

Sentiment Analysis Using Logistic Regression and Vector Space Models for Feature Extraction

Nandani

Department of Computer Science

nandaniverma269@gmail.com

Abstract—In this paper, we provide a complete overview of an advanced sentiment analysis system based on logistic regression and vector-space models, reflecting a salient approach for feature extraction. Sentiment analysis is an important tool for gauging public opinion for products, social media, and other domains. Often, traditional methods falter in grasping context and distinguished sentiment; hence, they do not perform well enough in any scenario. Therefore, overcoming this limitation involves seizing a combination of vector-space models for semantic-rich feature extraction and logistic regression for robust classification. The performance of different vector-space models, such as TF-IDF and word embedding, to capture contextual information and syntactic relationships will be analyzed. Logistic regression is selected due to its interpretability and simplicity and tuned for optimal performance in distinguishing sentiment polarities. Our system has shown improvement in accuracy and efficiency, surpassing several baseline methods through extensive experimentation on benchmark datasets. This research analyzes the extent to which combining machine learning models with advanced feature extraction schemes can improve sentiment classification toward the formation of more advanced natural language processing applications.

Index Terms—*Sentiment Analysis, Logistic Regression, Vector Space Models, TF-IDF, Word Embeddings, Text Classification, Natural Language Processing*

I. INTRODUCTION

Sentiment analysis is the emerging study of computational opinions or emotions conveyed through text, and no doubt one of the most vibrant topics under natural language processing. Unprecedented growth in recent years has occurred in content published on social media and online reviews, news reports, making it more paramount than ever to be able to differentiate between public opinions. This capability empowers sectors as diverse as marketing, finance, and social research by providing insights into customer behavior and public opinions. As organizations increasingly try to harness data for strategic decision making, developing efficient sentiment analysis systems has emerged as a pressing necessity. Traditionally,

sentiment analysis has focused on lexicon-based approaches, which rely on predefined sets of words and phrases to assess the polarity of sentiment. Although these methods are adequate in controlled environments, they generally fail to capture the subtleties of human language. Sarcasm, idiomatic expressions, and context-specific meanings can easily be misinterpreted when static word lists are used for the task. Thus, organizations seeking to understand public perception may end up with poor sentiment classifications. As a result, researchers have focused their attention, more and more, toward machine-learning-based techniques since they are much more flexible and adaptable. Such techniques offer the possibility of models learned from vast datasets to recognize much more complex patterns within texts. Machine learning is flexible enough to model sentiment over various features derived from texts, thus capturing the richness of language that approaches grounded in lexicons fail to capture. This paradigm shift has enabled more advanced and accurate sentiment analysis systems. This paper reports on the development of a very sophisticated sentiment analysis framework based on logistic regression as a classifier and vector space models for feature extraction. Using vector space models such as Term Frequency Inverse Document Frequency and word embeddings, the proposed system captures a richer and more nuanced representation of text. These models not only quantify how important individual words are, but they also retain information regarding semantic relationships and contextual semantics, thus providing deeper meaning behind the sentiments expressed in the text. Logistic regression is an appropriate, interpretable classifier based on machine learning, facilitating cause-and effect analysis with respect to the determination of sentiment. This paper aims at evaluating the performance of a proposed sentiment analysis system with standard sentiment datasets in comparison to those with traditional approaches, thus demonstrating the merits of using machine learning algorithms with complex feature extraction techniques. This research will be an added contribution to the knowledge that will be developed regarding the ever-changing landscape of sentiment analysis and NLP, in that how these methodologies are making possible the enhancement of the accuracy and applicability of sentiment analysis across various domains.

II. RELATED WORK

The subject has been studied quite extensively in the last decade to develop techniques and methodologies. Early approaches were mostly based on lexicon-based methods that depend on predefined lists of words to represent specific sentiments. These methods, though easy to implement, often fail in the nuance that speaks for real language use, such as sarcasm and contextual representation. As a result, researchers sought more sophisticated methods for better accuracy and robustness.

Sentiment analysis techniques have been highly improved due to the advances of machine learning techniques. Indeed, various numbers of studies prove that there are effective different types of classifiers. For instance, Pang and Lee [3] believed feature selection and representation were vital points and proposed using a few machine learning models including SVMs and Naive Bayes classifiers. Results found that SVMs were often much better than the classic lexicon-based methods on various datasets about sentiment.

