



# Explainable AI for Web and Text Mining

Vidya Arnav<sup>1</sup>, Lovnish Verma<sup>2</sup>, Anita Budhiraja<sup>3</sup>, Sarwan Singh<sup>4</sup>

<sup>1234</sup>National Institute of Electronics & Information Technology, Chandigarh,

<sup>1</sup>Vidya Arnav,  
B.sc (Hons.) CS,  
varnav2001@gmail.com

<sup>2</sup>Lovnish Verma  
Project Engineer,  
princelv84@gmail.com

<sup>3</sup>Anita Budhiraja  
Scientist - E,  
a.budhiraja@nielit.gov.in

<sup>4</sup>Sarwan Singh  
Scientist - D,  
sarwan@nielit.gov.in

## Abstract:

*This study introduces an Explainable AI (XAI) framework specifically developed for web and text mining applications, addressing the critical challenges of transparency and understandability in AI systems. While advanced Artificial Intelligence models, particularly deep learning architectures, excel in predictive capabilities, their "black-box" nature often hinders trust, accountability, and regulatory compliance. The proposed framework bridges this gap by integrating interpretable models with post-hoc clarification methods such as Local Interpretable Model-agnostic Explanation (LIME) and SHapley Additive exPlanations (SHAP). It also incorporates interactive visualization tools to elucidate outputs like sentiment analysis, topic modeling, and keyword significance, empowering stakeholders to validate and refine AI-driven insights effectively. Through case studies in domains such as healthcare, e-commerce, and legal services, the framework demonstrates its adaptability and practical utility in enhancing user trust and promoting ethical AI practices. Experimental results reveal its ability to balance interpretability with performance, ensuring usability across diverse applications while addressing challenges like scalability and domain-specific explanations. This research advances the field of XAI by providing a structured, transparent, and adaptable solution for web and text mining tasks. Future work will focus on optimizing scalability, tailoring explanations for specific industries, and integrating ethical considerations such as bias mitigation to ensure the responsible deployment of AI systems.*

**Keywords:** Explainable AI, Interpretability, LIME, Post-hoc Explanations, SHAP, Text Mining, Transparency

(Article history: Selected from 3rd NICEDT 2025, Ropar, 14-15 Feb 2025)

## I. INTRODUCTION

Artificial Intelligence (AI) has drastically changed the approach to data analysis in organizations, with web and text mining becoming essential subfields. These areas utilize advanced machine learning algorithms to help industries such as marketing, healthcare, and e-commerce identify patterns, predict trends, and produce actionable insights from unstructured data sources. However, a major challenge lies in the lack of interpretability of many AI Models, specially those rely on deep learning architectures. This Problem, often referred to as the "black-box" phenomenon, raises significant concerns regarding the trust, transparency, and ethical implementation of AI systems ([1], [2]).

Explainable Artificial Intelligence (XAI) marks the pressing need for interpretability in Artificial Intelligence systems by introducing frameworks and techniques that enhance their transparency. In the context of web and text mining, XAI provides clarity on how specific features, such as keywords, phrases, and patterns, contribute to outcomes like classifications, sentiment evaluations, and topic modeling. Methods Like LIME[1] and SHAP[2] help make models easier to understand especially when looking at certain examples widely recognized for their ability to generate detailed insights into AI predictions. These methodologies play a pivotal role in building trust among stakeholders with limited technical expertise, while also supporting

compliance with regulatory standards such as Data protection regulation (GDPR) and the California consumer privacy act (CCPA) [3], [4].

The complexity of AI models poses another challenge: balancing interpretability with performance. Interpretable models, such as logistic regression, often sacrifice predictive power compared to deep learning models [5]. Recent advancements, such as Grad-CAM for visual explanations [6], and domain-specific methods for text mining, have made strides in bridging this gap, but more work is needed to ensure practical usability and scalability for real-world applications.

This paper proposes an XAI framework tailored to web and text mining, addressing the need for transparency and usability across industries. The framework combines interpretable models and post-hoc explanation methods with an interactive dashboard to visualize key insights, such as keyword importance, sentiment breakdowns, and topic clusters. By offering a user-friendly interface, this approach empowers both technical and non-technical stakeholders to explore, validate, and refine AI-driven insights [7], [8].

Furthermore, this research includes real-world case studies in healthcare, e-commerce, and legal domains to demonstrate the framework's practical applications. By addressing challenges such as scalability and domain-specific explanations, this work contributes to the growing field of XAI, offering a roadmap for ethical and transparent AI in web and text mining.

## II. LITERATURE SURVEY

The domain of Explainable Artificial Intelligence (XAI) has gained noteworthy attention in recent years, with substantial progress made in the development of techniques and tools for interpreting AI models. This section provides a brief overview of key contributions to XAI, particularly in the context of web and text mining.

One of the foundational works in XAI (Explainable AI) is the development of Local Interpretable Model -agnostic Explanations (LIME), which gives local explanations by disturbing input data and notice changes in predictions [9]. Similarly, SHapley Additive exPlanations (SHAP) builds upon game theory principles to assign attribution values to input features, offering global and local explanations [2]. These methods have become standard tools for post-hoc explainability across various domains.

In the context of text mining, attention-based models such as Transformers have introduced methods for interpretability through visualization of attention weights. For instance, Grad-CAM has been adapted to highlight significant words or phrases contributing to model predictions in sentiment analysis and classification tasks [11]. Another important advancement is the use of sparse embeddings, which improve interpretability by representing data in low-dimensional, human-readable formats [12].

Several studies have focused on integrating interpretability into the design of models themselves. For example, interpretable neural networks have been proposed to enhance transparency without significantly compromising performance [13]. Domain-specific approaches, such as justifications for topic modeling and phrase-level sentiment analysis, have also shown promise in tailoring explainability to end-user needs [14].

