

```
In [524]: 1 import tweepy
2 import pandas as pd
3 import json
4 import re
5 import warnings
6 warnings.filterwarnings("ignore")
7 from nltk.corpus import stopwords
8 from nltk.stem.wordnet import WordNetLemmatizer
9 import gensim
10 from gensim import corpora
11 import seaborn as sns
12 import matplotlib.pyplot as plt
13 from textblob import TextBlob
14 from sklearn.feature_extraction.text import TfidfVectorizer
15 from sklearn.datasets import make_blobs
16 from sklearn.decomposition import PCA
17 from sklearn.preprocessing import normalize
18 from sklearn.metrics import pairwise_distances
19 import nltk
20 import string
21 import numpy as np
22 from wordcloud import WordCloud, STOPWORDS
23 from sklearn.feature_extraction.text import TfidfVectorizer
24 from nltk import sent_tokenize, word_tokenize, pos_tag
25 from PIL import Image
26 import gensim
27 from gensim import corpora
28 from nltk.stem.wordnet import WordNetLemmatizer
29 from sklearn.feature_extraction.text import CountVectorizer
30 from sklearn.decomposition import LatentDirichletAllocation
```

```
In [65]: 1 pd.set_option('display.max_rows', None)
2 pd.set_option('display.max_columns', None)
3 pd.set_option('display.width', None)
4 pd.set_option('display.max_colwidth', None)
```

```
In [21]: 1 %run ./key.ipynb
2 auth = tw.OAuthHandler(consumer_key,consumer_secret)
3 auth.set_access_token(access_token,access_secret)
```

```
In [25]: 1 output_file = 'tweeter.csv'
2 tweets_to_capture = 20000
```

```
In [26]: 1 tweet_list=[]
2 class MyStreamListener(tweepy.StreamListener):
3     def __init__(self,api=None):
4         super(MyStreamListener,self).__init__()
5         self.num_tweets=0
6         self.file=open(output_file,"w")
7     def on_status(self,status):
8         tweet=status._json
9         self.file.write(json.dumps(tweet)+ '\n')
10        tweet_list.append(status)
11        self.num_tweets+=1
12        if self.num_tweets<= tweets_to_capture:
13            return True
14        else:
15            return False
16        self.file.close()
```

```
In [27]: 1 %%time
2 l = MyStreamListener()
3 stream =tweepy.Stream(auth,l)
4 #this line filters twitter streams to capture data by keywords
5 stream.filter(track=['Bangladesh'])
```

Wall time: 2h 36min 15s

```
In [ ]: 1
```

```
In [543]: 1 df=pd.read_csv('C:\\Users\\MESSI\\Documents\\tweeter.csv',encoding='latin1',
```

In [547]:

1 df.info()

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 17 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Unnamed: 0            20000 non-null  int64
 1   text                  20000 non-null  object
 2   favorited             20000 non-null  bool
 3   favoriteCount         20000 non-null  int64
 4   replyToSN            1225 non-null   object
 5   created               20000 non-null  object
 6   truncated             20000 non-null  bool
 7   replyToSID           1171 non-null   float64
 8   id                    20000 non-null  int64
 9   replyToUID           1225 non-null   float64
10   statusSource          20000 non-null  object
11   screenName            20000 non-null  object
12   retweetCount          20000 non-null  int64
13   isRetweet             20000 non-null  bool
14   retweeted             20000 non-null  bool
15   longitude             1 non-null      float64
16   latitude              1 non-null      float64
dtypes: bool(4), float64(4), int64(4), object(5)
memory usage: 2.1+ MB

```

In [548]:

1 df.describe()

Out[548]:

	Unnamed: 0	favoriteCount	replyToSID	id	replyToUID	retweetCount	longit
count	20000.000000	20000.000000	1.171000e+03	2.000000e+04	1.225000e+03	20000.000000	1.0
mean	10000.500000	0.582100	1.440683e+18	1.450593e+18	5.024448e+17	1088.191600	90.4
std	5773.647028	23.777835	8.254990e+16	3.279439e+13	5.857636e+17	1400.126043	
min	1.000000	0.000000	7.193971e+09	1.450543e+18	1.200000e+01	0.000000	90.4
25%	5000.750000	0.000000	1.450359e+18	1.450561e+18	1.407433e+08	47.000000	90.4
50%	10000.500000	0.000000	1.450518e+18	1.450591e+18	2.876471e+09	383.000000	90.4
75%	15000.250000	0.000000	1.450565e+18	1.450625e+18	1.084767e+18	2016.000000	90.4
max	20000.000000	2243.000000	1.450641e+18	1.450641e+18	1.450618e+18	9177.000000	90.4

In [549]:

1 len(df)

Out[549]: 20000

In [550]: 1 df.tail(4)

Out[550]:

Unnamed: 0		text	favorited	favoriteCount	replyToSN	created	trunc
19996	19997	RT @upadhyayabhii: Big--<U+092C> <U+0921><U+093C> <U+093E> <U+0938> <U+091A> <U+0906> <U+092F><U+093E> <U+0938><U+093E> <U+092E><U+0928> <U+0947><U+0964> <U+092C><U+093E> <U+0902><U+0917> <U+094D><U+0932> <U+093E><U+0926> <U+0947><U+0936> <U+0915><U+0947> <U+0917><U+0943> <U+0939> <U+092E> <U+0902><U+0924> <U+094D><U+0930> <U+0940> <U+0915> <U+093E> <U+092C> <U+0921><U+093C> <U+093E> <U+092C> <U+092F><U+093E> <U+0928>- <U+0926> <U+0941><U+0930> <U+094D><U+0917> <U+093E> <U+092A> <U+0942><U+091C> <U+093E> <U+092A> <U+0902><U+0921> <U+093E><U+0932> <U+094B><U+0902> <U+092A><U+0930> <U+0939><U+092E> <U+0932><U+0947> <U+0915><U+0947> <U+0932><U+093F> <U+090F> <U+092A> <U+0939><U+0932> <U+0947> <U+0938> <U+0947> <U+0930> <U+091A><U+0940> <U+0917><U+0908> <U+0925><U+0940> <U+0938><U+093E>	False	0	NaN	2021-10-19 19:22:28	f
19997	19998	Breaking news. @SuPriyoBabul is planning to go to Bangladesh, he may have a better opportunity for a cabinete minis https://t.co/1XFmYjnKqd	False	51	dasgobardhan	2021-10-19 19:22:28	

Unnamed: 0		text	favorited	favoriteCount	replyToSN	created	trunc
19998	19999	RT @ManMundra: It has gone beyond a point of no return. You can't stop it. It has happened in Afganistan, happening in Pakistan and will ha	False	0	NaN	2021-10-19 19:22:27	f
19999	20000	RT @AskAnshul: After 'Rohingya football club', here is 'Bangladesh Youth'. And, later they changed it to 'Miya Bhai Youth' due to outrage.	False	0	NaN	2021-10-19 19:22:27	f

Cleaning

