

Twitter Sentiment Analysis Using Machine Learning with Python

¹Venati Meganadh Reddy²N Koteswara Rao

¹Final Year B.Tech, CSE, Narayana Engineering College, Gudur

²Associate Professor, CSE Department, Narayana Engineering College, Gudur

Abstract: *The paper aims to develop a robust sentiment analysis system tailored for Twitter data using advanced machine learning techniques implemented in Python. With the exponential growth of social media platforms like Twitter, understanding and analyzing public sentiment has become increasingly crucial for various applications including marketing strategies, brand management, and public opinion monitoring. The project leverages Natural Language Processing (NLP) techniques to preprocess the Twitter data, addressing challenges such as noise, slang, and contextual understanding. Various machine learning algorithms including supervised machine learning algorithms like logistic regression algorithm and evaluated for sentiment classification accuracy. Furthermore, techniques such as feature engineering, hyperparameter tuning, and model ensembling are employed to enhance the performance of sentiment analysis.*

Keywords: *Twitter; NLP; Machine Learning Algorithms; python*

I.INTRODUCTION

Twitter sentiment analysis is a real-time automated machine-learning technique that determines and categorizes the subjective context in tweets. Sentiment analysis of Twitter data involves Opinion Mining to analyze the psychological intent in a tweet – positive, negative, or neutral. Then, based on the patterns identified during text mining, it predicts the subsequent thread of texts., This analysis can be done for individual tweets or a larger dataset related to a particular topic or event. Twitter sentiment analysis has applications across various domains, including marketing, customer service, brand management, financial markets, and political analysis. It provides valuable insights into public opinion and sentiment trends, enabling organizations to make informed decisions and strategies based on Twitter data. Twitter is a free social networking site where users broadcast short posts known as tweets. These tweets can contain text, videos, photos, or links. Twitter is a service for friends, family, and coworkers to communicate and stay connected through the exchange of quick, frequent messages.

Twitter sentiment analysis using machine learning can be effectively performed with Python, employing logistic regression for classification. The process begins with data collection, typically using the Twitter API to gather tweets. These tweets are then preprocessed by cleaning text, removing stop words, and tokenizing. After preprocessing, feature extraction is done using techniques like TF-IDF (Term Frequency-Inverse Document Frequency) to convert text data into numerical form. The dataset is then split into training and testing sets. A logistic regression model is trained on the training data, learning to classify sentiments based on the features. Once trained, the model is evaluated on the test set to determine its accuracy and performance. This involves predicting the sentiment of test tweets and comparing these predictions to the actual sentiments. Fine-tuning the model and preprocessing steps may be necessary to improve accuracy. Finally, the model can be used to analyze the sentiment of new, unseen tweets, providing insights into public opinion and trends. This approach leverages the simplicity and effectiveness of logistic regression for binary classification tasks like sentiment analysis.

II.FUNCTIONAL OVERVIEW

In recent years, research efforts in Twitter sentiment analysis have leveraged a diverse array of methodologies and technologies to comprehensively understand and address the dynamics of public opinion. The integration of machine learning techniques with textual data from Twitter has emerged as a prominent approach, enhancing predictive models for sentiment analysis. By using algorithms like logistic regression, researchers can effectively analyze tweet data to classify sentiments as positive or negative with high precision.

The process begins with data collection via the Twitter API, followed by preprocessing steps such as cleaning text, removing stop words, and tokenizing. Feature extraction techniques like TF-IDF transform textual data into numerical vectors, making it suitable for machine learning models. Logistic regression, due to its simplicity and effectiveness, is then employed to train a sentiment classifier on labeled data. The model's performance is evaluated using metrics such as accuracy, precision, recall, and F1 score, ensuring its reliability.

Additionally, deep learning methodologies have demonstrated potential in capturing more complex patterns and dependencies in textual data, further improving sentiment analysis accuracy. This multifaceted approach, combining advanced modelling techniques, real-time monitoring, and community engagement, provides valuable insights into public opinion, aiding decision-making and fostering a deeper understanding of social dynamics on platforms like Twitter.

III.METHODOLOGY

The methodology for Twitter sentiment analysis involves several key steps. Initially, data is collected from Twitter using the Twitter API, targeting specific keywords or hashtags. This raw data undergoes preprocessing, which includes cleaning text by removing URLs, mentions, special characters, and stop words, followed by tokenization. Feature extraction is then performed using the TF-IDF method to convert text data into numerical vectors. A logistic regression model is trained on the processed data, learning to classify sentiments as positive or negative. The model is evaluated on a separate test set using metrics such as accuracy, precision, recall, and F1 score. Finally, the trained model is deployed to predict the sentiment of new tweets, providing real-time insights into public opinion.