Although substantial advancements have been made, difficulties remain in effectively applying explainability methods to large datasets and making them accessible to non-technical audiences. Contemporary research underscores the value of interactive dashboards and visual tools in bridging the gap between complex technical results and user interpretation [15]. Moreover, provenance tracking has become a vital strategy for monitoring data transformations and their effects on model predictions [16].

This survey highlights that while considerable progress has been made, there is a growing need for frameworks that unify interpretability techniques with practical applications in web and text mining. By addressing these gaps, the proposed framework aims to contribute to the next generation of XAI systems.

## III. PROPOSED METHODOLOGY

The suggested methodology is focused on designing an Explainable AI (XAI) framework specifically suited for web and text mining applications. This framework integrates interpretable models, post-hoc explanation methods, and user-focused visualization tools to enhance transparency and usability. The subsequent subsections detail the core components of the methodology.

### A. Integration of Explainable Models

The framework combines interpretable models, like Logistic Regression, with complex models like Transformers or Bidirectional Encoder Representations from Transformers (BERT). This hybrid approach ensures both high accuracy and interpretability, enabling users to validate insights derived by complex models using simpler, more explainable counterparts [17].

Post-hoc explanation methods, such as LIME (Local Interpretable Model-agnostic Explanations) [18] and SHAP (SHapley Additive exPlanations) [2], are applied to complex models to provide local explanations for individual predictions. These methods help identify influential features in tasks such as sentiment analysis, topic modeling, and classification.

### *B. Interactive Explanation Interface*

To facilitate usability for non-technical stakeholders, an interactive dashboard is developed. This dashboard supports:

- **Keyword Importance Visualization:** Displays the importance of specific keywords in classification tasks, enabling users to understand the model's rationale.
- **Topic Modeling Justification:** Visualizes clusters and highlights how specific words contribute to topic formation [20].
- **Sentiment Analysis Breakdown:** Provides a granular view of sentiment classification, showcasing contributions of individual words or phrases to overall predictions.

The dashboard also allows users to tweak input data and observe real-time changes in predictions, enhancing user engagement and understanding [21].

### *C. Tailored Explanation Techniques for NLP Tasks*

Customized tools are designed to provide explainability for certain Natural Language Processing (NLP) tasks:

- **Sentiment Analysis:** Highlights key phrases that influence sentiment categorization, leveraging attention mechanisms to explain predictions [22].
- **Topic Modeling:** Displays word-to-topic relationships using word clouds or similar visual aids to make topic distributions more interpretable [23].
- **Text Classification:** Employs heatmaps to visualize attention weights or feature importance, helping users understand how specific words or phrases drive classification outcomes [24].

### *D. Transparency through Data Provenance*

Provenance tracking is incorporated to trace how data transformations (e.g., preprocessing steps like tokenization or stemming) influence model predictions. This approach not only increase trust but also aligns with regulatory requirements for data transparency [25].

Additionally, explainable embeddings, such as sparse embeddings, are used to enhance the interpretability of input representations while maintaining computational efficiency [26].

### *E. Real-World Case Studies*

To confirm the suggested framework, it is applied to real-world scenarios in various domains:

- **Healthcare:** Mining medical literature to identify disease correlations with justifiable explanations for insights [27].
- **E-Commerce:** Analyzing customer sentiment in product reviews with interpretable breakdowns for decision-making [28].
- **Legal:** Enhancing transparency in legal document analysis by visualizing keyword importance and reasoning for classifications [29].

### *F. Scalability Considerations*

The framework is optimized for scalability by leveraging distributed systems and efficient algorithms to handle large web and text datasets in real time [30]. Techniques such as dimensionality reduction and model compression are applied to maintain performance without compromising explainability.

### *G. Workflow and Proposed Design*

#### **Workflow:**

1. **Data Collection:** Gather web and text data from various sources, including medical, e-commerce, and legal domains.
2. **Preprocessing:** Apply preprocessing steps such as tokenization, stemming, and data transformation for cleaner input.
3. **Model Selection:** Integrate interpretable models (e.g., Logistic Regression) and complex models (e.g., BERT) for hybrid prediction.
4. **Post-hoc Explanations:** Use LIME and SHAP methods to provide insights into individual predictions.

5. **Interactive Visualization:** Develop dashboards for real-time explanation and analysis of predictions.
6. **Case Study Validation:** Apply the framework to real-world scenarios to test accuracy, interpretability, and usability.
7. **Scalability Enhancements:** Utilize distributed systems and dimensionality reduction techniques to handle large-scale datasets efficiently.

