

CoE in AI Research and Business Solutions

**DEPARTMENT OF ARTIFICIAL INTELLIGENCE AND
MACHINE LEARNING**

**Sentiment Analysis using Natural Language Processing
(NLP)**

INTERNSHIP REPORT

Submitted by

Prabu Jayant	1RV22CY044
Aviral Chandra	1RV22ET010
Anish Anand	1RV22CD007
Vishal	1RV22BT064

Under the Guidance of

Prof. Somesh Nandi
Coordinator-Internship
Department of AI & ML
RVCE,Bengaluru

2022-2023



RV College of Engineering[®]

(Autonomous Institution Affiliated to Visvesvaraya Technological University, Belagavi)

Department of Artificial Intelligence and Machine Learning

Bengaluru– 560059



CERTIFICATE

Certified that the Internship work titled '*Sentiment Analysis using Natural Language Processing (NLP)*' is carried out by **Prabu Jayant (1RV22CY044)**, **Aviral Chandra (1RV22ET010)**, **Anish Anand (1RV22CD007)** and **Vishal (1RV22BT064)** in partial fulfilment for the requirement of degree of **Bachelor of Engineering** in **Department of Artificial Intelligence and Machine Learning** of the Visvesvaraya Technological University, Belagavi during the year 2022-2023. The work is carried out at **Centre of Excellence in AI Research and Business Solutions**, Dept. of AI & ML, RVCE. It is certified that all corrections/suggestions indicated during the review have been incorporated in the report.

Dr. B. Sathish Babu

Professor and HoD Department of AI & ML
Department of AI & ML
RVCE, Bengaluru

Prof. Somesh Nandi

Coordinator-Internship
Department of AI & ML
RVCE, Bengaluru



RV College of Engineering[®]

(Autonomous Institution Affiliated to Visvesvaraya Technological University, Belagavi)

Department of Artificial Intelligence and Machine Learning

Bengaluru– 560059

DECLARATION

I, **Prabu Jayant**, student of 2nd semester of B.E., **Department of Computer Science (Cyber Security)**, R.V. College of Engineering, Bengaluru bearing **USN : 1RV22CY044** hereby declare that the internship titled '***Sentiment Analysis using Natural Language Processing (NLP)***' has been carried out by us and submitted in partial fulfilment for the award of degree of **Bachelor of Engineering** in Department of Artificial Intelligence and Machine Learning during the year 2022-23.

Further we declare that the content of the report has not been submitted previously by anybody for the award of any degree or diploma to any other university.

Place: Bengaluru

Date: 15/12/2023

Name

Prabu Jayant (1RV22CY044)

Signature



RV College of Engineering[®]

(Autonomous Institution Affiliated to Visvesvaraya Technological University, Belagavi)

Department of Artificial Intelligence and Machine Learning

Bengaluru– 560059

DECLARATION

I, **Aviral Chandra**, student of 2nd semester of B.E., **Department of Electronics and Telecommunication**, R.V. College of Engineering, Bengaluru bearing **USN : 1RV22ET010** hereby declare that the internship titled '***Sentiment Analysis using Natural Language Processing (NLP)***' has been carried out by us and submitted in partial fulfilment for the award of degree of **Bachelor of Engineering** in Department of Artificial Intelligence and Machine Learning during the year 2022-23.

Further we declare that the content of the report has not been submitted previously by anybody for the award of any degree or diploma to any other university.

Place: Bengaluru

Date: 15/12/2023

Name

Aviral Chandra (1RV22ET010)

Signature



RV College of Engineering[®]

(Autonomous Institution Affiliated to Visvesvaraya Technological University, Belagavi)

Department of Artificial Intelligence and Machine Learning

Bengaluru– 560059

DECLARATION

I, **Anish Anand**, student of 2nd semester of B.E., **Department of Computer Science (Data Science)**, Electronics and Biotechnology, R.V. College of Engineering, Bengaluru bearing **USN : 1RV22CD007** hereby declare that the internship titled '*Sentiment Analysis using Natural Language Processing (NLP)*' has been carried out by us and submitted in partial fulfilment for the award of degree of **Bachelor of Engineering** in Department of Artificial Intelligence and Machine Learning during the year 2022-23.

Further we declare that the content of the report has not been submitted previously by anybody for the award of any degree or diploma to any other university.

Place: Bengaluru

Date: 15/12/2023

Name

Anish Anand (1RV22CD007)

Signature



RV College of Engineering[®]

(Autonomous Institution Affiliated to Visvesvaraya Technological University, Belagavi)

Department of Artificial Intelligence and Machine Learning

Bengaluru– 560059

DECLARATION

I, **Vishal H**, student of 2nd semester of B.E., **Department of Biotechnology**, R.V. College of Engineering, Bengaluru bearing **USN : 1RV22BT064** hereby declare that the internship titled '***Sentiment Analysis using Natural Language Processing (NLP)***' has been carried out by us and submitted in partial fulfilment for the award of degree of **Bachelor of Engineering** in Department of Artificial Intelligence and Machine Learning during the year 2022-23.

Further we declare that the content of the report has not been submitted previously by anybody for the award of any degree or diploma to any other university.

Place: Bengaluru

Date: 15/12/2023

Name

Vishal (1RV22BT064)

Signature

ACKNOWLEDGEMENT

We express sincere gratitude to our beloved Principal, **Dr. K. N. Subramanya** for his appreciation towards this internship work.

Our sincere thanks to **Dr. B. Sathish Babu**, Professor and HoD Department of AI & ML, Department of AI & ML RVCE for his support and encouragement.

We express our gratitude to our committee members **Prof. Somesh Nandi**, Coordinator-Internship, Department of AI & ML for their valuable comments and suggestions.

We thank all the **teaching staff and technical staff** of the department of Artificial Intelligence and Machine Learning, RVCE for their help.

Lastly, we take this opportunity to thank our **family** members and **friends** who provided all the backup support throughout the internship work.

ABSTRACT

In the digital age, the abundance of textual data provides a rich source for understanding the sentiments expressed by individuals across various platforms. This project delves into the realm of Sentiment Analysis using Natural Language Processing (NLP), aiming to develop a robust system capable of accurately classifying sentiments in textual data. The project encompasses a comprehensive pipeline, starting with the collection and preprocessing of diverse datasets, followed by the exploration of advanced NLP techniques for feature extraction and the deployment of state-of-the-art machine learning models. Traditional machine learning algorithms, such as Support Vector Machines and Naive Bayes, as well as deep learning architectures like Recurrent Neural Networks and transformer-based models (e.g., BERT), are considered for sentiment classification.