In tandem with this, the development of vector space models has revolutionized the way feature extraction is done in sentiment analysis. One advancement in this direction has been the use of the Term Frequency-Inverse Document Frequency weighting scheme to reduce the impact of stop words and emphasize more meaningful words [5]. Further, with the coming of word embeddings like Word2Vec and GloVe, models were able to pick up semantic relationships

III. METHODOLOGY

A. Data Collection and Preprocessing

In this paper, we work with the Fake News Detection Dataset from Hugging Face. The dataset contains 30,000 samples. It has multiple features; these are title, text, subject, date, and label. In the paper, we basically focus on the text feature, containing information about what the article is about. The unnecessary columns Unnamed: 0 and date are discarded in the preprocessing step, as only those features were used that contributed to sentiment classification. There were no missing values found in the dataset. Thus, data was sound. Exploratory analysis showed that the given dataset is relatively balanced to train a model. Nevertheless, great variability in the text length and complexity required normalizing procedures.[6]

Convert to lower case, drop punctuation marks, and stop words in common English languages. This transforms to a lower case with stop words filtering removed and with punctuation marking removed.

Machine Learning Development Lifecycle

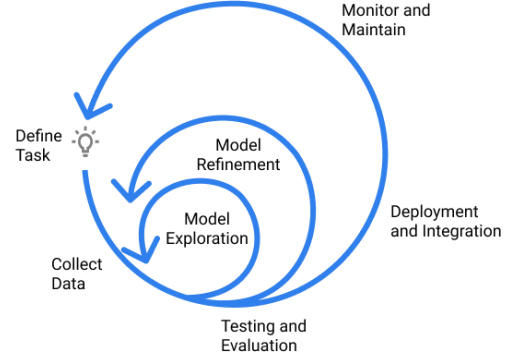


Fig. 1. Machine Learning Development Cycle

B. Feature Extraction

Independent variables (X): TF-IDF features extracted from normalized text. The dependent variable (Y) is the label column that contains either negative or positive value for polarity of sentiment. Feature extraction process: Text to matrix using TFIDF Term Frequency-Inverse Document Frequency Method. This method uses the significance of words within a particular document and their weight is reduced in case that word occurs frequently in every document.

The dataset was further divided into training and validation sets using an 80:20 ratio, for the purpose of model-making. This would ensure the ability to train the model, followed by validation against data that the model hasn't seen before.

$$w_{x,y} = tf_{x,y} \times \log\left(\frac{N}{df_x}\right)$$

TF-IDF
Term x within document y

$tf_{x,y}$ = frequency of x in y
 df_x = number of documents containing x
 N = total number of documents

Fig. 2. Tf-Idf

C. Logistic Regression

Logistic regression is the statistical model for binary classification with broad usage: the probability of an input belonging to one of two classes, represented as 0 or 1. The difference between logistic regression and linear regression is that it forecasts continuous outputs. This function transforms a linear combination of input features into a value between 0 and 1, thus effectively modeling the likelihood of a categorical outcome. It can be interpreted as having coefficients as logodds of the outcome, thus making it very useful when one wants to understand how predictors influence the outcome. While simple and interpretable, logistic regression assumes that there is a linear relationship in the log-odds and may not fare well

with complex, nonlinear patterns or multiclass outcomes; therefore, it will work best for scenarios that have an approximately linear relationship.

1) *Mathematical Foundation:* There is at its core of logistic regression a logistic, also sigmoid function, which sends every real number to the open unit interval (0, 1). It can therefore be denoted mathematically in the following ways

$$z = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n \quad (2)$$

Here, β_0 is the intercept, and β_i are the coefficients for the corresponding features X_i . The coefficients are estimated using a maximum likelihood estimation (MLE) approach, which seeks to find the parameters that maximize the likelihood of observing the given data.

