# Capstone Project (Python): Sentiment Analysis of Reviews

# &

# Model Evaluation

# BY

# THARAKH GEORGE CHACKO

# A Project Report Submitted

# in

# Partial Fulfillment of the

# Requirements for the Data Science and Artificial Intelligence (DSAI) Certification Course

# at

# Intellipaat

# CONTENTS

# LIST OF FIGURES

# PROBLEM STATEMENT:

The primary problem addressed by the project is to effectively analyze customer reviews to understand the sentiment and quality perception of products based on user-generated content. Specifically, the project aims to identify patterns and trends in the data that provide insights into customer satisfaction and product quality. Additionally, the project seeks to classify each review based on the sentiment expressed, aiding in the qualitative assessment of feedback.

This analysis is crucial for businesses as it helps:

- Understand customer sentiment at a granular level, enabling targeted product improvements.
- Identify potential issues in products that might not be evident through traditional quality assurance methods.
- Enhance customer experience by leveraging feedback for product and service refinement.

# PROJECT OBJECTIVE:

1. **Detailed Analysis of Reviews Data:**
   - **Trend Analysis**: Identify key trends over time in review sentiments and ratings. This includes analyzing how sentiments and product ratings change over different periods and potentially identifying any seasonal variations or impact due to specific events.
   - **Pattern Recognition**: Detect common themes and keywords in reviews through techniques such as word clouds and frequency analysis. Understanding these patterns can help pinpoint what aspects of a product stand out to consumers, whether positively or negatively.
   - **Rating and Helpfulness Analysis**: Examine the distribution of product ratings and the helpfulness of reviews. This includes understanding the relationship between the helpfulness of a review and its content, which can indicate what type of information is most valuable to consumers.
2. **Sentiment Classification:**
   - **Development of a Sentiment Analysis Model**: Utilize natural language processing (NLP) techniques to classify reviews into sentiment categories (positive, neutral, negative). This involves preprocessing of text data, feature extraction using techniques like TF-IDF, and the application of machine learning models to predict sentiment.
   - **Model Evaluation and Selection**: Compare the performance of different machine learning models such as Naive Bayes, Logistic Regression, and possibly more sophisticated algorithms if needed. The evaluation will focus on accuracy, precision, and recall ensuring the model reliably interprets the sentiment.
   - **Integration of Sentiment Analysis**: Implement the sentiment analysis as a functional component that can systematically categorize new reviews, providing ongoing insights into customer sentiment.
3. **Business Application of Insights:**
   - **Feedback Loop for Product Improvement**: Utilize insights from sentiment analysis and trend detection to inform product development teams about potential areas for improvement.
   - **Customer Service Enhancement**: Identify critical issues in customer feedback that need immediate attention, potentially improving customer satisfaction and loyalty.
   - **Market Strategy Adaptation**: Adjust marketing strategies based on the insights derived from sentiment analysis and product ratings to better align with customer expectations and improve sales performance.

The project's outcome will assist the business in harnessing the power of customer reviews to fine-tune products and services, ultimately leading to improved customer satisfaction and potentially increased revenue. By systematically analyzing and acting on customer feedback, the company can foster a positive brand image and strengthen customer relationships.

# DATA DESCRIPTION:

The detailed overview of the dataset based on the provided columns and their descriptions are given below:

1. **Id**:
   - **Type**: int64
   - **Description**: It is a unique identifier for each review. This column is typically used as an index to uniquely identify each entry in the dataset.
2. **ProductId**:
   - **Type**: object
   - **Description**: It is the unique identifier for the product being reviewed. This allows for aggregation of reviews by product, which can be crucial for analyzing product-specific feedback.
3. **UserId**:
   - **Type**: object
   - **Description**: It is the identifier for the user who wrote the review. Analyzing user-level data can help in identifying trends among different user demographics or identifying frequent reviewers.
4. **ProfileName**:
   - **Type**: object
   - **Description**: It is the name of the user profile. This could be used for displaying data in a more personalized way or for linking reviews to specific user profiles, though care must be taken to maintain user privacy.
5. **HelpfulnessNumerator**:
   - **Type**: int64
   - **Description**: It represents the number of users who found the review helpful. This is part of the mechanism through which users can interact with reviews, offering insights into the perceived utility of reviews.
6. **HelpfulnessDenominator**:
   - **Type**: int64
   - **Description**: It represents the number of users who indicated whether they found the review helpful or not. This denominator provides context for the helpfulness numerator, giving a ratio that indicates the overall helpfulness of the review.
7. **Score**:
   - **Type**: int64
   - **Description**: It is the rating given to the product by the reviewer, typically on a scale (e.g., 1-5). This is a critical measure as it directly reflects the reviewer's sentiment towards the product.
8. **Time**:
   - **Type**: int64
   - **Description**: It is the timestamp when the review was posted. This can be useful for analyzing trends over time and for understanding any temporal factors affecting reviews.
9. **Summary**:
   - **Type**: object
   - **Description**: It is the summary of the review. This text data can be crucial for quick scans of sentiment and for algorithms that need shorter text inputs.

10. **Text**:

- **Type**: object
- **Description**: It is the full text of the review. This is the most substantial piece of text data in the dataset and is key for detailed text analysis and NLP tasks.

