

SENTIMENTAL ANALYSIS

```
import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

import os
for dirname, __, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

Index

- Importing Libraries
- Reading csv file
- Simplify the data
- analysis
- dropping unnecessary columns
- EDA
- Word cloud positive data
- Word cloud negative data
- TFIDF for sentiment analysis
- Train Test Split
- LogisticRegression
- Accuracy Score: 85%

▾ Importing Libraries

```
import matplotlib.pyplot as plt
import seaborn as sns
from string import punctuation
from nltk.tokenize import word_tokenize
from nltk.stem import LancasterStemmer
from string import punctuation
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem import LancasterStemmer
from nltk.stem.wordnet import WordNetLemmatizer
import re
import warnings
warnings.filterwarnings('ignore')
```

▾ Reading csv file

```
df = pd.read_csv('/content/training.1600000.processed.noemoticon.csv',delimiter=',', encoding='ISO-8859-1')
```

```
df.head()
```

	polarity of tweet	id of the tweet	date of the tweet	query	user	text of the tweet
0	0	1467810672	Mon Apr 06 22:19:49 PDT 2009	NO_QUERY	scotthamilton	is upset that he can't update his Facebook by ...
1	0	1467810917	Mon Apr 06 22:19:53 PDT 2009	NO_QUERY	mattycus	@Kenichan I dived many times for the ball. Man...
2	0	1467811184	Mon Apr 06 22:19:57 PDT 2009	NO_QUERY	ElleCTF	my whole body feels itchy and like its on fire

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 53164 entries, 0 to 53163
Data columns (total 6 columns):
#   Column          Non-Null Count  Dtype
---  -
#   Column          Non-Null Count  Dtype
```

```
0    polarity of tweet    53164 non-null    int64
1    id of the tweet      53164 non-null    int64
2    date of the tweet    53164 non-null    object
3    query                53164 non-null    object
4    user                 53164 non-null    object
5    text of the tweet    53164 non-null    object
dtypes: int64(2), object(4)
memory usage: 2.4+ MB
```

```
df.isnull().sum()

polarity of tweet    0
id of the tweet      0
date of the tweet    0
query                0
user                 0
text of the tweet    0
dtype: int64
```

▼ Simplify the data

```
df.columns=['sentiment','id','date','query','username','text']

df.head()

      sentiment      id      date      query  username      text
0            0  1467810672  Mon Apr 06 22:19:49 PDT 2009  NO_QUERY  scotthamilton  is upset that he can't update his Facebook by ...
1            0  1467810917  Mon Apr 06 22:19:53 PDT 2009  NO_QUERY      mattycus  @Kenichan I dived many times for the ball. Man...
2            0  1467811184  Mon Apr 06 22:19:57 PDT 2009  NO_QUERY      EliaCOTE  my whole body feels

df.shape

(53164, 6)
```

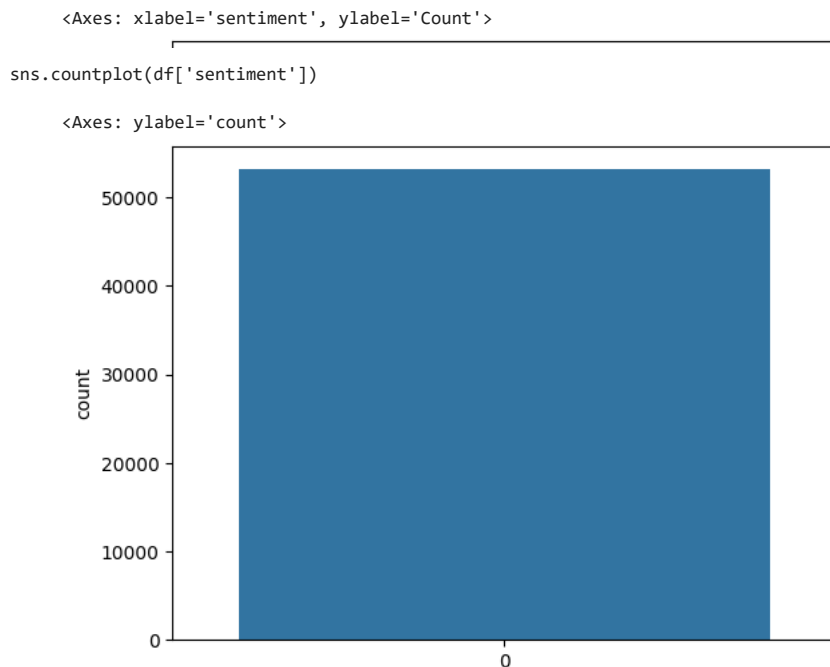
▼ Analysis

```
df['sentiment'].value_counts()

0    53164
Name: sentiment, dtype: int64

import seaborn as sns
import matplotlib.pyplot as plt

sns.histplot(df['sentiment'],kde=True)
```



```
df['query'].value_counts()
```

```
NO_QUERY    53164
Name: query, dtype: int64
```

