

Hate Speech Recognizer

Yashita Vajpayee
Masters in Data science
Stevens Institute of Technology
Hoboken, USA
yvajpaye@stevens.edu

Sharven Subhash Rane
Masters in Data science
Stevens Institute of Technology
Hoboken, USA
srane8@stevens.edu

Kalyan Varma Patchamatla
Masters in Data science
Stevens Institute of Technology
Hoboken, USA
kpatcham@stevens.edu

Abstract—The aim of this project is to detect hate speech from twitter data using logistic regression, Random forest and Naive Bayes methods and report a comparison between results. In this mid-term report we have done pre-processing of the data and implemented Logistic regression

I. INTRODUCTION

In recent years, the advent of social media platforms has led users to freely express their opinions on various subjects, including politics, society, health, education, finance, and even business-related issues. However, this widespread usage of social media has also increased the risk of its misuse by some groups resulting in spreading hate speeches or offensive language. This new problem of hate speech on social media has been addressed by recent studies that utilized a number of feature engineering techniques and machine learning algorithms.

II. OUR SOLUTION

In this project, we are going to do sentiment analysis and hate speech recognition. The data from Twitter, a popular microblogging social media platform for sharing short digital content, was used for this project. Preprocessing of data involves basic regex operations, data cleaning and lemmatization. Three machine learning methods will be investigated: logistic regression (LR), random forest (RF), naive Bayes (NB). We use hyperparameter tuning for optimum learning rate and increase accuracy of the model

III. DESCRIPTION OF DATASET

Dataset comprises of three columns : id, label, tweet. Id is the id assigned to a particular tweet and label. Given a training sample of tweets and labels, where label '1' denotes the tweet is racist/sexist and label '0' denotes the tweet is not racist/sexist. We say a tweet contains hate speech if it has a racist or sexist sentiment associated with it. Full tweet texts are provided with their labels for training data. Mentioned users' username is replaced with @user.

Data pre-processing is done as following: Firstly we converted upper case to lower case. Following this we removed URL and hashtags and punctuations. We also removed all the stop words from english language stop words. Then we proceeded to remove any duplicate data by dropping one of them. Finally we perform lemmatization on the processed data.

```
In [18]: #creating a function to process the data
def data_processing(tweet):
    tweet = tweet.lower()
    tweet = re.sub(r"https?|www?|http?", '', tweet, flags = re.MULTILINE)
    tweet = re.sub(r"@w+|\#", '', tweet)
    tweet = re.sub(r"[^a-zA-Z]", '', tweet)
    tweet = re.sub(r'0+', '', tweet)
    tweet_tokens = word_tokenize(tweet)
    filtered_tweets = [w for w in tweet_tokens if not w in stop_words]
    return " ".join(filtered_tweets)

In [19]: tweet_df.tweet = tweet_df['tweet'].apply(data_processing)

In [20]: tweet_df = tweet_df.drop_duplicates('tweet')

In [21]: lemmatizer = WordNetLemmatizer()
def lemmatizing(data):
    tweet = [lemmatizer.lemmatize(word) for word in data]
    return data

In [22]: tweet_df['tweet'] = tweet_df['tweet'].apply(lambda x: lemmatizing(x))
```

IV. DATA VISUALIZATION

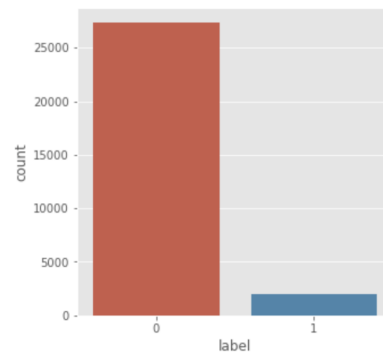
The processed data after lemmatization is visualized using count plot and pie-chart as following:

A. Count plot:

Tells the count of both the labels in the form of bar graph

```
In [26]: fig = plt.figure(figsize=(5,5))
sns.countplot(x='label', data = tweet_df)

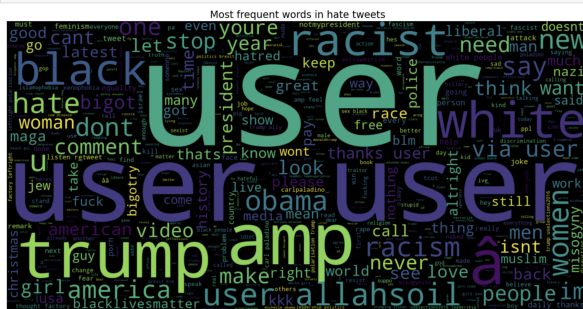
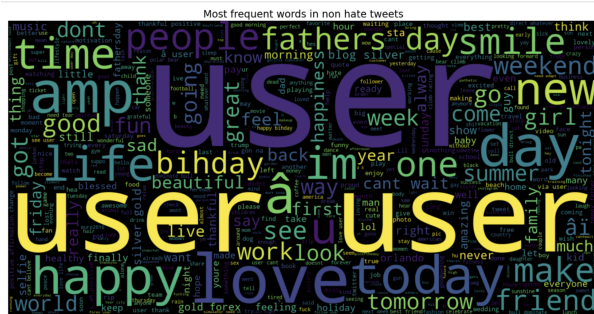
Out[26]: <AxesSubplot:xlabel='label', ylabel='count'>
```