## General Observations:

- The dataset has 5,68,411 rows and 10 columns after removing the null values and duplicates.
- The dataset is quite comprehensive, covering both quantitative and qualitative aspects of customer reviews.
- Text data from the "Summary" and "Text" columns are particularly important for NLP to analyze the sentiment expressed by the reviewer.
- The "Score" alongside the helpfulness metrics provides a quantitative measure of customer sentiment and the perceived utility of each review.
- Timestamps can be used to analyze trends and changes in customer sentiment over time.

```
data.info()

<class 'pandas.core.frame.DataFrame'>
Int64Index: 568411 entries, 0 to 568453
Data columns (total 10 columns):
 #   Column                  Non-Null Count    Dtype
---  ------                  --------------    -----
 0   Id                      568411 non-null   int64
 1   ProductId               568411 non-null   object
 2   UserId                  568411 non-null   object
 3   ProfileName             568411 non-null   object
 4   HelpfulnessNumerator    568411 non-null   int64
 5   HelpfulnessDenominator  568411 non-null   int64
 6   Score                   568411 non-null   int64
 7   Time                    568411 non-null   int64
 8   Summary                 568411 non-null   object
 9   Text                    568411 non-null   object
dtypes: int64(5), object(5)
memory usage: 47.7+ MB
```

Figure 1 : Data Description

## Potential Use Cases:

- Trend analysis in customer sentiment relative to product launches or other events.
- Identifying key words or phrases in positive or negative reviews using text analysis.
- Predicting product ratings based on the textual content of reviews.
- Assessing the impact of reviewer credibility on the perceived helpfulness of reviews.

The detailed breakdown of each column helps in understanding the structure and potential uses of the dataset in various analytical tasks related to sentiment analysis.

# DATA PREPROCESSING STEPS AND INSPIRATION:

Data processing and preprocessing are essential to transforming raw review data into a structured format that can be analyzed and utilized in sentiment analysis models. The detailed explanation of these steps is given below:

## 1. Data Loading and Initial Inspection:

- **Loading the Data**: The dataset is loaded into a pandas DataFrame from a CSV file containing the reviews. This is the starting point for all subsequent operations.
- **Initial Inspection**: Basic checks include viewing the first few rows of the dataset to understand the structure and types of data available (data.head()). Additionally, checking for the total number of entries and data types (data.info()), and summarizing the statistics of numerical fields (data.describe()) are crucial for initial assessments.

## 2. Data Cleaning:

- **Handling Missing Values**: Any missing values are identified and removed (data.dropna()). This is important because missing data can introduce bias or errors in the analysis if not addressed properly.
- **Removing Duplicates**: Duplicate entries are checked and removed (data.duplicated().sum()), ensuring the uniqueness of data points for accurate analysis.
- **Consistency Checks**: Ensuring that helpfulness numerators do not exceed denominators and standardizing text data for uniformity (e.g., converting all text to lowercase).

## 3. Feature Engineering:

- **Helpfulness Ratio**: Creating a new feature by dividing the helpfulness numerator by the helpfulness denominator to get a ratio that reflects how helpful users found the review. This feature is particularly useful for filtering out noise in the data, as reviews with very low helpfulness scores might not be reliable.
- **Text Length Metrics**: New features such as text length and summary length are generated to analyze the length of reviews and summaries, providing insights into whether longer texts correlate with more helpful or more positive/negative reviews.

## 4. Exploratory Data Analysis (EDA):

- **Visualizing Distributions**: Using histograms and count plots to visualize the distributions of scores, text lengths, and the newly created helpfulness ratio. This step is crucial for understanding the underlying patterns and ensuring that the data is suitable for further analysis.
- **Demand Analysis**: Identifying which products have the highest number of reviews, which can be an indicator of their popularity or visibility on the platform.

```
data.isnull().sum()
```

```
Id                        0
ProductId                 0
UserId                    0
ProfileName              16
HelpfulnessNumerator      0
HelpfulnessDenominator    0
Score                     0
Time                      0
Summary                  27
Text                      0
dtype: int64
```

```
data=data.dropna()
```

```
data.duplicated().sum()
```

```
0
```

```
data.shape
```

```
(568411, 10)
```

Figure 2 : Data Cleaning

```
data.ProductId.value_counts()
```

```
B007JFMH8M    913
B002QWHJOU    632
B002QWP8H0    632
B002QWP89S    632
B0026RQTGE    632
              ...
B004DSPTTM      1
B008C9QWU8      1
B007O5A6BM      1
B003Q4TZ08      1
B001LR2CU2      1
Name: ProductId, Length: 74258, dtype: int64
```

Figure 3: ProductId Details and Length

```
data.UserId.value_counts()

A3OXHLG6DIBRW8        448
A1YUL9PCJR3JTY        421
AY12DBB0U420B         389
A281NPSIMI1C2R        365
A1Z54EM24Y40LL        256
                      ...
A1C6KXG47GAQ7B          1
A1TU5DS89D9OVD          1
A3N005JS5FG5FI          1
AQ8W157G7F6I2           1
A3LGQPJCZVL9UC          1
Name: UserId, Length: 256047, dtype: int64
```

```
data.ProfileName.value_counts()

C. F. Hill "CFH"                           451
O. Brown "Ms. O. Khannah-Brown"            421
Gary Peterson                              389
Rebecca of Amazon "The Rebecca Review"     365
Chris                                      363
                                           ...
zinbc                                        1
Steven Wolff                                 1
joycomeau                                    1
Lizz                                         1
srfell17                                     1
Name: ProfileName, Length: 218413, dtype: int64
```

