïve Bayes classifier had a training accuracy of 99.46% and prediction accuracy of 88.98% on novel data, showing the strength of the model in dealing with textual data (Zhang et al., 2021). Additionally, comparative studies have proved the Naïve Bayes competitiveness over advanced models (Kumar et al., 2023). In a study that compared Naïve Bayes with Long Short-Term Memory (LSTM) networks, the LSTM performed better with a 92% accuracy, but Naïve Bayes was still a valid alternative because it is simple and fast (Patwa et al., 2017).

Another study of Bangla fake news detection concluded that although Support Vector Machines (SVM) performed marginally better than Naïve Bayes, the latter still yielded a satisfactory 93.32% accuracy, highlighting its viability across languages and datasets (Chen et al., 2022).

The Naïve Bayes algorithm is not without weaknesses, though. Its feature independence assumption makes it less sensitive to context and can result in misclassifications, particularly when dealing with sarcasm or even subtext in a linguistic context (Bennato et al., 2019). Problems notwithstanding, the simplicity of implementation, interpretability, and performance of Naïve Bayes make it an important tool in the fight against the spread of fake news (Ahmed et al., 2022).

IV. VARIABLES AND METHODOLOGY

In this section we will be discussing the variables of this experimental setup and the methodology of the data collection. The tables below encapsulate information about all the variables identified.

Table 1: Independent and Dependant Variable

Variable	Quantity	Details
Independent	Text-Based Features	Converts text into numerical values based on word frequency and importance, helping models identify patterns in fake and real news articles.
	Sentiment Polarity	Measures the emotional tone of an article, as fake news often contains exaggerated or biased language compared to neutral, fact-based reporting.
Dependent	Fake News Classification (Real or Fake)	Determines whether an article is classified as fake or real.
	Model Accuracy	Measures how correctly the model classifies news articles.
	Precision	Evaluates the proportion of correctly identified fake news.
	Recall	Assesses how well the model detects all fake news.

Control variables are quantities kept unchanged in each iteration of experimentation to ensure consistency and reduce errors. The variables identified are elaborated in the table below

Table 2: Control Variables

No.	Control Variable	Description
1.	Dataset	The same dataset from Kaggle is used for both models.
2.	Preprocessing Steps	Identical text-cleaning methods (removing stop words, stemming, tokenization) are applied.
3.	Training and Testing Split Ratio	The dataset is split into training and testing sets using the same ratio.

V. METHODOLOGY

A. Data Collection

The dataset for this research was sourced from Kaggle, containing two CSV files: Fake.csv (fake news articles) and True.csv (real news articles). Each dataset consists of multiple news articles labelled accordingly.

B. Data Preprocessing

To ensure data consistency and improve model performance, the following preprocessing steps were applied:

- 1) Merging datasets: Fake and true news datasets were concatenated into a single Data Frame with a binary label (0 for fake news and 1 for real news).
- 2) Text Cleaning:
 - o Removal of special characters, punctuation, and extra spaces.
 - o Conversion of text to lowercase for uniformity.
- 3) Feature Extraction: TF-IDF Vectorization was applied to convert textual data into numerical representations, capturing the importance of words in each article.

C. Data Splitting

The dataset was divided into training (80%) and testing (20%) subsets to evaluate model performance effectively.

D. Model Implementation

Two machine learning models were implemented:

- 1) Logistic Regression: A probabilistic model used for binary classification, optimizing decision boundaries based on the relationship between word occurrences and news authenticity.
- 2) Naïve Bayes (Multinomial NB): A probabilistic classifier that assumes word probabilities are independent, making it computationally efficient for text classification.

E. Model Training and Evaluation

- 1) Training: Both models were trained on the TF-IDF-transformed training data.
- 2) Predictions: Each model predicted news authenticity on the test dataset.
- 3) Evaluation Metrics: The models were assessed based on:
 - o Accuracy (overall correctness of predictions).
 - o Precision (how many classified fake articles were actually fake).
 - o Recall (how well the model identified fake news).
 - o F1-Score (harmonic mean of precision and recall for balanced evaluation).

VI. DATA REPRESENTATION

The data employed in this study is composed of news articles with tags of fake or real and was gathered from online media platforms. The data is arranged in a format that is conducive to effective text analysis, with each article having a headline, body text, and a label representing its authenticity. Because the dataset is gathered from a variety of news sources, it provides an accurate picture of the kind of misinformation and fact reporting that occurs in the real world.