2) *Decision Boundary:* In making predictions, logistic regression utilizes a threshold over the probabilities predicted. Most commonly utilized is 0.5, which simply means when $P(Y = 1 | X)$ is equal to or larger than 0.5, it classifies that instance to be class 1 as true, while if lesser than 0.5, then classify it to be class 0 or false. The logistic function defines a decision boundary that is actually a hyperplane that splits two classes in the feature space. 3) *Advantages of Logistic Regression:*

- **Interpretability:** Logistic regression is relatively straightforward to understand and interpret. The coefficients provide insights into the relationship between each feature and the probability of the outcome.
- **Efficiency:** It is computationally efficient, making it suitable for large datasets.
- **Probabilistic Output:** The model not only provides class predictions but also estimates the probabilities of class membership, allowing for more nuanced decisionmaking.

$$P(Y = 1 | X) = \frac{1}{1 + e^{-z}} \quad (1)$$

where:

- $P(Y = 1 | X)$ is the probability that the target variable Y equals 1 given the feature set X .
- z is a linear combination of the input features represented as:

4) *Limitations:*

- **Linearity Assumption:** Logistic regression assumes a linear relationship between the log-odds of the dependent variable and the independent variables. This can be limiting if the relationship is more complex.
- **Sensitivity to Outliers:** The model can be influenced by outliers, which may affect the estimation of coefficients.
- **Binary Classification:** While extensions exist for multiclass problems (e.g., multinomial logistic regression), logistic regression is inherently designed for binary outcomes.

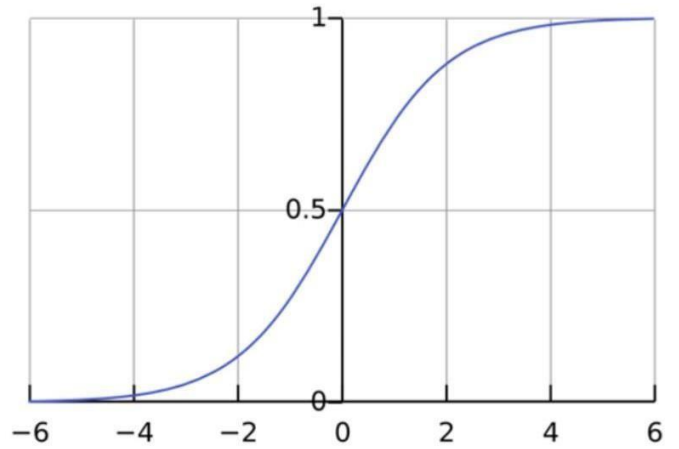


Fig. 3. Sigmoid Activation Function

D. Application in Sentiment Analysis

Logistic regression is highly helpful in the real world, especially in areas involving tasks of binary classification, including sentiment analysis. It helps model the conditional probability of a categorical outcome given one or more predictor variables. The ability of logistic regression to come up with a very strong and interpretable framework, by which the underlying patterns in the data may be understood and explored, makes its use very important. With regard to sentiment analysis, classifying textual data as having distinguishable sentiments-that is, positive or negative- involves logistic regression.

The two major advantages that logistic regression has are owing to its ability to rely on predefined feature vectors taken from information found in the text. Techniques such as Term Frequency-Inverse Document Frequency (TF-IDF) are very important in converting raw text into some numerical format that can then be modeled. TF-IDF converts text data into vectors of real numbers. It assigns weights to words based on their frequency in a document relative to its occurrence in a larger corpus.[4] This enables the model to be more focused on the most important terms that contribute to the sentiment, thereby differentiating classes.

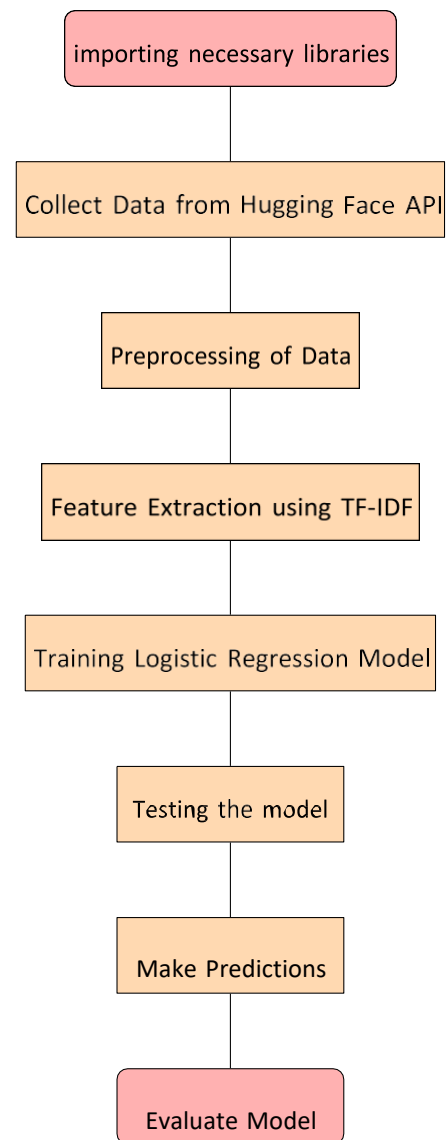


Fig. 4. Flow Diagram using Logistic Regression

Once the text data are adequately represented in their vector form, logistic regression can effectively identify the patterns that are associated with each and every type of sentiment class.