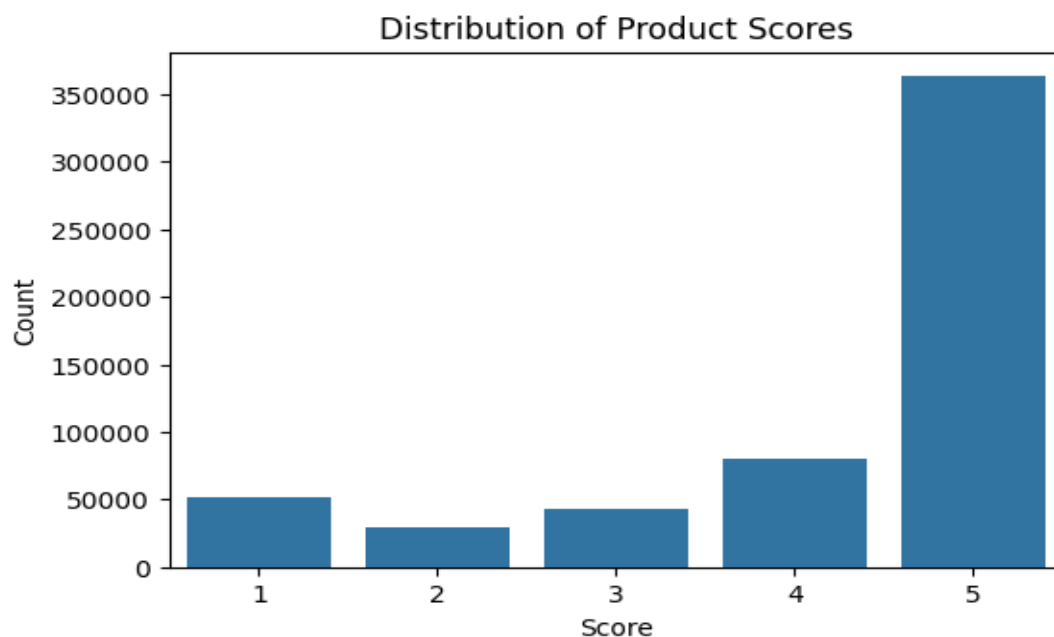Figure 4: Details and Length for UserId and ProfileName



Figure 5: Distribution of Product Scores

```
data['HelpfulnessRatio'] = data['HelpfulnessNumerator'] / data['HelpfulnessDenominator']
plt.figure(figsize=(6,4))
sns.histplot(data['HelpfulnessRatio'], bins=30)
plt.title('Distribution of Helpfulness Ratio')
plt.xlim(0,1.2)
plt.xlabel('Helpfulness Ratio')
plt.ylabel('Count')
plt.show()
```
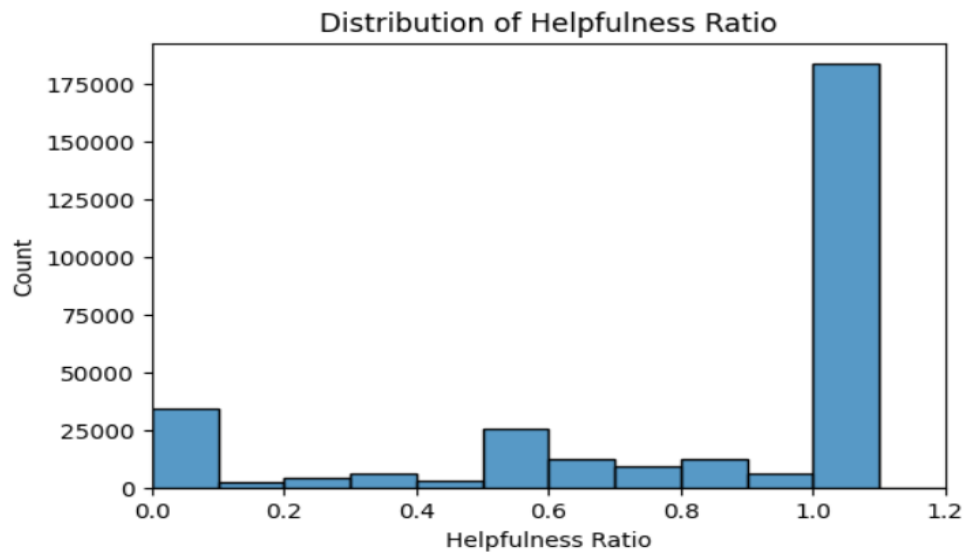