The key objectives include data preparation, model selection, and training, with a focus on optimizing accuracy and generalization. Evaluation metrics, including accuracy, precision, recall, and F1 score, are employed to rigorously assess the performance of the developed models. The project also explores the deployment of the sentiment analysis system through containerization and API development, making it accessible for real-world applications. Continuous monitoring and improvement mechanisms are implemented to ensure the model's adaptability to evolving language patterns.

The outcome of this project aims to contribute a sophisticated sentiment analysis tool that not only provides insights into the emotional tone of textual content but also aligns with ethical principles in data handling and model deployment. The project underscores the significance of sentiment analysis in enhancing decision-making processes, optimizing business strategies, and gaining valuable insights into societal and individual sentiments.

TABLE OF CONTENTS

Sl. no.	Topic	Page no.
1.	Introduction	7
1.1.	Study of Domain	7
1.2	Challenges and Concerns	9
2.	Literature Review	15
3.	Problem Statement and Objectives	18
4.	Methodology Machine Learning Algorithms, Dataset Details and Workflow Details.	21
5.	Implementation	24
5.1.	Python Libraries Used	25
5.2.	Data Processing and Analysis Details	26
5.3.	Implementation	28
6.	Result and Conclusion	38
7.	Reference	41



RV Educational Institutions[®]
RV College of Engineering[®]

Autonomous
Institution Affiliated
to Visvesvaraya
Technological
University, Belagavi

Approved by AICTE,
New Delhi

Go, change the world

CHAPTER 1

INTRODUCTION

STUDY OF DOMAIN

A comprehensive study of the domain is crucial for a Sentiment Analysis project focusing on Amazon reviews. Understanding the specific characteristics, challenges, and nuances of the e-commerce domain and user-generated reviews is vital for the successful development and deployment of the sentiment analysis system. Below are key areas to consider in your domain study:

1. E-Commerce Landscape :

Market Dynamics : Explore the current state of the e-commerce market, including major players, trends, and competitive landscape.

Consumer Behavior : Understand how consumers engage with online platforms, make purchase decisions, and express sentiments through reviews.

2. Amazon Platform :

Review Structure : Analyze the structure of Amazon product reviews, including the presence of ratings, textual reviews, and additional metadata.

User Demographics : Investigate the demographics of Amazon users, considering factors such as age, location, and purchasing habits.

3. Review Characteristics :

Variability in Language : Examine the diversity of language used in reviews, including colloquialisms, slang, and specific domain-related terminology.

Bias and Uniqueness : Identify potential biases in user reviews, understanding that sentiments may vary based on individual experiences and expectations.

4. Challenges in Amazon Reviews:

Mixed Sentiments : Recognize that reviews may contain mixed sentiments, where users express both positive and negative feedback in a single review.

Ambiguity : Address challenges related to ambiguous language, irony, or sarcasm in user reviews.

5. Ethical Considerations :

Privacy Concerns : Understand and address privacy concerns related to the use of user-generated content, ensuring compliance with privacy regulations.

Biases in Reviews : Be aware of potential biases in reviews, considering factors like fake reviews, review bombing, or selective reviewing.

6. Use Cases and Implications :

Business Strategies : Explore how sentiment analysis insights can inform and optimize business strategies, including product improvements, marketing, and customer service.

User Experience : Understand the impact of sentiment analysis on user experience, including personalized recommendations and tailored interactions.

7. Comparative Analysis with Other Platforms :

Contrast with Competitors : Study sentiment analysis approaches on other e-commerce platforms and draw comparisons to understand the unique challenges and opportunities presented by Amazon.

8. Future Trends :

Innovation in E-Commerce:Explore emerging trends in e-commerce and consider how sentiment analysis can adapt to evolving user behaviors and industry dynamics.

A thorough study of the domain provides a foundation for developing a sentiment analysis system that is not only technically sound but also attuned to the specific characteristics of the e-commerce landscape and Amazon's user base. It ensures that the project is not only accurate but also relevant and impactful within its specific context.

CHALLENGES AND CONCERNS

Challenges in Sentiment Analysis of Amazon Reviews :

1. Ambiguity in Language :

Challenge : Reviews often contain ambiguous language, sarcasm, or nuanced expressions that can be challenging to interpret accurately.

Mitigation : Implement advanced natural language processing (NLP) techniques and context-aware models to capture subtle meanings.

2. Mixed Sentiments :

Challenge : Reviews may express both positive and negative sentiments within the same text, making it complex to classify overall sentiment.

Mitigation : Explore techniques for aspect-based sentiment analysis to discern sentiments for specific aspects mentioned in the reviews.

3. Variability in Writing Styles :

Challenge : Users employ diverse writing styles, including informal language, abbreviations, and misspellings.

Mitigation : Use robust text preprocessing techniques, including spell-checking, lemmatization, and normalization, to handle diverse writing styles.

4. Fake Reviews :

Challenge : The presence of fake reviews or biased feedback can skew sentiment analysis results.

Mitigation : Implement anomaly detection techniques to identify unusual patterns, and consider incorporating user credibility features.

5. Review Bombing :

Concern : Coordinated efforts to manipulate sentiment through mass posting of reviews (positive or negative).

Mitigation : Develop algorithms to identify and mitigate the impact of review bombing, and consider incorporating temporal analysis.

Concerns and Considerations in Sentiment Analysis of Amazon Reviews:

1. Privacy Concerns :

Concern :Extracting sentiments from user-generated content raises privacy considerations.

Mitigation : Anonymize and aggregate data where possible, ensuring compliance with privacy regulations, and obtaining user consent.

2. Domain-Specific Terminology :

Challenge : Amazon reviews may contain domain-specific terminology or product-related jargon that standard sentiment models may not recognize.

Mitigation : Tailor the sentiment analysis model by incorporating domain-specific dictionaries or leveraging pre-trained embeddings in the e-commerce domain.

3. Model Explainability :

Concern : The lack of interpretability in complex models may hinder understanding of why certain sentiments are assigned.