```
In [551]: 1 def cleaner(tweet):
2         tweet = re.sub("@[A-Za-z0-9]+", "", tweet) #Remove @ sign
3         tweet = re.sub(r"(?:\@|http?\:\/\/|https?\:\/\/|www)\S+", "", tweet) #Remove
4         tweet = " ".join(tweet.split())
5         #tweet = ''.join(c for c in tweet if c not in emoji.UNICODE_EMOJI) #Remo
6         tweet = tweet.replace("#", "").replace("_", " ") #Remove hashtag sign bu
7         tweet = re.sub(r'[RT:]+', '', tweet) #replace RT-tags
8         tweet = re.sub('<[>]+>', '', tweet)
9
10        return tweet
11 df['text'] = df['text'].map(lambda x: cleaner(x))
```

In [552]:

1df.head(6)

Out[552]:

Unnamed: 0		text	favorited	favoriteCount	replyToSN	created	truncated	replyToSID
0	1	BIG At the request of the Bangladesh Govt. witter deletes Bangladesh Hindu Unity Council's twitter ha	False	0	NaN	2021-10-20 01:52:15	False	NaN
1	2	Hundreds protest in Bangladesh over religious violence	False	0	NaN	2021-10-20 01:52:14	False	NaN
2	3	We are appalled by recent reports of deadly attacks on the Hindu community in Bangladesh. All, including members of religious	False	0	NaN	2021-10-20 01:52:13	False	NaN
3	4	Now, protests in USA against the violence on emples & Hindus in Bangladesh.	False	0	NaN	2021-10-20 01:52:12	False	NaN
4	5	in	False	0	NaN	2021-10-20 01:52:11	False	NaN
5	6	Bangladesh 27 20WorldCup2021 Mahmudullah	False	0	NaN	2021-10-20 01:52:11	True	NaN

In [553]:

1df = df.drop(['replyToSN','longitude','latitude','replyToSID','replyToUID',

In [554]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 20000 entries, 0 to 19999
Data columns (total 7 columns):
 #   Column          Non-Null Count  Dtype
---  -
 0   text            20000 non-null  object
 1   created         20000 non-null  object
 2   id              20000 non-null  int64
 3   statusSource    20000 non-null  object
 4   screenName      20000 non-null  object
 5   retweetCount    20000 non-null  int64
 6   isRetweet       20000 non-null  bool
dtypes: bool(1), int64(2), object(4)
memory usage: 957.2+ KB
```

In [555]:

```
1 def clear(tweet):
2     tweet = re.sub('<[>]+>', '', tweet)
3     return tweet
4 df['statusSource'] = df['statusSource'].map(lambda x: clear(x))
```

```
In [556]: 1 df.head(4)
```

Out[556]:

	text	created	id	statusSource	screenName	retweetCount	isRetwe
0	BIG At the request of the Bangladesh Govt. witter deletes Bangladesh Hindu Unity Council's twitter ha	2021-10-20 01:52:15	1450640978625236992	Twitter for iPhone	arvind291	292	Tr
1	Hundreds protest in Bangladesh over religious violence	2021-10-20 01:52:14	1450640973856317445	Twitter Web App	gojharan	280	Tr
2	We are appalled by recent reports of deadly attacks on the Hindu community in Bangladesh. All, including members of religious	2021-10-20 01:52:13	1450640971083907073	Twitter for Android	amodbhardwaj	898	Tr
3	Now, protests in USA against the violence on emples & Hindus in Bangladesh.	2021-10-20 01:52:12	1450640963878096901	Twitter for Android	VamsiKandula2	4260	Tr

```
In [557]: 1 df['text'] = df['text'].str.encode('ascii', 'ignore').str.decode('ascii')
```


In [558]:

```
1 df.text
46
47
48
49
50
51
52
53
54
55
```

In [559]:

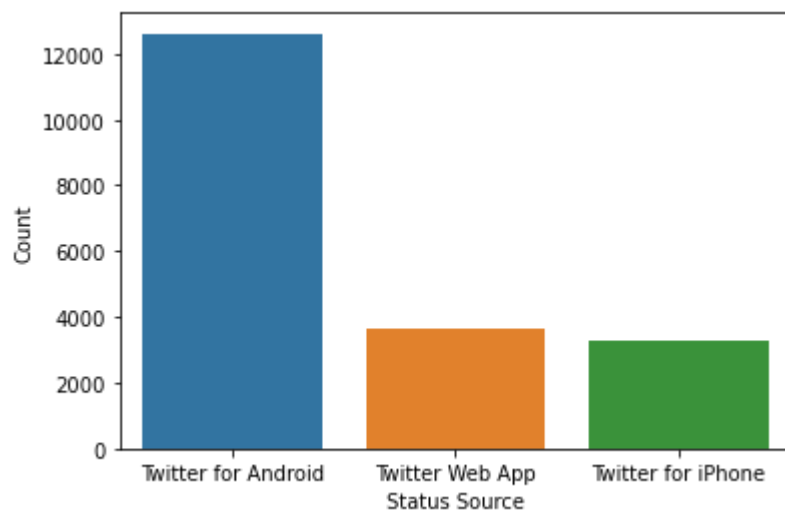
```
1 df["statusSource_count"] = 1
2 df_statusSource = df.groupby(['statusSource'], as_index=False, sort=False)[[
3 df_statusSource = df_statusSource.sort_values("statusSource_count", axis = 0
4 df_statusSource
```