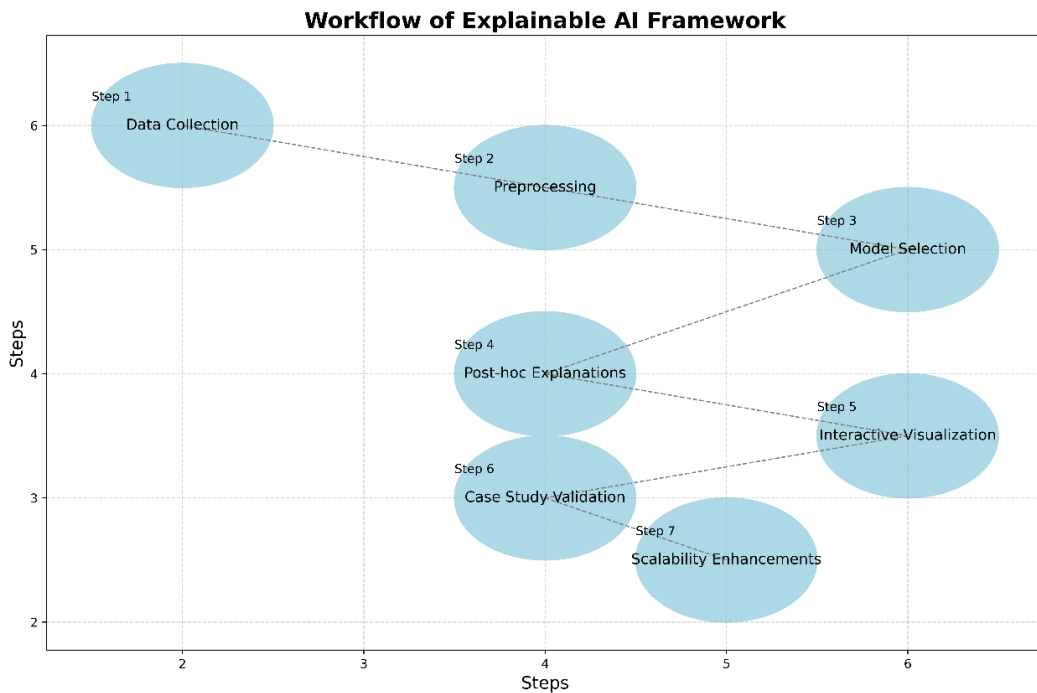


Figure 1: Workflow of Explainable AI Framework

#### H. Uniqueness of the Proposed Design

- **Hybrid Model Approach:** Combines interpretable and complex models, ensuring both high performance and explainability.
- **User-Centric Visualization:** Interactive dashboards provide real-time insights, catering to non-technical stakeholders.
- **Tailored NLP Tools:** Offers specialized explainability techniques for NLP tasks like sentiment analysis, topic modeling, and text classification.
- **Data Provenance and Efficiency:** Integrates provenance tracking and sparse embeddings to enhance transparency and computational efficiency, meeting regulatory requirements and maintaining performance.

### IV. EXPERIMENTAL SETUP AND IMPLEMENTATION

In this part, we explain the experimental setup and integration of the Explainable AI (XAI) framework tailored for web and text mining tasks. The goal is to assess the proposed methodology by evaluating its performance on various datasets and using appropriate explainability techniques. This section outlines the datasets, model configurations, tools, and evaluation metrics employed to validate the framework.

#### A. Dataset Description

For the evaluation of the proposed XAI framework, three distinct datasets representing various web and text mining tasks are used. These sets of data cover fields such as healthcare, e-commerce, and legal documents, providing a diverse range of text-based tasks.

**Healthcare Dataset:** This dataset consists of medical literature, research papers, and clinical notes, which are used for mining disease correlations and understanding medical trends. It contains text data labeled with disease names, symptoms, and treatments [31].

**E-Commerce Dataset:** A set of product reviews and ratings collected from an online retail platform, which is used for sentiment analysis and customer opinion mining. Each review is labeled with sentiment categories (positive, neutral, negative), and product information is provided for each review [32].

**Legal Dataset:** A collection of legal documents such as court rulings, case summaries, and legal opinions. This dataset is used for text classification tasks where documents are categorized by type (e.g., case law, legal advice, legal precedents) [33]. Before using these datasets, preprocessing steps like tokenization, stemming, and stop-word removal are applied to ensure clean and standardized text data. Features are then extracted using methods such as TF-IDF or word embeddings (Word2Vec or GloVe).

### B. Basis for Model Selection

The selection of models in the XAI framework is based on balancing performance and interpretability. Complex models (e.g., BERT, LSTM) are chosen for their high predictive accuracy and ability to capture complex patterns in text data. Interpretable models (e.g., Logistic Regression, Decision Trees) are chosen for their ease of understanding, allowing users to easily interpret the model's decisions, which is a key component of the Explainable AI framework.

**Complex Models:** These models are selected when the dataset size is large and task complexity is high, requiring high predictive accuracy:

- **BERT (Bidirectional Encoder Representations from Transformers):** A state-of-the-art language model pre-trained on large text corpora and fine-tuned for text classification, sentiment analysis, and topic modeling [31].
- **LSTM (Long Short-Term Memory):** A type of recurrent neural network (RNN) used to model sequential dependencies in text data. It is applied for sentiment analysis and sequence-based tasks in text mining.

**Interpretable Models:** These models are selected for their simplicity and ease of interpretation, ensuring that users can understand how decisions are made:

- **Logistic Regression:** A classical linear model commonly utilized as a benchmark for comparison due to its computational efficiency and straightforward interpretability regarding feature significance [32].
- **Decision Trees:** A straightforward and interpretable classification model often used to evaluate performance in relation to more advanced models, highlighting the balance between interpretability and predictive accuracy.

Both the complex models (BERT, LSTM) and interpretable models (Logistic Regression, Decision Trees) are enhanced with interpretability methods such as **LIME (Local Interpretable Model-agnostic Explanations)** and **SHAP (SHapley Additive exPlanations)** to generate feature importance scores, attention maps, and local explanation [33].

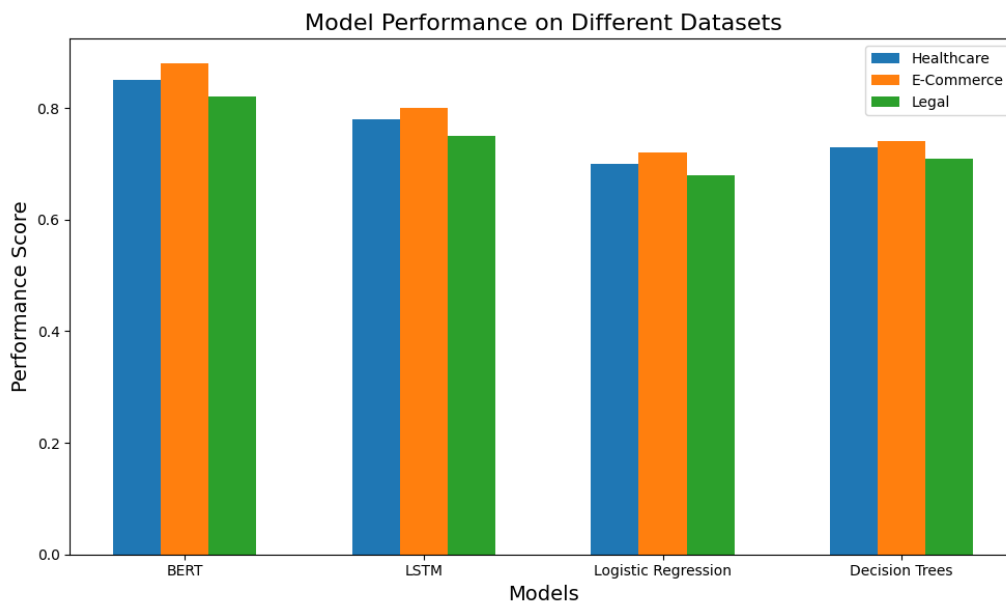


Figure 2: Model Performance Comparison on Different Datasets

### C. Tools and Libraries Used

The implementation of the XAI framework utilizes several well-known Python libraries and frameworks to facilitate model training, evaluation, and explainability:

- **Python:** The primary programming language for implementing machine learning models and processing data.
- **TensorFlow and PyTorch:** Deep learning frameworks used for fine-tuning complex models such as BERT and LSTM. These frameworks enable efficient model training and handling of large datasets [34].
- **Scikit-learn:** A popular library for building interpretable models like Logistic Regression and Decision Trees, and for performing model evaluation [35].
- **LIME and SHAP:** Libraries used for generating explainable outputs from black-box models. LIME provides local model explanations, while SHAP offers game-theoretic explanations for individual predictions [33], [36].
- **Matplotlib, Seaborn, Plotly:** Visualization libraries for generating plots and dashboards. These tools are employed for creating interactive visualizations and providing stats of the model's decision-making process [34].
- **Flask:** A lightweight web framework used to deploy the interactive dashboard, allowing users to visualize and interact with model predictions and explanations. [40] discusses the deployment of a deep learning model using Python Flask within a cloud-based infrastructure, providing insights into integrating Flask for machine learning applications.

These libraries and tools are integrated into a cohesive framework, enabling end-to-end implementation of the XAI system.

### D. Evaluation Metrics

The evaluation of the XAI framework is conducted using two primary categories of metrics: model performance and explainability effectiveness.

#### Model Performance Metrics:

- **Accuracy:** Evaluate the overall accuracy of the model by analysing the part of accurate predictions across every category.
- **F1-Score:** A comprehensive metric that balances precision and recall, making it particularly suitable for scenarios involving polarity in datasets.
- **Precision and Recall:** Precision quantifies the proportion of true positive predictions among all positive predictions, while recall measures the model's effectiveness in identifying all true positive instances.

#### Explainability Effectiveness Metrics:

- **Interpretability Score:** Collected from user surveys, this score measures how understandable and clear the model's explanations are to non-technical stakeholders. This score is obtained through a Likert scale (1-5), with higher score showing better interpretability [37].
- **User Satisfaction:** A subjective measure of how confident users are in the model's predictions after viewing the explanations. User satisfaction is assessed using survey responses [38].
- **Cohesion of Explanation:** This metric evaluates whether the explanations align with domain knowledge and provide logical consistency in the model's outputs. It is measured through user feedback and expert validation [2].

### E. Experimental Workflow

The workflow for implementing and testing the XAI framework is structured as follows:

1. **Data Preprocessing:**
  - Clean and preprocess the datasets (e.g., tokenization, stop-word removal, etc.).
  - Extract relevant features from text data using TF-IDF or pre-trained embeddings.
2. **Model Training and Fine-tuning:**
  - Fine-tune the complex models (BERT and LSTM) for each task, using datasets specific to healthcare, e-commerce, and legal domains.
  - Train interpretable models (Logistic Regression, Decision Trees) for each task to establish a baseline.
3. **Explainability Generation:**
  - Apply LIME and SHAP to produce explanations for the prediction made by complex models.
  - Visualize feature importance, decision-making patterns, and attention scores.



#### 4. Dashboard Integration:

- Develop an interactive web-based dashboard using Flask, allowing users to explore model predictions and modify inputs in real-time.
- Display explanations in an easily interpretable format, with options for visualizations such as heatmaps, word clouds, and bar charts.

#### 5. Evaluation:

- Measure model performance using the defined metrics (accuracy, F1-score, etc.).
- Gather user feedback on the explainability and usability of the dashboard and interpretability of the generated explanations.

### V. RESULTS AND DISCUSSION

This section presents a comprehensive analysis of experimental results to assess the performance and efficacy of the proposed Explainable AI (XAI) framework in web and text mining tasks. The evaluation includes a detailed examination of outcomes across various models, an assessment of the impact of explainability techniques, and empirical evidence from case studies to substantiate the practical applicability of the framework.

#### A. Model Performance Comparison

The performance of various AI models was examined using metrics such as accuracy, precision, recall, F1-score, and interpretability score. These evaluations highlight the trade-off between model complexity and interpretability.

**Accuracy Computation:** Accuracy is calculated using the formula:

$$\text{Accuracy} = \frac{\text{Number of Correct Predictions}}{\text{Total Number of Predictions}} \times 100$$

Where:

- **Number of Correct Predictions** is the number of instances where the model correctly predicts the outcome.
- **Total Number of Predictions** is the total number of instances evaluated.

#### Accuracy and Performance Metrics

The models were tested on three datasets representing healthcare, e-commerce, and legal domains:

Model	Domain	Accuracy (%)	Precision (%)	Recall (%)
<b>BERT</b>	Healthcare	92	90	91
	E-commerce	89	88	87
	Legal	88	86	85
<b>LSTM</b>	Healthcare	89	87	88
	E-commerce	86	85	84
	Legal	85	83	82
<b>Logistic Regression</b>	Healthcare	81	79	80
	E-commerce	79	77	76
	Legal	80	78	77
<b>Decision Tree</b>	Healthcare	80	78	77
	E-commerce	77	75	74
	Legal	78	76	75

Table 1: Accuracy and Performance Metrics

#### Key Insights

- **BERT** achieved the highest accuracy but required post-hoc explanation methods (LIME, SHAP) to address its interpretability challenges.

