# Depression Detection on Twitter

PROJECT REPORT

Anto Francis

Candidate ID: C0825095

Lambton College, Toronto

Canada

Omer Volkan Guney

Candidate ID: C0831373

Lambton College, Toronto

Canada

Sachin Sreekumar

Candidate ID: C0825096

Lambton College, Toronto

Canada

Rupesh Chandran

Candidate ID: C0826779

Lambton College, Toronto

Canada

*Abstract* — **The project focuses on identifying depressive tweets posted on Twitter. Twint, a python library, is used for extracting tweets from Twitter without using Twitter API and authentication. Using Twint, three types of tweets are extracted: depressed, neutral, and positive. All the ML techniques and preprocessing techniques are performed on the dataset, which is the combined tweets of the three types. Multivariate classification models are performed on the dataset, and a comparison of the models is illustrated in the project.**

## I. INTRODUCTION

Depression is a common form of disability that negatively affects the way someone thinks and act. It can result in losing interest in those activities you once enjoyed. It leads to various emotional and physical problems and can affect your ability to do things. Above 400 million people suffer from depression, and out of that, a fraction of people needs medical treatment to cure it. It is becoming the leading cause of suicide and death, even in teens.

Most youngsters use social media. Hence social media is one of the platforms to identify the early stages of depression, especially in youth. Every second, more than 6K tweets are generated on Twitter, which results in around 250 billion tweets per year. This project identifies and analyses textual markers associated with depression symptoms in order to create an algorithm that effectively predicts depression. Hence, by creating an algorithm that can analyze tweets and can be used by medical professionals.

## II. RESOURCES

Link to the colab file:
https://colab.research.google.com/drive/1kKlydS6AliSaFlQbk NqBwCV0kJxKE28i#scrollTo=_IiZ-n497_2W

Link to one drive account: CBD_Final_Project_resources

## III. DATA SCRAPPING

Data is scraped from Twitter using the python library twint. Unlike Twitter API, no authentication is required to scrap tweets from Twitter. And there is no limit to amount of tweets that can be extracted compared to the 3,200 tweet limit of Twitter API. Using twint, three types of tweets are scrapped, namely depressed, neutral and positive. Depressed tweets are scrapped with keywords depressed or lonely; neutral tweets are scrapped with keywords unbiased or impartial, and positive tweets are scrapped with keywords happy or love. There are about 10,000 tweets for each category, so a total of 30,000 tweets are extracted from Twitter.

## IV. ANALYSIS

### i. Loading the Dataset

The extracted tweets are loaded into the respective dataset such that depressive tweets are loaded into the depression dataset, neutral tweets are loaded into the neutral dataset, and positive tweets are loaded into the positive dataset.

### ii. Cleaning the Dataset

After loading the dataset, an extra attribute is added to each dataset according to its type. Target is the new attribute added to each dataset in such a way that for the depressed dataset, the target attribute is '0', for the neutral dataset target variable is '2' and for the positive dataset is '4', respectively. Each dataset contains about 40 attributes, including the newly assigned target attribute. Most of these features are not required for analysis and are dropped. Only the required attributes are maintained.

```
df3.columns

Index(['id', 'conversation_id', 'created_at', 'date', 'timezone', 'place',
       'tweet', 'language', 'hashtags', 'cashtags', 'user_id', 'user_id_str',
       'username', 'name', 'day', 'hour', 'link', 'urls', 'photos', 'video',
       'thumbnail', 'retweet', 'nlikes', 'nreplies', 'nretweets', 'quote_url',
       'search', 'near', 'geo', 'source', 'user_rt_id', 'user_rt',
       'retweet_id', 'reply_to', 'retweet_date', 'translate', 'trans_src',
       'trans_dest'],
      dtype='object')
```

Attributes 'id', 'date', 'tweet', 'username', 'target' are the required attributes for analysis. After selecting these attributes, each dataset is checked for null values, there are no null values, but as a handler, the rows with null values are dropped for higher accuracy.
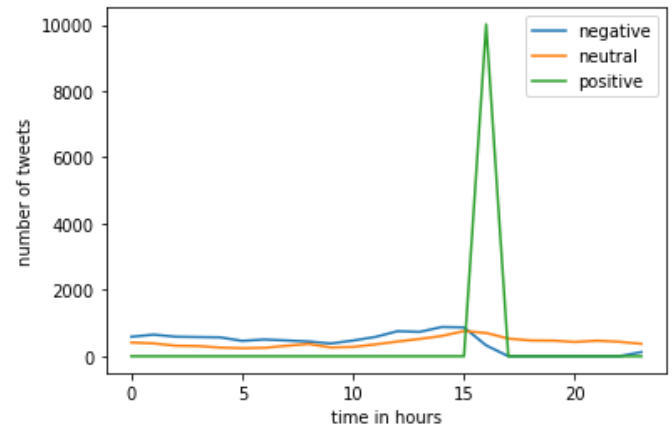
### iii. Final Dataset

The final dataset is obtained by concatenating the three datasets, and the row values are also shuffled.

```
[ ] df
```

| | id | date | tweet | username | target |
|---|---|---|---|---|---|
| 6801 | 1516452499082858500 | 2022-04-19 16:23:44 | @spideyferrari16 Happy birthday! | formulanhlcarly | 4 |
| 6250 | 1515730751018389512 | 2022-04-17 16:35:46 | @Sobana_Blupaw @gopisdirty @FECKFECK611 @CodeO... | Columbro1 | 2 |
| 5846 | 1516312760002674688 | 2022-04-19 07:08:28 | @SchizophrenicNY I have depression | Enders_Pattern | 0 |
| 7233 | 1515603464671285250 | 2022-04-17 08:09:59 | @trusty_reviews @SinghGreg @ohhnoitsco Unbiase... | HoneyIsVegan | 2 |
| 4294 | 1516366470179078149 | 2022-04-19 10:41:53 | @saintorochimaru I like to think of it as "Mis... | KDsMomFIRST | 0 |

### iv. Timeline of tweets in a day

The behavior of tweets is visualized over the hours of a day. For this visualization, the date must be converted into date-time format, which is taken care of by the date-time format library in python. With the help of the target attribute, the behavior of the tweet can be identified over time.



