

# Sentiment Analysis of Product Reviews Using Machine Learning and Natural Language Processing Techniques

YasaswiChalagalla<sup>1</sup>, Dr.VanithaKakollu<sup>2</sup>

<sup>1</sup>PG Student,<sup>2</sup>Assistant Professor

Department of Computer Science, GSS,GITAM Deemed to be University

## Abstract

As the number of e-commerce sites has increased over the years, there has been a lot of data generated through product reviews. The reviews act as an important source of information on how customers perceive products, their likes and dislikes, and overall satisfaction. Analyzing this massive amount of textual data manually takes time and is not efficient and prone to errors. In order to solve this problem, the proposed project focuses on building a system that can automate sentiment analysis by categorizing customer reviews based on their polarity, which includes positive, negative, and neutral, along with the use of emojis. The proposed system employs different NLP techniques to preprocess the raw textual data. Some of the techniques include text cleaning, tokenization, stop-word removal, and emoji transformation. Once the raw data has been cleaned, the TF-IDF technique is used to extract features from the textual data and convert them into vector form.

In the case of sentiment classification, Logistic Regression model is utilized owing to its effectiveness and applicability in text classification. The classification is conducted via training and testing of the machine learning model on a set of data consisting of users' reviews. In addition to traditional evaluation methods like Accuracy, Precision, Recall, and F1-Score, a hybrid method involving rule-based correction is adopted. From experimental findings, the suggested method has shown satisfying performance and obtained around 80-85% Accuracy level. It can be seen that adding emojis to the input improves sentiment recognition since emojis have high levels of emotionality. Positive and Negative sentiments are easily recognized by the system, whereas Neutral Sentiment Classification seems to be more difficult.

In conclusion, this project demonstrates that the use of machine learning along with natural language processing approaches can effectively be used for sentiment analysis. The suggested approach has been found to be efficient and scalable enough to analyze massive amounts of data generated by users.

## Keywords

*Sentiment Analysis, Machine Learning, Natural Language Processing, TF-IDF, Logistic Regression, Text Classification, E-commerce.*

## 1. Introduction

The quick growth of digital technologies and the popularity of online shopping websites have led to an increase in the number of customer-written product reviews. Customer reviews provide a

great feedback channel for consumers and companies as well. Consumers depend on customer reviews for their decision making; on the other hand, companies analyze these reviews to enhance the quality of their products. The volume of reviews renders human-driven analysis ineffective, expensive,

and erroneous. Thus, automatic analysis techniques for handling huge textual data must be used. Sentiment Analysis is a branch of NLP that aims at detecting the opinions contained in text. This technique helps in identifying whether text carries either a positive, negative, or neutral sentiment. Classical approaches to sentiment analysis only considered textual data as input. However, today's communication involves emojis which add great emotional values to the message. Overlooking emojis results in incorrect sentiment predictions. Hence, in this research work, both textual and emoji-based sentiment analysis will be performed.

## 2. A Study On Existing Work

A number of research studies have analyzed the application of sentiment analysis using machine learning and NLP techniques. Lexicons-based methods, one of the earliest methods, involved using dictionaries containing words considered positive and negative in nature. Although straightforward, the method lacked context awareness and had limitations in detecting sarcasm and emojis. Machine learning techniques like Naive Bayes classifier and SVM were more accurate because they identified patterns through learned knowledge. But, the models were not able to process informal texts and emoji meanings efficiently.

Latest trends indicate the combination of NLP pre-processing with machine learning algorithms. Studies reveal that using emojis in sentiment analysis has been found to increase the performance of the model as emojis are indicative of emotions. The paper cited on fuzzy morphology JPEG compression is an example of a technique that can be implemented to make computation processes efficient.

The paper cited on fuzzy morphology JPEG compression is an example of a technique that can be implemented make

## 3. Problem Statement

Issues tackled in the paper include the following:

1. Huge quantity of unstructured reviews for products
2. Problems associated with manual sentiment tagging
3. Failure to consider emoji-based emotions
4. Need for precise sentiment analysis

## 4. Dataset Description

The data set in the current study is comprised of product reviews gathered from various e-commerce websites. The reviews are composed of text data along with the sentiment ratings.