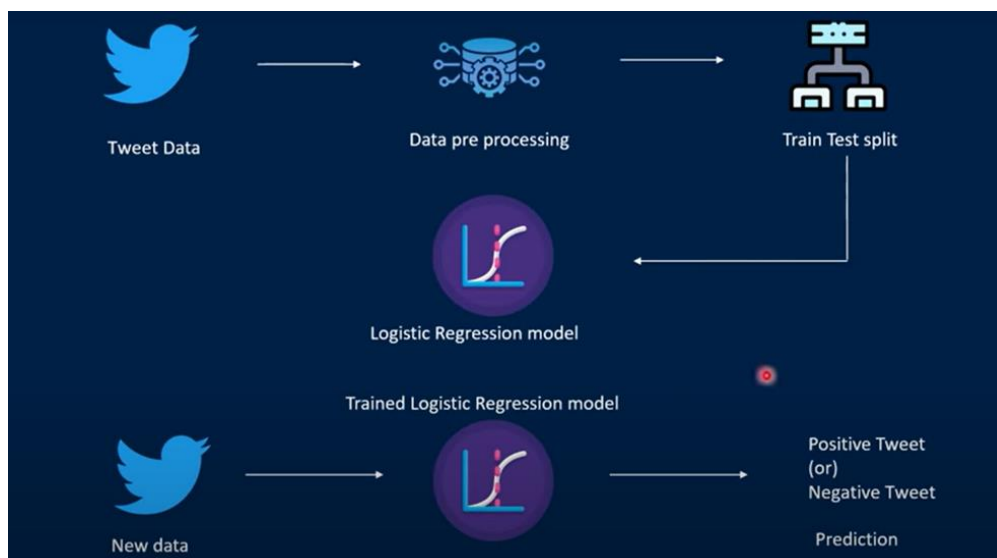


Fig.1 Architecture of Twitter Sentiment Analysis

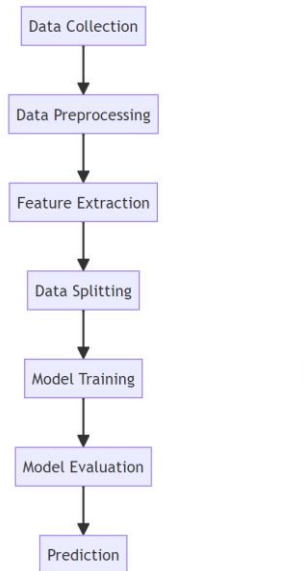


Fig.2 Flow Representation of the System

1. Data Collection:

Data collection for the Twitter sentiment analysis project involves integrating with the Twitter API to fetch real-time tweets based on search queries or streaming topics. Collected tweets undergo preprocessing to clean and tokenize text for subsequent sentiment analysis. Data is stored in a database or file format, ensuring scalability and compliance with privacy regulations.

Upload Your Kaggle.json File

```
[ ] !mkdir -p ~/.kaggle
!cp kaggle.json ~/.kaggle/
!chmod 600 ~/.kaggle/kaggle.json
```

Importing twitter dataset

```
[ ] !kaggle datasets download -d kazanova/sentiment140
```

```
Dataset URL: https://www.kaggle.com/datasets/kazanova/sentiment140
License(s): other
Downloading sentiment140.zip to /content
99% 80.0M/80.9M [00:01<00:00, 90.6MB/s]
100% 80.9M/80.9M [00:01<00:00, 70.4MB/s]
```

```
[ ] # extracting the compressed dataset

from zipfile import ZipFile
dataset='/content/sentiment140.zip'

with ZipFile(dataset, 'r') as zip:
    zip.extractall()
    print('The dataset is extracted')
```

Fig.3 Data Collection

2.Data Preprocessing:

Data preprocessing is a critical step in Twitter sentiment analysis, ensuring that raw tweet data is clean and suitable for model training. Here's an overview of the preprocessing steps:

