# Sentiment Analysis with Twitter API

# Feedback and Reviews matter!

- Companies need reviews and feedback on products during development stage
- Governments need to know popular opinion over different subjects
- Brands need to know general preferences before launching new products
- Services providers need to track clients satisfaction

But where to find all this information?

# On the Internet of course!

Social Media platforms provide tons of thoughts, reviews, opinions and feedback on the most diverse topics and fields. On Twitter, more than 340 million people engage on the platform providing an infinite source of data.

# Analysing the US elections of 2020

The project:

- Twitter API
- TextBLob
- SKlearn
- MatplotLib
- Decision Tree Algorithm
- Tags on twitter: #DonaldTrump, #JoeBiden and #USElections2020

# Data Gathering

Twitter API

```
1   import tweepy as tw
2
3   #Authentication protocol to access the twitter developer account
4   consumer_key = 'consumer_key'
5   consumer_secret = 'consumer_secret'
6   access_key = 'access_key'
7   access_secret = 'access_secret'
8
9   auth = tw.OAuthHandler(consumer_key, consumer_secret)
10  auth.set_access_token(access_key, access_secret)
11  api = tw.API(auth, wait_on_rate_limit=True)
12
13  # Define the search term and the date_since date as variables
14  search_words = ["#JoeBiden"]
15  #user = "JoeBiden"
16  date_since = "2020-11-01"
17
18  # Collect tweets
19  tweets = tw.Cursor(api.search, q=search_words, lang="en", since=date_since, tweet_mode="extended").items(1000)
20  #tweets = tw.Cursor(api.user_timeline, id=user, lang="en", since=date_since, tweet_mode="extended").items(200)
```

# Data Cleaning

```python
def clean_tweet(tweet):
    '''
    Utility function to clean tweet text by removing links, special characters
    using simple regex statements.
    '''
    return ' '.join(re.sub("(@[A-Za-z0-9]+)|([^0-9A-Za-z \t])|(\w+:\/\/\S+)", " ", tweet).split())
```

# Getting the Sentiment

**TextBlob and Natural Language Processing (NLP)**

```python
1   def get_tweet_sentiment(tweet):
2       '''
3       Utility function to classify sentiment of passed tweet
4       using textblob's sentiment method
5       '''
6       # create TextBlob object of passed tweet text
7       analysis = TextBlob(clean_tweet(tweet))
8       # set sentiment
9       if analysis.sentiment.polarity > 0:
10          return 'positive'
11      elif analysis.sentiment.polarity == 0:
12          return 'neutral'
13      else:
14          return 'negative'
15
16  # Collect a list of tweets
17  tweet_list = [clean_tweet(tweet.full_text) for tweet in tweets]
18  #print(tweet_list)
19  sentiment_list = [get_tweet_sentiment(tweet) for tweet in tweet_list]
20  #print(sentiment_list)
21
22  tweets_and_sentiments = pd.DataFrame({"Text":tweet_list, "Sentiment":sentiment_list})
23  tweets_and_sentiments.to_csv("Coronavirus.csv", index=False)
24
25  #print(tweets_and_sentiments)
```