Out[559]:

	statusSource	statusSource_count
2	Twitter for Android	12617
1	Twitter Web App	3661
0	Twitter for iPhone	3288

In [560]:

```
1
2 sns.barplot(df_statusSource.statusSource,df_statusSource.statusSource_count)
3 plt.xlabel("Status Source")
4 plt.ylabel("Count")
5 plt.show()
```

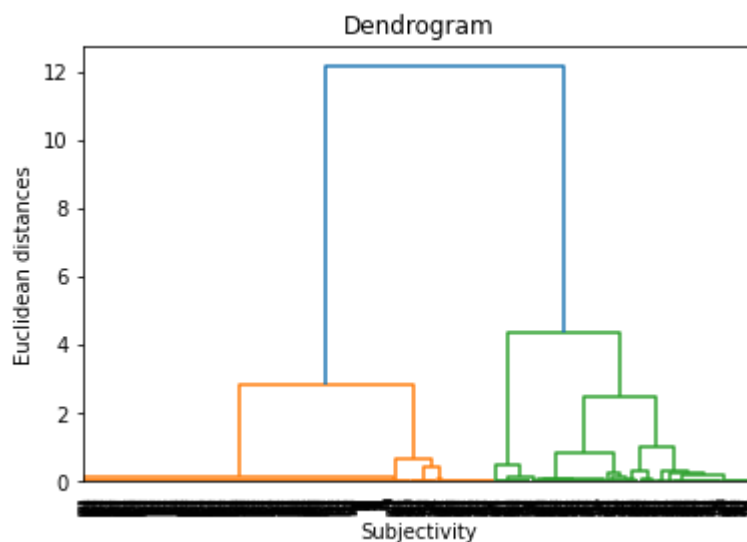


Clustering

In [561]:

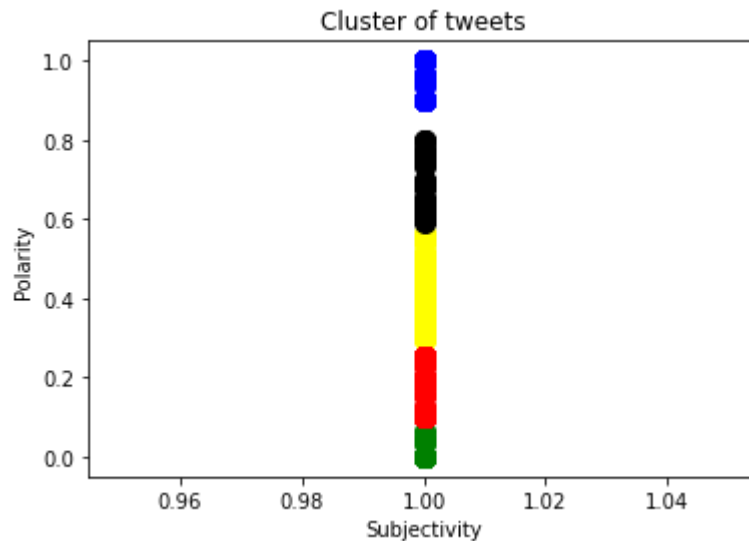
```
1 def subjectivity_check(tweet):
2     return TextBlob(tweet).sentiment.subjectivity
3 def polarity_check(tweet):
4     return TextBlob(tweet).sentiment.polarity
5
6 #Applying the values of subjectivity_check and polarity_check to our newly a
7 df['subjectivity']=df['text'].apply(subjectivity_check)
8 df['polarity']=df['text'].apply(polarity_check)
```

```
In [562]: 1 data=df.iloc[0:1000, [7,8]].values
2 import scipy.cluster.hierarchy as sch
3 dendrogram = sch.dendrogram(sch.linkage(data, method = 'ward'))
4 plt.title('Dendrogram')
5 plt.xlabel('Subjectivity')
6 plt.ylabel('Euclidean distances')
7 plt.show()
```



```
In [563]: 1 from sklearn.cluster import AgglomerativeClustering
2 clusters = AgglomerativeClustering(n_clusters=5, affinity='euclidean', linka
3 p_clusters=clusters.fit_predict(data)
```

```
In [564]: 1 plt.scatter(data[p_clusters == 0,0], data[p_clusters == 0,1], s=100, c='yellow')
2 plt.scatter(data[p_clusters == 1,0], data[p_clusters == 1,1], s=100, c='black')
3 plt.scatter(data[p_clusters == 2,0], data[p_clusters == 2,1], s=100, c='blue')
4 plt.scatter(data[p_clusters == 3,0], data[p_clusters == 3,1], s=100, c='green')
5 plt.scatter(data[p_clusters == 4,0], data[p_clusters == 4,1], s=100, c='red')
6 plt.title("Cluster of tweets")
7 plt.xlabel('Subjectivity')
8 plt.ylabel('Polarity')
9 plt.show()
```



K Means Clustering

```
In [565]: 1 data = df['text']
2 data.head()
```

```
Out[565]: 0 BIG At the request of the Bangladesh Govt. witte
r deletes Bangladesh Hindu Unity Council's twitter ha
1 H
undreds protest in Bangladesh over religious violence
2 We are appalled by recent reports of deadly attacks on the Hindu communi
ty in Bangladesh. All, including members of religious
3 Now, protests in USA again
st the violence on emples & Hindus in Bangladesh.
4
in
Name: text, dtype: object
```

Converting text to TF-IDF

```
In [566]: 1 tf_idf_vectorizer = TfidfVectorizer(stop_words = 'english',
2                                           max_features = 5000)
3 %time tf_idf = tf_idf_vectorizer.fit_transform(data.values.astype('U')) #ast
4 tf_idf_norm = normalize(tf_idf)
5 tf_idf_array = tf_idf_norm.toarray()
```

Wall time: 385 ms

```
In [567]: 1 pd.DataFrame(tf_idf_array, columns=tf_idf_vectorizer.get_feature_names()).he
```

Out[567]:

	000	007	01	034	04	07	070	0774	08	0mar	10	100	1000	1000s	101	101st	105	10
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0