1. **Data Cleaning:** Remove URLs, mentions (usernames), hashtags, special characters, and punctuation from the tweet text to reduce noise and focus on meaningful content.
2. **Lowercasing:** Convert all text to lowercase to maintain consistency and avoid treating the same word in different cases as different tokens.
3. **Stop Words Removal:** Eliminate common stop words (e.g., "and," "the," "is") that carry little semantic value and do not contribute significantly to sentiment analysis.
4. **Tokenization:** Split the tweet text into individual words or tokens. This helps in analyzing the tweet at a granular level.
5. **Lemmatization/Stemming:** Reduce words to their base or root form (e.g., "running" to "run") to group similar words together and improve the efficiency of the model.
6. **Handling Emoticons and Emojis:** Optionally, convert emoticons and emojis to text representations, as they often convey sentiment.
7. **Removing Duplicates and Irrelevant Data:** Remove duplicate tweets and any irrelevant content that does not add value to sentiment analysis.

```
import nltk
nltk.download('stopwords')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
True

[] # printing the stopwords in English
print(stopwords.words('english'))

['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves',
```

Fig.4 Stop Words Removal

```
def stemming(content):
    stemmed_content = re.sub('[^a-zA-Z]', ' ', content)
    stemmed_content = stemmed_content.lower()
    stemmed_content = stemmed_content.split()
    stemmed_content = [port_stem.stem(word) for word in stemmed_content if not word in stopwords.words('english')]
    stemmed_content = ' '.join(stemmed_content)

    return stemmed_content

twitter_data['stemmed_content'] = twitter_data['text'].apply(stemming)

twitter_data.head()
```

	target	id	date	flag	user	text	stemmed_content
0	0	1467810369	Mon Apr 06 22:19:45 PDT 2009	NO_QUERY	_TheSpecialOne_	@switchfoot http://twitpic.com/2y1zl - Awww, t...	switchfoot http twitpic com zl awww bumper sho...
1	0	1467810672	Mon Apr 06 22:19:49 PDT 2009	NO_QUERY	scotthamilton	is upset that he can't update his Facebook by ...	upset updat facebook text might cri result sch...
2	0	1467810917	Mon Apr 06 22:19:53 PDT 2009	NO_QUERY	mattycus	@Kenichan I dived many times for the ball. Man...	kenichan dive mani time ball manag save rest g...
3	0	1467811184	Mon Apr 06 22:19:57 PDT 2009	NO_QUERY	ElleCTF	my whole body feels itchy and like its on fire	whole bodi feel itchi like fire

Fig.5 Stemming

3.Feature Extraction:

Feature extraction in the context of Twitter sentiment analysis involves transforming raw tweet data into numerical representations that machine learning algorithms can process effectively. This process includes techniques like TF-IDF vectorization or word embeddings to capture semantic meaning and contextual information from text, enabling accurate sentiment classification based on extracted features.

1. **Term Frequency (TF):** Measures how frequently a term appears in a document. It is calculated as the ratio of the number of times a term occurs in a document to the total number of terms in the document.

$$TF(t, d) = \frac{\text{Number of times term } t \text{ appears in document } d}{\text{Total number of terms in document } d}$$

2. **Inverse Document Frequency (IDF):** Measures the rarity of a term across documents in the corpus. It is calculated as the logarithm of the ratio of the total number of documents to the number of documents containing the term.

$$IDF(t) = \log \left(\frac{\text{Total number of documents}}{\text{Number of documents containing term } t} \right)$$

The TF-IDF score for a term in a document is obtained by multiplying its TF by its IDF:

$$TF\text{-}IDF(t, d) = TF(t, d) \times IDF(t)$$

Fig.6 TF-IDF vectorization

4.Data Splitting:

In the context of machine learning with scikit-learn (sklearn), `train_test_split` is a function used to split a dataset into training and testing subsets for model evaluation. Here's a brief explanation of how to use `train_test_split`

`train_test_split` Function:

The `train_test_split` function is part of the `model_selection` module in scikit-learn. It splits arrays or matrices into random train and test subsets. It divides the dataset into training and testing sets (e.g., 80% training, 20% testing)

Splitting the data to training data and test data

```
[27] 1 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.2)
[28] 1 print(X.shape, X_train.shape, X_test.shape)
      (1600000,) (1280000,) (320000,)
[29] 1 print(X_train)
      ['watch saw iv drink lil wine' 'hatermagazin'
      'even though favourit drink think vodka coke wipe mind time think im gon
      ... 'eager monday afternoon'
      'hope everyon mother great day wait hear guy store tomorrow'
      'love wake folger bad voic deeper']
```

Fig.7 Splitting dataset into train and test data

5. Model Training:

Logistic Regression Model: Use scikit-learn's Logistic Regression module to train the model on the training data.

```
Logistic Regression

[34] 1 model = LogisticRegression(max_iter=1000)

1 model.fit(X_train, Y_train)

LogisticRegression(max_iter=1000)
```