# Visualization Tools

MatplotLib

```python
1   #imports and reads the dataset
2   data = "USElections2020.csv"
3   tweets_list = pd.read_csv(data)
4   #print(tweets_list)
5
6
7   #set the plot size
8   plot_size = plt.rcParams["figure.figsize"]
9
10  plot_size[0] = 8
11  plot_size[1] = 6
12  plt.rcParams["figure.figsize"] = plot_size
13
14
15  #plt.subplot(2, 1, 2) #divides the window in two plots
16  tweets_list.Sentiment.value_counts().plot(kind='pie', autopct='%1.0f%%')#creates a pie graph
17
18  plt.show()#displays the plot window
```
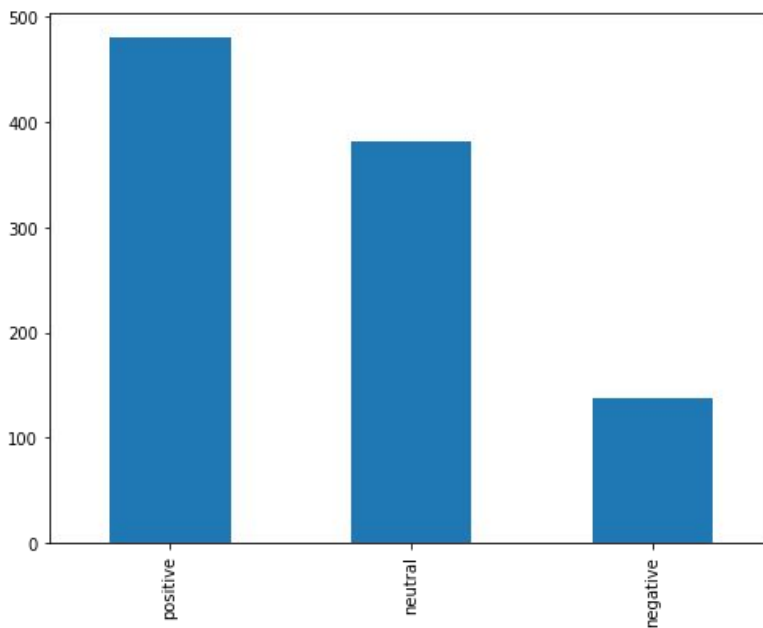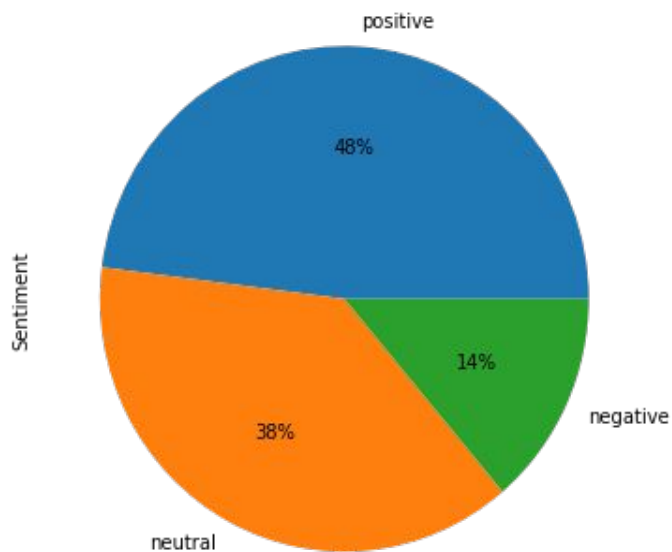
# Prediction Model