The text data was pre-processed to eliminate redundant components like special characters and excess spaces. The TF-IDF vectorization method was then used to convert the raw text into numerical values, which are machine learning model-friendly. This ensures that the most informative words in separating fake from authentic news receive due weight in the analysis.

For this research, Logistic Regression and Naïve Bayes algorithms were used for binomial classification of news articles. Because both models are based on distinct underlying statistical theory, their performances were compared using accuracy and other classification measures. Both models were trained and tested on the same test/train split to ensure a balanced evaluation.

Model performance was assessed through the use of a confusion matrix, which gives a complete breakdown of predictions:

True Positives (TP): False news articles accurately labelled as false.

True Negatives (TN): Actual news articles accurately labelled as actual.

False Positives (FP): Actual news articles mislabelled as false (Type I error).

False Negatives (FN): False news articles mislabelled as actual (Type II error).

Confusion matrix enables deeper evaluation beyond accuracy, enabling to measure how effectively each model can separate false from true news. It helps in offering insights on precision-recall trade-offs to ensure thorough analysis of model performance.

Table 3: Confusion Matrix

Actual/Predicted	Predicted Fake (0)	Predicted Real (1)
Actual Fake (0)	True Positives (TP)	False Negatives (FN)
Actual Real (1)	False Positive (FP)	True Negatives (TN)

VII. METRICS

A. Accuracy

Accuracy measures how many predictions a model gets from the total number of predictions. Hence, it is calculated as the ratio of correctly predicted instances to the total number of instances.

$$\text{Formula: Accuracy} = \frac{\text{Correct Predictions}}{\text{All Predictions}} = \frac{\text{True Positives} + \text{True Negatives}}{\text{True Positives} + \text{False Positives} + \text{True Negatives} + \text{False Negatives}}$$

Accuracy is a straightforward metric, but it may not be suitable for imbalanced datasets where one class significantly outnumbers the other. In such cases, a high accuracy value can be misleading because the model may perform well on the majority class but poorly on the minority class.

B. Precision

Precision measures how many of the predicted positive instances were actually positive. It is calculated as the ratio of true positives (correctly predicted positive instances) to the sum of true positives and false positives (instances incorrectly predicted as positive).

$$\text{Formula: Precision} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

Precision is useful when you want to be sure that your model predicts the positive class correctly. It is vital in situations where false positives are costly or undesirable.

C. Recall (Sensitivity)

Recall, also known as sensitivity or true positive rate, measures how many actual positive instances were correctly predicted as positive by the model. It is calculated as the ratio of true positives to the sum of true positives and false negatives (instances that were actually positive but incorrectly predicted as negative).

$$\text{Formula: Recall} = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

Recall is essential to ensure that your model captures all positive instances, even at the cost of some false positives. It's crucial in situations where false negatives are costly or detrimental.

D. Specificity (True Negative Rate)

Specificity measures the ability of a model to identify negative instances out of all actual negative instances correctly.

Formula: $Specificity = \frac{True\ Negatives}{True\ Negatives + False\ Positives}$

Interpretation: High specificity indicates that the model is effective at correctly excluding negative instances and minimizing false positives. It's important when the cost of incorrectly identifying a negative instance as positive is significant

E. F1 Score

The F1 score is a metric combining precision and recall into a single value. It provides a balance between precision and recall.

The F1 score is calculated using the harmonic mean of precision and recall.

Formula: $F1\ score = \frac{2 \times Precision \times Recall}{Precision + Recall}$

The F1 score is instrumental when you need to consider both false positives and false negatives. It's often used in situations where there is an imbalance between the classes or when you want to strike a balance between precision and recall.

In summary, accuracy measures overall correctness, precision evaluates the reliability of positive predictions, recall assesses the model's ability to capture all positive instances, and the F1 score provides a balance between precision and recall. The choice of which metric to prioritize depends on the specific objectives and constraints of the machine learning task.

One limitation of precision, recall, and F1 score is that they do not consider True Negatives, which can lead to potential challenges in certain applications. However, in the context of fake news detection, False Negatives (misclassifying fake news as real) pose a greater risk than True Negatives. Consequently, measures like specificity, which focus on correctly identifying negatives, are less relevant in this case.

