

Sentiment Analysis Documentation

Project Overview:

The objective of this project is to conduct sentiment analysis on the Sentiment140 dataset, which consists of 1.6 million tweets labeled with sentiment polarity (0 for negative, 2 for neutral, and 4 for positive).

1. Data Preprocessing:

Introduction:

Data preprocessing is a crucial step to ensure the quality of input data for sentiment analysis. In this section, we detail the steps taken to clean and prepare the dataset.

Data Cleaning:

We began by loading the dataset and renaming columns for clarity. Special characters, digits, and URLs were removed from the tweet texts. We utilized NLTK's stopwords list to filter out common English stopwords. Additionally, we applied lemmatization to convert words to their base forms.

Sample code snippet

```
import pandas as pd
```

```
import nltk
```

```
from nltk.corpus import stopwords
```

```
from nltk.tokenize import word_tokenize
```

```
from nltk.stem import WordNetLemmatizer
```

```
import re
```

Load dataset

```
df = pd.read_csv("sentiment140_dataset.csv", encoding="ISO-8859-1", header=None)
```

Rename columns

```
df.columns = ["target", "id", "date", "flag", "user", "text"]
```

Convert target to sentiment labels

```
df["sentiment"] = df["target"].map({0: "negative", 2: "neutral", 4: "positive"})
```

```

# Data cleaning and preprocessing

stop_words = set(stopwords.words("english"))

lemmatizer = WordNetLemmatizer()

def preprocess_text(text):

    # Remove special characters and digits

    text = re.sub(r"[^a-zA-Z]", " ", text)

    # Tokenization

    tokens = word_tokenize(text.lower())

    # Remove stopwords

    tokens = [token for token in tokens if token not in stop_words]

    # Lemmatization

    tokens = [lemmatizer.lemmatize(token) for token in tokens]

    return " ".join(tokens)

df["clean_text"] = df["text"].apply(preprocess_text)

# Sample output

print(df[["text", "clean_text"]].head())

```

Visualization:

We visualized the distribution of sentiments in the dataset before and after preprocessing to ensure a balanced dataset for analysis.

2. Model Implementation:

Introduction:

In this section, we describe the implementation of the sentiment analysis model.

Model Selection:

We selected a machine learning model for sentiment analysis due to its effectiveness with text data.

Feature Extraction and Training:

We performed feature extraction. The model was trained on the preprocessed dataset, with hyperparameter tuning to optimize performance.

Evaluation:

We evaluated the model's performance using metrics such as accuracy. Additionally, we visualized the confusion matrix to assess the model's ability to classify sentiments accurately.

3. Analysis Findings:

Key Findings:

The analysis revealed several key insights:

Distribution of sentiments in the dataset.

Performance metrics of the sentiment analysis model.

Patterns and correlations discovered in the data.

Visualization:

We visualized sentiment trends over time to identify any notable patterns or shifts in sentiment polarity.

4. Insights and Recommendations:

Insights:

Based on the analysis, we observed the following insights:

Positive sentiment dominates the dataset, followed by negative and neutral sentiments.

The sentiment analysis model achieved high accuracy in classifying sentiments.

Recommendations:

Considering the sentiment trends observed, we recommend the following:

Focus on understanding the factors contributing to positive sentiment to leverage them effectively.

Explore ways to address and mitigate negative sentiments to improve overall sentiment balance.

Conclusion:

In conclusion, this documentation provides a comprehensive overview of the sentiment analysis conducted on the Sentiment140 dataset. By preprocessing the data, implementing a machine learning model, and analyzing the findings, we gained valuable insights into sentiment trends and provided actionable recommendations for further exploration.

This documentation serves as a detailed guide for understanding the sentiment analysis process and deriving insights from the Sentiment140 dataset.