



IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 12 Issue: XI Month of publication: November 2024 DOI: https://doi.org/10.22214/ijraset.2024.65670

www.ijraset.com

Call: 🕥 08813907089 🔰 E-mail ID: ijraset@gmail.com



Explainable AI for News Classification

Devang Chavan¹, Shrihari Padatare²

Department Computer Science and Engineering, DY Patil International University

Abstract: The proliferation of news content across digital platforms necessitates robust and interpretable machine learning models to classify news into predefined categories effectively. This study investigates the integration of Explainable AI (XAI) techniques within the context of traditional machine learning models, including Naive Bayes, Logistic Regression, and Support Vector Machines (SVM), to achieve interpretable and accurate news classification. Utilizing the News Category Dataset, we preprocess the data to focus on the top 15 categories while addressing class imbalance challenges. Models are trained using Term Frequency-Inverse Document Frequency (TF-IDF) vectorization, achieving an acceptable classification accuracy of 67% across all models despite the complexity introduced by the high number of classes.

To elucidate the decision-making processes of these models, we employ feature importance visualizations derived from model coefficients and feature log probabilities, complemented by local interpretability techniques such as LIME (Local Interpretable Model-agnostic Explanations). These methodologies enable granular insights into word-level contributions to predictions for each news category. Comparative heatmaps across models reveal significant consistencies and divergences in feature reliance, highlighting nuanced decision-making patterns.

The integration of explainability into news classification provides critical interpretive capabilities, offering transparency and mitigating the risks associated with algorithmic opacity. The findings demonstrate how XAI enhances stakeholder trust by aligning model predictions with human interpretability, particularly in ethically sensitive domains. This work emphasizes the role of XAI in fostering responsible AI deployment and paves the way for future advancements, including deep learning integration and multilingual news classification with inherent interpretability frameworks.

Keywords: Explainable AI, Responsible AI, AI Ethics, News Classification, LIME.

I. INTRODUCTION

The emergence of Artificial Intelligence (AI) has significantly altered the landscape of the media industry, particularly in the realm of news classification. With the increasing volume of digital content available, automated systems have become essential for efficiently categorizing articles across various topics such as politics, technology, and health. These AI-driven classification systems not only enhance the speed and accuracy of content delivery but also facilitate personalized news feeds for users. However, despite the advantages of these technologies, there are critical concerns regarding their transparency, ethical implications, and the potential for bias.

AI models often function as "black boxes," leading to a lack of clarity about how decisions are made. This opacity is particularly concerning in the context of news media, where biased classifications can distort public understanding and influence societal narratives. For instance, if a news classification system is trained predominantly on articles from a specific political perspective, it may inadvertently reinforce existing biases, further polarizing public opinion. The ability to understand and trust AI outputs is crucial, especially when misinformation can have serious repercussions.

To address these challenges, Explainable AI (XAI) has gained traction as a solution that promotes interpretability and transparency in AI systems. XAI techniques such as SHAP (SHapley Additive exPlanations) and LIME (Local Interpretable Model-Agnostic Explanations) provide valuable insights into the decision-making processes of AI models. By elucidating which features contribute most significantly to classifications, these techniques help users, including journalists and editors, understand the rationale behind AI-generated outputs, thereby fostering trust and accountability.

In this research, the focus is on the integration of XAI into news classification systems. The goal is to evaluate how these techniques can mitigate bias, improve interpretability, and enhance the ethical deployment of AI in journalism. By utilizing various machine learning models—including Logistic Regression, Support Vector Machines, and Multinomial Naive Bayes—this study will analyze a curated dataset of news articles. The findings aim to demonstrate that incorporating XAI not only enhances model transparency but also supports more equitable AI practices in the media landscape.



As digital news consumption continues to evolve, ensuring that AI systems not only provide efficiency but also maintain the integrity of information dissemination is imperative. This research seeks to contribute to a deeper understanding of how XAI can play a pivotal role in promoting ethical AI in journalism, ultimately serving the public's right to reliable and unbiased news.