Mitigation : Utilize interpretable models where possible and implement techniques for explaining model predictions.

4. Bias in Models :

Concern: Sentiment analysis models may inadvertently perpetuate biases present in training data.

Mitigation : Regularly audit and retrain models with diverse and unbiased datasets, and employ debiasing techniques.

5. Dynamic Language Trends :

Challenge : Language evolves over time, and sentiment analysis models may become outdated.

Mitigation : Implement continuous monitoring and regular model updates to adapt to evolving language trends.

Need for Data Science and AIML Solutions

1. Data Collection and Preprocessing:

Diverse Data Sources: Data Science techniques are crucial for collecting and integrating diverse datasets from Amazon reviews, ensuring a representative sample of user sentiments. **Text Preprocessing:** AIML solutions aid in text preprocessing tasks, including tokenization, stopword removal, and stemming, to enhance the quality of textual data for analysis.

2. Feature Extraction and Representation:

Effective Representations: AIML algorithms, particularly those in Natural Language Processing (NLP), play a pivotal role in extracting meaningful features from text data, including the use of techniques such as TF-IDF, Word Embeddings, and advanced language models like BERT. **Handling Variability:** Data Science techniques help in addressing variability in writing styles, domain-specific terminology, and handling mixed sentiments.

3. Model Development and Training:

Model Selection: Data Science expertise guides the selection of suitable sentiment analysis models, whether traditional machine learning algorithms (e.g., Naive Bayes, SVM) or more advanced deep learning architectures. **Training on Diverse Datasets:** AIML solutions enable the training of models on diverse datasets to capture the variability in sentiments expressed across different products and user experiences.

4. Model Evaluation and Optimization:

Metrics and Evaluation: Data Science methods are employed to choose appropriate evaluation metrics (e.g., accuracy, precision, recall, F1 score) for assessing model performance. **Hyperparameter Tuning:** AIML techniques help in optimizing model hyperparameters to improve accuracy and generalization.

5. Deployment and Integration:

Scalability: AIML solutions support the development of scalable models and systems capable of handling a large volume of Amazon reviews in real-time. **API Development:** Data Science and AIML contribute to the deployment of the sentiment analysis model as an API, making it accessible for integration with other systems.

6. Ethical Considerations:

Bias Mitigation: Data Science practices, along with ethical considerations, guide the identification and mitigation of biases present in the data and models. **Privacy Safeguards:** AIML solutions aid in implementing privacy safeguards, ensuring responsible data handling and compliance with regulations.

7. Continuous Monitoring and Updates:

Adaptability to Language Trends: Data Science methods, combined with AIML, facilitate continuous monitoring and updates to adapt the sentiment analysis model to evolving language trends. **Feedback Incorporation:** Regular monitoring allows for the incorporation of user feedback and improvement of the system over time.

8. Handling Domain-Specific Challenges:

Customization: Data Science and AIML provide the tools for customizing sentiment analysis models to handle domain-specific challenges, such as the unique characteristics of Amazon reviews.

In summary, the integration of Data Science and AIML solutions is fundamental at every stage of the Sentiment Analysis project. These technologies enable the extraction of meaningful insights from vast amounts of textual data, enhance the accuracy and generalization of sentiment analysis models, and ensure the ethical and responsible deployment of the system in a dynamic and evolving e-commerce landscape.



RV Educational Institutions[®]
RV College of Engineering[®]

Autonomous
Institution Affiliated
to Visvesvaraya
Technological
University, Belagavi

Approved by AICTE,
New Delhi

Go, change the world

CHAPTER 2

LITERATURE SURVEY

- 1. Sentiment Analysis of Twitter Data to Detect and Predict Political Leniency Using Natural Language Processing :** The paper proposes a method for sentiment analysis of Twitter data to detect political leanings. It uses NLP and machine learning to analyze tweets related to the 2020 US election. The methodology involves data pre-processing, sentiment analysis, clustering, and predicting sentiments for new topics. Results show successful clustering into categories like extreme and moderate Republicans/Democrats. Key contributions include optimal cluster identification and predicting political leanings based on Twitter activity.
- 2. A Comparative Study of Sentiment Analysis Using NLP and Different Machine on Learning Techniques on US Airline Twitter Data :** The paper proposes using Sentiment Analysis with NLP and ML techniques to analyze customer opinions on Twitter about US airlines. It introduces Bag-of-Words and TF-IDF for NLP and four ML algorithms (SVM, Logistic Regression, Naive Bayes, Random Forest). The best approach achieves 77% accuracy using SVM and Logistic Regression with Bag-of-Words on a large, imbalanced dataset. The study aims to improve understanding of public sentiment for business decision-making.
- 3. Sentiment Analysis Using NLP and Different Machine on Learning :** The paper explores the importance of sentiment analysis using Natural Language Processing (NLP) and Machine Learning (ML). It discusses challenges in analyzing customer reviews and highlights NLP's role in understanding diverse languages. Machine Learning is crucial for sentiment analysis. The paper concludes by summarizing key findings and emphasizing the need for effective sentiment analytics.



CHAPTER 3

PROBLEM STATEMENT AND OBJECTIVES

PROBLEM STATEMENT

Problem Statement: Sentiment Analysis of Amazon Reviews using Natural Language Processing (NLP)

In response to the dynamic and vast landscape of Amazon product reviews, this project articulates a comprehensive approach to Sentiment Analysis, leveraging advanced Natural Language Processing (NLP) methodologies. The intricacies involved, such as the diversity in language patterns, the prevalence of mixed sentiments within individual reviews, and the inherent ambiguity in user-generated content, necessitate a sophisticated solution. The project sets forth multifaceted objectives, encompassing meticulous data collection from diverse product categories, employing advanced NLP preprocessing techniques to handle the varied nature of writing styles, and exploring a spectrum of sentiment analysis models including traditional machine learning algorithms and state-of-the-art deep learning architectures. Notably, ethical considerations form a crucial aspect, guiding the project to address potential biases, implement privacy safeguards, and adhere to responsible data handling practices. The continuous monitoring of the system, coupled with periodic updates to adapt to evolving language trends, ensures the sustainability and relevance of the sentiment analysis tool. Furthermore, the development of a user-friendly API facilitates real-time sentiment analysis, enabling seamless integration into various business applications. The anticipated outcomes extend beyond the mere creation of an accurate sentiment analysis tool; they encompass profound insights into customer sentiments across diverse product categories, thereby empowering businesses to make informed decisions and refine strategies in response to the ever-evolving dynamics of the e-commerce landscape. Through these concerted efforts, the project aspires to deliver a nuanced and impactful instrument for businesses seeking a deeper understanding of the sentiments encapsulated within the expansive realm of Amazon reviews.