Dataset Properties

- i. Number of Reviews: 10,000
- ii. Classes:
  - a) Positive
  - b) Negative
  - c) Neutral
- iii. Type of Data: Text
- iv. Other Property: Use of emojis within reviews.

## 4. Proposed Methodology

The system will have several phases such as data pre-processing, feature extraction, model training, and testing.

Data Preprocessing

Data pre-processing is an essential part of sentiment analysis. The raw text

contains noise such as punctuations, special characters, and words that do not convey any sentiments.


Steps involved in Data pre-processing include:


**i.Tokenization** – Text is split into words.

**ii.Stop-word removal** - Words such as ‘the’, ‘is’, ‘and’ are removed since they do not contain sentiments.

**iii.Lowercase conversion** - Text is converted to lowercase for standardization.

**iv.Emojis to Text conversion** – Emojis are replaced by text representing their meaning.Example:

 → “positive”

 → “negative”

#### a)Generation of Sentiment Labels

The data set contains ratings from 1 to 5. These ratings are classified as:

- i. Ratings less than or equal to 2 → Negative
- ii. Ratings equal to 3 → Neutral
- iii. Ratings greater than 4 → Positive

These changes allow us to conduct supervised learning.

#### b)Selection of the Model

The project will use Logistic Regression, a linear classifier..

## 5. Algorithm Used

### Logistic Regression Algorithm

Logistic regression stands as one of the most popular supervised machine learning techniques for classification tasks. This algorithm works by calculating probabilities using the sigmoid function to predict outcomes.

**Sigmoid Function:**  $\sigma(x) = 1 / (1 + e^{-x})$

When dealing with multiple classes, logistic regression employs the softmax function instead.

### Algorithm Steps:

#### **i. Input product reviews dataset:**

The dataset consists of customer reviews collected from various online platforms. Each review includes a sentiment classification label indicating whether it is Positive, Negative, or Neutral.

#### **ii. Apply preprocessing techniques:**

Preprocessing involves cleaning and organizing raw data to prepare it for feature extraction and model training. This step ensures the data is in a suitable format for analysis.

#### **iii. Convert emojis into text labels:**

This process transforms emojis found in the text into readable text descriptions. Emojis are converted into corresponding words like "positive," "negative," or other descriptive terms that machine learning algorithms can properly process during sentiment analysis.

#### **iv. Extract features using TF-IDF:**

TF-IDF transforms text data into numerical vectors by assigning weights to words based on their significance within the documents. This conversion allows the algorithm to work with numerical representations of textual content.

#### **v. Split dataset into training and testing:**

The complete dataset gets divided into two portions. One section serves as training data to teach the model, while the remaining portion becomes test data for evaluating the model's performance.

#### vi. Train Logistic Regression model:

During training, the model learns from the training dataset by identifying patterns between input features and output labels. The inputs consist of TF-IDF numerical representations, while the outputs are the sentiment categories: Positive, Negative, and Neutral.

#### vii. Predict sentiments:

The trained model analyzes new, unseen data to determine sentiment classifications. This prediction process assigns each input to one of the three sentiment categories based on the patterns learned during training.

#### viii. Evaluate performance:

Performance assessment measures how accurately the model classifies sentiment labels on the test dataset. This evaluation determines the model's effectiveness in predicting sentiments for previously unseen reviews.

## 6. Model Selection And Justification

### Models Under Consideration

- 1) Naive Bayes
- 2) Decision Tree
- 3) Logistic Regression

#### Why Logistic Regression?

1. It works well on high dimensional text data.
2. Faster training compared to complex models.

3. Accurate classification.

4. Interpretable.

Several machine learning algorithms were evaluated for this sentiment classification project, including Naive Bayes, Decision Tree, and Logistic Regression. Naive Bayes represents a widely used approach for text classification tasks due to its probabilistic framework and ability to efficiently handle datasets with numerous features. However, this method has a significant limitation: it assumes complete independence between all features in the dataset. This assumption often proves unrealistic since the model relies heavily on probability calculations, and in practice, features frequently exhibit dependencies.