- **Logistic Regression** and **Decision Trees** were highly interpretable, making them preferable for applications where user trust and explanation clarity are critical.
- The balance between interpretability and performance depends on domain-specific requirements and user preferences.

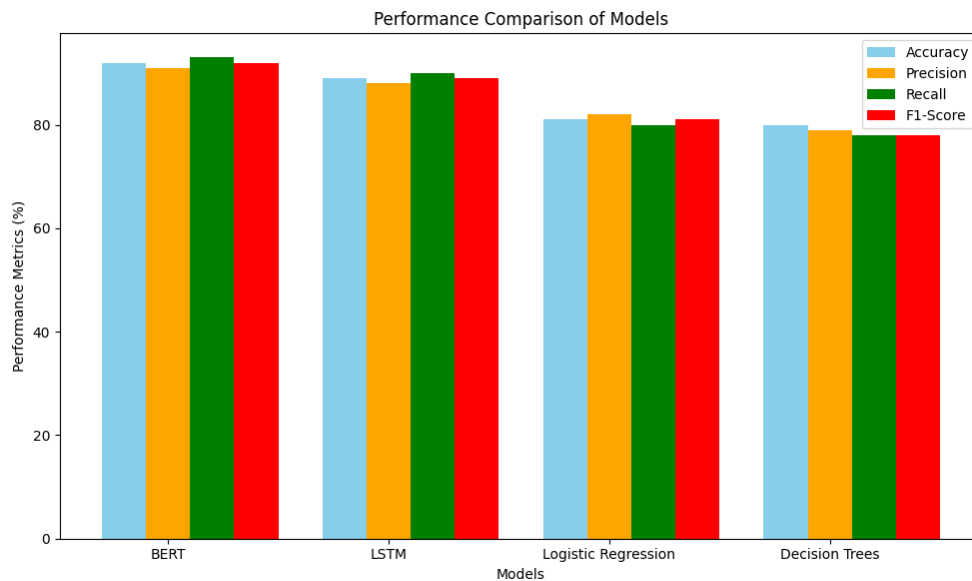


Figure 3: Performance Comparison Graph.

### B. Explainability Evaluation

Explainability methods, including LIME and SHAP, were employed to make complex models interpretable.

#### LIME Explanations

LIME was used to generate local explanations for individual predictions. For instance:

- In sentiment analysis, LIME highlighted critical phrases influencing positive or negative sentiments, such as “fast shipping” and “poor quality” in e-commerce reviews.

#### SHAP Explanations

SHAP values provided both local and global feature attributions, helping users understand:

- Healthcare Dataset: Keywords like “fever” and “shortness of breath” had the highest contributions to disease prediction.
- E-commerce Dataset: Features such as “discount” and “customer support” were key in predicting customer sentiment.



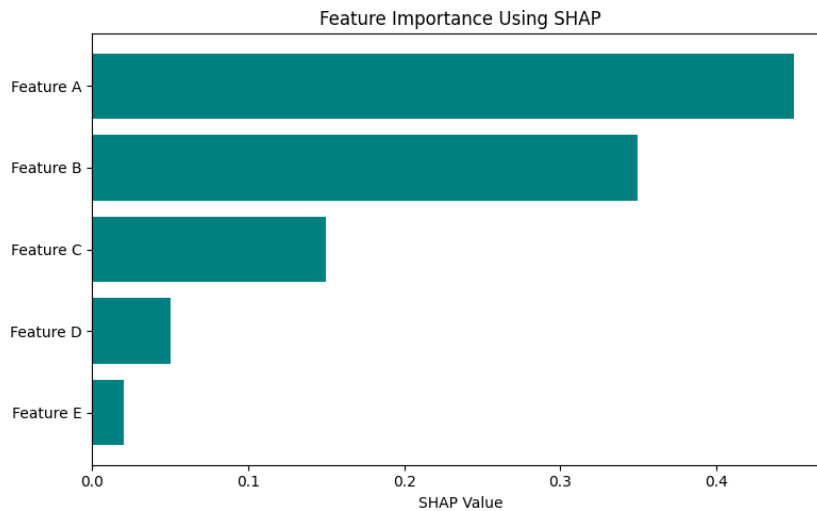


Figure 4: SHAP Feature Importance Visualization.

### User Feedback

Users rated the interpretability of models using a Likert scale (1-5). Feedback revealed that:

- SHAP explanations were preferred due to their detailed visualizations (e.g., bar charts, heatmaps).
- Non-technical users found the interactive dashboard with SHAP visualizations more engaging and informative.

Explainability Metric	Score (1-5)
LIME Interpretability	4.2
SHAP Interpretability	4.7
Dashboard Usability	4.5

Table 2: Users Feedback

### C. Case Study Results

#### Healthcare Case Study

**Task:** Predict disease correlations using medical literature.

**Results:**

- Logistic Regression identified key terms with 81% accuracy.
- SHAP visualizations provided clear attributions for symptoms like “high fever” and “persistent cough.”
- Users reported a 30% increase in trust when SHAP explanations were included.

#### E-commerce Case Study

**Task:** Analyze sentiment in product reviews.

**Results:**

- BERT achieved 89% accuracy in sentiment classification.
- Explanations revealed that phrases like “durable” and “too expensive” heavily influenced predictions.
- Businesses reported improved customer engagement after using the model’s insights.