Figure 6: Distribution of Helpfulness Ratio

```
data['TextLength'] = data['Text'].apply(len)
plt.figure(figsize=(10,6))
sns.histplot(data['TextLength'], bins=30)
plt.title('Distribution of Text Length')
plt.xlabel('Text Length')
plt.xlim(0,6000)
plt.ylabel('Count')
plt.show()
```
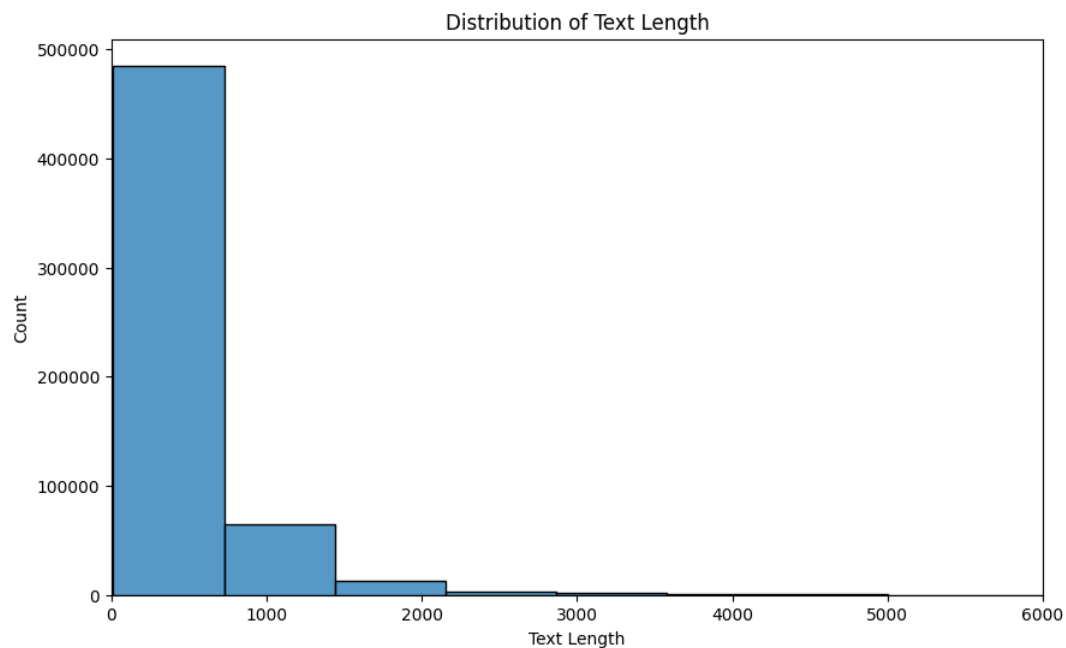


Figure 7 : Distribution of Text Length

```
data['SummaryLength'] = data['Summary'].apply(len)
plt.figure(figsize=(10,6))
sns.histplot(data['SummaryLength'], bins=30)
plt.title('Distribution of Summary Length')
plt.xlabel('Summary Length')
plt.ylabel('Count')
plt.show()
```



Figure 8: Distribution of Summary Length



Figure 9: Scatter Plot of Demand and Rating

Figure 10: Word Cloud for Good Reviews



Figure 11: Word Cloud for Bad Reviews

## 5. Text Preprocessing:

- **Tokenization**: Breaking down the text into individual words or tokens. This is a fundamental step in text analysis to prepare the input for vectorization.
- **Stop Words Removal**: Eliminating common words (e.g., "and", "the") that might appear frequently in texts but offer little value for analysis.
- **Lemmatization**: Converting words into their base form, aiding in the normalization of the text data.

- **Vectorization**: Transforming text data into numerical format using techniques like Count Vectorization and TF-IDF. These transformed features are used by machine learning models since they cannot directly process raw text data.

```python
def preprocess_text(text, remove_numbers=True, stemming=False):
    text = text.lower()
    text = text.replace('\n', ' ')
    text = text.translate(str.maketrans('', '', string.punctuation))
    stop_words = set(stopwords.words('english'))
    text_tokens = word_tokenize(text)
    text = ' '.join([word for word in text_tokens if word not in stop_words])
    if remove_numbers:
        text = ''.join([i for i in text if not i.isdigit()])
    if stemming:
        stemmer = SnowballStemmer("english")
        text = ' '.join([stemmer.stem(word) for word in text.split()])
    else:
        lemmatizer = WordNetLemmatizer()
        text = ' '.join([lemmatizer.lemmatize(word) for word in text.split()])
    return text
```

Figure 12: Text Preprocessing Command

## 6. Model Preparation:

- **Training and Test Split**: Dividing the data into training and test sets to ensure that the model can be trained and then evaluated on unseen data.
- **Model Training**: Using different machine learning algorithms to train sentiment analysis models. This includes comparing various classifiers like Naive Bayes and Logistic Regression to determine which provides the best performance based on the training data.

## Inspiration for Data Preprocessing

The inspiration for the specific preprocessing steps in this project comes from typical challenges encountered in natural language processing and sentiment analysis tasks, particularly:

- **Noise Reduction**: Text data from reviews is inherently noisy with variations in formatting, slang, typos, and grammatical errors. Preprocessing steps such as lemmatization and stop words removal help reduce this noise, making the data cleaner for analysis.
- **Dimensionality Reduction**: Text data can become extremely high-dimensional due to the vast number of unique words. Techniques like TF-IDF help in reducing dimensionality by giving importance to more informative words.
- **Bias Removal**: By creating features like the helpfulness ratio and analyzing text lengths, biases in the data (e.g., longer reviews being perceived as more helpful regardless of content) can be identified and mitigated.

# CHOOSING THE ALGORITHM FOR THE PROJECT:

When selecting algorithms for a sentiment analysis project, the choice largely depends on the nature of the text data, the complexity of the task, and the desired output. The detailed explanation on various algorithms and text vectorization techniques before delving into specific model considerations for the project are given below:

## Text Vectorization Techniques

1. **TF-IDF (Term Frequency-Inverse Document Frequency)**:
   o **Purpose**: Weighs words based on how unique they are in the corpus, thus highlighting important words that could be more influential in sentiment analysis.
   o **Utility**: Useful in reducing the weight of words that appear frequently across documents and are not unique, helping to focus on words that offer more significance in context.
2. **Count Vectorizer**:
   o **Purpose**: Converts text documents into a matrix of token counts, simply counting the number of times each word appears in the document.
   o **Utility**: Effective for models where the mere presence of words (regardless of their relative importance or rarity) is sufficient to build a predictive model.