II. LITERATURE REVIEW

A. Introduction to Explainable AI (XAI)

Artificial Intelligence (AI) has revolutionized many fields, including journalism, where it is used to automate tasks such as categorizing news articles. However, the decision-making processes of many AI models remain opaque, leading to concerns about transparency, bias, and accountability. Explainable AI (XAI) seeks to address these issues by providing users with interpretable insights into how models make decisions. This transparency is crucial in sensitive applications like news classification, where trust in AI systems plays a vital role in public perception [1] [2].

Arrieta et al. (2019) emphasized the importance of XAI for ensuring ethical and responsible AI practices. Their work highlighted the potential of XAI to make AI systems more accessible and interpretable for a broad audience, particularly in domains where fairness is a key concern. Doshi-Velez and Kim (2017) further advocated for interpretable machine learning models, noting that clear explanations can enhance trust and usability in AI systems, especially in high-stakes scenarios [1] [2].

The use of AI in news classification has revolutionized information dissemination by automating tasks like categorizing articles into topics such as politics, sports, or entertainment. Despite this, traditional AI models are often criticized for their lack of transparency, as they function like "black boxes," offering little clarity on how decisions are made. This opacity has raised significant concerns about bias, misinformation, and the ethical use of AI in media. Explainable Artificial Intelligence (XAI) addresses these issues by making model behavior interpretable and accessible, ensuring trust and accountability among users.

B. Key Techniques in Explainable AI

Numerous techniques have been developed to enhance the interpretability of AI models. These methods enable researchers and practitioners to understand the internal workings of models applied to news classification.

- SHAP (SHapley Additive exPlanations): SHAP utilizes cooperative game theory to quantify the contribution of each feature in a prediction. Developed by Lundberg and Lee (2017), SHAP uses principles from cooperative game theory to attribute importance scores to individual features. This approach is particularly valuable for identifying which factors most influence a model's predictions, providing insights at both the global and local levels [3] [5]. For instance, in classifying a political news article, SHAP may highlight keywords like "election" or "campaign" as significant determinants. This method excels in providing both global insights (overall model behavior) and local explanations (specific predictions).
- 2) LIME (Local Interpretable Model-Agnostic Explanations): LIME generates simplified, interpretable models that approximate the predictions of complex classifiers for specific data points. Ribeiro et al. (2016) introduced LIME as a method to locally approximate the behavior of complex models. By perturbing input data and observing changes in outputs, LIME generates interpretable explanations for specific predictions, helping users understand localized model behavior [4]. In the context of news, it can show how minor changes in an article's content might alter its classification outcome.
- 3) Attention Mechanisms: Common in natural language processing (NLP), attention mechanisms help identify the parts of a text that a model focuses on during classification. Yang et al. (2016) demonstrated the role of attention mechanisms in highlighting key parts of text that influence model predictions. These mechanisms, widely used in transformer-based models like BERT, provide visual cues that make the decision-making process more intuitive and comprehensible [5] [6]. For example, these mechanisms can visually highlight phrases or sentences in an article that influenced its categorization.
- 4) Integrated Gradients: This method quantifies the contribution of each input feature to a model's decision by comparing its baseline to the actual prediction. It has been applied effectively to textual data to explain classification outcomes.

C. Applications in News Classification

The integration of XAI into news classification systems addresses several critical challenges:

1) Bias Identification: By revealing the features that influence decisions, XAI helps detect and mitigate biases in training data or model architectures.SHAP enables the identification of biases by pinpointing features that disproportionately affect classifications.



This capability is crucial for ensuring fairness in news classification, particularly when datasets contain imbalances or skewed representations (1) (3). For instance, SHAP can expose political biases by showing which terms predominantly affect classifications.

- 2) Trust and Transparency: Journalists and editors often rely on AI-driven tools for content curation. Providing clear, interpretable outputs ensures confidence in these systems. Interpretability fosters confidence among users, such as journalists and editors, who rely on AI systems for content curation. By explaining how specific predictions are made, tools like LIME and attention mechanisms ensure that users can validate the fairness and accuracy of classifications [4] [6].
- 3) Misinformation Control: XAI enables stakeholders to verify the accuracy and integrity of AI-driven classifications, reducing the risk of misinformation spreading through automated systems. Transparent AI systems allow stakeholders to verify the reliability of model outputs, reducing the likelihood of misinformation. XAI techniques help ensure that classification models promote accurate and unbiased reporting [2] [3] [5].