The x-axis represents the hours in a day, and the y-axis represents the number of tweets in the hours of the day. Three different behavior of the tweets is shown in the figure. It must be noted that positive or happy tweets are happening around 15:00 to 17:00.

### v. Behavior-based on tweet handles

Tweets containing '@' are taken for detecting the behavior of the tweet. The count of positive, neutral, and depressive tweets is taken by comparing the values of target attributes. Only the tweets which contain the symbol '@' are taken for the count.

The x-axis represents the types of tweet behavior, and the y-axis represents the number of tweets. It is evident from the graph that the number of neutral tweets is highest based on the Twitter handlers.

## vi.    Tweet Preprocessing

Tweets in the dataset need to be preprocessed before fitting the data into the model. The following operations are performed on the tweets to obtain cleaned tweets for further analysis:

- Removing Html tags
- Removing twitter handles.
- Removing Punctuations
- Removing Stopwords
- Tokenization
- Stemming using PorterStemmer()
- Combining words together and removing small words whose length is less than3.
- Check for null values in the cleaned tweet. If present, drop the rows respectively.

## vii.    Word Cloud Visualization

After preprocessing of tweets, the cleaned tweets are visualized in the word cloud to see the most common words. World cloud is visualized for each dataset: depressed, neutral, happy, and the combined dataset.


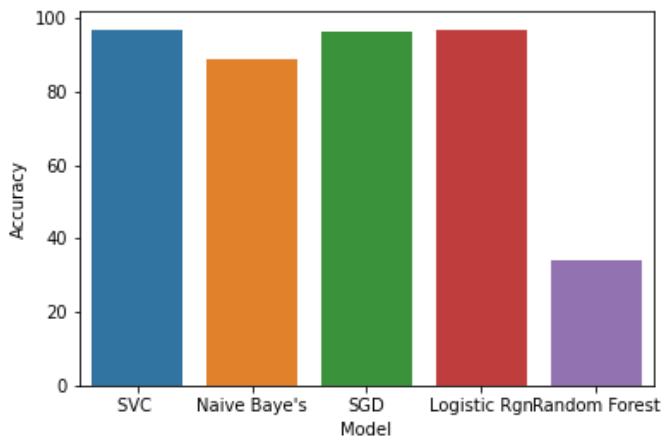
## viii.    Model Building

Various multivariate classifier models are used in the analysis for classifying the depressing tweets, neutral tweets, and positive tweets. For feeding the input cleaned tweet into various models, the tweets are vectorized using TF-IDF, and the input data is split into train and test data where the model trains the data. Test data is given to the model, and the accuracy of the model is determined using a classification report. The accuracy of each model is described below in the chart.

```
+-----------------------------------------------------+
|                 Model Comparision                   |
+-------------------------+-----------+---------------+
|          Model          | Tokenizer |   Accuracy    |
+-------------------------+-----------+---------------+
|           SVC           |   TFIDF   |     96.782    |
|  Naive Baye's Classifier|   TFIDF   |     88.693    |
|      SGD Classifier     |   TFIDF   |     96.405    |
|    Logistic Regression  |   TFIDF   |     96.871    |
| Random Forest Classifier|   TFIDF   |     34.166    |
+-------------------------+-----------+---------------+
```

The classification model accuracy is visualized in a bar graph which is shown below:

## V. CONCLUSION

For scrapping tweets, twint is easier to use than Twitter API. Twint is free and can extract more tweets from Twitter compared to Twitter. By using twint, data is almost unskewed such that all categories of tweets are equal in number. Various multivariate classifiers are used for detecting depressive tweets from the dataset. It is found that most of the classifiers are giving excellent accuracy results except for the random forest classifier. The logistic Regression classifier achieves the maximum accuracy of 96.871%.

## VI. REFERENCES

[1] Redaabdou. (2019, November 15). Depression on social media. Kaggle. Retrieved April 20, 2022, from https://www.kaggle.com/code/redaabdou/depression-on-social-media

[2] Yonebayashi. (2019, July 23). WIP: Detect early depression through tweets. Kaggle. Retrieved April 20, 2022, from https://www.kaggle.com/code/yonebayashi/wip-detect-early-depression-through-tweets/notebook

[3] Datetime - basic date and time types¶. datetime - Basic date and time types - Python 3.10.4 documentation. (n.d.). Retrieved April 20, 2022, from https://docs.python.org/3/library/datetime.html

[4] D, E. (2019, December 15). Scraping tweets off Twitter with TWINT. Medium. Retrieved April 20, 2022, from https://medium.com/@erika.dauria/scraping-tweets-off-twitter-with-twint-a7e9d78415bf