Fig.8 Training the model

6. Model Evaluation:

Model evaluation is a crucial step in the sentiment analysis workflow. It helps to assess the performance of your trained logistic regression model on unseen data. Here's a detailed breakdown of the evaluation process, including common metrics and methods:

Accuracy: The proportion of correctly classified instances among the total instances.

```
Model Evaluation

Accuracy score

[36] 1 # accuracy score on the training data
      2 X_train_prediction = model.predict(X_train)
      3 training_data_accuracy = accuracy_score(Y_train, X_train_prediction)

1 print('Accuracy score on the training data :', training_data_accuracy)

Accuracy score on the training data : 0.81021484375
```

Fig.9 Accuracy Score

7. Prediction:

Definition: Prediction involves applying a trained machine learning model to input data to generate an output or prediction. This output can be in the form of class labels (classification), continuous values (regression), or other types of predictions depending on the task.

Process:

1. **Trained Model:** After training a machine learning model using historical data (training set), the model learns patterns and relationships between input features (independent variables) and target labels or values (dependent variable).
2. **Input Data:** New, unseen data (test set or real-world data) is provided to the trained model for prediction.
3. **Output:** The model processes the input data and generates predictions based on the learned patterns. For classification tasks, it assigns class labels to each instance. For regression tasks, it predicts continuous values.

IV.RESULTS AND ANALYSIS

This prediction process allows you to analyze the sentiment of new tweets in real-time, providing valuable insights into public opinion or user sentiment on Twitter.

We provide a random index from the csv file, then our trained model predicts the output as either positive tweet or negative tweet.

Prediction Process:

Input: A random index is provided as input.

Tweet Retrieval: The tweet corresponding to the provided index is retrieved from the dataset.

Model Prediction: The pre-trained logistic regression model predicts the sentiment of the tweet.

Output: The model returns 0 for a negative tweet and 1 for a positive tweet.

Example Prediction:

Input:

Random Index: 123456

Retrieval:

Tweet at Index 123456: "Just finished a great workout! Feeling fantastic. #fitness"

Prediction:

Predicted Sentiment: 1 (Positive)

```
1 X_new = X_test[200]
2 print(Y_test[200])
3
4 prediction = model.predict(X_new)
5 print(prediction)
6
7 if (prediction[0] == 0):
8     print('Negative Tweet')
9
10 else:
11     print('Positive Tweet')
```

1
[1]
Positive Tweet

Fig.10 Positive Tweet

```
1 X_new = X_test[3]
2 print(Y_test[3])
3
4 prediction = loaded_model.predict(X_new)
5 print(prediction)
6
7 if (prediction[0] == 0):
8     print('Negative Tweet')
9
10 else:
11     print('Positive Tweet')
```

0
[0]
Negative Tweet

Fig.11 Negative Tweet

V.CONCLUSION

In this project, we conducted sentiment analysis on Twitter data using machine learning techniques, with a focus on logistic regression implemented in Python. The goal was to analyze sentiments expressed in tweets and classify them as either positive or negative. The results of our sentiment analysis provide valuable insights into public opinion and sentiment trends related to our chosen topic on Twitter. This project lays the foundation for further research and applications in sentiment analysis, such as sentiment monitoring, trend analysis, and brand reputation management. Future directions could involve exploring more advanced machine learning techniques, incorporating domain-specific features, and integrating real-time data streaming for continuous sentiment analysis.

VI.REFERENCES

- [1] Pak, A., & Paroubek, P. (2010). Twitter as a Corpus for Sentiment Analysis and Opinion Mining. In Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10) (pp. 1320-1326). European Language Resources Association (ELRA).
- [2] Go, A., Bhayani, R., & Huang, L. (2009). Twitter Sentiment Classification using Distant Supervision. Stanford University. Retrieved from <http://cs.stanford.edu/people/alecmgo/papers/TwitterDistantSupervision09.pdf>
- [3] Kouloumpis, E., Wilson, T., & Moore, J. (2011). Twitter Sentiment Analysis: The Good the Bad and the OMG! In Proceedings of the Fifth International Conference on Weblogs and Social Media (ICWSM-11) (pp. 538-541). The AAAI Press.
- [4] Zimbra, D., Abbasi, A., Zeng, D., & Chen, H. (2018). The State-of-the-Art in Twitter Sentiment Analysis: A Review and Benchmark Evaluation. ACM Transactions on Management Information Systems (TMIS), 9(2), 5:1-5:29. ACM.
- [5] Saif, H., He, Y., Fernandez, M., & Alani, H. (2016). Contextual Semantics for Sentiment Analysis of Twitter. Information Processing & Management, 52(1), 5-19. Elsevier.