A notable case study demonstrated how integrating SHAP into a news classification system uncovered over-reliance on specific politically charged terms, leading to refinements in the dataset and improved system fairness.

D. Challenges and Research Gaps

While XAI holds great promise, its practical application in news classification is not without challenges:

- Computational Intensity: Many XAI methods, such as SHAP, require significant computational resources, making them less suitable for real-time systems, particularly for large datasets or real-time applications. Optimizing these techniques for efficiency is essential for their broader adoption [3] [8].
- 2) Subjectivity in Interpretation: Different stakeholders may interpret model explanations differently, leading to inconsistencies in assessing fairness and accuracy.
- 3) Data Quality Dependence: The effectiveness of XAI techniques is heavily influenced by the diversity and quality of training datasets. Biases in the data can propagate through explanations, compromising their utility.
- 4) Scalability: Explanations provided by XAI methods can be interpreted differently depending on the user's perspective. A standardized framework for assessing interpretability could address this issue, ensuring consistency across applications [1]

[2].

Future research must focus on optimizing these methods to reduce their computational cost while maintaining robustness and reliability.

Sr	Technique	Strengths	Limitations
No.			
1	SHAP	Detailed local and global explanations	Computationally expensive
2	LIME	Simple and user- friendly	Sensitive to input perturbations
3	Attention Mechanisms	Intuitive and visually interpretable	Limited to specific neural architectures
4	Integrated Gradients	Robust for deep learning models	Requires baseline comparison

E. Comparative Analysis of Techniques

Each technique contributes uniquely to the field, and combining them often results in a more comprehensive understanding of model behavior.

F. Conclusion

The literature on Explainable AI highlights its transformative potential in making AI-driven systems transparent and ethical. By leveraging techniques like SHAP, LIME, and attention mechanisms, researchers have addressed key concerns such as bias, misinformation, and user trust in news classification systems. However, challenges related to scalability, computational efficiency, and interpretability persist.



Addressing these gaps will ensure the broader adoption of XAI in media and beyond, ultimately leading to fairer and more accountable AI systems. The integration of XAI into news classification systems is a critical step toward creating more ethical and transparent AI applications, addressing issues of bias and fostering trust. As AI continues to play a significant role in shaping public discourse, ensuring that these systems operate transparently will be vital for their success and societal acceptance. Future research should prioritize scalability, real-time application, and adaptability across diverse languages and cultures to maximize the impact of XAI in the media domain.

III. METHODOLOGY

Our methodology was designed to ensure a systematic, interpretable, and ethical approach towards solving the problem of news category classification. This paper explains the four stages of our approach: acquisition and preprocessing of datasets, model building, performance measurements, and finally, a comprehensive explainability analysis, along with both the technical merit and ethical undertones associated with our study.

A. Data Acquisition and Preprocessing

We utilized the News Category Dataset, sourced from Kaggle, which contains over 200,000 news headlines categorized into 42 unique classes. An initial exploration of the dataset revealed significant class imbalance, with certain categories dominating the distribution. For the purpose of this study, we focused on the top 15 categories, which were selected based on their frequency of occurrence. This selection ensures both statistical robustness and relevance to real-world applications.

To counterbalance this imbalance, we did undersampling, so that all classes would be of the same size as the minor class. This preprocessing ensures our models do not bias against the major classes and might lead to biased predictions.

B. Feature Representation

Text was converted into numerical representations that the machine learning algorithms could take in. We tried multiple approaches including:

CountVectorizer, which converts text into token frequencies, Word2Vec, a distributed representation method capturing semantic relationships, and TF-IDF (Term Frequency-Inverse Document Frequency).

After these experiments, TF-IDF was chosen as the best feature extraction technique for this task. The TF-IDF not only performs well but also gives a weighted representation of terms by importance, which was crucial for downstream explainability analyses.