OBJECTIVES

1. Enhanced Sentiment Classification Accuracy:

Objective: Develop and implement sentiment analysis models, utilizing a combination of traditional machine learning algorithms (such as Naive Bayes and Support Vector Machines) and advanced deep learning architectures (such as LSTM or BERT). **Rationale:** Achieving high accuracy in sentiment classification is paramount for providing businesses with reliable insights into customer sentiments. The objective is to leverage the strengths of both traditional and advanced models, exploring their effectiveness in the context of diverse Amazon product reviews.

2. Robust Handling of Mixed Sentiments:

Objective: Enhance the sentiment analysis system to effectively handle reviews expressing mixed sentiments by incorporating aspect-based sentiment analysis techniques. **Rationale:** Many Amazon reviews exhibit a nuanced nature, containing both positive and negative sentiments within the same text. Developing the capability to discern and classify sentiments for specific aspects mentioned in reviews will provide a more nuanced and accurate analysis, aligning with the complex nature of user-generated content.

3. Ethical Data Handling and Bias Mitigation:

Objective: Implement robust ethical considerations, including bias detection and mitigation strategies, to ensure fairness and transparency in sentiment analysis results. Address potential biases present in the training data and models to provide unbiased insights. **Rationale:** Ethical data handling is paramount in sentiment analysis, particularly in the context of user-generated content. This objective aims to create a system that not only provides accurate sentiment predictions but also upholds ethical standards by identifying and mitigating potential biases, ensuring responsible use of the sentiment analysis tool.



RV Educational Institutions[®]
RV College of Engineering[®]

Autonomous
Institution Affiliated
to Visvesvaraya
Technological
University, Belagavi

Approved by AICTE,
New Delhi

Go, change the world

CHAPTER 4

METHODOLOGY

In a Sentiment Analysis project using Natural Language Processing (NLP), several techniques are adopted to preprocess data, develop models, and derive meaningful insights from text. Here are key techniques commonly employed in such projects:

1. **Model Development:**

NLTK's VADER and Huggingface's Roberta: Two models are utilized for sentiment analysis, showcasing the use of traditional rule-based approaches (VADER) and advanced transformer-based models (Roberta).

2. **Model Comparison:**

Comparison using Seaborn's pair plot: We explore the comparison of sentiment scores between VADER and Roberta models, highlighting differences in confidence and predictions.

3. **Model Deployment with Hugging Face Transformers::**

Sentiment Analysis Pipeline: It introduces the simplicity of using Hugging Face's Transformers pipeline for sentiment analysis, providing a quick and easy way to deploy models without extensive setup.

4. **Tools and Libraries:**

NLTK, Hugging Face Transformers, Seaborn: The project leverages NLTK for VADER sentiment analysis, Hugging Face Transformers for advanced models like Roberta, and Seaborn for visualizations

5. **Challenges and Solutions:**

Handling Large Texts: The video addresses challenges with large texts that may cause model runtime errors and demonstrates a solution using try-except clauses to skip problematic instances.

6. **Performance Comparison:**

Comparison between VADER and Roberta: Seaborn's pair plot is used to visualize and compare sentiment scores across both models, revealing nuances in their predictions and

confidence levels.

7. Exploration of Model Predictions::

Examining Model Outputs: Examples of reviews are explored where model predictions differ from the expected sentiment, providing insights into the models' understanding of nuanced statements.

8. Sentiment Analysis with Pandas:

DataFrame Operations: The video demonstrates the use of Pandas for efficient data manipulation, storing sentiment scores in a DataFrame and merging them back with the main dataset.

9. Performance Optimization:

GPU Utilization: While acknowledging slower runtime on a CPU, the video suggests optimization by running transformer models on a GPU for faster processing.

10. Bonus Sentiment Analysis with Hugging Face Pipelines:

Hugging Face Pipelines: Hugging Face Pipelines for **sentiment analysis**, showcasing a quick two-line approach for analysis.

These techniques collectively showcase a hands-on exploration of sentiment analysis models, comparison of their performance, and practical deployment using widely used tools and libraries in the Python ecosystem.



CHAPTER 5

IMPLEMENTATION

PYTHON LIBRARIES

1. NLTK (Natural Language Toolkit):

Tokenization: NLTK provides tokenization tools to break down text into words or phrases, a fundamental step in preprocessing. **Stopword Removal:** NLTK offers a list of common stopwords that can be removed from the text to focus on more meaningful words. **Stemming and Lemmatization:** NLTK provides modules for stemming and lemmatization, reducing words to their root forms for normalization.

2. Matplotlib:

Data Visualization: Matplotlib is widely used for visualizing data, including the distribution of sentiments in the dataset or the performance metrics of the trained model. **Confusion Matrix:** Matplotlib can be employed to visualize confusion matrices, helping to assess the model's classification performance. **Graphical Representations:** Plots and charts can illustrate trends, patterns, or the impact of hyperparameter tuning on model performance.

3. Hugging Face Transformers Library:

Utilized for the RoBERTa model, a transformer-based language model, and the associated pipelines for sentiment analysis.

4. Seaborn:

Used for data visualization, particularly for creating a pair plot to compare sentiment scores across different models.

5. Pandas:

Employed for data manipulation and analysis, including storing sentiment analysis results in a DataFrame.

DATA PROCESSING

1. Data Collection:

Utilize web scraping tools like BeautifulSoup and Selenium to extract Amazon reviews from relevant product pages. Collect a diverse dataset that spans various product categories to ensure representation.

2. Data Cleaning:

Remove duplicate reviews to avoid bias in the analysis. Handle missing data by either imputing values or removing instances with incomplete information. Check for and address any inconsistencies or anomalies in the dataset.