## Machine Learning Algorithms

1. **Logistic Regression**:
   o **Nature**: A statistical model that predicts binary outcomes based on the linear relationships between dependent and independent variables.
   o **Suitability**: Excellent for binary classification tasks, such as determining if a review is positive or negative.
   o **Performance**: Typically, strong with high-dimensional sparse data, like that from text vectorization.
2. **Naive Bayes Classifier**:
   o **Nature**: A probabilistic classifier based on applying Bayes' theorem, assuming independence between predictors.
   o **Suitability**: Particularly good for large datasets and can be used effectively with text data due to its assumption about feature independence.
   o **Variants**: Multinomial Naive Bayes is specifically tailored for text classification, making it a robust choice for this domain.

## Sentiment Analysis Tools

1. **SIA (Sentiment Intensity Analyzer)**:
   o **Nature**: Part of the NLTK library, it provides a way to calculate text sentiment by categorizing it into positive, negative, and neutral sentiments based on the polarity scores.
   o **Utility**: Quick and useful for preliminary sentiment analysis, especially when processing requirements are minimal and speed is valued over depth of insight.

## Considered Models for the Project

1. **Logistic Regression - TFIDF**:
   - Combines Logistic Regression with TFIDF to prioritize important but less frequent words, potentially improving the model's ability to distinguish between nuanced sentiments.
2. **Naïve Bayes – Count Vectorizer**:
   - Leverages the simplicity of Count Vectorizer with the probabilistic approach of Naive Bayes, suitable for datasets where prevalence of specific words strongly indicates sentiment.
3. **Logistic Regression – Count Vectorizer**:
   - Uses Logistic Regression with a basic form of text representation, which can be very effective when the presence of certain words strongly suggests sentiment, regardless of their frequency across documents.
4. **Naïve Bayes - TFIDF**:
   - Applies Naive Bayes to a TFIDF-transformed dataset, offering a balance between word frequency and their contextual importance, which can be particularly potent for datasets with diverse vocabulary.
5. **NLTK SIA Polarity Scores**:
   - Utilizes a built-in tool for a quick assessment of sentiments, providing a baseline understanding of the general sentiment trends within the data.

Each algorithm and vectorization technique offers unique advantages and potential drawbacks. The choice of the model should align with the specific characteristics of the data and the project's goals. For instance, Logistic Regression with Count Vectorizer might be chosen for its performance and interpretability when the presence of certain words strongly correlates with sentiment, while TFIDF techniques paired with Naive Bayes could be better for understanding the relative importance of words in a larger corpus. The NLTK SIA is particularly useful for rapid, real-time sentiment analysis where complex model training is infeasible.

# MOTIVATION AND REASON FOR CHOOSING THE ALGORITHM:

The motivation behind selecting these algorithms revolves around a few key considerations:

## Applicability to Text Data:

- Both Logistic Regression and Multinomial Naive Bayes have proven track records in handling text data. Naive Bayes has been a traditional choice for spam detection and similar text classification tasks due to its basis in probability theory and its effectiveness in dealing with high-dimensional data.

## Scalability and Performance:

- Considering the dataset size (over half a million reviews), it is crucial to choose algorithms that scale well to large datasets and do not require extensive computational resources. Logistic Regression and Naive Bayes are both known for their efficiency in this regard.

## Baseline Establishing:

- Logistic Regression serves as an excellent baseline model with which to compare other more complex algorithms. It provides a clear benchmark for performance improvements and helps validate the effectiveness of more sophisticated approaches.

## Flexibility and Ease of Implementation:

- The chosen algorithms are supported by popular Python libraries like scikit-learn, making them easy to implement and integrate into a data processing pipeline. Their flexibility and the availability of pre-built functions for training and evaluation streamline the development process.

In summary, the choice of algorithms is driven by the nature of the task (text classification), the characteristics of the data (large and high-dimensional), and practical considerations like computational efficiency, ease of implementation, and the ability to handle the specifics of text data effectively. These factors ensure that the project not only achieves its analytical goals but does so in a computationally efficient and theoretically sound manner.

# ASSUMPTIONS:

When deploying machine learning models, particularly for NLP tasks such as sentiment analysis, several key assumptions are made. These assumptions can affect both the approach to the problem and the interpretation of the results:

## 1. Independence of Features:

- **Multinomial Naive Bayes** relies heavily on the assumption that features (words or tokens in this case) are independent of each other given the output class. This is rarely true in natural language, as words often depend on context, but Naive Bayes can still perform well despite this simplification.

## 2. Linear Relationships:

- **Logistic Regression** assumes that the decision boundary between classes can be defined by a linear combination of the input features. In the context of sentiment analysis, this implies that the sentiment can be linearly separable based on the weighted presence of words, which may not always capture the nuances and context dependencies in language.

## 3. Text Preprocessing Decisions:

- The effectiveness of the models depends significantly on the text preprocessing steps. It is assumed that steps like removing stopwords, applying lemmatization, and vectorizing text (e.g., using TF-IDF or Count Vectorizer) adequately capture the important features from the text for the models to learn effectively.

## 4. Quality and Completeness of the Data:

- It is assumed that the dataset accurately represents the population of interest (all reviews) and contains all relevant attributes needed to predict sentiments accurately. The assumption that the reviews are not biased and that the sample is random and representative can heavily influence the outcomes.