▼ dropping unnecessary columns

```
df=df.drop(columns=['query'])
```

```
df.head()
```

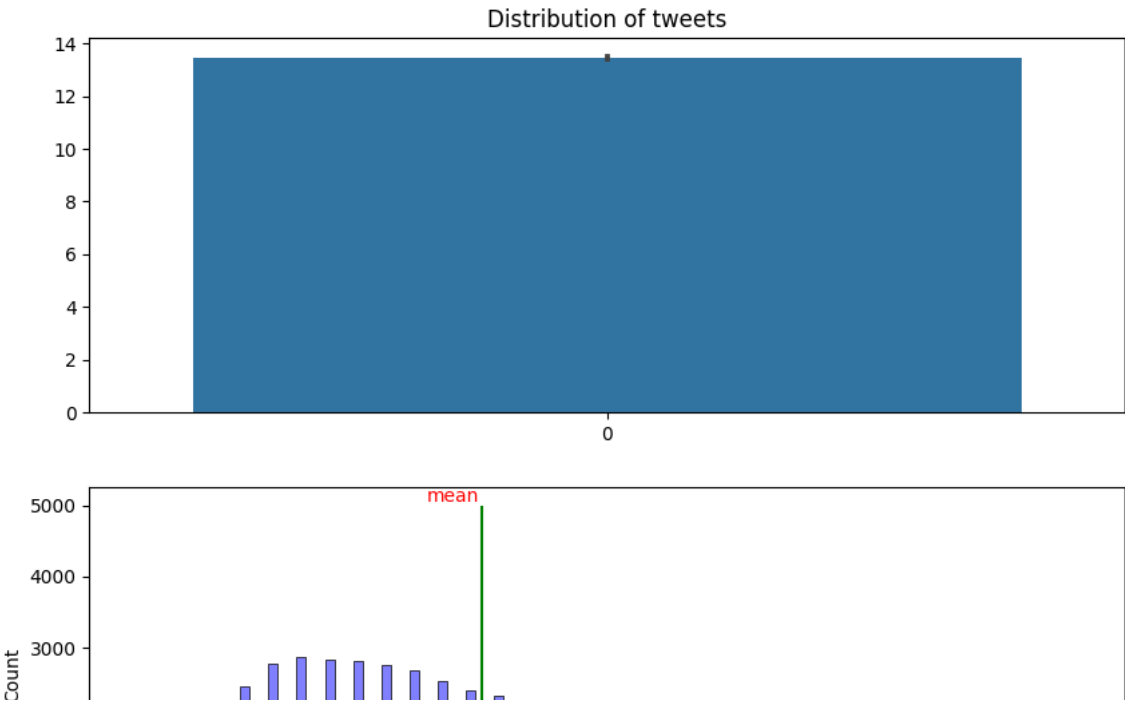
	sentiment	id	date	username	text
0	0	1467810672	Mon Apr 06 22:19:49 PDT 2009	scotthamilton	is upset that he can't update his Facebook by ...
1	0	1467810917	Mon Apr 06 22:19:53 PDT 2009	mattycus	@Kenichan I dived many times for the ball. Man...
2	0	1467811184	Mon Apr 06 22:19:57 PDT 2009	ElleCTF	my whole body feels itchy and like its on fire
3	0	1467811400	Mon Apr 06 22:19:57 PDT	...	@nationwideclass no, it's not behaving at

```
texts = df['text']
```

```
text_lens = [len(t.split()) for t in texts.values]
len_mean = np.mean(text_lens)
```

▼ EDA

```
fig, axes = plt.subplots(2,1, figsize=(10, 8))
axes[0].set_title('Distribution of tweets')
sns.barplot(text_lens, ax=axes[0])
sns.histplot(text_lens,bins=100, kde=True, ax=axes[1],color='blue')
axes[1].vlines(len_mean, 0, 5000, color = 'g')
plt.annotate("mean", xy=(len_mean, 5000), xytext=(len_mean-2, 5050),color='r')
plt.show()
```



```
import nltk
nltk.download('stopwords')
stuff_to_be_removed = list(stopwords.words('english'))+list(punctuation)
stemmer = LancasterStemmer()
corpus = df['text'].tolist()
print(len(corpus))
print(corpus[0])
```

```
53164
is upset that he can't update his Facebook by texting it... and might cry as a result School today also. Blah!
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```