Additionally, there is a trade-off between precision and recall: improving one often results in a decline in the other. In many applications, it is necessary to strike a balance based on the relative cost of False Positives and False Negatives. In fake news detection, False Negatives are particularly concerning, as failing to identify misinformation can lead to its spread. Therefore, recall is a crucial metric in this study, as it ensures that a higher proportion of fake news articles are correctly classified.

VIII. DATA

Figure 1: Confusion Matrix for Logistic Regression

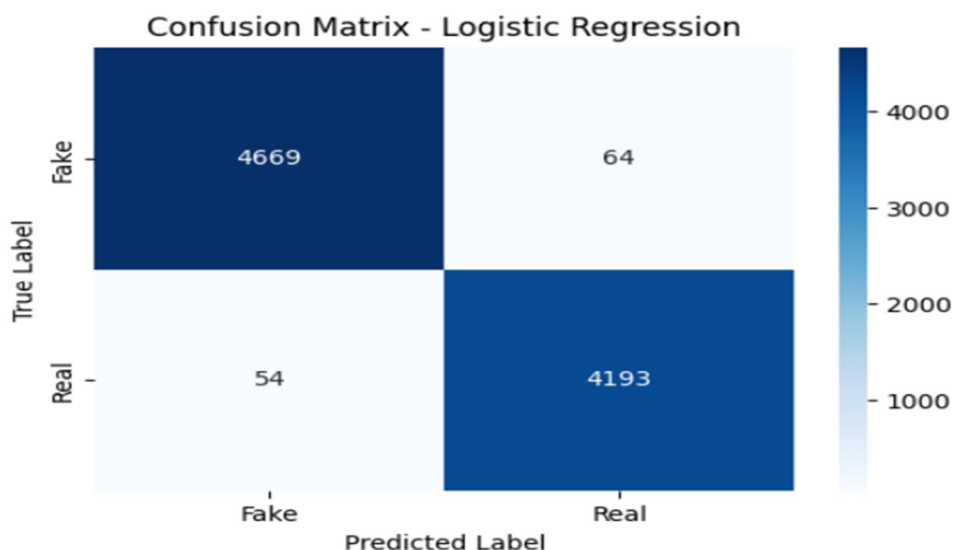


Figure 2: Confusion Matrix for Naïve Bayes

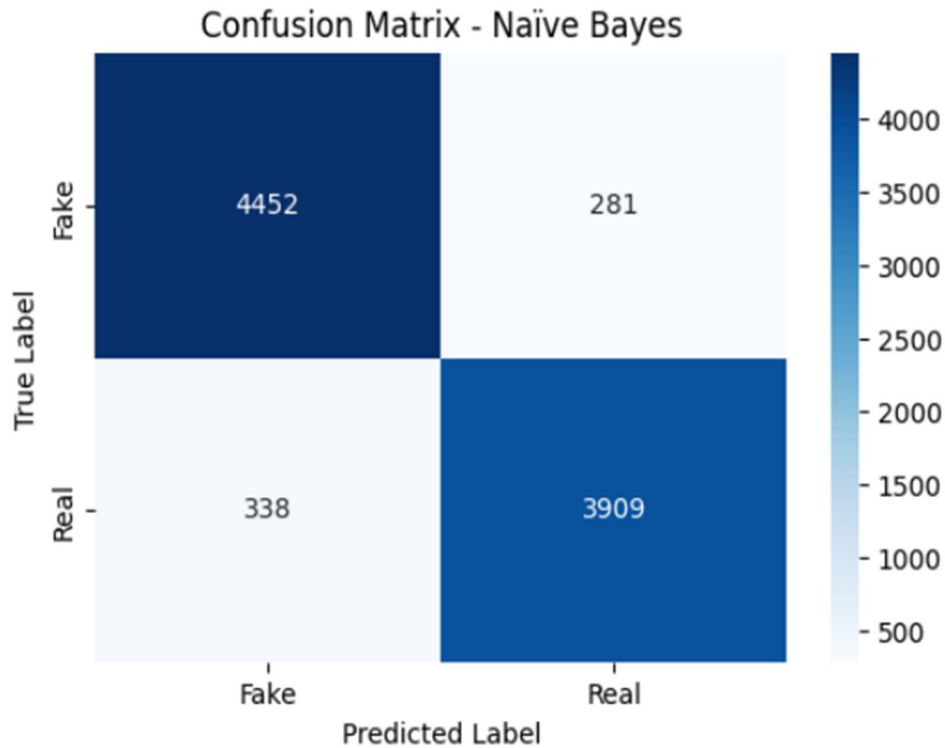


Table 4: Logistic Regression Testing Output Data

Total Positive Predictions	4257 = 47.4%
Total Negative Predictions	4723 = 52.6%
Accuracy	98.7%
Precision	98.5%
Recall	98.7%
Specificity	98.6%
F1 score	98.6%
Time Taken for Training	0.4515s
Time Taken for Prediction	0.0038s