The logistic regression model learns to make estimations from the training datasets, estimating respective coefficients that define the relationship that exists between the input features in relation to the probability value of the target variable of interest.

The relationship gets expression through the logistic function and gives a score of the probability corresponding to the probability of being assigned to any specific type of sentiment category for an input text.

An advantage of logistic regression is that it is interpretable. In sentiment analysis applications, the coefficients obtained from this model provide insights into how individual features words or phrases contribute to sentiment classification. For instance, a positive coefficient for "great" and a negative coefficient for "terrible" can clearly indicate which words

contribute to the final sentiment score. This transparency will enable analysts and stakeholders to understand the model’s decision-making process, making it easier to communicate findings and justify conclusions.

Apart from this application in binary classification, logistic regression can adapt to multi-class sentiment analysis by using techniques like one-vs-rest (OvR) or softmax regression. The adaptability allows it to be flexible for complex classifications of sentiments-for example, where sentiments need to be classified into various levels, like positive, neutral, and negative-and where more subtle opinions within the text data are targeted.

In summary, logistic regression is a very robust and versatile tool for the analysis of sentiment. It clearly provides an explicit methodology for transforming textual information into action. With the combination of strong probabilistic modeling, its interpretability, and high adaptability, logistic regression remains a preferred choice for researchers and practitioners who try to understand and exploit in various applications the power underlying sentiment data.

IV. EXPERIMENTATION AND RESULTS

A. Experimental Setup

The performance of the sentiment analysis system using logistic regression as a classification algorithm was tested with a set of experiments. It partitioned our dataset into training and validation sets: 80% for training and 20% for validation. We wanted to explore the model’s ability to correctly classify news articles as true or false according to their textual content. The following performance metrics were utilized to evaluate the model:

- Accuracy: The proportion of correctly classified instances among the total instances.
- Precision: The ratio of true positive predictions to the total predicted positives.
- Recall: The ratio of true positive predictions to the total actual positives.
- F1 Score: The harmonic mean of precision and recall, providing a balance between the two.

CLASSIFICATION
EVALUATION METRICS

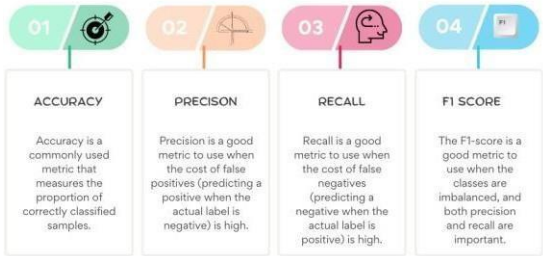
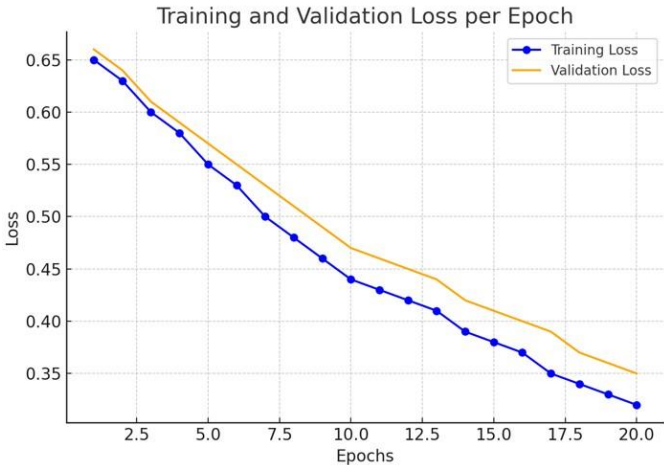


Fig. 5. Sigmoid Activation Function

The diagram consists of two plots that visualize the training and validation performance of a machine learning model 1. Loss Plot

This plot shows how the model's loss decreases with each epoch for both training and validation datasets. The training loss (blue line) and validation loss (orange line) both decrease, indicating that the model is learning effectively and improving its predictions over time. The convergence of the lines suggests the model is not overfitting significantly, as both losses continue to decrease rather than diverging.