```
import nltk
nltk.download('wordnet')
final_corpus = []
final_corpus_joined = []
for i in df.index:
    text = re.sub('[^a-zA-Z]', ' ', df['text'][i])
    text = text.lower()
    text=re.sub("</?.*?>"," " &lt;&gt; ",text)
    text=re.sub("(\d|\\W)+"," ",text)
    text = text.split()
    lem = WordNetLemmatizer()
    text = [lem.lemmatize(word) for word in text if not word in stuff_to_be_removed]
    text1 = " ".join(text)
    final_corpus.append(text)
    final_corpus_joined.append(text1)
```

```
[nltk_data] Downloading package wordnet to /root/nltk_data...
```

```
data_cleaned = pd.DataFrame()
data_cleaned["text"] = final_corpus_joined
data_cleaned["sentiment"] = df["sentiment"].values
```

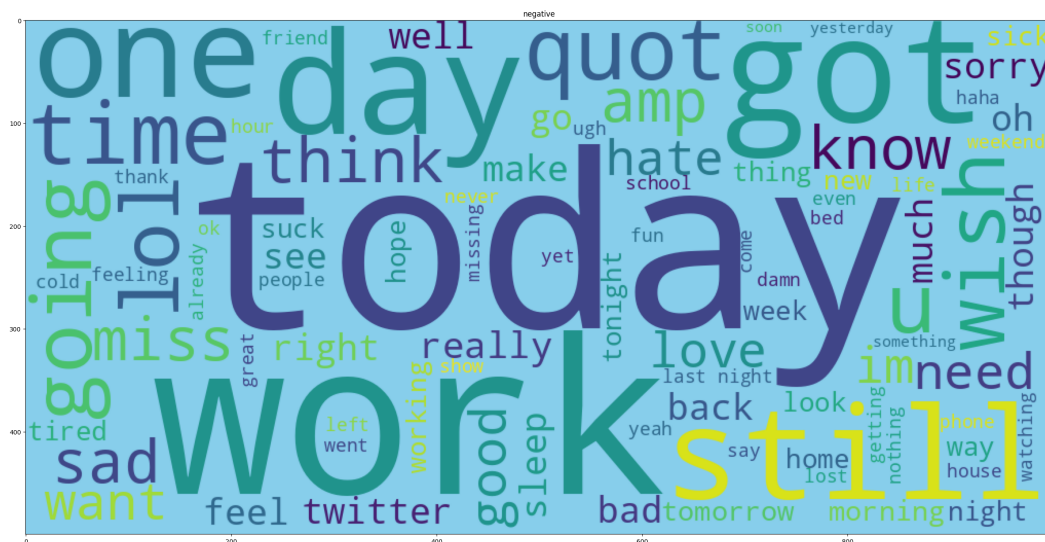
```
data_eda = pd.DataFrame()
data_eda['text'] = final_corpus
data_eda['sentiment'] = df['sentiment'].values
data_eda.head()
```

	text	sentiment	
0	[upset, update, facebook, texting, might, cry,...	0	
1	[kenichan, dived, many, time, ball, managed, s...	0	
2	[whole, body, feel, itchy, like, fire]	0	
3	[nationwideclass, behaving, mad, see]	0	
4	[kwesidei, whole, crew]	0	

Word cloud positive data

[illegible]

5/7



- TFIDF for sentiment analysis

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer()
xt = tfidf.fit_transform(data_cleaned['text'])
y = data_cleaned['sentiment']
```

- ▼ Train Test Split

[illegible]

- ▼ LogisticRegression

```
from sklearn.metrics import accuracy_score, confusion_matrix, classification_report
from sklearn.linear_model import LogisticRegression
```

```
lr = LogisticRegression()  
lr.fit(X_train,y_train)
```

```
▼ LogisticRegression
LogisticRegression()
```

- Accuracy Score: 85%

```
y_train_pred = lr.predict(X_train)
y_test_pred = lr.predict(X_test)
accuracy score(y train,y train pred)*100
```

85.72699464659102

✓ 0s completed at 12:12 AM

● ×