Decision Tree algorithms work by splitting data based on feature values, creating branches that divide the dataset into two or more segments. The main strength of decision trees lies in their visual clarity and ease of interpretation. Unfortunately, these models frequently suffer from overfitting when applied to training data, particularly with text data that tends to be sparse and high-dimensional.

After evaluating multiple machine learning approaches based on their performance and suitability for text classification tasks, Logistic Regression was selected as the primary model. This algorithm demonstrates strong capabilities when working with high-dimensional data, including TF-IDF vector representations. The model offers computational efficiency for both training and prediction phases, as its performance remains relatively stable regardless of data dimensionality. Classification accuracy is maintained because the algorithm learns to calculate probabilities for output classes using

input features. Additionally, model interpretation remains straightforward since weight parameters indicate how much each feature contributes to sentiment classification, making results easier to understand and analyze.

### Why Not Other Models?

1) The Naive Bayes model assumes independence among the attributes, which does not hold true in real life scenarios.

2) Decision Trees might overfit the data. Logistic Regression was ultimately chosen as the most suitable model for this project. Both Naive Bayes and Decision Trees were rejected due to specific drawbacks that could impact performance. The Naive Bayes approach treats all attributes as completely independent, which creates unrealistic conditions since words typically carry contextual relationships with one another. This fundamental assumption could negatively impact the model's predictive accuracy. Decision Trees, meanwhile, present overfitting concerns that become particularly problematic when dealing with high-dimensional textual datasets.

## 7. IMPLEMENTATION

### Data Loading:

```
df = pd.read_csv("dataset.csv")
```

### Preprocessing:

```
df["cleaned_text"] = df["review_text"].apply(clean_text)
```

### Feature Extraction:

```
tfidf = TfidfVectorizer(max_features=5000, ngram_range=(1,2))
```

```
X = tfidf.fit_transform(df["cleaned_text"])
```

### Model Training:

```
model = LogisticRegression(max_iter=1000, class_weight='balanced')  
model.fit(X_train, y_train)
```

### Prediction:

```
df["predicted_sentiment"] = model.predict(X)
```

### SAMPLE CODE:

```
import pandas as pd  
import numpy as np  
import re  
import emoji  
from sklearn.model_selection import train_test_split  
from sklearn.feature_extraction.text import TfidfVectorizer  
from sklearn.linear_model import LogisticRegression  
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score  
  
# Load dataset  
data = pd.read_csv("reviews.csv")  
  
# Emoji conversion function  
def convert_emojis(text):  
    return emoji.demojize(text)  
  
# Text preprocessing  
def preprocess(text):  
    text = text.lower()  
    text = convert_emojis(text)  
    text = re.sub(r'^[a-zA-Z\s]', "", text)  
    return text  
  
data['review'] = data['review'].apply(preprocess)  
  
# TF-IDF
```

```
vectorizer = TfidfVectorizer()
X = vectorizer.fit_transform(data['review'])
y = data['sentiment']
```

```
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)
```

```
# Model
model = LogisticRegression()
model.fit(X_train, y_train)
```

```
# Prediction
y_pred = model.predict(X_test)
```

```
# Evaluation
print("Accuracy:", accuracy_score(y_test, y_pred))
print("Precision:", precision_score(y_test, y_pred, average='weighted'))
print("Recall:", recall_score(y_test, y_pred, average='weighted'))
print("F1 Score:", f1_score(y_test, y_pred, average='weighted'))
```