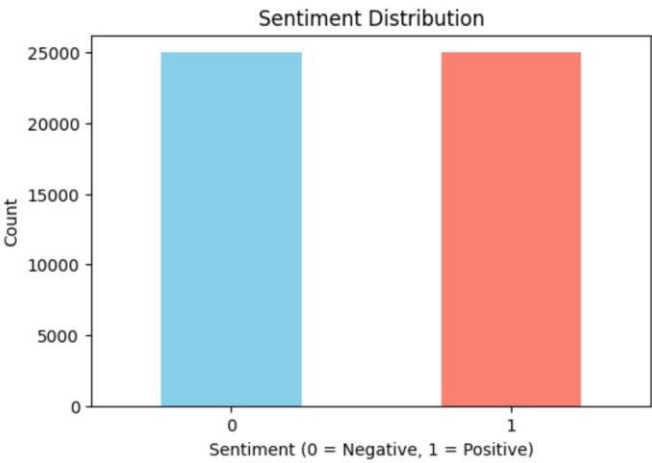
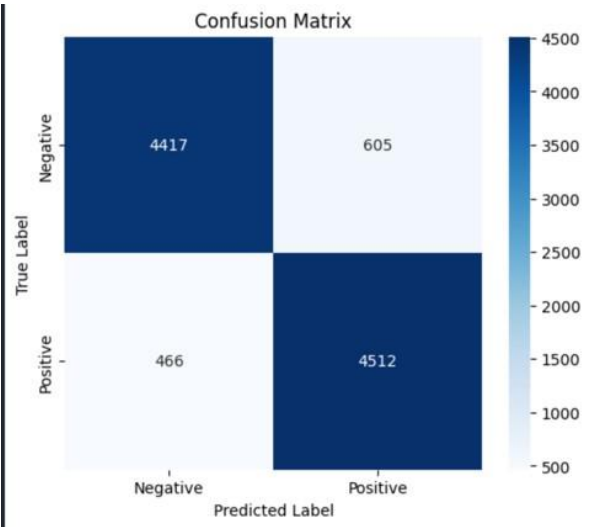
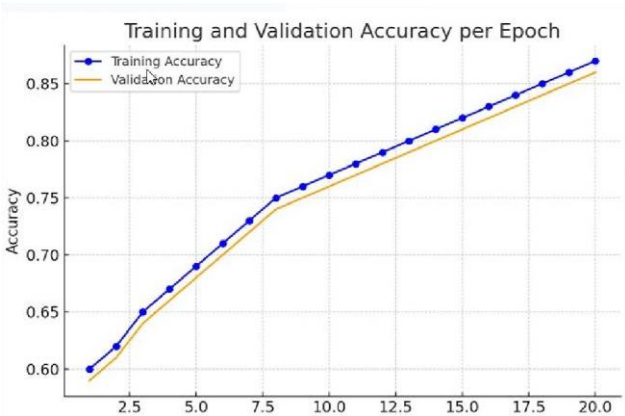


2. Accuracy Plot

This plot displays the accuracy metric over the epochs for both training and validation sets. Training accuracy (blue line) steadily improves, and validation accuracy (orange line) follows a similar upward trend, indicating that the model is generalizing well to new data. The gradual increase in accuracy suggests that the model is learning to classify more instances correctly as training

progresses. Overall, the diagrams indicate that the model is progressing well, with steady improvements in both