## 5. Sentiment Labelling Accuracy:

- The sentiment labels (derived from the review scores) are assumed to be correct. For instance, reviews with a score of 4 or 5 are considered positive, if the numeric score perfectly aligns with the textual sentiment, which might not always be the case due to subjectivity in how users rate products.

# MODEL EVALUATION TECHNIQUES:

In the sentiment analysis, the importance of each metric—accuracy, precision, recall, and sensitivity—can vary based on the specific goals and context of the analysis. Below explanation shows how each metric applies and why they might be important:

## 1. Accuracy

- **Definition**: Accuracy measures the proportion of total predictions that were correct.
- **Importance**: It provides a straightforward indication of overall model performance across all classes. It is particularly useful when the classes are balanced. If the dataset has an equal number of positive and negative reviews, accuracy gives a quick snapshot of model effectiveness.

## 2. Precision

- **Definition**: Precision (also known as positive predictive value) measures the accuracy of positive predictions. It is the ratio of true positives to the sum of true and false positives.
- **Importance**: Precision is crucial when the cost of a false positive is high. In sentiment analysis, high precision means that when a review is labelled as positive, it is very likely to be genuinely positive. This is important in scenarios where businesses act based on positive reviews, such as using them in promotional materials or prioritizing them in product improvements.

## 3. Recall (Sensitivity)

- **Definition**: Recall measures the ability of a model to find all the relevant cases within a dataset. It is the ratio of true positives to the sum of true positives and false negatives.
- **Importance**: Recall is essential when it is critical to capture as many positives as possible. For sentiment analysis, high recall means the model captures most of the positive reviews, which is crucial for not missing out on feedback that could indicate important improvements. In customer service contexts, high recall can ensure that most customer complaints or praises are captured and addressed.

## 4. Sensitivity

- **Definition**: Sensitivity is another term for recall.
- **Importance**: As described under recall, sensitivity is important for capturing all potential true positive cases, ensuring comprehensive coverage of positive sentiments.

## Determining the Most Important Metrics

- **Scenario-Dependent**: The importance of each metric can depend on the specific business objectives and operational requirements. For instance:

- o **Customer Feedback Monitoring**: If the primary goal is to capture as much feedback as possible to inform product development or customer service improvements, **recall** (sensitivity) might be prioritized.
- o **Marketing or Promotional Use**: If positive reviews are used for promotional purposes, where falsely advertising a product based on incorrect sentiment analysis could harm reputation, **precision** becomes crucial.
- **Trade-offs**: Often, there is a trade-off between precision and recall (known as the precision-recall trade-off). Increasing precision typically reduces recall and vice versa. The choice of which to prioritize may depend on the specific costs associated with false positives and false negatives in the context of the project.
- **F1 Score**: Given the trade-offs between precision and recall, the F1 score (the harmonic mean of precision and recall) is often used as a balanced metric that can be particularly useful when you need to balance the importance of precision and recall.

In the context of this project, while accuracy provides a general sense of model performance, precision and recall (sensitivity) offer deeper insights into the model's capabilities in specific operational scenarios. Depending on the business implications of false positives and false negatives, the emphasis might shift between these metrics to align with strategic objectives.

# INFERENCES:

To determine which model performs best among the five cases, we need to examine each model's metrics — precision, recall, f1-score, and accuracy — considering both classes (positive and negative reviews) and overall effectiveness. Each scenario might prioritize different metrics based on the specific application, but for a general assessment, we will focus on balanced performance across classes.

## Breakdown of Each Model's Performance:

1. **Logistic Regression - TFIDF**
   - o **Accuracy**: 91.29%
   - o **Precision**: 0.84 (negative), 0.93 (positive)
   - o **Recall**: 0.73 (negative), 0.96 (positive)
   - o **F1-Score**: 0.79 (negative), 0.95 (positive)
2. **Naïve Bayes – Count Vectorizer**
   - o **Accuracy**: 89.42%
   - o **Precision**: 0.77 (negative), 0.93 (positive)
   - o **Recall**: 0.73 (negative), 0.94 (positive)
   - o **F1-Score**: 0.75 (negative), 0.93 (positive)
3. **Logistic Regression – Count Vectorizer**
   - o **Accuracy**: 91.64%
   - o **Precision**: 0.84 (negative), 0.93 (positive)
   - o **Recall**: 0.76 (negative), 0.96 (positive)
   - o **F1-Score**: 0.80 (negative), 0.95 (positive)
4. **Naïve Bayes - TFIDF**
   - o **Accuracy**: 85.38%
   - o **Precision**: 0.90 (negative), 0.85 (positive)
   - o **Recall**: 0.37 (negative), 0.99 (positive)
   - o **F1-Score**: 0.52 (negative), 0.91 (positive)
5. **NLTK SIA Polarity Scores**
   - o **Accuracy**: 81.97%
   - o **Precision**: 0.74 (negative), 0.83 (positive)
   - o **Recall**: 0.26 (negative), 0.97 (positive)
   - o **F1-Score**: 0.39 (negative), 0.89 (positive)

## Analysis and Recommendation:

- **Balanced Performance**: Looking at the balance between precision, recall, and f1-scores across both classes, **Logistic Regression – Count Vectorizer** stands out as the best model. It has the highest accuracy of 91.64%, a very good balance between precision and recall for both classes, and the highest f1-score for negative reviews among all models.
- **Precision and Recall**: This model does not just offer high accuracy; it also maintains high precision and recall across both classes. Its ability to effectively identify negative reviews (higher recall for negative class than other models except NLTK and

MultinomialNB - TFIDF) while keeping a high precision indicates fewer false positives and false negatives, crucial for applications where misclassification costs are significant.