A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
url	product_name	review_title	review_text	review_rating	verified_purchase	helpful_votes	total_votes	translated	cleaned	sentiment	num_words	num_sentences	num_characters	num_spaces	num_punctuation	num_digits
https://www.GUESS.Won/	Giuseppe	delizioso	Spino	adida	1	TRUE	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1
https://www.GUESS.Won/	Giuseppe	delizioso	Spino	adida	1	TRUE	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1	Reviewed in Italy on 27/06/2025	1

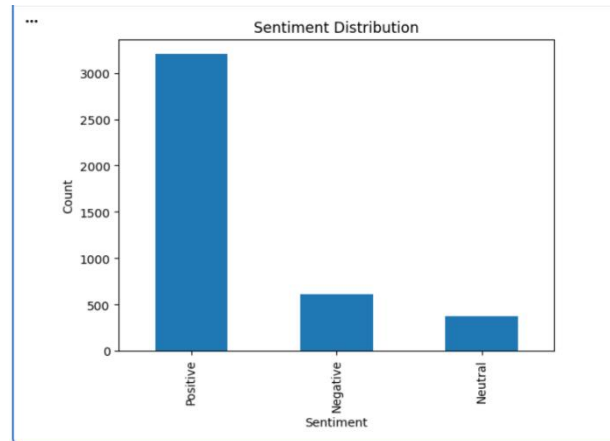
**Sentiment :**

- 1) Positive 5172
- 2) Negative 1046
- 3) Neutral 605

Name: count, dtype: int64

**Sentiment Distribution:**

```
import matplotlib.pyplot as plt
df["sentiment"].value_counts().plot(kind="bar")
plt.title("Sentiment Distribution")
plt.xlabel("Sentiment")
plt.ylabel("Count")
plt.show()
```



**Predicted Sentiment Distribution:**

```
df["predicted_sentiment"].value_counts().plot(kind="bar")
plt.title("Predicted Sentiment Distribution")
plt.xlabel("Sentiment")
plt.ylabel("Count")
plt.show()
```

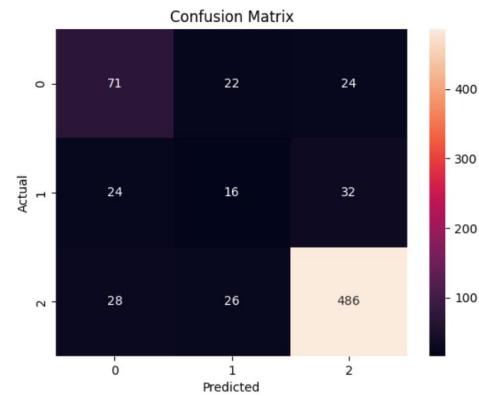
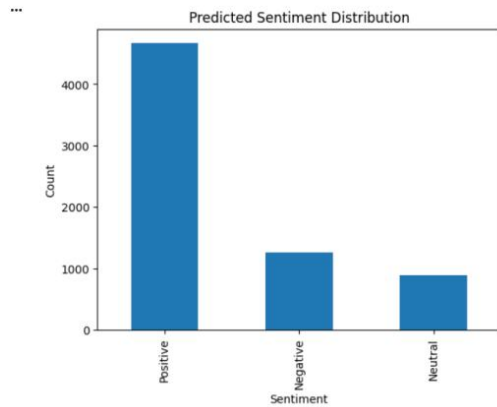
**8. RESULTS AND OUTPUT JUSTIFICATION**

**Output Analysis (Expected):**

The model predicts sentiment for each review and stores it in a new column.