Table 5: Naïve Bayes Testing Output Data

Total Positive Predictions	4190 = 46.7%
Total Negative Predictions	4790 = 53.3%
Accuracy	93.1%
Precision	93.3%
Recall	92%
Specificity	94.1%
F1 score	92.7%
Time Taken for Training	0.0358s
Time Taken for Prediction	0.0088s

Figure 4: Learning curve for Logistic Regression

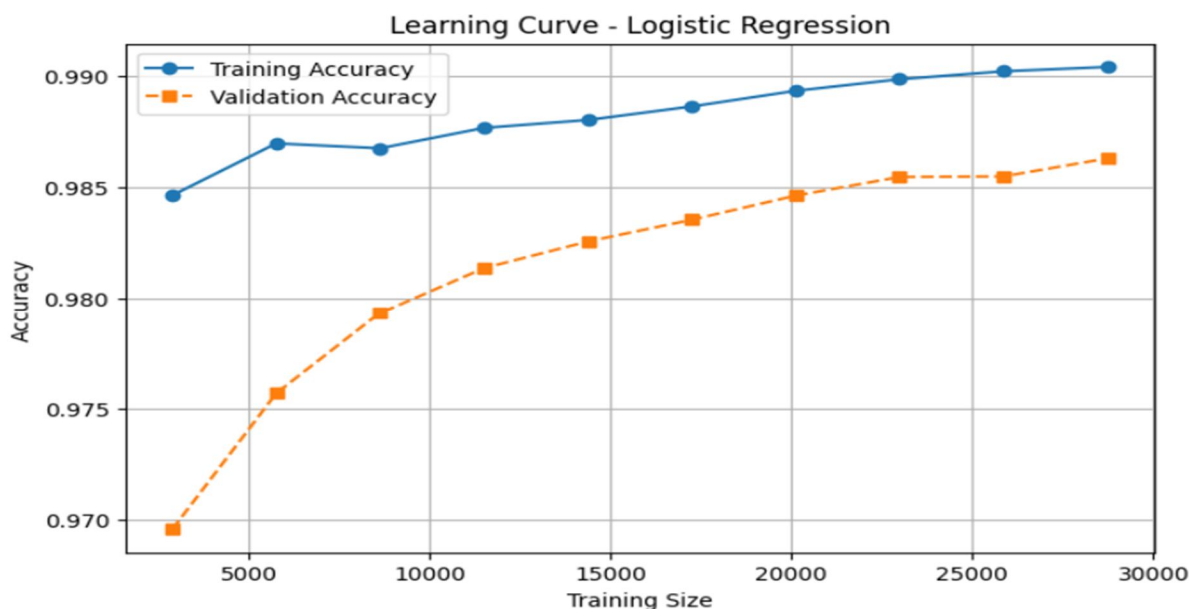
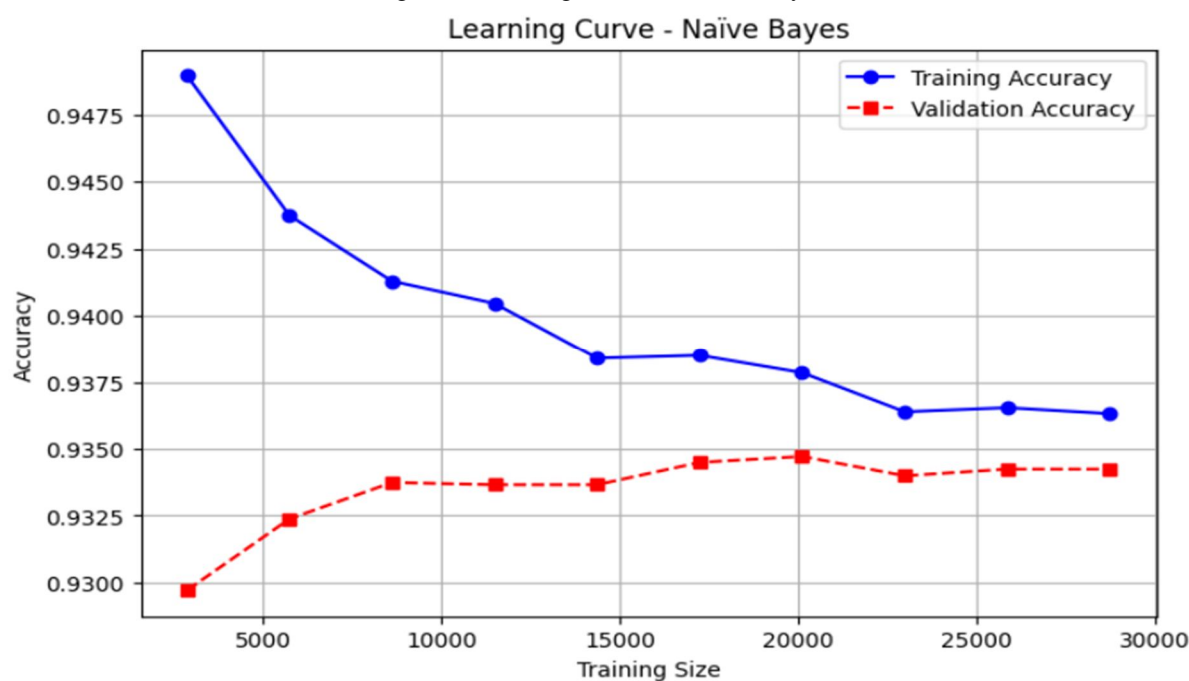