- **Application Consideration**: For applications which were missing out on negative sentiments is costly (e.g., customer service scenarios where all complaints need addressing), the higher recall for the negative class in Logistic Regression – Count Vectorizer is beneficial. For marketing or product analysis purposes where precision in identifying true positive sentiments is more critical, this model also stands strong.

Given the balanced performance across precision, recall, and overall accuracy, the **Logistic Regression – Count Vectorizer** model is recommended as the best choice among the presented options. This model provides robust predictive capability, making it suitable for a wide range of applications, from automated review processing to in-depth sentiment analysis for strategic decision-making.

| | Model | Accuracy |
|---|---|---|
| **0** | Logistic Regression - CountVectorizer | 0.916372 |
| **1** | MultinomialNB - TFIDF | 0.912871 |
| **2** | MultinomialNB - CountVectorizer | 0.894205 |
| **3** | Logistic Regression - TFIDF | 0.853794 |
| **4** | NLTK SIA Polarity Scores | 0.819699 |

Figure 13: Model Accuracy Table

We were able to divide the products into 5 segments based on the reviews. The graph on the next page provides insights into how many products fall into each category of satisfaction based on predictive scores. This segmentation reflects overall customer sentiment towards these products as assessed through sentiment analysis or similar predictive modeling.

1. **Very Satisfied (45,611 products)**:
   - **Interpretation**: The majority of products are in this segment, indicating that a large number of products have received high satisfaction scores from predictive modeling. This suggests that these products consistently meet or exceed customer expectations.
   - **Implication**: Products in this category are likely well-received by the market. It could be beneficial to analyze the characteristics or features of these products that contribute to high satisfaction to replicate this success in future products.
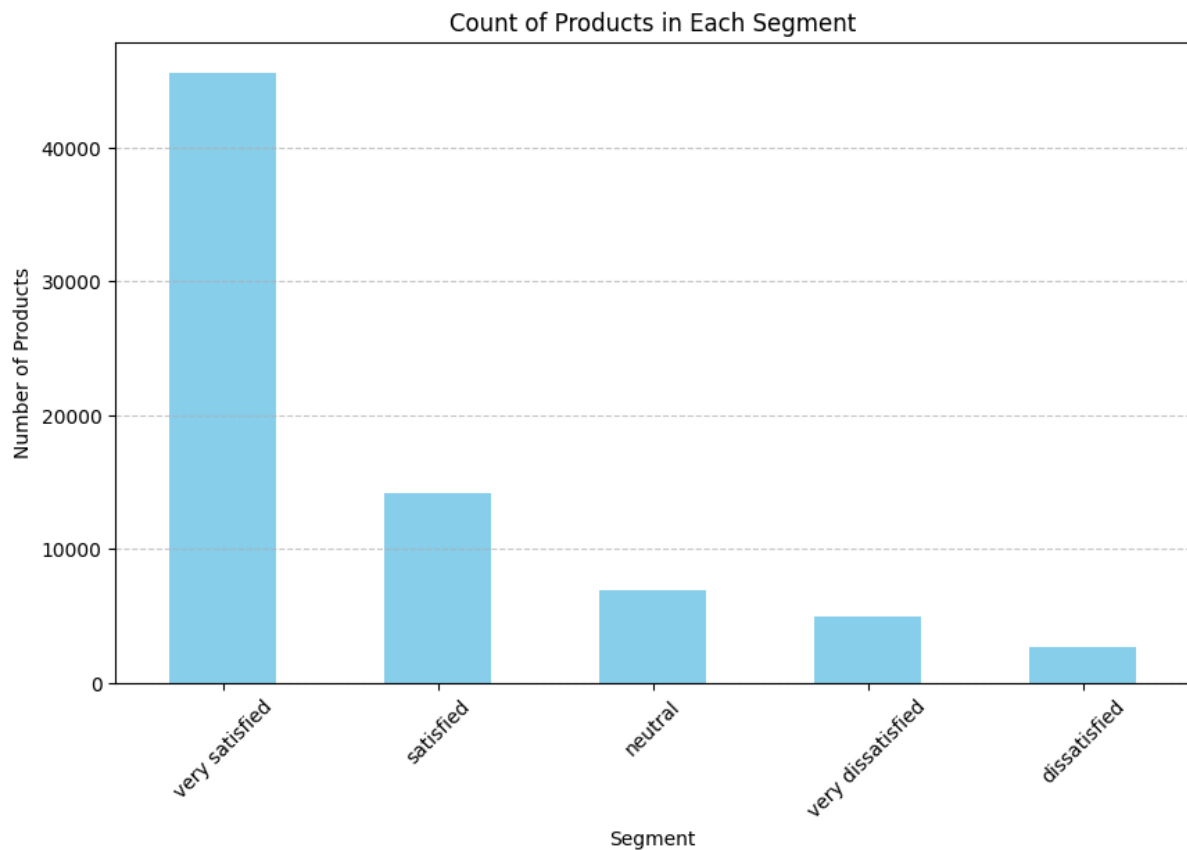
Figure 14 : Count of Products in Each Segment