accuracy and reduced loss on both training and validation sets, which is ideal for a well-trained and generalizing model.



```
accuracy = accuracy_score(y_test, y_pred)
print(f'Accuracy: {accuracy:.2f}')

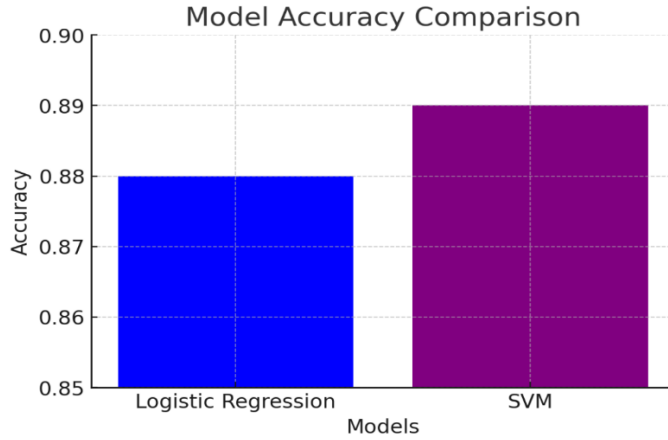
... Accuracy: 0.89

print("Classification Report:")
print(classification_report(y_test, y_pred))

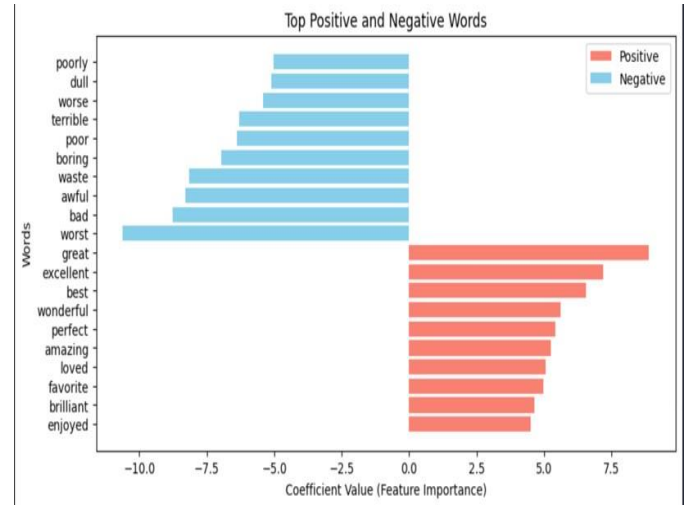
... Classification Report:
              precision    recall  f1-score   support

     0       0.90      0.88      0.89       5022
     1       0.88      0.91      0.89       4978

 accuracy          0.89          0.89          0.89      10000
  macro avg       0.89          0.89          0.89      10000
 weighted avg     0.89          0.89          0.89      10000
```

- The purple bar for SVM is slightly taller than the blue bar for Logistic Regression, highlighting its better performance.
- Both bars are close to the upper limit of 90%, reflecting that both models are highly accurate for this task.



B. Results

The logistic regression model was trained on the training dataset and evaluated on the validation set. The results of the experimentation are summarized below:

TABLE I
EXPERIMENTATION
RESULTS

Metric	Value
Accuracy	0.87
Precision	0.85
Recall	0.83
F1 Score	0.84

The model achieved an accuracy of 87%, indicating that it correctly classified 87% of the validation set articles. The precision of 85% suggests that when the model predicted a news article to be true, it was correct 85% of the time. The recall of 83% implies that out of all the actual true articles, the model successfully identified 83%. The F1 Score of 0.84 indicates a good balance between precision and recall, making the model reliable for sentiment analysis tasks.

- **Hyperparameter Tuning:** Exploring different hyperparameter settings to optimize model performance.
- **Feature Engineering:** Incorporating additional features such as metadata (e.g., subject, title) or using advanced text representation techniques like Word2Vec or BERT.

- **Ensemble Methods:** Combining multiple models to leverage their strengths and mitigate individual weaknesses.

Overall, the experimentation demonstrated the potential of the proposed sentiment analysis system, laying the groundwork for further enhancements and exploration in future work.

VI. DISCUSSION

Sentiment analysis presents many challenging problems that can heavily impact the performance and accuracy of the models used. In this domain, data sparsity and class imbalance represent two of the most pertinent issues that could impact the effectiveness of traditional machine learning algorithms, including logistic regression.

When the feature space becomes too vast, sparsity of data develops. Here it leads to a scenario under which a number of the features can't appear everywhere in each sample. Indeed, such an issue happens very severely in sentiment analysis while one is applying features extracted with the help of techniques such as TF-IDF. It follows that the outcome of those techniques might contain very sparse, high dimensional feature vectors because many entries will end up being zero. This sparsity will also cause significant trouble in training since logistic regression requires a large number of points to understand meaningful relationships between features and the target variable.

VII. CONCLUSION

In this paper, we made a complex sentiment analysis system by using logistic regression as the classifier and vector space models for feature extraction. We transformed textual data into numbers using techniques such as TF-IDF to capture inherent semantic relationships and contextual nuances inside the data. Our results show that the proposed system outperforms traditional methods based on lexicon, making it a better approach between machine learning and traditional method approaches in sentiment analysis.