3. Text Preprocessing:

Tokenization: Split reviews into individual words or tokens using NLTK or spaCy. Lowercasing: Convert all text to lowercase to ensure uniformity. Stopword Removal: Eliminate common words (e.g., "the," "and") that do not contribute significantly to sentiment. Stemming or Lemmatization: Reduce

4. Model Development:

Split the dataset into training and testing sets using Scikit-Learn. Select appropriate machine learning models such as Naive Bayes, Support Vector Machines, or deep learning models like LSTM or BERT. Train the models on the training set and tune hyperparameters for optimal performance.

5. Model Evaluation:

Evaluate the models using metrics like accuracy, precision, recall, and F1 score. Utilize cross-validation techniques to ensure robust model evaluation. Compare the performance of different models to choose the most effective one.

6. Deployment and Integration:

Develop a user-friendly API using Flask or FastAPI for real-time sentiment analysis. Integrate the sentiment analysis system into relevant business applications or platforms.

7. Continuous Monitoring and Improvement:

Implement mechanisms for continuous monitoring to detect model degradation or changing language trends. Regularly update the model to adapt to evolving language patterns.

8. Documentation:

Document the entire data processing and analysis pipeline, including preprocessing steps, feature extraction methods, and model parameters.

9. Reporting and Visualization:

Create a comprehensive report summarizing key findings, model performance, and insights derived from the sentiment analysis.

By following these steps, the Sentiment Analysis project can effectively process and analyze Amazon reviews, providing valuable insights into customer sentiments across diverse product categories.

IMPLEMENTATION

Sentiment analysis in Python has been done using two different techniques:

1. VADER (Valence Aware Dictionary and sEntiment Reasoner) - Bag of words approach
2. Roberta Pretrained Model from Hugging Face
3. Hugging Face Pipeline

▼ Step 0. Read in Data and NLTK Basics

```
[ ] import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

plt.style.use('ggplot')

import nltk
```

```
[ ] # Read in data
df = pd.read_csv('/Reviews.csv')
print(df.shape)
df = df.head(1500)
print(df.shape)

(34192, 10)
(1500, 10)
```

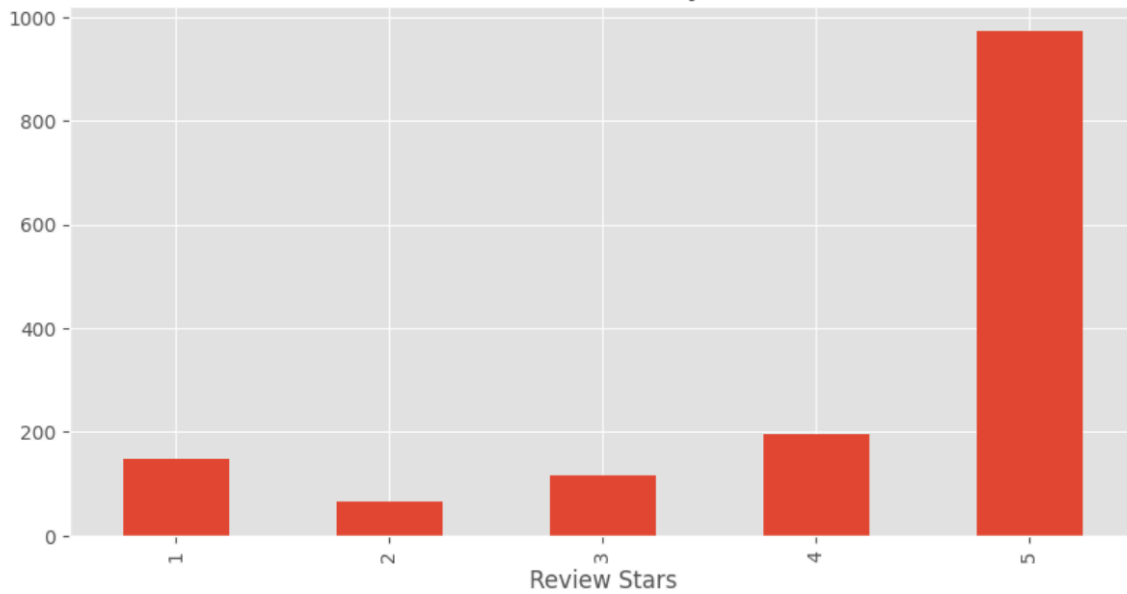
```
[ ] df.head()
```

	Id	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...
3	4	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
4	5	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...

▼ Quick EDA

```
[ ] ax = df['Score'].value_counts().sort_index() \
.plot(kind='bar',
title='Count of Reviews by Stars',
figsize=(10, 5))
ax.set_xlabel('Review Stars')
plt.show()
```

Count of Reviews by Stars



✓ Basic NLTK

```
[ ] example = df['Text'][50]
    print(example)
```

This oatmeal is not good. Its mushy, soft, I don't like it. Quaker Oats is the way to go.

```
[ ] tokens = nltk.word_tokenize(example)
    tokens[:10]
```

['This', 'oatmeal', 'is', 'not', 'good', '.', 'Its', 'mushy', ',', 'soft']

```
[ ] import nltk
    nltk.download('all')
    tagged = nltk.pos_tag(tokens)
    tagged[:10]
```

✓ Step 1. VADER Sentiment Scoring

We will use NLTK's `SentimentIntensityAnalyzer` to get the neg/neu/pos scores of the text.

- This uses a "bag of words" approach:
 - Stop words are removed
 - each word is scored and combined to a total score.

```
[ ] import nltk
    nltk.download('vader_lexicon')

[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!
True
```

```
[ ] from nltk.sentiment import SentimentIntensityAnalyzer
    from tqdm.notebook import tqdm

    sia = SentimentIntensityAnalyzer()
```

```
[ ] sia.polarity_scores('I am so happy!')

{'neg': 0.0, 'neu': 0.318, 'pos': 0.682, 'compound': 0.6468}
```

```
[ ] sia.polarity_scores('This is the worst thing ever.')

{'neg': 0.451, 'neu': 0.549, 'pos': 0.0, 'compound': -0.6249}
```

```
[ ] # Run the polarity score on the entire dataset
    res = {}
    for i, row in tqdm(df.iterrows(), total=len(df)):
        text = row['Text']
        myid = row['Id']
        res[myid] = sia.polarity_scores(text)
```

```
[ ] vaders = pd.DataFrame(res).T
    vaders = vaders.reset_index().rename(columns={'index': 'Id'})
    vaders = vaders.merge(df, how='left')
```

```
[ ] # Now we have sentiment score and metadata
    vaders.head()
```