C. Model Training

To classify news headlines, we used a diverse set of machine learning models. These included: Logistic Regression, a robust baseline for multi-class classification, Multinomial Naive Bayes, a probabilistic model suitable for text data, and Support Vector Machines (SVM), a high-margin classifier that can work with high-dimensional data.

All the models were trained with stratified k-fold cross-validation, which ensures consistent evaluation across folds and reduces biases from dataset partitioning. Hyperparameters were tuned systematically with grid search to optimize performance metrics such as accuracy, precision, recall, and F1-score.

D. 3.4 Explainability and Interpretability

A critical part of our approach was the embedding of explainability into the model evaluation pipeline. This aligns with the tenets of transparent and responsible AI in that predictions made are interpretable, hence promoting accountability.

1) Local Explanation Using LIME

We utilized the LIME (Local Interpretable Model-Agnostic Explanations) library for inspecting individual predictions. LIME generates interpretable explanations by perturbing input data and analyzing the model predictions resulting from such a perturbation. For every instance, LIME shows key features (words) which influence the prediction, making possible granular insights into decision-making. This technique made it possible to check whether models were relying on semantically meaningful terms rather than artifacts or noise.

2) Global Explanation by Coefficient Analysis

In case of linear models such as Logistic Regression and SVM, we carried out coefficient analysis to determine how features in general contribute.



Every coefficient of the linear model will represent how a feature is contributing towards some class, and the features for every class were selected along with corresponding validation for matching the ones from the knowledge domain.

This coefficient analysis is important because it uncovers systemic biases or over-reliance on specific features, thus ensuring the model is fair and ethical in its soundness. For example, a disproportionately high weight for sensitive terms would suggest unintended biases, which would then be addressed iteratively.

3) Probabilistic Insights in Naive Bayes

In the case of Multinomial Naive Bayes, we analyzed log probabilities for each feature. This probabilistic perspective provided additional transparency, offering a view of how strongly a feature contributes to class likelihoods. By examining these probabilities, we ensured that the model's assumptions were interpretable and aligned with expected patterns in the data.

E. Performance Evaluation

To rigorously compare model performance, we employed a comprehensive suite of evaluation metrics, including:

Accuracy: Global accuracy of predictions,

Precision and Recall: Class-specific performance measures, and

F1-Score: Harmonic mean of precision and recall, useful for imbalanced datasets.

We also developed confusion matrices to visualize errors in classification and identify patterns of misclassification. Such analyses provided both quantitative and qualitative insights into the relative strengths of each model.

F. Ethical Considerations and Broader Implications

Our methodological choices were guided by ethical imperatives, especially in the domains of transparency, accountability, and fairness. In doing so, we have integrated explainability techniques such as LIME and coefficient analysis into our models, making sure that our models are both accurate and interpretable, which is important in real-world applications where the influence of automated systems can have a bearing on important decisions.

These insights from explainability analyses also help us identify and mitigate any potential biases in accordance with the principles of responsible AI. We also believe that this makes our models' decision-making processes interpretable to stakeholders, thereby helping in fostering trust and increasing acceptability of AI systems in sensitive domains.

IV. RESULTS AND DISCUSSION

A. Model Performance

The performance of the three classification models—Naive Bayes, Logistic Regression, and Support Vector Machine (SVM)—was evaluated over 15 categories. Despite the complexity introduced by the large number of classes, all models demonstrated a similar level of performance, achieving an accuracy of 67%, which is acceptable for a multi-class classification task of this nature.

To provide a deeper understanding of the models' capabilities, additional metrics including precision, recall, and F1-score were calculated. These metrics, visualized using a heatmap (Figure 1), illustrate that all three models maintain comparable weighted scores across these metrics, with only minor variations. This consistency suggests that the models perform similarly in balancing precision and recall while accounting for the dataset's multi-class structure.

Overall, the heatmap highlights Logistic Regression and SVM as slightly more consistent in their performance metrics compared to Naive Bayes, which exhibits a marginally lower precision. However, given the similarity in accuracy, all three models are viable options for this classification task.