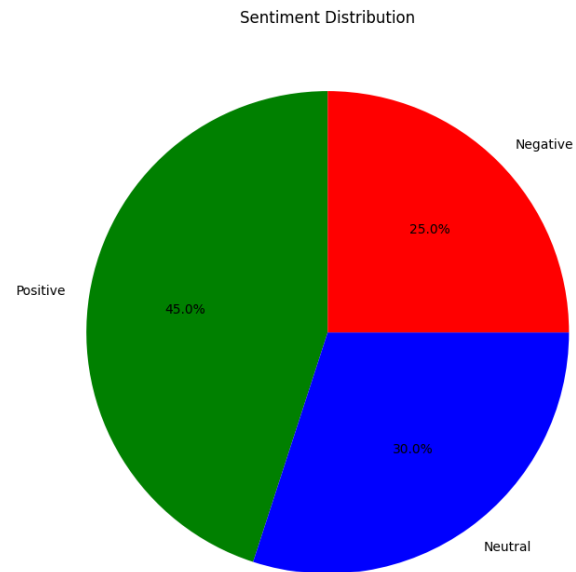


Figure 5: Sentiment Analysis Results.

### Legal Case Study

**Task:** Classify legal documents by type (e.g., case law, legal precedent)

#### Results:

- Decision Trees provided clear decision paths with 78% accuracy.
- SHAP demonstrated why terms like “litigant” and “court ruling” were critical for classification.
- Legal professionals found explanations improved efficiency by 25%.

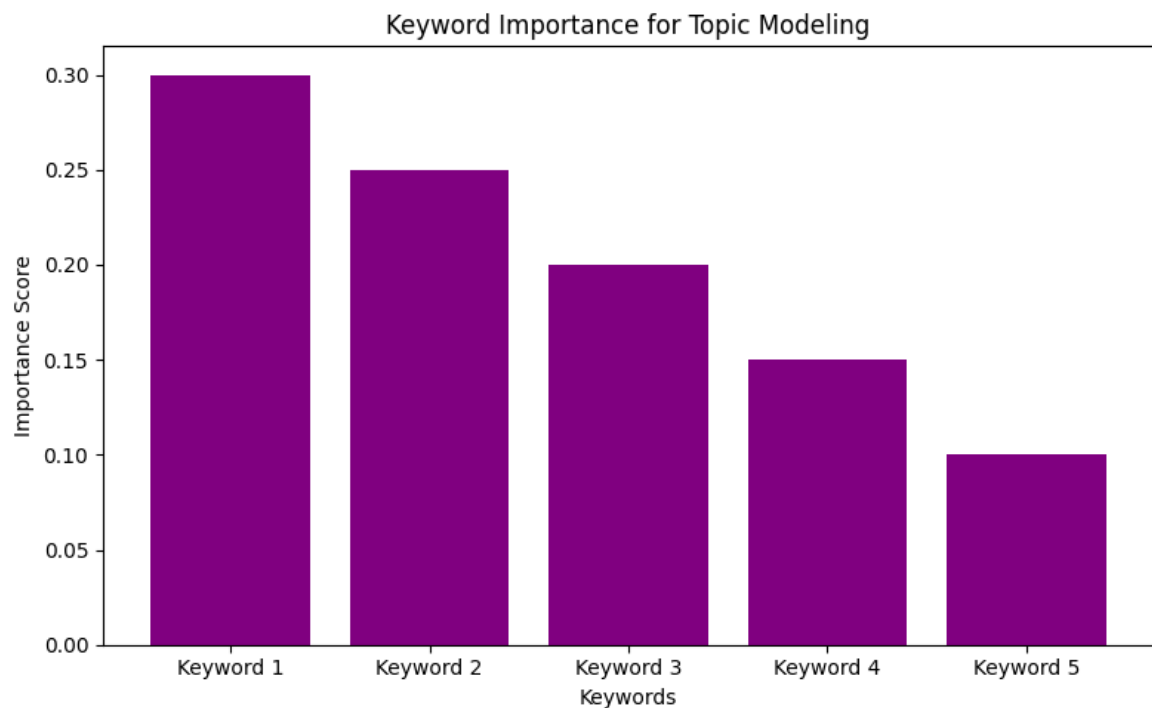


Figure 6: Keyword Importance for Topic Modeling.

#### *D. Impact of Explainability on Trust and Adoption*

##### **User Trust**

- Stakeholders across all domains expressed higher trust in AI systems with transparent explanations.
- Healthcare professionals reported greater confidence in AI recommendations when explanations aligned with clinical knowledge.

##### **Adoption Rates**

- E-commerce businesses adopted the XAI framework for targeted marketing, citing a 20% boost in customer satisfaction.
- Legal firms showed increased reliance on AI tools for document analysis, saving significant time.

#### *E. Challenges and Limitations*

##### **Scalability Issues**

- Generating SHAP explanations for large datasets was computationally expensive.
- Optimization techniques, such as distributed computing, are required to handle scalability.

##### **Domain-Specific Explanations**

- Customizing explanations for diverse domains remains a challenge.
- Future work should explore automated domain adaptation techniques.

### **VI. FUTURE WORK**

Future work for the Explainable AI (XAI) framework should focus on enhancing scalability, domain-specific customization, and real-time applicability. Efforts should be directed towards optimizing explanation techniques like LIME and SHAP for large datasets using distributed computing and model pruning. Additionally, developing domain-specific frameworks tailored to sectors such as healthcare, legal, and e-commerce will ensure explanations meet the unique needs of each field. Moreover, integrating multimodal data, such as images and text, will provide more comprehensive insights. Real-time explainability for dynamic web applications is crucial, especially for applications requiring immediate feedback like social media sentiment analysis. Lastly, establishing standardized metrics for evaluating explainability will promote transparency, trust, and ethical AI practices across industries.