In [568]:

```

1 class Kmeans:
2     def __init__(self, k, seed = None, max_iter = 200):
3         self.k = k
4         self.seed = seed
5         if self.seed is not None:
6             np.random.seed(self.seed)
7         self.max_iter = max_iter
8
9     ##Randomly initialising centroids which returns array of k centroids chosen
10    def initialise_centroids(self, data):
11        initial_centroids = np.random.permutation(data.shape[0])[:self.k]
12        self.centroids = data[initial_centroids]
13
14        return self.centroids
15
16    ##Compute distance of data from clusters and assign data point to closest cl
17    def assign_clusters(self, data):
18        if data.ndim == 1:
19            data = data.reshape(-1, 1)
20
21        dist_to_centroid = pairwise_distances(data, self.centroids, metric
22        self.cluster_labels = np.argmin(dist_to_centroid, axis = 1)
23
24        return self.cluster_labels
25
26    ##Computes average of all data points in cluster and assigns new centroids a
27
28    def update_centroids(self, data):
29        self.centroids = np.array([data[self.cluster_labels == i].mean(axis
30        return self.centroids
31    ##to Predict which cluster data point belongs to
32    def predict(self, data):
33        return self.assign_clusters(data)
34
35    ## contains the main loop to fit the algorithm Implements initialise centri
36    ## Returns instance of kmeans class
37
38    def fit_kmeans(self, data):
39        self.centroids = self.initialise_centroids(data)
40
41        for iter in range(self.max_iter):
42
43            self.cluster_labels = self.assign_clusters(data)
44            self.centroids = self.update_centroids(data)
45            if iter % 100 == 0:
46                print("Running Model Iteration %d " %iter)
47        print("Model finished running")
48        return self

```

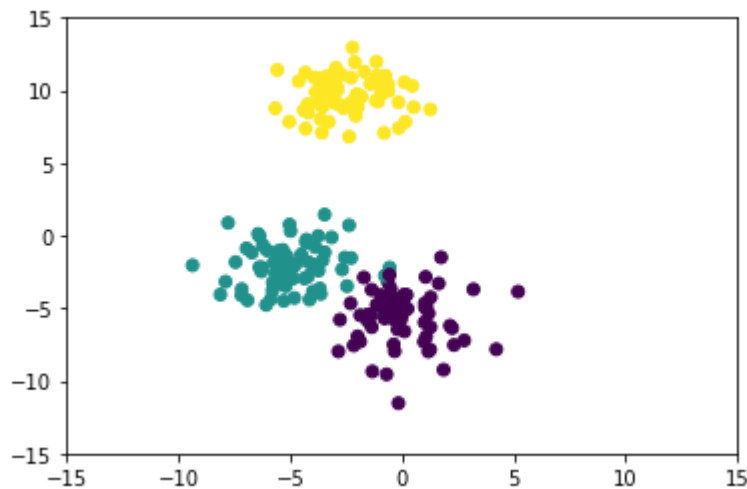
In [569]:

```

1  # create blobs
2  data = make_blobs(n_samples=200, n_features=2, centers=3, cluster_std=1.6, r
3  # create np array for data
4  points = data[0]
5  # create scatter plot
6  plt.scatter(data[0][:,0], data[0][:,1], c=data[1], cmap='viridis')
7  plt.xlim(-15,15)
8  plt.ylim(-15,15)
9
10
11 X = data[0]
12 X[2]

```

Out[569]: array([-3.58040006, 7.08578225])



In [570]:

```

1  temp_k = Kmeans(3, 1, 600)
2  temp_fitted = temp_k.fit_kmeans(X)
3  new_data = np.array([[1.066, -8.66],
4                        [1.87876, -6.516],
5                        [-1.59728965,  8.45369045],
6                        [1.87876, -6.516]])
7  temp_fitted.predict(new_data)

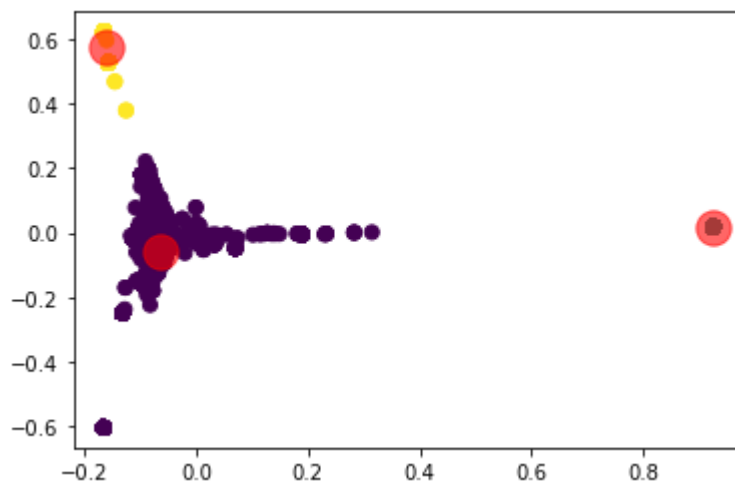
```

Running Model Iteration 0
Running Model Iteration 100
Running Model Iteration 200
Running Model Iteration 300
Running Model Iteration 400
Running Model Iteration 500
Model finished running

Out[570]: array([0, 0, 1, 0], dtype=int64)

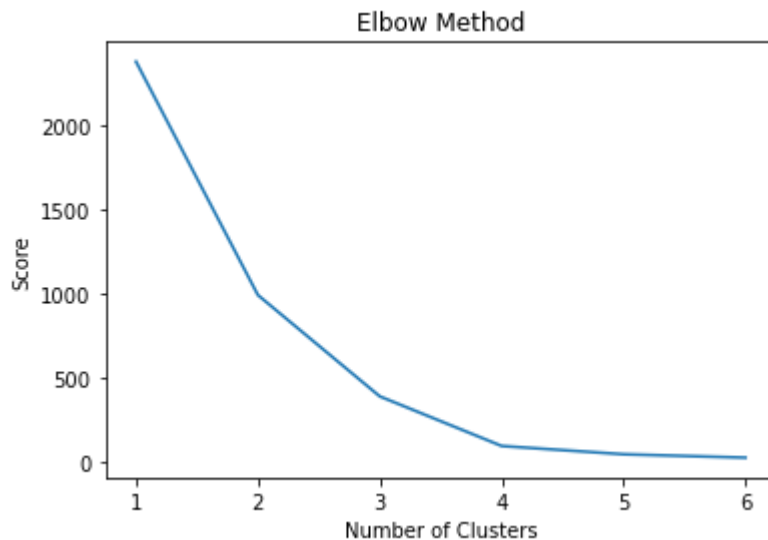
```
In [571]: 1 from sklearn.cluster import KMeans
2 n_clusters = 3
3 sklearn_pca = PCA(n_components = 2)
4 Y_sklearn = sklearn_pca.fit_transform(tf_idf_array)
5 kmeans = KMeans(n_clusters=n_clusters, max_iter=600, algorithm = 'auto')
6 %time fitted = kmeans.fit(Y_sklearn)
7 prediction = kmeans.predict(Y_sklearn)
8
9 plt.scatter(Y_sklearn[:, 0], Y_sklearn[:, 1],c=prediction ,s=50, cmap='virid
10
11 centers2 = fitted.cluster_centers_
12 plt.scatter(centers2[:, 0], centers2[:, 1],c='red', s=300, alpha=0.6);
```