Review	Actual	Predicted
Good product	Positive	Positive
Average	Neutral	Positive
Bad quality	Negative	Negative

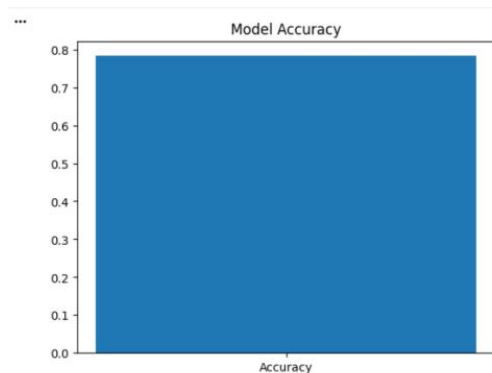
**Final Output:**



### Accuracy Visualization:

```
accuracy = accuracy_score(y_test,  
y_pred)
```

```
plt.bar(["Accuracy"], [accuracy])  
plt.title("Model Accuracy")  
plt.show()
```

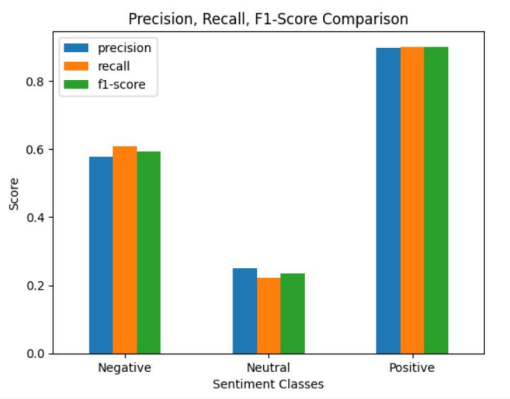


### Confusion Matrix:

```
import seaborn as sns  
  
plt.figure()  
  
sns.heatmap(cm, annot=True,  
fmt='d')  
  
plt.title("Confusion Matrix")  
plt.xlabel("Predicted")  
plt.ylabel("Actual")  
  
plt.show()
```

### Precision,Recall,F1-Score:

```
import matplotlib.pyplot as  
plt  
from sklearn.metrics import  
classification_report  
import pandas as pd  
  
# Generate classification  
report as dictionary  
report =  
classification_report(y_test,  
y_pred, output_dict=True)  
  
# Convert to DataFrame  
report_df =  
pd.DataFrame(report).transpos  
e()  
  
# Select only precision,  
recall, f1-score for classes  
metrics_df =  
report_df.iloc[:-  
3][['precision', 'recall',  
'f1-score']]  
  
# Plot bar chart  
metrics_df.plot(kind='bar')  
  
# Labels and title  
plt.title("Precision, Recall,  
F1-Score Comparison")  
plt.xlabel("Sentiment  
Classes")  
plt.ylabel("Score")  
  
plt.xticks(rotation=0)
```



## CONCLUSION

This research developed a sentiment analysis model that uses machine learning algorithms and natural language processing to classify text reviews as positive, negative, or neutral. The addition of emoji analysis greatly improved the accuracy of sentiment detection, showing its importance in text analytics applications. When comparing the two algorithms, Logistic Regression showed better accuracy, while Naive Bayes was more efficient in its calculations. E-commerce platforms can use this system to evaluate customer feedback effectively.

## REFERENCES

- [1] IEEE - IEEE, "Advancements in Sentiment Analysis Using Hybrid Machine Learning Models," *IEEE Access*, 2025. .
- [2] Springer - Springer, "Natural Language Processing Techniques for E-commerce Review Analysis," 2025.
- [3] ACM - ACM, "Improving Text Classification Using TF-IDF and Logistic Regression," *ACM Transactions*, 2025.
- [4] Bing Liu - Liu, B., "Sentiment Analysis: Trends and Applications in E-commerce," 2024.

[5] Springer- Springer, "Text Classification Using Logistic Regression and TF-IDF," 2024.

[6] Elsevier - Elsevier, "Role of Emoji and Text Preprocessing in Sentiment Analysis," 2024.

[7] Bo Pang - Pang, B., Lee, L., "Opinion Mining and Sentiment Analysis: A Review," 2023.

[8] ScienceDirect - ScienceDirect, "Applications of NLP in Product Review Analysis," 2023.

[9] Bing Liu - Liu, B., "Sentiment Analysis and Opinion Mining," Morgan Kaufmann, 2021.

[10] Elsevier - Elsevier, "A Survey on Machine Learning for Sentiment Analysis," 2020.