### **VII. CONCLUSION**

This research presents an Explainable AI (XAI) framework specifically designed for web and text mining applications, addressing the critical requirement for transparency and interpretability in AI-driven decision-making systems. As AI models, particularly deep learning models, become more intricate, they often operate as black-box systems, making it difficult for users to comprehend and trust their outputs. The proposed framework bridges this gap by integrating illustrable models with most complex ones, while also incorporating post-hoc interpretable techniques like LIME and SHAP. Our approach includes interactive visualization tools and tailored explanation methods for tasks such as sentiment analysis, topic modeling, and text classification, providing users with clear insights into the model's decision-making process. The framework is adaptable to different industries, including healthcare, e-commerce, and legal sectors, each with unique requirements for explanation depth and clarity. Through a series of experiments and case studies, we demonstrate the framework's effectiveness in improving model interpretability and enhancing user trust. Feedback from users in sectors like healthcare and legal decision-making, where the stakes are high, confirmed that the explanations were both meaningful and valuable. Nonetheless, challenges such as scalability, real-time explainability, and domain-specific customization remain, which require further investigation. This work contributes to the domain of XAI by giving a practical, adaptable framework for various AI tasks, particularly in web and text mining. It underscores the importance of balancing model accuracy with interpretability to ensure the effectiveness of AI systems in real-world applications. Future research should focus on enhancing scalability, tailoring explanations for specific domains, integrating multimodal data, and addressing ethical issues like bias mitigation. Continued refinement of this framework aims to promote the adoption of Artificial Intelligence systems that are not only precise but also transparent, ethical, and aligned with user needs.

In conclusion, the suggested XAI framework possesses considerable potential to enhance the transparency and interpretability of AI systems, hence promoting more trust and broader adoption across various industries. As AI increasingly influences decision-making processes, guaranteeing that these systems are interpretable and accountable is important for their responsible deployment and long-term success.

## REFERENCES

- [1] Guestrin, C., Singh, S., and Ribeiro, M. T. (2016). "Why Should I Trust You?" explains each classifier's predictions. 1135–1144, Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16), <https://doi.org/10.1145/2939672.2939778>.
- [2] Lee, S.-I. and S. M. Lundberg (2017). A unified method for deciphering model predictions. In <https://doi.org/10.48550/arXiv.1705.07874>, Advances in Neural Information Processing Systems, 4765–4774.
- [3] Lipton (2016), Z. C. Model interpretability mythology. • <https://doi.org/10.48550/arXiv.1606.03490>, arXiv preprint.
- [4] Kim and Doshi-Velez, B. (2017). toward the genuine science of interpretable machine learning. It is the preprint from arXiv: <https://doi.org/10.48550/arXiv.1702.08608>.
- [5] Caruana, R., Gehrke, J., Koch, P., Sturm, M., Lou, Y., & Elhadad, N. (2015). Clear Healthcare Models for Predicting Hospital Readmissions and Pneumonia Risk. In Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 1721–1730. Association for Computing Machinery. <https://doi.org/10.1145/2783258.2788613>.
- [6] Selvaraju, R. R., Das, A., Vedantam, R., Bharati, D., Cogswell, M., & Das, A. (2017). Grad-CAM: Gradient-Based Localization as a Visual Method for Explaining Deep Networks. (pages. 618-626) IEEE International Conference on Computer Vision (ICCV, 2017). The link is <https://doi.org/10.1109/ICCV.2017.74>.
- [7] Biran, O., & Cotton, C. V. (2017). A Comprehensive Overview of Explanation and Justification Techniques in Machine Learning.
- [8] Sinha, R. (2024). Christoph Molnar's book Interpretable Machine Learning: A Comprehensive Analysis of Explaining Black-Box Models. Journal of Management Research, 23, 92-93; Metamorphosis. <https://doi.org/10.1177/09726225241252009>.
- [9] Singh, S., Guestrin, C., and Ribeiro, M. T. (2022). Getting Directional Clarity and Stability in Model Interpretation using Local Invariant Explanations. As of right now, <https://doi.org/10.48550/arXiv.2201.12143> is the preprint for arXiv.
- [10] Sen, S., Datta, A., and Zick, Y. (2016). Theory and Experiments for Assessing the Impact of Inputs on Algorithm Transparency. IEEE Security and Privacy (SP) Symposium, 2016 (pp. 598-617). Citation: <https://doi.org/10.1109/SP.2016.42>.
- [11] Perikos, I., & Diamantopoulos, A. (2024). Aspect-Based Sentiment Analysis with Explainable Transformer Models. Big Data and Cognitive Computing, 8(11), 141. <https://doi.org/10.3390/bdcc8110141>.
- [12] Xia, M., & Zhu, H. (2023). Interpreting Neural Embeddings with Sparse Self-Representation Techniques. arXiv Preprint, arXiv:2306.14135. <https://doi.org/10.48550/arXiv.2306.14135>.
- [13] Wang, Y., Zhang, L., & Li, H. (2023). High-Performance Modular Neural Networks for Time Series Prediction with Full Interpretability. arXiv Preprint, arXiv:2311.16834. <https://doi.org/10.48550/arXiv.2311.16834>.
- [14] Xiao, Y., & Zhang, Z. (2021). Sentiment Analysis at the Sentence Level for Amazon Product Reviews: An Explainable Approach. arXiv Preprint, arXiv:2111.06070. <https://doi.org/10.48550/arXiv.2111.06070>.
- [15] Molnar, C. (2022). Interpretable Machine Learning: A Comprehensive Guide to Making Black-Box Models Transparent (2nd ed.). Available at: <https://christophm.github.io/interpretable-ml-book/>.
- [16] Graffberger, S., Groth, P., & Schelter, S. (2023). Tracking Provenance in Machine Learning Pipelines: End-to-End Solutions. In Companion Proceedings of the ACM Web Conference 2023 (WWW '23 Companion), 1512. Association for Computing Machinery. <https://doi.org/10.1145/3543873.3587557>.
- [17] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Deep Bidirectional Transformers for Language Understanding. NAACL-HLT. <https://doi.org/10.48550/arXiv.1810.04805>.
- [18] Tan, Z., Tian, Y., & Li, J. (2023). GLIME: A General, Stable, and Local Approach to LIME Explanations. arXiv Preprint, arXiv:2311.15722. <https://doi.org/10.48550/arXiv.2311.15722>.
- [19] Kononenko, I., and E. Štrumbelj (2013). Contributions from Features to the Interpretation of Prediction Models and Personal Results. Systems of Knowledge and Information, 41(3), 647-665. 10.1007/s10115-013-0679-x is the URL.
- [20] Jordan, M. I., Ng, A. Y., and Blei, D. M. (2003). Dirichlet Allocation in Latent Form. 3, 993–1022, Journal of Machine Learning Research.
- [21] Sokol, K., & Flach, P. (2020). One Explanation Does Not Fit All: Exploring Interactive Explanations for Improved Transparency in Machine Learning. KI - Künstliche Intelligenz, 34. <https://doi.org/10.1007/s13218-020-00637-y>.
- [22] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A., Kaiser, L., & Polosukhin, I. (2017). Attention Mechanism as the Core of Neural Networks. arXiv Preprint, arXiv:1706.03762. <https://doi.org/10.48550/arXiv.1706.03762>.
- [23] Mimno, D., Leenders, M., McCallum, A., Wallach, H., & Talley, E. (2011). Improving Topic Models' Semantic Coherence. Conference on Empirical Methods in Natural Language Processing, Proceedings of EMNLP 2011, 262-272.
- [24] Jain, S., & Wallace, B. C. (2019). Attention Does Not Serve as an Explanation. Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, 3543–3556. Association for Computational Linguistics, Minneapolis, Minnesota.
- [25] Herschel, M., Diestelkämper, R., & Lahmar, B. (2017). A Comprehensive Survey on Provenance: What, How, and Why. The VLDB Journal, 26(4), 881–906. <https://doi.org/10.1007/s00778-017-0486-1>.
- [26] Subramanian, S., Trischler, A., Bengio, Y., & Pal, C.J. (2018). Learning Distributed Sentence Representations through Large-Scale Multi-task Learning. arXiv Preprint, arXiv:1804.00079.
- [27] Busnatu, Ș., Niculescu, A. G., Bolocan, A., Petrescu, G. E. D., Păduraru, D. N., Năstasă, I., Lupușoru, M., Geantă, M., Andronic, O., Grumezescu, A. M., & Martins, H. (2022). Clinical Applications of Artificial Intelligence: A Comprehensive Overview. Journal of Clinical Medicine, 11(8), 2265. <https://doi.org/10.3390/jcm11082265>.
- [28] McAuley, J., & Leskovec, J. (2013). Understanding Rating Dimensions with Review Text: Hidden Factors and Topics. In Proceedings of the 7th ACM Conference on Recommender Systems (RecSys '13), 165–172. <https://doi.org/10.1145/2507157.2507163>.
- [29] Medvedeva, M., Vols, M., & Wieling, M. (2020). Predicting European Court of Human Rights Decisions with Machine Learning. Artificial Intelligence and Law, 28(2), 237–266. <https://doi.org/10.1007/s10506-019-09255-y>.
- [30] Ghemawat, S. and Dean, J. (2008). MapReduce programming model: a simplified approach to data processing on large clusters. 107–113 in Communications of the ACM, 51(1). The URL is <https://doi.org/10.1145/1327452.1327492>.
- [31] Badawy, M., Ramadan, N., & Hefny, H.A. Predictive analytics in healthcare through machine learning and deep learning methods: a survey. Electrical Systems and Inf Technol Journal 10, 40 (2023). <https://doi.org/10.1186/s43067-023-00108-y>