B. Pie-chart:

Pie-chart showing both the labels. It can be seen that label 1 i.e. hate speech is 6.8 per cent and non-hate speech is 93.2 per cent.

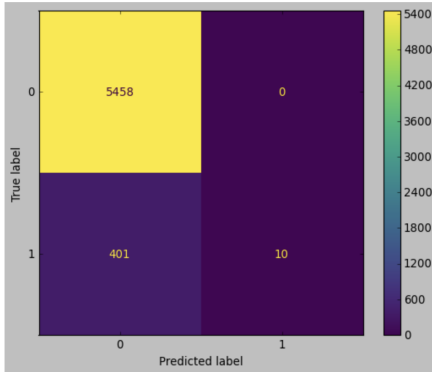
```
In [27]: fig = plt.figure(figsize=(7,7))
         colors = ("red", "gold")
         wp = plt.subplot(1,2, 'edgecolor':'black')
         tags = tweet_df['label'].value_counts()
         explode = (0.1, 0.1)
         tags.plot(kind='pie', autopct = '%1.1%', shadow=True, colors = colors, startangle = 90,
         wedgeprops = wp, explode = explode, label=wp)
         plt.title('Distribution of sentiments')
```



We also printed confusion matrix for performance analysis of our model as below:

```
style.use('classic')
cm = confusion_matrix(y_test, logreg_predict, labels=logreg.classes_)
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=logreg.classes_)
disp.plot()

<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x1f17be5f2e0>
```



B. Hyperparameter Tuning

Next we performed hyperparameter tuning to improve the performance of our model. We used GridSearchCV for this task and for best cross-validation score and best parameter. Following this we predicted our test values. This can be seen as below:

```
from sklearn.model_selection import GridSearchCV
import warnings
warnings.filterwarnings('ignore')

param_grid = {'C': [100, 10, 1.0, 0.1, 0.01], 'solver': ['newton-cg', 'lbfgs', 'liblinear']}
grid = GridSearchCV(LogisticRegression(), param_grid, cv = 5)
grid.fit(x_train, y_train)
print("Best Cross validation score: {:.2f}".format(grid.best_score_))
print("Best parameters: ", grid.best_params_)

Best Cross validation score: 0.95
Best parameters: {'C': 100, 'solver': 'newton-cg'}

y_pred = grid.predict(x_test)
```

Lastly we checked the accuracy of the model after hyperparameter tuning and printed the confusion matrix for performance analysis

```
logreg_acc = accuracy_score(y_pred, y_test)
print("Test accuracy: {:.2f}%".format(logreg_acc*100))
```

Test accuracy: 94.89%

```
print(confusion_matrix(y_test, y_pred))
print("\n")
print(classification_report(y_test, y_pred))
```

```
[[5450   8]
 [ 292 119]]
```

	precision	recall	f1-score	support
0	0.95	1.00	0.97	5458
1	0.94	0.29	0.44	411
accuracy			0.95	5869
macro avg	0.94	0.64	0.71	5869
weighted avg	0.95	0.95	0.94	5869

It can be seen that the model accuracy improved to 94.89 per cent.

VII. FUTURE DIRECTIONS

Results will be compared with two other machine learning methods: Random Forest (RF) and Naive Bayes (NB), trying to obtain the most accurate classification. In the future, this

can be implemented for larger datasets and our final goal is to implement this project in any social media platform and confront the current hate speech issues.

VIII. CONCLUSION

The extensive experiments and analysis show logistic regression to be a good algorithm for hate speech detection. More efficient algorithms can be used alongside deep learning and neural networks techniques to further enhance the model. This is to provide important insights into their detection accuracy, computational efficiency, capability in using pre-trained models, and domain generalizability for their deployment in real-world applications.

REFERENCES

- [1] <https://www.kaggle.com/datasets/arkhoshghalb/twitter-sentiment-analysis-hatred-speech?select=train.csv>
- [2] <https://towardsdatascience.com/logistic-regression-detailed-overview-46c4da4303bc>
- [3] <https://realpython.com/logistic-regression-python/>