SKlean and Decision Trees

```python
1   from sklearn.feature_extraction.text import CountVectorizer
2   from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
3   from sklearn.ensemble import RandomForestClassifier
4   from sklearn.model_selection import train_test_split
5
6   corpus = tweet_list
7
8   vectorizer = CountVectorizer(max_features=2500)
9   #print( vectorizer.fit_transform(corpus).toarray() )
10  processed_features = vectorizer.fit_transform(corpus).toarray()
11  #print(processed_features)
12  labels = sentiment_list
13
14  #print(features)
15  print(labels)
16  print(len(labels))
17
18  X_train, X_test, y_train, y_test = train_test_split(processed_features, labels, test_size=0.2, random_state=0)
19
20  print(len(X_train),len(X_test), len(y_train), len(y_test))
21
22  text_classifier = RandomForestClassifier(n_estimators=200, random_state=0)
23  text_classifier.fit(X_train, y_train)
24
25  predictions = text_classifier.predict(X_test)
26  #print(predictions)
27
28  print(confusion_matrix(y_test,predictions))
29  print(classification_report(y_test,predictions))
30  print('ACCURACY = ', accuracy_score(y_test, predictions))
```
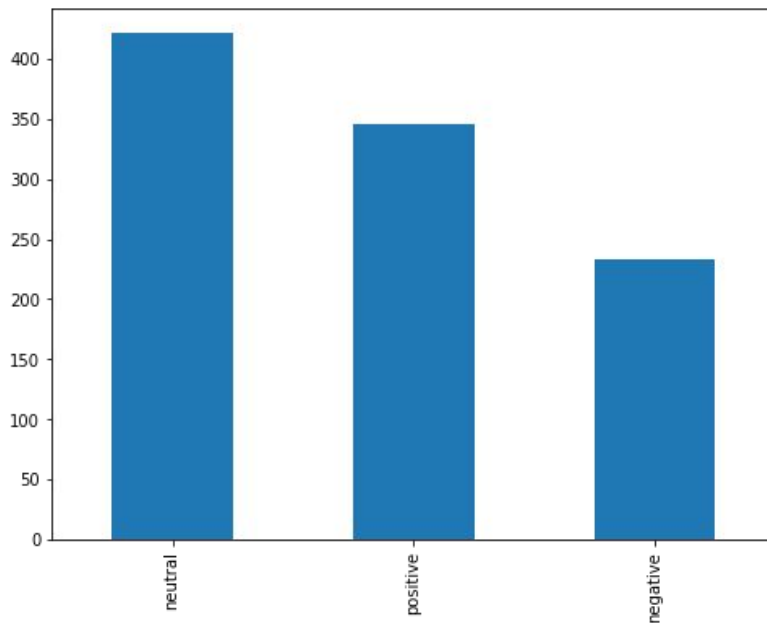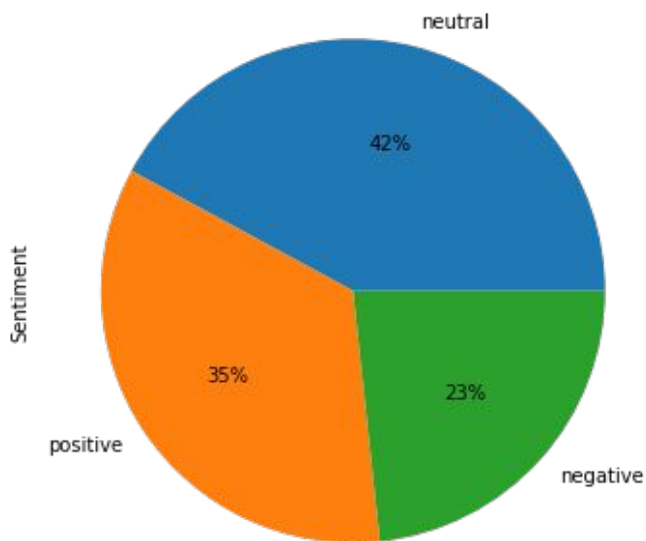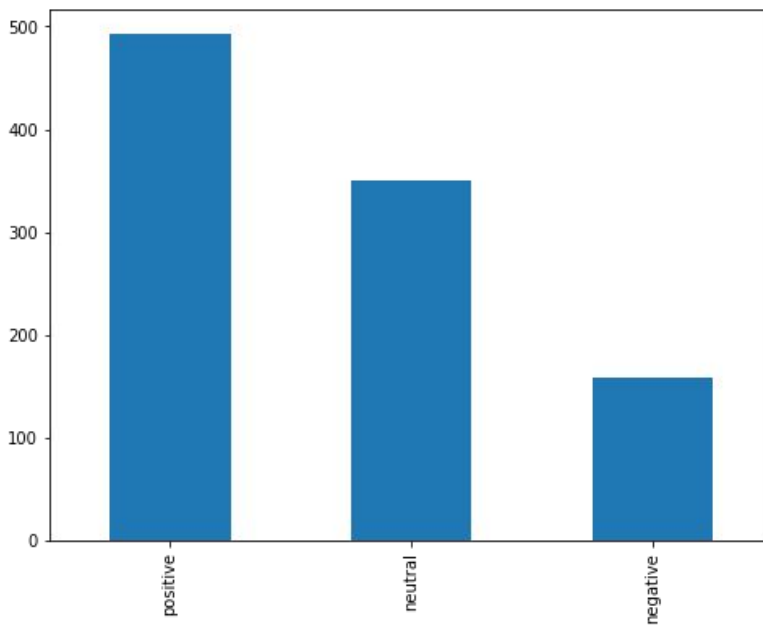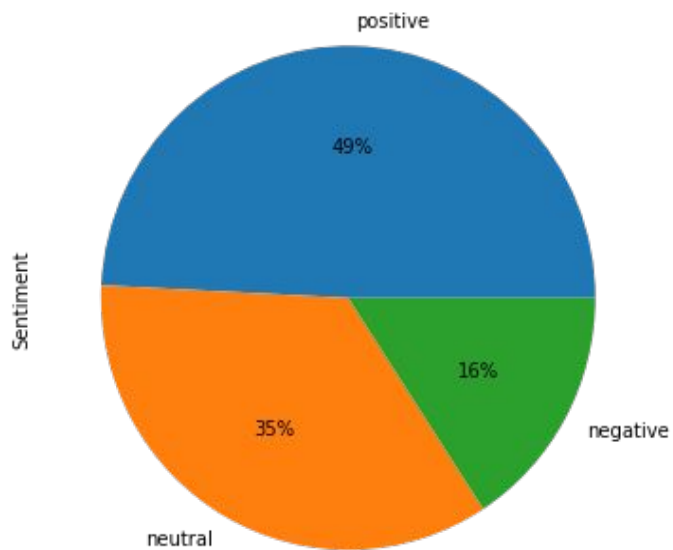
# #USElections2020 Results

# #DonaldTrump Results

# #JoeBiden Results

# Prediction Model Metrics

ACCURACY = 0.82

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| negative | 0.91 | 0.36 | 0.51 | 28 |
| neutral | 0.80 | 0.93 | 0.86 | 80 |
| positive | 0.83 | 0.87 | 0.85 | 92 |
| accuracy |  |  | 0.82 | 200 |
| macro avg | 0.85 | 0.72 | 0.74 | 200 |
| weighted avg | 0.83 | 0.82 | 0.81 | 200 |

# Conclusion and Further Improvements

- Good Results with simple and cheap implementation
- More Data = More Accuracy
- Fields Flexibility
- NLP is still a developing technology
- Improvements on the specific vocabulary
- Possibilities are limitless!!

# Thank you!!