2. **Satisfied (14,124 products)**:
    o **Interpretation**: A significant number of products are simply 'satisfied', indicating moderate positive feedback from customers.
    o **Implication**: While these products are generally favourable, there may be room for improvement. Understanding the limitations or less favourable aspects could guide enhancements to boost customer satisfaction.
3. **Neutral (6,934 products)**:
    o **Interpretation**: These products did not elicit strong positive or negative reactions, suggesting that while they meet basic expectations, they do not impress.
    o **Implication**: Identifying features that could be improved or added might help to enhance customer perceptions and move these products into higher satisfaction categories.
4. **Very Dissatisfied (4,900 products)**:
    o **Interpretation**: A smaller segment of products falls under very dissatisfied, indicating significant issues or shortcomings that lead to poor customer reviews.
    o **Implication**: It is critical to investigate the reasons behind such negative feedback. Addressing these issues could be crucial for preventing future dissatisfaction and improving the product's reception.
5. **Dissatisfied (2,689 products)**:
    o **Interpretation**: These products are somewhat unsatisfactory but not to the extent of the very dissatisfied category.

    o **Implication**: Minor improvements and quick fixes could potentially elevate the perception of these products. Understanding specific customer grievances would be key in targeting these improvements effectively.

## Overall Strategic Considerations

- **Enhancing Product Quality**: Leveraging insights from both the very satisfied and dissatisfied segments could inform balanced product development that focuses on maintaining strengths while addressing weaknesses.
- **Focused Marketing and Product Development**: Products that fall into lower satisfaction segments might benefit from targeted marketing strategies that address perceived shortcomings or from product development efforts aimed at redesigning or improving these products.
- **Customer Feedback Integration**: Integrating detailed customer feedback for products especially in neutral to dissatisfied segments can guide more specific improvements and help better align products with customer expectations.

This detailed segmentation allows businesses to tailor their strategies more precisely, ensuring that products not only meet the market needs but are also continuously improved based on customer feedback.

# FUTURE POSSIBILITIES OF THE PROJECT:

Looking forward, the project opens several avenues for enhancement and application:

1. **Deep Learning Integration**: Integrating deep learning models such as CNNs or LSTMs could potentially improve model accuracy and the ability to capture semantic nuances in text data. These models excel in understanding context and sequence in text, which might provide better insights into complex customer sentiments.
2. **Real-time Analysis Capability**: Developing the capability to perform sentiment analysis in real-time would allow businesses to react instantly to customer feedback, enhancing customer service and ensuring immediate resolution of issues.
3. **Multilingual Analysis**: Expanding the model to include multilingual capabilities would cater to a global customer base, making the analysis more inclusive and expanding its applicability to non-English speaking regions.
4. **Integration with Other Data Sources**: Combining review data with other forms of customer data, such as purchase history or demographic information, could lead to more personalized insights and enable more targeted marketing and product development strategies.
5. **Automated Response Generation**: Leveraging the insights from sentiment analysis to not only analyze but also respond to customer queries and reviews through automated, intelligent responses could enhance customer interaction, reducing the workload on human agents.

# CONCLUSION:

The sentiment analysis project leveraging machine learning models to classify customer reviews has significantly showcased its potential in enhancing customer interactions and informing product development strategies. The utilization of models like Logistic Regression with Count Vectorizer and TFIDF has proven highly effective, with Logistic Regression with Count Vectorizer standing out for its superior performance in balancing precision and recall across different sentiments.

## Key Project Insights:

- **Model Effectiveness**: Logistic Regression with Count Vectorizer emerged as particularly potent, offering the best balance of accuracy, precision, and recall. This model is well-suited for practical deployment in sentiment analysis tasks, addressing both positive and negative reviews effectively.
- **Importance of Vectorization**: The comparison between Count Vectorizer and TFIDF underlined the profound impact of text preprocessing and feature engineering in NLP. Count Vectorizer, with its straightforward approach, has shown commendable success in capturing relevant features without the complexity introduced by TFIDF's normalization.
- **Robustness and Scalability**: The models demonstrated robustness and scalability, efficiently handling large datasets which is vital for real-world applications where data volume and velocity can be substantial.

## Strategic Model Selection and Algorithm Suitability:

- **Model Selection**: For applications that prioritize high accuracy in sentiment classification, Logistic Regression with Count Vectorizer is recommended due to its outstanding performance. However, for scenarios where interpretability of word importance is crucial (e.g., in feature analysis), TFIDF with Multinomial Naive Bayes could provide deeper insights.
- **Algorithm Suitability**: The effectiveness of an algorithm can vary considerably based on the method of text vectorization, highlighting the importance of experimenting with different preprocessing techniques in NLP tasks.
- **Use of Heuristic Methods**: Quick heuristic methods like NLTK's SIA, although not as precise for detailed sentiment analysis in specific domains, are better suited for applications where speed and computational efficiency are prioritized over precision.

These insights reinforce the need for a nuanced approach to model selection in NLP tasks, considering both the nature of the text data and the specific requirements of the sentiment analysis application. The project's outcomes encourage ongoing refinement and adaptation of models to keep pace with evolving data characteristics and business needs, ensuring that the sentiment analysis remains a potent tool for data-driven decision-making in customer-centric industries.

## REFERENCES:

• **Bird, S., Klein, E., & Loper, E. (2009). Natural Language Processing with Python. O'Reilly Media, Inc.** - This book offers comprehensive coverage of using Python for NLP tasks.

• **Jurafsky, D., & Martin, J. H. (2019). Speech and Language Processing (3rd ed.).** - A fundamental resource for advanced NLP methodologies including deep learning approaches.

• **Kaggle Notebooks:** https://www.kaggle.com/datasets/snap/amazon-fine-food-reviews