Wall time: 80.8 ms



Choosing optimal clusters using Elbow method


```
In [572]: 1 number_clusters = range(1, 7)
2
3 kmeans = [KMeans(n_clusters=i, max_iter = 600) for i in number_clusters]
4 kmeans
5
6 score = [kmeans[i].fit(Y_skllearn).score(Y_skllearn) for i in range(len(kmeans)
7 score = [i*-1 for i in score]
8
9 plt.plot(number_clusters, score)
10 plt.xlabel('Number of Clusters')
11 plt.ylabel('Score')
12 plt.title('Elbow Method')
13 plt.show()
```



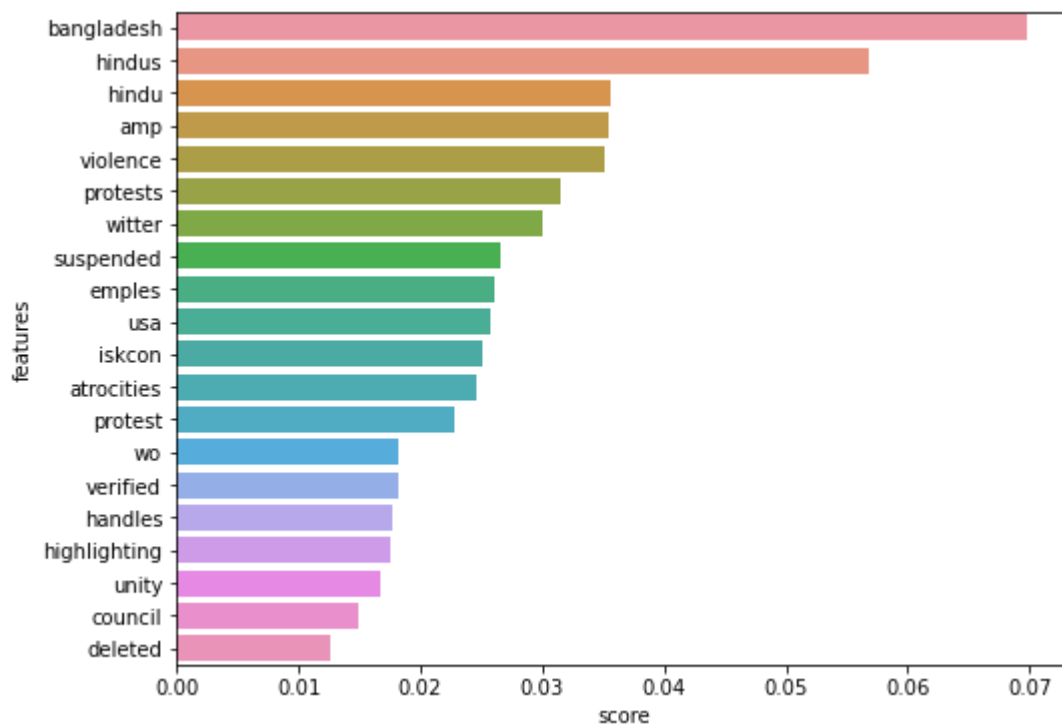
Extracting top features

```
In [573]: 1 def get_top_features_cluster(tf_idf_array, prediction, n_feats):
2     labels = np.unique(prediction)
3     dfs = []
4     for label in labels:
5         id_temp = np.where(prediction==label) # indices for each cluster
6         x_means = np.mean(tf_idf_array[id_temp], axis = 0) # returns average
7         sorted_means = np.argsort(x_means)[:n_feats] # indices with top
8         features = tf_idf_vectorizer.get_feature_names()
9         best_features = [(features[i], x_means[i]) for i in sorted_means]
10        df1 = pd.DataFrame(best_features, columns = ['features', 'score'])
11        dfs.append(df1)
12    return dfs
13    dfs = get_top_features_cluster(tf_idf_array, prediction, 20)
```

Cluster 1

```
In [574]: 1 import seaborn as sns
          2 plt.figure(figsize=(8,6))
          3 sns.barplot(x = 'score' , y = 'features', orient = 'h' , data = dfs[0][:20])
```