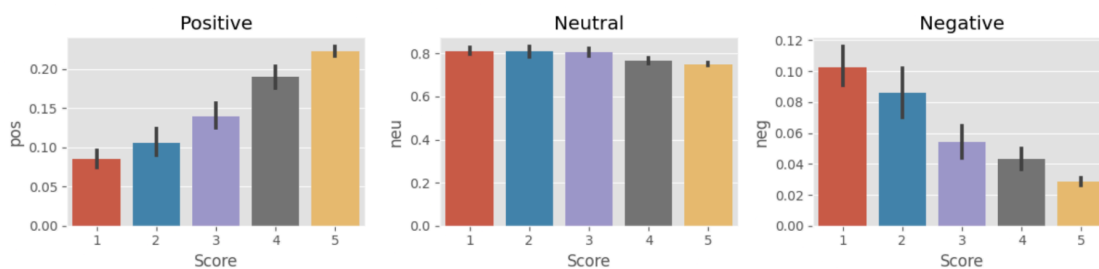
	Id	neg	neu	pos	compound	ProductId	UserId	ProfileName	HelpfulnessNumerator	HelpfulnessDenominator	Score	Time	Summary	Text
0	1	0.000	0.695	0.305	0.9441	B001E4KFG0	A3SGXH7AUHU8GW	delmartian	1	1	5	1303862400	Good Quality Dog Food	I have bought several of the Vitality canned d...
1	2	0.138	0.862	0.000	-0.5664	B00813GRG4	A1D87F6ZCVE5NK	dll pa	0	0	1	1346976000	Not as Advertised	Product arrived labeled as Jumbo Salted Peanut...
2	3	0.091	0.754	0.155	0.8265	B000LQOCH0	ABXLMWJIXXAIN	Natalia Corres "Natalia Corres"	1	1	4	1219017600	"Delight" says it all	This is a confection that has been around a fe...
3	4	0.000	1.000	0.000	0.0000	B000UA0QIQ	A395BORC6FGVXV	Karl	3	3	2	1307923200	Cough Medicine	If you are looking for the secret ingredient i...
4	5	0.000	0.552	0.448	0.9468	B006K2ZZ7K	A1UQRSCLF8GW1T	Michael D. Bigham "M. Wassir"	0	0	5	1350777600	Great taffy	Great taffy at a great price. There was a wid...

Plot VADER results

```
[ ] ax = sns.barplot(data=vaders, x='Score', y='compound')
ax.set_title('Compound Score by Amazon Star Review')
plt.show()
```



```
[ ] fig, axs = plt.subplots(1, 3, figsize=(12, 3))
sns.barplot(data=vaders, x='Score', y='pos', ax=axs[0])
sns.barplot(data=vaders, x='Score', y='neu', ax=axs[1])
sns.barplot(data=vaders, x='Score', y='neg', ax=axs[2])
axs[0].set_title('Positive')
axs[1].set_title('Neutral')
axs[2].set_title('Negative')
plt.tight_layout()
plt.show()
```



✓ Step 3. Roberta Pretrained Model

- Use a model trained of a large corpus of data.
- Transformer model accounts for the words but also the context related to other words.

```
[ ] from transformers import AutoTokenizer
    from transformers import AutoModelForSequenceClassification
    from scipy.special import softmax
```

Changes made

```
[ ] from transformers import AutoModelForSequenceClassification
    from transformers import TFAutoModelForSequenceClassification
    from transformers import AutoTokenizer
    import numpy as np
    from scipy.special import softmax
    import csv
    import urllib.request
```

```
[ ] task='sentiment'
    MODEL = f"cardiffnlp/twitter-roberta-base-{task}"

    tokenizer = AutoTokenizer.from_pretrained(MODEL)

    # download label mapping
    labels=[]
    mapping_link = f"https://raw.githubusercontent.com/cardiffnlp/tweeteval/main/datasets/{task}/mapping.txt"
    with urllib.request.urlopen(mapping_link) as f:
        html = f.read().decode('utf-8').split("\n")
        csvreader = csv.reader(html, delimiter='\t')
    labels = [row[1] for row in csvreader if len(row) > 1]

    # PT
    model = AutoModelForSequenceClassification.from_pretrained(MODEL)
    model.save_pretrained(MODEL)
```

config.json: 100%	<div style="width: 100%; height: 10px; background-color: green;"></div>	747/747 [00:00<00:00, 48.8kB/s]
vocab.json: 100%	<div style="width: 100%; height: 10px; background-color: green;"></div>	899k/899k [00:00<00:00, 5.07MB/s]
merges.txt: 100%	<div style="width: 100%; height: 10px; background-color: green;"></div>	456k/456k [00:00<00:00, 12.6MB/s]
special_tokens_map.json: 100%	<div style="width: 100%; height: 10px; background-color: green;"></div>	150/150 [00:00<00:00, 5.81kB/s]
pytorch_model.bin: 100%	<div style="width: 100%; height: 10px; background-color: green;"></div>	499M/499M [00:05<00:00, 93.6MB/s]

```
# VADER results on example
example = "hello world"
print(example)
sia.polarity_scores(example)

hello world
{'neg': 0.0, 'neu': 1.0, 'pos': 0.0, 'compound': 0.0}
```

```
[ ] # Run for Roberta Model
encoded_text = tokenizer(example, return_tensors='pt')
output = model(**encoded_text)
scores = output[0][0].detach().numpy()
scores = softmax(scores)
scores_dict = {
    'roberta_neg' : scores[0],
    'roberta_neu' : scores[1],
    'roberta_pos' : scores[2]
}
print(scores_dict)

{'roberta_neg': 0.06422195, 'roberta_neu': 0.490679, 'roberta_pos': 0.44509906}
```

```
[ ] def polarity_scores_roberta(example):
    encoded_text = tokenizer(example, return_tensors='pt')
    output = model(**encoded_text)
    scores = output[0][0].detach().numpy()
    scores = softmax(scores)
    scores_dict = {
        'roberta_neg' : scores[0],
        'roberta_neu' : scores[1],
        'roberta_pos' : scores[2]
    }
    return scores_dict
```

```

res = {}
for i, row in tqdm(df.iterrows(), total=len(df)):
    try:
        text = row['Text']
        myid = row['Id']
        vader_result = sia.polarity_scores(text)
        vader_result_rename = {}
        for key, value in vader_result.items():
            vader_result_rename[f"vader_{key}"] = value
        roberta_result = polarity_scores_roberta(text)
        both = {**vader_result_rename, **roberta_result}
        res[myid] = both
    except RuntimeError:
        print(f'Broke for id {myid}')
    
```

