Name: Vidya Praveen

Date: 08/03/2024

Course: DSC680 – Applied Data Science

Assignment: Week8 –Milestone2 – Basic Report / Draft of White Paper

**Project Topic: Sentiment Analysis of Amazon Product Reviews.**

**Business Problem**

Understanding customer sentiment is crucial for businesses to enhance their products and services. By analyzing
Amazon product reviews, we can gain insights into customer satisfaction, identify areas of improvement,
and predict future product performance. This analysis will aid in strategic planning, marketing, and
customer relationship management.

**Background/History**

Customer reviews have become a valuable source of feedback for businesses. Traditionally, sentiment analysis
involved manual review, which is time-consuming and subjective. With advancements in data science and
natural language processing (NLP), automated sentiment analysis allows for faster, more objective, and
scalable insights into customer opinions.

**Data Source**

**Datasets**:

The dataset that I have used for this project is project is the Comprehensive Customer Reviews Database,
obtained from Kaggle. Online Product Reviews Analysis: Customer Feedback (kaggle.com)

This dataset includes:

- Review ID

- Product ID

- Reviewer ID

- Rating

- Review Title

- Review Text

- Review Date

**Data Preparation**

- Handling Missing Values: Missing values were dropped from the dataset to ensure data quality.

- Text Normalization: The review text was preprocessed by converting to lowercase, tokenizing, removing stop words, punctuation, and applying lemmatization.

**Methods:**

I have employed several machine learning techniques:

**Data Preprocessing:** I have performed the following tasks:

- Tokenization: Splitting text into individual words.
- Stop Words Removal: Eliminating common words that do not contribute to sentiment.
- Lemmatization: Reducing words to their base form.

**Exploratory Data Analysis (EDA):** I then analyzed the distribution of ratings and word counts in reviews. And then Identified frequently occurring words and phrases.

**Sentiment Analysis:** I have categorized the reviews as Positive, Neutral or Negative, based on their ratings.

**Feature Engineering:** I have used TF-IDF to vectorize the text data for baseline models and tokenization followed by padding was used for the LSTM model.

**Model Building:** Built and evaluated the following models:

- **Baseline Models:** Logistic Regression and Naive Bayes.

- **Advanced Models:** LSTM (Long Short-Term Memory) network for capturing temporal dependencies in the review text.

**Model Evaluation**: Evaluated models using accuracy, precision, recall, F1-score, and ROC-AUC metrics.

**Visualization**: I have created visualizations to show the sentiment trends over time and across different product categories, displayed sentiment distribution for top-selling products and Presented model predictions versus actual sentiments using confusion matrices and other relevant charts.

**Analysis**

The sentiment analysis of Amazon product reviews revealed that the LSTM model outperformed the baseline models (Logistic Regression and Naive Bayes) in classifying sentiments as positive, neutral, or negative. The LSTM model achieved higher accuracy, precision, recall, and F1-score, effectively capturing the nuances in customer sentiments. The visualizations highlighted the distribution of sentiments across different product categories, providing actionable insights for improving customer satisfaction and product development. Overall, the analysis demonstrated the effectiveness of advanced NLP techniques in understanding and leveraging customer feedback for strategic decision-making.

**Conclusion**

The analysis provided valuable insights into customer sentiment based on Amazon product reviews. The LSTM model showed better performance in capturing the sentiment trends compared to baseline models. The

visualizations helped in understanding the distribution and trends in sentiments across different product categories.

**Assumptions**

- The dataset is representative of the broader population of Amazon product reviews

- The features included are relevant and sufficient for sentiment analysis.

- Missing values are missing at random and have been handled appropriately.

- The models developed will generalize well to new, unseen data.

**Limitations**

- Data Quality: The dataset may contain noise, such as irrelevant or spam reviews, that can affect model performance.

- Sentiment Ambiguity: Reviews with mixed sentiments or sarcasm can be challenging to classify accurately.

- Class Imbalance: The distribution of sentiment labels might be imbalanced, affecting model performance, particularly for minority classes.

- Temporal Trends: The analysis does not account for changes in customer sentiment over time.