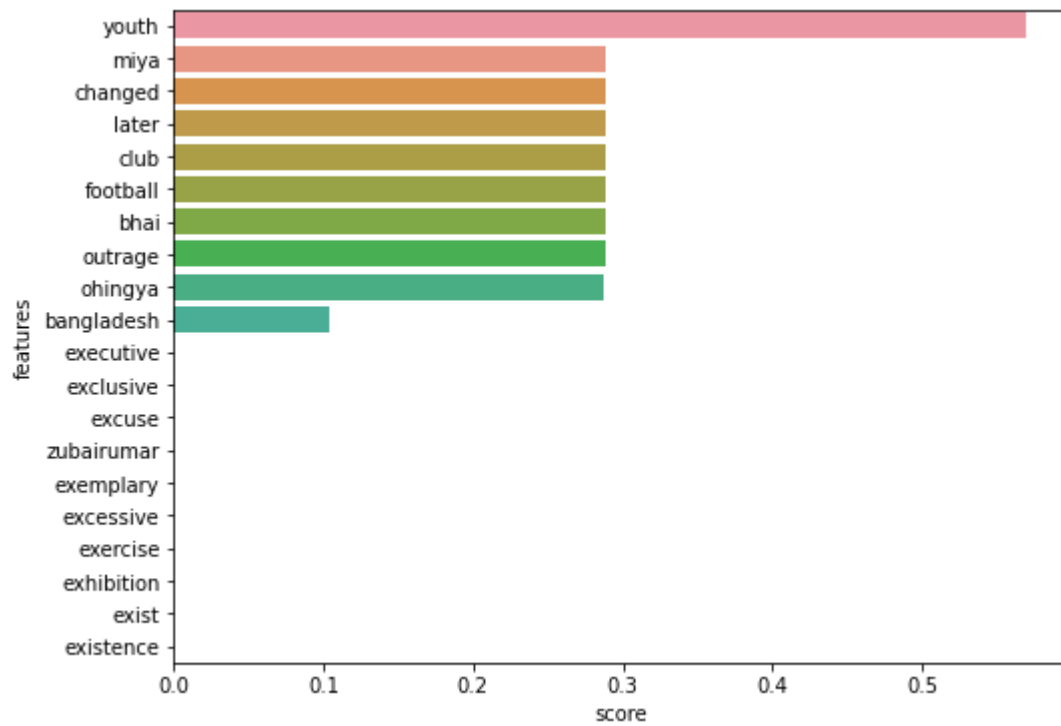
Out[574]: <AxesSubplot:xlabel='score', ylabel='features'>



Cluster 2

```
In [575]: 1 plt.figure(figsize=(8,6))  
          2 sns.barplot(x = 'score' , y = 'features', orient = 'h' , data = dfs[1][:20])
```

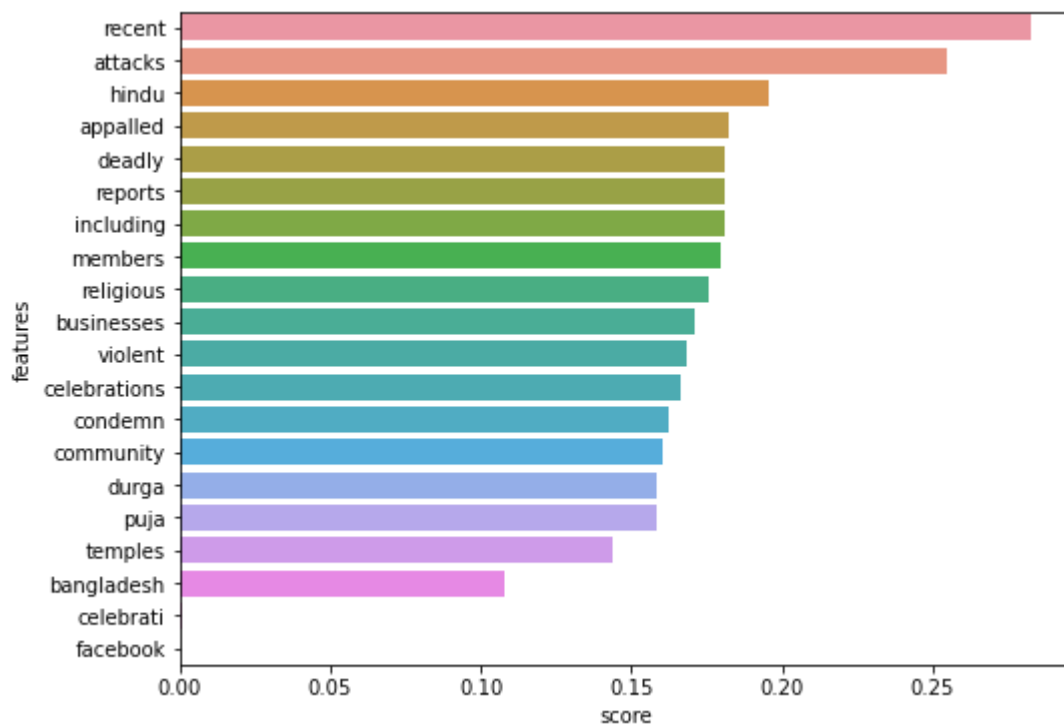
Out[575]: <AxesSubplot:xlabel='score', ylabel='features'>



Cluster 3

```
In [576]: 1 plt.figure(figsize=(8,6))
          2 sns.barplot(x = 'score' , y = 'features', orient = 'h' , data = dfs[2][:20])
```