B. Feature Importance Analysis

To further interpret the models' behavior, an analysis of feature importance was conducted. The significance of words within each category was determined for each model, and the results were visualized through a series of heatmaps (Figures 2). These heatmaps display the top 10 most important words for each category and their corresponding importance scores, offering valuable insights into the decision-making processes of the models.

- 1) Naive Bayes: This model relies on the frequency-based probabilities of words, and its heatmap highlights highly discriminative terms for each category. Notably, it tends to emphasize words with strong associations to specific categories, but its simplistic assumptions may lead to a narrower feature selection compared to the other models.
- 2) Logistic Regression: This model's feature coefficients reveal a more nuanced understanding of word importance, particularly in cases where terms may contribute positively or negatively to multiple categories. The heatmap underscores Logistic Regression's ability to identify features that are both highly relevant and context-sensitive.
- 3) SVM: Similar to Logistic Regression, SVM demonstrates an effective feature selection process, leveraging the magnitude of its coefficients to highlight impactful words. However, its emphasis on margin maximization occasionally leads to lower importance for features that might be more critical in overlapping categories.







International Journal for Research in Applied Science & Engineering Technology (IJRASET) ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538

Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

C. Explainability through LIME

The final step in the analysis involved generating explanations for model predictions using LIME (Local Interpretable Model-Agnostic Explanations). This explainability framework was applied to individual test instances to uncover how the models arrived at their predictions.

- For each model, LIME provided a breakdown of the most influential words and their contributions to the predicted category. The explanations align well with the feature importance heatmaps, further validating the models' reliance on category-specific terms.
- 2) Example instance explanations: LIME explanations for a test instance demonstrated how words like "technology," "politics," or "sports" contributed to the classification of corresponding categories. For instance, words such as "election" and "candidate" were heavily weighted in predictions related to the "politics" category across all three models.

This approach enhances the interpretability of the models, providing users with a clear understanding of the decision-making process. Furthermore, it bridges the gap between the statistical performance of the models and their practical application by elucidating how predictions are derived from the input features.







Figure 4: LIME Explanations for Logistic Regression







Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

D. Discussion

The analysis reveals several key insights:

- Model Performance: While the overall accuracy of 67% might appear moderate, it reflects the complexity of the multi-class nature of the problem. The consistent performance across models suggests that the TF-IDF vectorizer effectively captures relevant patterns in the text data, enabling comparable outcomes across algorithms with distinct methodologies.
- 2) Feature Importance: The heatmaps emphasize the models' reliance on distinct and meaningful words, with Logistic Regression and SVM displaying more nuanced feature selection compared to Naive Bayes. This aligns with the models' underlying mechanisms, as Logistic Regression and SVM optimize for finer decision boundaries.
- 3) Explainability: The incorporation of LIME adds a critical layer of interpretability, allowing us to validate model decisions and identify the contribution of individual words. This feature is particularly significant for applications where understanding model predictions is essential for trust and reliability.

Future work could involve experimenting with ensemble methods or deep learning models to improve accuracy while maintaining interpretability. Additionally, addressing data imbalance and exploring alternative feature engineering techniques may further enhance model performance.

Overall, this study highlights the importance of combining traditional evaluation metrics with explainability techniques to achieve both robust performance and transparency in text classification tasks

V. RESULTS AND DISCUSSION

In this study, we explored the application of Explainable AI (XAI) techniques in the domain of news classification, focusing on three traditional machine learning models: Naive Bayes, Logistic Regression, and Support Vector Machines. While these models demonstrated comparable performance with an overall accuracy of 67% across 15 news categories, the integration of explainability added a critical dimension to our analysis. Using feature importance visualizations and LIME-based explanations, we gained meaningful insights into the internal workings of these models, particularly their reliance on specific words to differentiate between news categories.

Explainable AI not only enhances our understanding of model predictions but also strengthens trust and accountability in AI systems. By providing interpretable justifications for classification decisions, XAI bridges the gap between the model's statistical performance and its real-world applicability. This interpretability is especially crucial in sensitive areas like news classification, where the consequences of algorithmic bias or erroneous predictions can be significant. The insights obtained from explainability frameworks enable stakeholders to audit and refine models, ensuring that their decisions align with ethical principles and societal expectations.