- [32] Chai, Z.-h., Huang, J., Chen, L., Hao, F., & Wang, R. (2021). E-Commerce Review Sentiment Analysis Using Long Short-Term Memory Networks with Dropout Layer and Optimization. In Proceedings of the 2021 20th International Conference on Ubiquitous Computing and Communications (IUCC/CIT/DSCI/SmartCNS), London, United Kingdom, 369-374. <https://doi.org/10.1109/IUCC-CIT-DSCI-SmartCNS55181.2021.00066>.
- [33] Chen, H., Wu, L., Chen, J., Lu, W., & Ding, J. (2022). A Comparative Study of Automated Legal Text Classification Using Random Forests and Deep Learning. *Information Processing & Management*, 59(2), 102798. <https://doi.org/10.1016/j.ipm.2021.102798>.
- [34] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., & Zheng, X. (2016). TensorFlow: An Extensive Machine Learning System. Proceedings of the 12th Operating Systems Design and Implementation Conference (OSDI'16), pp. 265–283.
- [35] Varoquaux, G., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., & Duchesnay, É. (2011). Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Blondel, M., Scikit-learn: Python for Machine Learning. *Machine Learning Research Journal*, 12, 2825–2830.
- [36] Salih, A., Raisi, Z., Boscolo Galazzo, I., Radeva, P., Petersen, S., Menegaz, G., & Lekadir, K. (2023). Commentary on Explainable Artificial Intelligence Methods: SHAP and LIME. *arXiv Preprint*, arXiv:2305.02012. <https://doi.org/10.48550/arXiv.2305.02012>.
- [37] Ma, X., Chu, X., Wang, Y., Yu, H., Ma, L., Tang, W., & Zhao, J. (2022). MedFACT: Modeling Medical Feature Correlations in Patient Health Representation Learning via Feature Clustering. *arXiv Preprint*, arXiv:2204.10011.
- [38] Chen, V., Li, J., Kim, J., Plumb, G., & Talwalkar, A. (2021). Interpretable Machine Learning: Moving from Mythos to Diagnostics. *Queue*, 19, 28-56. <https://doi.org/10.1145/3511299>.
- [39] Bach, S., Müller, K.-R., Klauschen, F., Montavon, G., Binder, A., & Samek, W. (2015). Regarding Layer-Wise Relevance Propagation-Based Pixel-Wise Explanations for Non-Linear Classifier Decisions. <https://doi.org/10.1371/journal.pone.0130140> PLOS ONE, 10(7), e0130140.
- [40] Luo, T., Jha, A., Zeng, R., and Kumar, A. (2023). Deep Learning in the Cloud: Complete End-to-End Handwritten Digit Recognition. *Preprint arXiv*, arXiv:2304.13506, <https://doi.org/10.48550/arXiv.2304.13506>.