Out[576]: <AxesSubplot:xlabel='score', ylabel='features'>



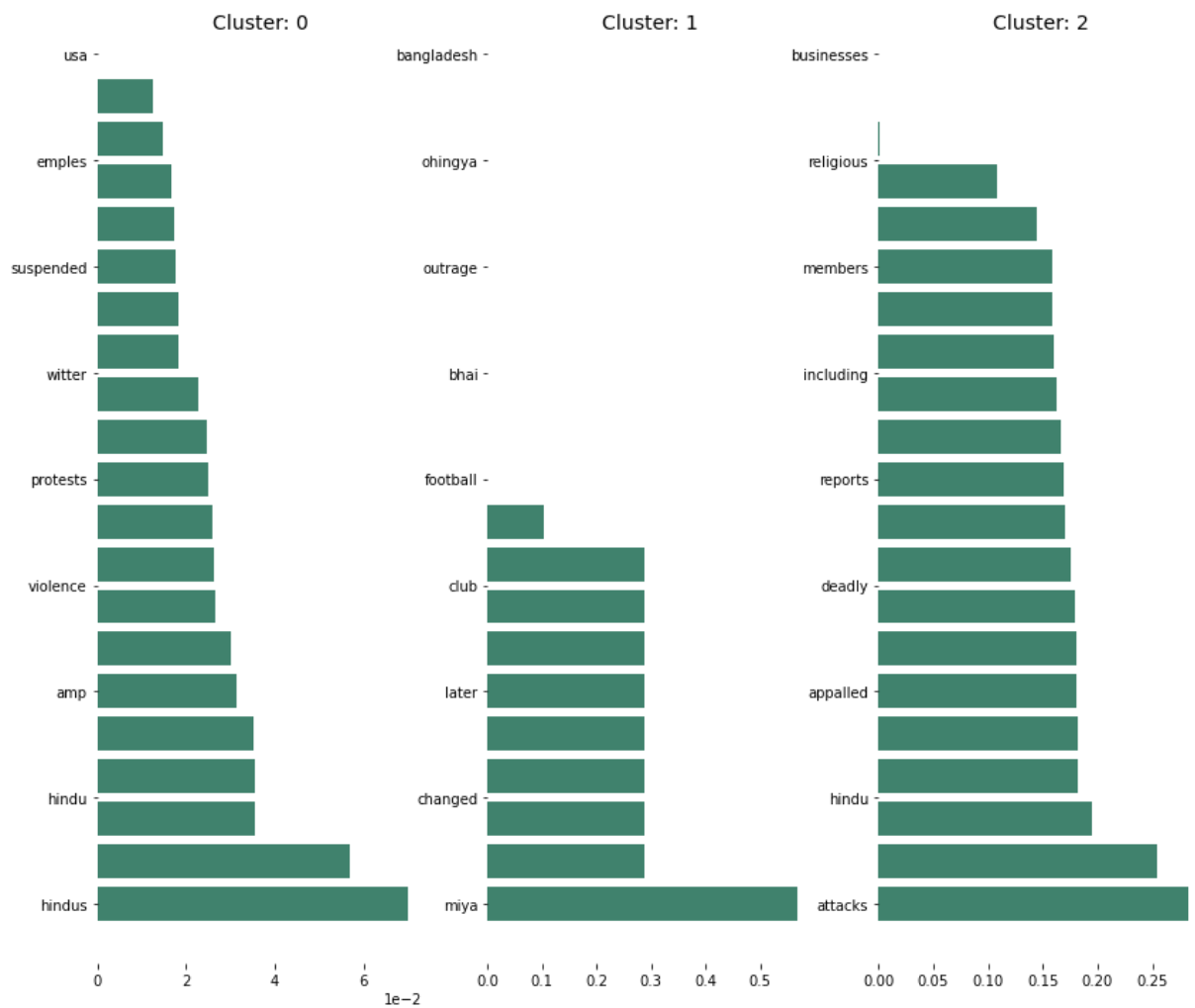
Cluster visualization

```
In [577]: 1 for i, df1 in enumerate(dfs):
          2     df1.to_csv('df_'+str(i)+'.csv')
```

```

In [578]: 1 def plot_features(dfs):
2     fig = plt.figure(figsize=(14,12))
3     x = np.arange(len(dfs[0]))
4     for i, df1 in enumerate(dfs):
5         ax = fig.add_subplot(1, len(dfs), i+1)
6         ax.set_title("Cluster: " + str(i), fontsize = 14)
7         ax.spines["top"].set_visible(False)
8         ax.spines["right"].set_visible(False)
9         ax.set_frame_on(False)
10        ax.get_xaxis().tick_bottom()
11        ax.get_yaxis().tick_left()
12        ax.ticklabel_format(axis='x', style='sci', scilimits=(-2,2))
13        ax.barh(x, df1.score, align='center', color='#40826d')
14        yticks = ax.set_yticklabels(df1.features)
15    plt.show();
16    plot_features(dfs)

```



Word Cloud

```
In [579]: 1 stopwords = set(STOPWORDS)
```

```
In [580]: 1 mask = np.array(Image.open("yoyo.png"))
```

```
In [581]: 1 mask
```

```
Out[581]: array([[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]],

                [[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]],

                [[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]],

                ...,

                [[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]],

                [[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]],

                [[255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255],
                 ...,
                 [255, 255, 255],
                 [255, 255, 255],
                 [255, 255, 255]]], dtype=uint8)
```

```
In [582]: 1 wordcloud = WordCloud(max_font_size=200,random_state = 42,
2 width = mask.shape[1],
3 height = mask.shape[0], background_color="white", max_words=2000000,
4 mask=mask,font_path = 'arial',colormap="Paired").generate(' '.join(df.text))
5 plt.figure( figsize=(100,100) )
6 plt.imshow(wordcloud, interpolation='bilinear')
7 plt.axis("off")
8 plt.show()
```


In [584]:

```
1
2 # Printing positive tweets
3 print("Printing positive tweets:\n")
4 j=1
5 sortedDF = df.sort_values (by=[ 'polarity']) #Sort the tweets
6 for i in range(0, sortedDF.shape[0] ):
7     if( sortedDF['Analysis'][i] == 'Positive'):
8         print(str(j) + ') ' + sortedDF['text'][i])
9     print()
10 j= j+1
```

1) gs Exactly a week ago, this picture triggered a wave of anti-hindu violence in Bangladesh (book planted by suspected Jamaat-e-Islami)

1) Protest going in many parts of the world. We are planning a one-day protest and prayer meetings for the victims in Bangladesh, on

1) com Sheikh Hasina is playing fast and loose on the genocide of Bangladeshis

```
In [585]: 1 # Printing negative tweets
2 print("Printing negative tweets:\n")
3 j=1
4 sortedDF = df.sort_values (by=[ 'polarity'],ascending=False) #Sort the tweet
5 for i in range(0, sortedDF.shape[0] ):
6     if( sortedDF['Analysis'][i] == 'Negative'):
7         print(str(j) + ' ' + sortedDF['text'][i])
8     print()
9     j= j+1
```

ity in Bangladesh. All, including members of religious

1) Isckon devote of Bangladesh has organized iftar party to Muslim brother and recently he was killed. Pray f

1) salam ights of minorities are inviolable. Sheikh Hasina citing Muslim p ersecutions in India as an excuse for the attack on Hindus

1) We are appalled by recent reports of deadly attacks on the Hindu commun ity in Bangladesh. All, including members of religious

1) Hindus are forced to migrate from Bangladesh and Kashmir. Understand th e importance of demography. Understand the ta

```
In [586]: 1 # Print the percentage of positive tweets
2 ptweets = df[df.Analysis == 'Positive']
3 ptweets = ptweets['text']
4 ptweets
5 round( (ptweets.shape[0] / df.shape[0]) * 100,1)
```