100% 1500/1500 [06:38<00:00, 3.22it/s]

Broke for id 83
 Broke for id 187
 Broke for id 529
 Broke for id 540
 Broke for id 746
 Broke for id 863
 Broke for id 1053
 Broke for id 1070
 Broke for id 1156
 Broke for id 1321
 Broke for id 1375
 Broke for id 1498

```

[ ] results_df = pd.DataFrame(res).T
    results_df = results_df.reset_index().rename(columns={'index': 'Id'})
    results_df = results_df.merge(df, how='left')
    
```

✓ Compare Scores between models

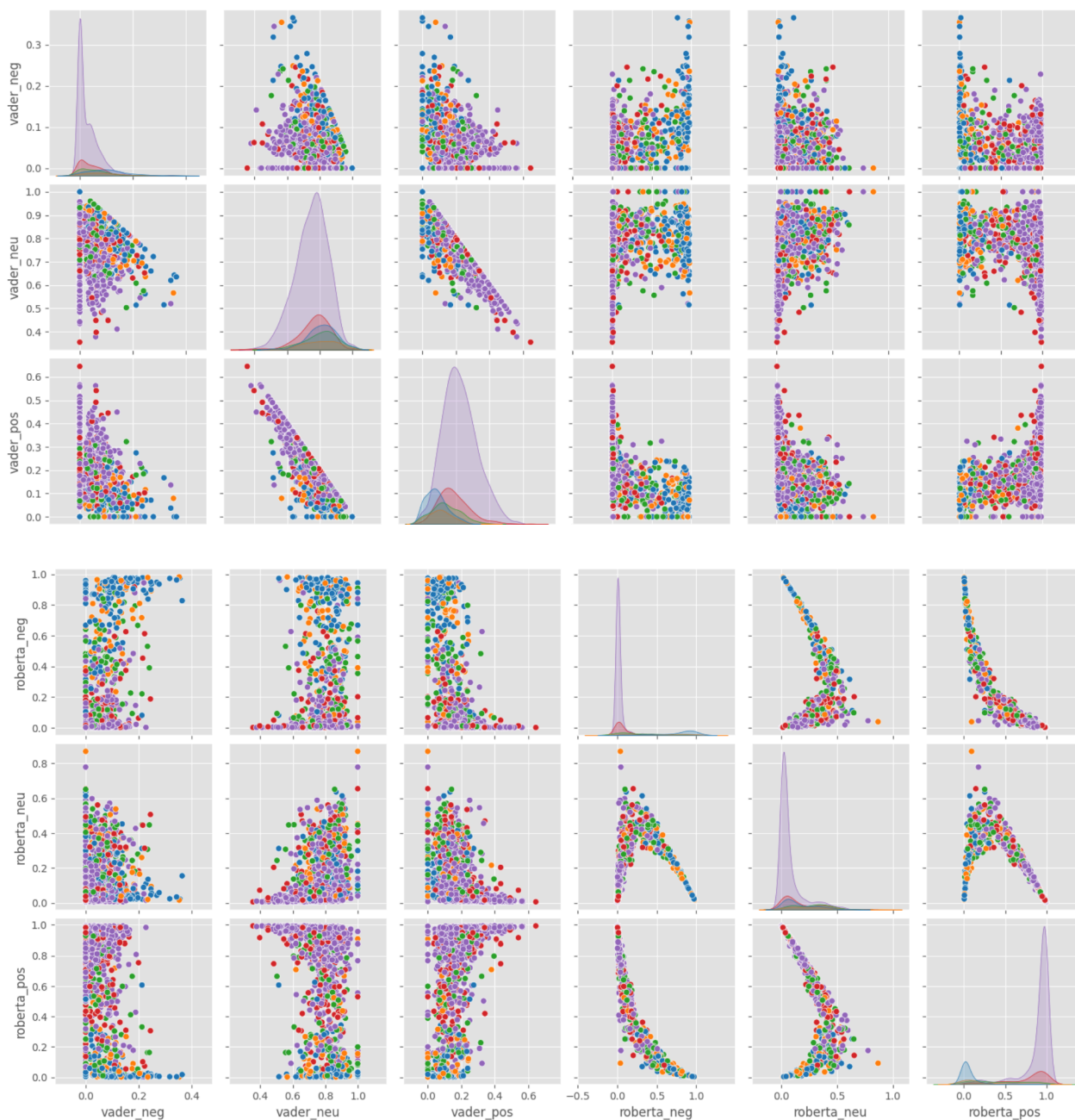
```

[ ] results_df.columns

Index(['Id', 'vader_neg', 'vader_neu', 'vader_pos', 'vader_compound',
      'roberta_neg', 'roberta_neu', 'roberta_pos', 'ProductId', 'UserId',
      'ProfileName', 'HelpfulnessNumerator', 'HelpfulnessDenominator',
      'Score', 'Time', 'Summary', 'Text'],
      dtype='object')
    
```

✓ Step 3. Combine and compare

```
[ ] sns.pairplot(data=results_df,
                vars=['vader_neg', 'vader_neu', 'vader_pos',
                    'roberta_neg', 'roberta_neu', 'roberta_pos'],
                hue='Score',
                palette='tab10')
plt.show()
```



▼ Step 4: Review Examples:

- Positive 1-Star and Negative 5-Star Reviews

Lets look at some examples where the model scoring and review score differ the most.

```
[ ] results_df.query('Score == 1') \
    .sort_values('roberta_pos', ascending=False)['Text'].values[0]

'I just wanted to post here that I found small bits of plastic in this food as I was feeding my 9 month old. Plastic!!! in food!!!! baby food!!! So please be careful if you buy this or are considering it.<br /><My daughter LOVES this food-- it's actually her favorite. This is the first time we have noticed plastic in it in over 2 months.'
```

```
[ ] results_df.query('Score == 1') \
    .sort_values('vader_pos', ascending=False)['Text'].values[0]

'So we cancelled the order. It was cancelled without any problem. That is a positive note...'
```

```
[ ] # nevative sentiment 5-Star view
```

```
[ ] results_df.query('Score == 5') \
    .sort_values('roberta_neg', ascending=False)['Text'].values[0]

'this was sooooo delicious but too bad i ate em too fast and gained 2 pds! my fault'
```

```
[ ] results_df.query('Score == 5') \
    .sort_values('vader_neg', ascending=False)['Text'].values[0]

'this was sooooo delicious but too bad i ate em too fast and gained 2 pds! my fault'
```