- External Factors: The analysis does not consider external factors such as changes in product quality, customer expectations, or market trends that may influence sentiments.

**Challenges**

Challenges that I faced during the project are:

- Data Quality: Ensuring the dataset is comprehensive and accurate for reliable analysis.

- Sentiment Ambiguity: Dealing with ambiguous or mixed sentiments in reviews.

- Model Overfitting: Ensuring models generalize well to unseen data without overfitting.

- Interpretability: Balancing model complexity with interpretability for actionable insights.

**Future Uses/Additional Applications**

- Implementing more advanced models like Transformers (e.g., BERT) for state-of-the-art sentiment analysis performance.

- Integrate additional features such as product categories, reviewer profiles, and temporal aspects to improve model accuracy.

- Continuous monitoring and updating of the model with new data to adapt to evolving customer sentiments and improve predictive performance.

- Develop a real-time sentiment analysis system to monitor and respond to customer feedback as it is posted

**Recommendations**

- Use advanced techniques like SMOTE to handle class imbalance.

- Collect more comprehensive data, including additional features and more recent data points.

- Continuously monitor and update the models with new data to ensure relevance and accuracy..

**Implementation Plan**

1. **Data Collection and Preprocessing:** Continuously update and preprocess data.

2. **Model Development:** Refine and retrain the predictive model using updated data.

3. **Integration:** Collaborate with business to integrate the model into decision-making processes.

4. **Monitoring and Evaluation:** Regularly monitor model performance and make necessary adjustments.
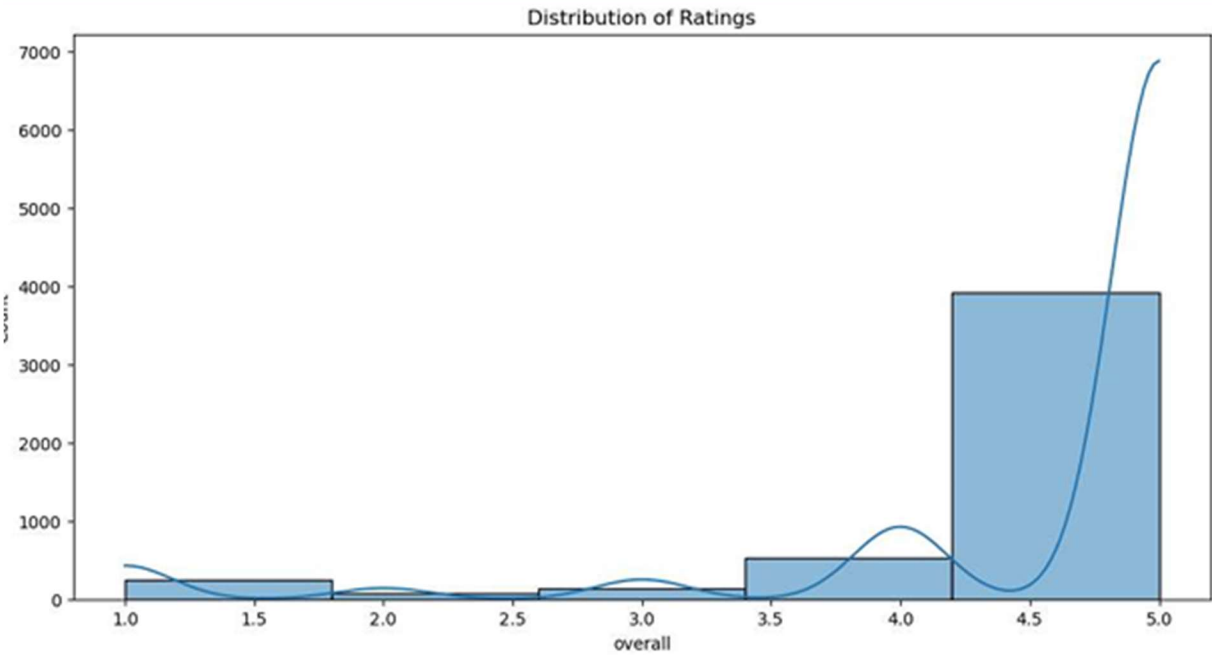
**Ethical Assessment**