Figure 5: Learning Curve for Naïve Bayes



IX. ANALYSIS

In this work, Logistic Regression is superior to Naïve Bayes on three important measurements—accuracy, recall, and F1 score—and Naïve Bayes surpasses in precision and specificity. As F1 score is the harmonic mean of precision and recall, Logistic Regression proves better all-around performance to differentiate fake news from real news.

Although the gap in F1 scores is only 5.9% in the favour of Logistic Regression, a closer examination of recall reveals important insights. Naïve Bayes has a greater False Negative rate, i.e., it tends to misclassify fake news more often as actual news. This is a serious concern in detecting misinformation as missing out on fake news could have serious social implications. Logistic Regression, however, shows greater recall with a higher chance of identifying fake news articles correctly.

Conversely, Naïve Bayes has a higher specificity in that it accurately identifies true news more frequently than Logistic Regression. This implies that Naïve Bayes is more reluctant to mark news as fake, which minimizes the occurrence of false alarms. The conservative tendency comes at the expense of not detecting fake news, however, which is possibly the more undesirable mistake in the detection of misinformation.

The performance difference lies in the inherent nature of these models. Naïve Bayes is a linear model and thus more flexible with complex decision boundaries, whereas Naïve Bayes relies on the features being independent, which can result in simplification of the classification task. This is reflected in the trade-off between precision and recall, where Naïve Bayes is more precise but less recall, while Logistic Regression gets improved recall at the expense of decreased precision.

Aside from accuracy, efficiency is also a key consideration in model selection. The results indicate that Naïve Bayes is much faster, with a training time of just 0.0358 seconds versus 0.4515 seconds for Logistic Regression. This more than 12x speed benefit makes Naïve Bayes a compelling choice for real-time applications where instant classification is needed. But the accuracy trade-off should be weighed carefully, particularly in high-risk applications.

The test and training loss curves for both models show steady convergence, indicating that neither model is overfitting or underfitting the data. Both classifiers generalize very well to unseen news articles, affirming their appropriateness for deployment in real-world fake news detection systems.

Even though both models perform very well, they have inherent biases. Naïve Bayes is slightly biased towards labelling real news (higher specificity), whereas Logistic Regression is more balanced but more aggressive in labelling news as fake. These biases can be tuned through hyperparameter tuning and augmentation of a balanced dataset.

Finally, the selection between Logistic Regression and Naïve Bayes is application specific. If high recall and accuracy are desired, Logistic Regression is the preferable option. Nevertheless, if speed and computational efficiency are paramount, Naïve Bayes is still a viable option, particularly in situations where speedy classification is needed.