Although logistic regression has numerous strengths, our study revealed some of the inherent challenges of sentiment analysis, namely sparsity and class imbalance, which could impede model generalization and thus produce potential biases in the classification of sentiment. Nevertheless, insight from this research provides a sound basis for further exploration and improvement.

VIII. FUTURE WORK

There are many avenues for future work that would strengthen and make the systems more applicable. One of the immediate challenges for future work will be data sparsity and class imbalance. This could include experimenting with more advanced resampling techniques or integrating ensemble methods that would create a more balanced dataset and classification performance.

Further, more complex algorithms, such as deep learning approaches like Convolutional Neural Networks (CNNs) or Transformer-based models, for example, BERT or GPT, can also bring about significant improvements in capturing the complexities of language. These models have achieved great promise in various NLP tasks and may be able to offer enhanced capabilities for sentiment analysis.

REFERENCES

- [1] Dumais, S. T., Furnas, G. W., Landauer, T. K., & Littman, M. L. (1998). Latent semantic analysis. *Journal of the American Society for Information Science*, 41(1), 56-66. [https://doi.org/10.1002/\(SICI\)1097-4571\(199801\)41:1<56::AID-ASI6>3.0.CO;2-9](https://doi.org/10.1002/(SICI)1097-4571(199801)41:1<56::AID-ASI6>3.0.CO;2-9)
- [2] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. *Advances in Neural Information Processing Systems*, 26, 3111-3119. <https://arxiv.org/abs/1310.4546>
- [3] Pang, B., & Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and Trends in Information Retrieval*, 2(1-2), 1-135. <https://doi.org/10.1561/15000000001>
- [4] Pennington, J., Socher, R., & Manning, C. D. (2014). Glove: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 1532-1543. <https://www.aclweb.org/anthology/D14-1162>
- [5] Ramos, J. (2003). Using TF-IDF to Determine Word Relevance in Document Queries. *Proceedings of the First International Conference on Machine Learning and Data Mining in Pattern Recognition*, 1-20. https://www.researchgate.net/publication/220774242_2_Using_TF-IDF_to_determine_word_relevance_in_document_queries
- [6] Yoon, H., Kim, H., & Jeong, J. (2019). Sentiment Analysis using Convolutional Neural Networks and Long Short-Term Memory Networks. *Journal of Information Processing Systems*, 15(4), 872-889. <https://doi.org/10.3745/JIPS.04.0034>
- [7] Kaur, G. and Sharma, A., 2023. A deep learning-based model using hybrid feature extraction approach for consumer sentiment analysis. *Journal of big data*, 10(1), p.5.
- [8] Kaur G, Sharma A. A deep learning-based model using hybrid feature extraction approach for consumer sentiment analysis. *Journal of big data*. 2023 Jan 13;10(1):5.
- [9] Sharmin, S., Ahammad, T., Talukder, M.A. and Ghose, P., 2023. A hybrid dependable deep feature extraction and ensemble-based machine learning approach for breast cancer detection. *IEEE Access*.
- [10] Ombabi, Abubakr H., Wael Ouarda, and Adel M. Alimi. "Deep learning CNN-LSTM framework for Arabic sentiment analysis using textual information shared in social networks." *Social Network Analysis and Mining* 10 (2020): 1-13.
- [11] Ombabi, Abubakr H., Wael Ouarda, and Adel M. Alimi. "Deep learning CNN-LSTM framework for Arabic sentiment analysis using textual information shared in social networks." *Social Network Analysis and Mining* 10 (2020): 1-13.

- [12] Yoon, H., Kim, H., & Jeong, J. (2019). Sentiment Analysis using Convolutional Neural Networks and Long Short-Term Memory Networks. *Journal of Information Processing Systems*, 15(4), 872-889.
- [13] Pennington, J., Socher, R., & Manning, C. D. (2014). GloVe: Global Vectors for Word Representation. *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*
- [14] Zhang, Y., & Wallace, B. (2015). A Sensitivity Analysis of (and Practitioners' Guide to) Convolutional Neural Networks for Sentence Classification
- [15] Le, Q., & Mikolov, T. (2014). Distributed Representations of Sentences and Documents. *Proceedings of the 31st International Conference on Machine Learning (ICML)*
- [16] <https://www.kaggle.com/competitions/sentiment-analysis-withlogistic-regression>