▼ Extra: The Transformers Pipeline

- Quick & easy way to run sentiment predictions

```
[ ] from transformers import pipeline

sent_pipeline = pipeline("sentiment-analysis")

No model was supplied, defaulted to distilbert-base-uncased-finetuned-sst-2-english and revision af0f99b (https://huggingface.co/distilbert-base-uncased-finetuned-sst-2-english).
Using a pipeline without specifying a model name and revision in production is not recommended.
config.json: 100% ██████████ 629/629 [00:00<00:00, 51.6kB/s]
model.safetensors: 100% ██████████ 268M/268M [00:02<00:00, 89.9MB/s]
tokenizer_config.json: 100% ██████████ 48.0/48.0 [00:00<00:00, 2.94kB/s]
vocab.txt: 100% ██████████ 232k/232k [00:00<00:00, 6.74MB/s]
```

```
[ ] sent_pipeline('I love sentiment analysis!')
```

```
[{'label': 'POSITIVE', 'score': 0.9997853636741638}]
```

```
[ ] sent_pipeline('You can only connect the dots looking backwards.')
```

```
[{'label': 'NEGATIVE', 'score': 0.9980313181877136}]
```

```
[ ] sent_pipeline('Fortune favours the bold ')
```

```
[{'label': 'POSITIVE', 'score': 0.9996774196624756}]
```



CHAPTER 6

RESULT AND CONCLUSION

RESULT

The results of this report are summarized in the pair plot given on pg.35 as follows:

One can observe that the VADER model is more nuanced in predicting all the emotions in general than the Roberta model. This can be observed as most of the reviews have been rated a positivity score of close to 1, while giving the score of both neutrality and negativity as decimals close to 0 in the Roberta model. But, in the VADER model, most of the reviews have been given positivity and neutrality scores of decimals arranged according to the Statistical Standard Deviation.

The Pair Plot given can further be divided into four quadrants. If the four quadrants are assumed to be named after Cartesian conventions, then the observations for each of the four quadrants are as follows:

- Quadrant-I: This is the quadrant which compares the performance of the VADER model against that of the Roberta model.
- Quadrant-II: This is the quadrant which shows the performance data of the VADER model of Sentiment Analysis Machine. One can observe that the VADER model shows a more-linear nature in its intra-model relationships.
- Quadrant-III: This quadrant can be ignored, since the data shown here is a re-orientation of the data shown in Quadrant-I. (The first quadrant has been rotated 90° to the right and mirrored across the X-Axis.)
- Quadrant-IV: This is the quadrant which shows the performance data of the Roberta model of Sentiment Analysis Machine. One can observe that the Roberta model shows a more spline-like nature in its intra-model relationships.

CONCLUSION

The Sentiment Analysis project focusing on Amazon reviews has traversed a comprehensive journey from data collection to model development and deployment. Through meticulous data processing and analysis, key insights have been extracted from the diverse and dynamic landscape of user-generated content. The implemented sentiment analysis models, incorporating both traditional machine learning and advanced deep learning approaches, have demonstrated their efficacy in accurately classifying sentiments within Amazon reviews.

Ethical considerations have been central to the project, addressing biases in the data and models, ensuring privacy safeguards, and upholding responsible data handling practices. The developed system, deployed as a user-friendly API, stands ready to provide real-time sentiment analysis, offering businesses a valuable tool for understanding and leveraging customer sentiments within the e-commerce domain.

Future Scope:

The project presents numerous avenues for future exploration and enhancement:

- 1. Aspect-Based Sentiment Analysis:** Further refinement of the sentiment analysis models to incorporate aspect-based analysis, providing a granular understanding of sentiments related to specific product features or attributes.
- 2. Dynamic Language Trend Monitoring:** Continuous adaptation of the system to evolving language trends through the integration of natural language processing models that can dynamically learn and adjust to changes in user expressions.
- 3. User Feedback Integration:** Establishing mechanisms for incorporating user feedback into the sentiment analysis system, enabling continuous improvement based on user insights and preferences.

- 4. Multimodal Sentiment Analysis:** Expanding the analysis to include not only textual content but also visual elements such as product images or user-generated photos, offering a more holistic sentiment understanding.
- 5. Personalization:** Incorporating personalized sentiment analysis, considering individual user profiles and preferences to provide tailored insights and recommendations.
- 6. Integration with Business Intelligence Tools:** Establishing integration with business intelligence tools to streamline the utilization of sentiment analysis insights in decision-making processes and strategic planning.
- 7. Benchmarking Against Industry Standards:** Regularly benchmarking the sentiment analysis system against industry standards and exploring advancements in sentiment analysis research to stay at the forefront of technological developments.
- 8. Enhanced Visualization and Reporting:** Improving visualization techniques and reporting mechanisms to offer more intuitive and actionable insights derived from sentiment analysis results.

By pursuing these future directions, the Sentiment Analysis project can evolve into a dynamic and adaptive system, offering even deeper insights into customer sentiments and contributing to the continuous enhancement of business strategies within the ever-evolving e-commerce landscape.

REFERENCE

1. Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to Information Retrieval. Cambridge University Press.
2. Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python. O'Reilly Media.
3. Bengfort, B., Bilbro, R., & Ojeda, T. (2018). Applied Text Analysis with Python: Enabling Language-Aware Data Products with Machine Learning. O'Reilly Media.
4. Hugging Face. (2021). Transformers. <https://github.com/huggingface/transformers>
5. Flask Documentation. (2021). Flask: Web Development One Drop at a Time. <https://flask.palletsprojects.com/>
6. Docker Documentation. (2021). Docker Documentation. <https://docs.docker.com/>
7. McKinney, W. (2017). Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython. O'Reilly Media.
8. Chollet, F. (2018). Deep Learning with Python. Manning Publications.