The implications of this work are far-reaching. From an AI ethics perspective, XAI empowers users by offering transparency and reducing the "black box" nature of machine learning models. It supports fairness by revealing potential biases in feature importance and helps mitigate harm by providing clarity around how predictions are made. For instance, in news classification, understanding the reasons behind a model's decision can help identify biases in content curation or prevent the spread of misinformation.

Furthermore, the integration of XAI aligns with the broader goals of responsible AI development, where systems are designed not only to perform well but also to be interpretable, accountable, and inclusive. As AI continues to be applied in high-stakes domains, the need for explainability will become increasingly critical. The methodologies demonstrated in this work—such as feature analysis and instance-level explanations—serve as foundational tools for designing AI systems that are both effective and ethical.

Future directions include extending this work to more advanced models, such as deep learning architectures, which often face greater challenges in interpretability. Additionally, addressing data imbalance and exploring multilingual datasets could enhance the generalizability of the findings. Ultimately, this research underscores the value of explainable AI in fostering trust, improving transparency, and ensuring the ethical deployment of AI systems in news classification and beyond.

REFERENCES

- [1] Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., et al. (2019). Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities, and Challenges Toward Responsible AI. IEEE Access, 7, 121205-121223.
- [2] This paper discusses the concepts and various frameworks within XAI, emphasizing the need for transparency in AI applications across diverse domains.
- [3] Doshi-Velez, F., & Kim, P. (2017). Towards a rigorous science of interpretable machine learning. Proceedings of the 34th International Conference on Machine Learning, 70, 3961-3970.
- [4] The authors outline key principles for developing interpretable machine learning models, focusing on their importance in enhancing user understanding and trust.



International Journal for Research in Applied Science & Engineering Technology (IJRASET)

ISSN: 2321-9653; IC Value: 45.98; SJ Impact Factor: 7.538

Volume 12 Issue XI Nov 2024- Available at www.ijraset.com

- [5] Lundberg, S. M., & Lee, S. I. (2017). A unified approach to interpreting model predictions. In Advances in Neural Information Processing Systems (Vol. 30).
- [6] This paper introduces SHAP (SHapley Additive exPlanations), a method that interprets complex models by assigning importance values to features based on cooperative game theory.
- [7] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining the predictions of any classifier. In Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1135-1144.
- [8] This study presents LIME, an approach that generates local, interpretable explanations for individual predictions made by any classification model.
- [9] BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. (2019). arXiv:1810.04805.
- [10] The BERT model represents a breakthrough in natural language processing by employing a transformer-based architecture that allows for nuanced understanding of text.
- [11] Yang, Z., et al. (2016). Hierarchical Attention Networks for Document Classification. In Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies.
- [12] This paper introduces attention mechanisms in deep learning, providing insights into how models focus on specific parts of input data, relevant for text classification tasks.
- [13] Zhang, Y., & Wang, H. (2020). A Review on Explainable Artificial Intelligence (XAI) in Medical Imaging. Journal of Healthcare Engineering, 2020, Article ID 8858774.
- [14] Although focused on medical imaging, this review discusses various XAI techniques applicable across domains, highlighting their significance in enhancing interpretability.
- [15] Chen, J., Song, L., et al. (2019). Learning to Explain: An Information-Theoretic Perspective on Model Interpretation. In Proceedings of the 36th International Conference on Machine Learning.
- [16] This research investigates the theoretical foundations of model interpretability, proposing methods to enhance the clarity of AI-driven predictions.
- [17] Gilpin, L. H., et al. (2018). Explaining Explanations: An Overview of Interpretability of Machine Learning. ACM Computing Surveys (CSUR), 51(5), 1-42.
- [18] This survey provides an extensive overview of various interpretability methods, their applications, and the implications for AI systems across different fields.
- [19] Friedler, S. A., et al. (2019). A Comparative Study of the Effect of Class Imbalance on Machine Learning Models. Proceedings of the 27th ACM International Conference on Information and Knowledge Management.











45.98



IMPACT FACTOR: 7.129







INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 🕓 (24*7 Support on Whatsapp)