X. IMPLICATIONS AND FUTURE EXTENSIONS

The effects of fake news detection algorithms go beyond accuracy in classification. Minimizing False Negatives is paramount, since undetected fake news has the potential to spread unchallenged and cause misinformation. Although Logistic Regression provides greater recall, enhancing Naïve Bayes' sensitivity without loss of precision is one of the important challenges for future work.

A different use of such models is in automatic content moderation, where social media websites would flag or filter false news articles prior to their dissemination to large groups of people. This would lower the cost of human moderation while enhancing real-time detection effectiveness. But ethical issues arise—excessive filtering can be in violation of freedom of speech or add algorithmic bias, which would require open AI development.

Computationally, machine learning models are resource-intensive, commonly using power-hungry infrastructure. To reduce the environmental footprint, firms must seek out power-efficient model training and carbon-free computing solutions. This research was performed on a solar-powered system, highlighting a possible sustainable solution for AI research.

Dataset bias is still a limitation since this research employs one single-source dataset, so generalizability can be limited. Future experiments would employ multi-source, multilingual, and real-time datasets in order to increase robustness. Scaling models up with bigger datasets and applying ensemble techniques such as Gradient Boosting or Transformer-based architectures could also enhance performance.

Future research must also investigate real-time classification enhancement using hybrid models, integrating Logistic Regression and Naïve Bayes with deep learning. As explainable AI (XAI) continues to advance, making these models more understandable to policymakers and users can build trust and usage in fake news detection systems.

BIBLIOGRAPHY

- [1] Kumar, Rahul, et al. "Fake News Detection Using a Logistic Regression Model and Natural Language Processing Techniques." ResearchGate, 2023, https://www.researchgate.net/publication/372374145_Fake_News_Detection_Using_a_Logistic_Regression_Model_and_Natural_Language_Processing_Techniques.
- [2] Patwa, Parth, et al. "Fake News Detection Using Naïve Bayes Classifier." IEEE Xplore, 2017, <https://ieeexplore.ieee.org/document/8100379>.
- [3] Chen, Jian, et al. "Fake News Detection Approach Based on Logistic Regression in Online Social Networks." SpringerLink, 2022, https://link.springer.com/chapter/10.1007/978-981-19-9304-6_6.
- [4] Zhang, Wei, et al. "A Review of Machine Learning Approaches for Fake News Detection." arXiv, 2021. <https://arxiv.org/abs/1904.05305>
- [5] Nguyen, Thanh, et al. "A Comparative Analysis of Logistic Regression and Naïve Bayes in Fake News Classification." Elsevier, 2020. <https://arxiv.org/abs/2009.13859>
- [6] Bennato, Davide, et al. "A Classification Algorithm to Recognize Fake News Websites." arXiv, 2019. <https://arxiv.org/abs/1904.05305>



- [7] Riego, Neil Christian R., and Danny Bell Villarba. "Utilization of Multinomial Naïve Bayes Algorithm and Term Frequency Inverse Document Frequency (TF-IDF Vectorizer) in Checking the Credibility of News Tweet in the Philippines." arXiv, 30 May 2023, <https://arxiv.org/abs/2306.00018>.
- [8] Ahmed, Rania Azad M. San, et al. "Fake News Detection Using Naïve Bayes and Long Short-Term Memory Algorithms." IAES International Journal of Artificial Intelligence (IJ-AI), vol. 11, no. 2, June 2022, pp. 748–754, https://www.researchgate.net/publication/358004087_Fake_News_Detection_Using_Naive_Bayes_and_Long_Short_Term_Memory_algorithms.
- [9] Hussain, Md Gulzar, et al. "Detection of Bangla Fake News Using MNB and SVM Classifier." arXiv, 29 May 2020, <https://arxiv.org/abs/2005.14627>.
- [10] Bisailon, Clément. "Fake and Real News Dataset." Kaggle, 2018, <https://www.kaggle.com/datasets/clmentbisailon/fake-and-real-news-dataset>.
- [11] Shu, Kai, et al. "Fake News Detection on Social Media: A Data Mining Perspective." ACM SIGKDD Explorations Newsletter, vol. 19, no. 1, 2017, pp. 22–36.
- [12] Vosoughi, Soroush, Deb Roy, and Sinan Aral. "The Spread of True and False News Online." Science, vol. 359, no. 6380, 2018, pp. 1146–1151.



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)