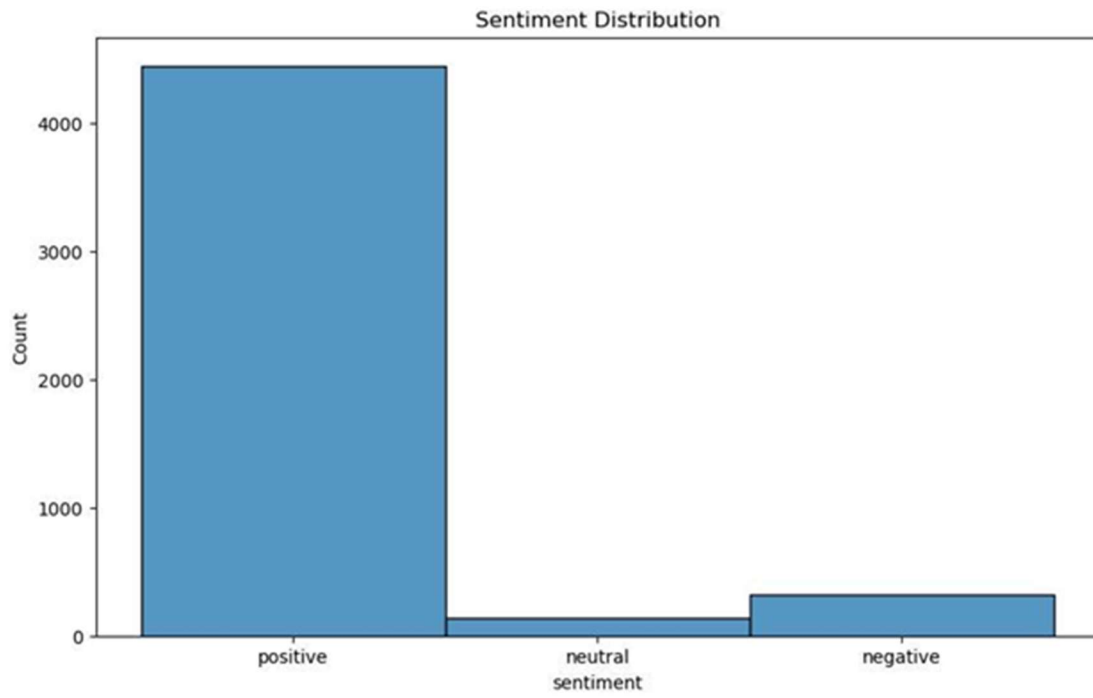Potential ethical concerns that I considered while gathering data were:

- Data Privacy: Ensured that any sensitive information is anonymized and protected.

-  Bias and Fairness: Looked for data from reliable sources where there is no biases in the data that may lead to unfair predictions.

**Appendix**

**Data Dictionary**

| Feature | Description |
|---|---|
| Review ID | Unique identifier for each review |
| Product ID | Identifier for the product |
| Reviewer ID | Identifier for the reviewer |
| Rating | Rating given by the reviewer (1-5) |
| Review Title | Title of the review |
| Review Text | Full text of the review |
| Review Date | Date when the review was written |



Distribution of Ratings

Sentiment Distribution

# Questions

1. Did you explore alternative techniques for encoding text data apart from TF-IDF and tokenization?

2. What influenced your decision to use LSTM for advanced sentiment analysis?

3. How did you address class imbalance in the sentiment labels?

4. What strategies did you use to optimize the hyperparameters of the LSTM model?

5. Have you considered other deep learning models like Transformers, and if so, how do they compare in terms of performance?

6. How can the recall for the neutral sentiment class be improved?

7. Besides accuracy, precision, recall, and F1-score, which other metrics did you consider for model evaluation?

8. Are there any external data sources that could be introduced to enhance the model's performance?

9. What measures were taken to prevent model overfitting?

10. How do you envision the integration of this sentiment analysis model into actual business workflows?