Out[586]: 19.9

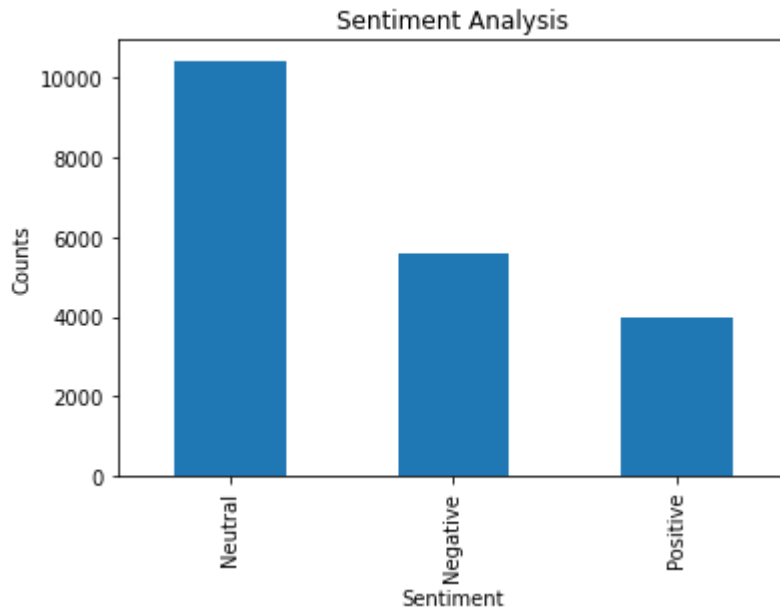
```
In [587]: 1 # Print the percentage of negative tweets
2 ntweets = df [df.Analysis == 'Negative']
3 ntweets = ntweets['text']
4 ntweets
5 round ((ntweets.shape[0] / df.shape[0]) * 100, 1)
```

Out[587]: 27.9

```
In [588]: 1 df['Analysis'].value_counts()
```

```
Out[588]: Neutral      10430
Negative      5587
Positive      3983
Name: Analysis, dtype: int64
```

```
In [589]: 1 # Plotting and visualizing the counts
2 plt.title('Sentiment Analysis')
3 plt.xlabel('Sentiment')
4 plt.ylabel('Counts')
5 df['Analysis'].value_counts().plot(kind = 'bar')
6 plt.show()
```



Topic Model

```
In [590]: 1 cv = CountVectorizer(max_df=0.95, min_df=2, stop_words = 'english')
2 dtm = cv.fit_transform(df.text)
3 dtm
```

Out[590]: <20000x4494 sparse matrix of type '<class 'numpy.int64'>' with 186894 stored elements in Compressed Sparse Row format>

```
In [591]: 1 LDA = LatentDirichletAllocation(n_components=7, random_state=42)
2 LDA.fit(dtm)
```

Out[591]: LatentDirichletAllocation(n_components=7, random_state=42)

```
In [592]: 1 len(cv.get_feature_names())
```

Out[592]: 4494

```
In [593]: 1 import random
```

```
In [594]: 1 for i in range(10):
          2     random_word_id = random.randint(0,4494)
          3     print(cv.get_feature_names()[random_word_id])
```

```
iss
cricfreak
manik
foreign
event
alam
acadmy
cases
shopee
maqsood
```

```
In [595]: 1 len(LDA.components_)
```

Out[595]: 7

```
In [596]: 1 LDA.components_
```

```
Out[596]: array([[ 0.14295631,  0.14297732,  0.14285716, ...,  0.14285718,
                  0.14285717,  0.14285716],
                [ 0.14285725,  0.14285728,  0.14285717, ...,  0.1428572 ,
                  0.14285718,  0.14285717],
                [ 0.14287411,  0.14285728,  0.14285717, ...,  0.1428572 ,
                  3.14273036,  0.14285717],
                ...,
                [ 0.14286319,  0.14296415,  0.14285719, ...,  0.14285724,
                  0.14298375,  0.1428572 ],
                [ 3.26392242,  0.14285727,  0.14285717, ...,  0.14294099,
                  0.14285718,  3.14285695],
                [46.02147272,  0.14290996,  0.14285716, ...,  0.14287758,
                  0.14285717,  0.14285717]])
```

```
In [597]: 1 len(LDA.components_[0])
```

Out[597]: 4494

```
In [598]: 1 single_topic = LDA.components_[0]
```

```
In [599]: 1 single_topic.argsort()
          2 single_topic[4094]
          3 single_topic[4456]
          4
          5 single_topic.argsort()[-10:]
          6 top_word_indices = single_topic.argsort()[-10:]
          7
          8
```

```
In [600]: 1 for index in top_word_indices:  
          2     print(cv.get_feature_names()[index])
```

```
miya  
changed  
later  
club  
football  
bhai  
outrage  
ohingya  
youth  
bangladesh
```

```
In [601]: 1 for index, topic in enumerate(LDA.components_):  
2         print(f'THE TOP 15 WORDS FOR TOPIC #{index}')  
3         print([cv.get_feature_names()[i] for i in topic.argsort()[-15:]])  
4         print('\n')
```

THE TOP 15 WORDS FOR TOPIC #0

['handles', 'wo', 'witter', 'suspended', 'hindus', 'miya', 'changed', 'later', 'club', 'football', 'bhai', 'outrage', 'ohingya', 'youth', 'bangladesh']

THE TOP 15 WORDS FOR TOPIC #1

['worker', 'bablu', 'salauddin', 'need', 'help', 'know', 'district', 'south', 'usa', 'emples', 'violence', 'amp', 'protests', 'hindus', 'bangladesh']

THE TOP 15 WORDS FOR TOPIC #2

['october', 'incidents', 'today', 'council', 'protests', 'iskcon', 'unity', 'atrocities', 'taken', 'protest', 'community', 'genocide', 'hindus', 'hindu', 'bangladesh']

THE TOP 15 WORDS FOR TOPIC #3

['victims', 'working', 'going', 'offer', 'ban', 'ussia', 'amp', 'devotees', 'prayers', 'violence', 'iskcon', 'hindus', 'india', 'protest', 'bangladesh']

THE TOP 15 WORDS FOR TOPIC #4

['hindus', 'simultaneous', 'asserted', 'australian', 'body', 'planned', 'com', 'attacks', 'https', 'temples', '17', 'amp', 'like', 'bangladesh', 'hindu']

THE TOP 15 WORDS FOR TOPIC #5

['unity', 'iskcon', 'witter', 'appalled', 'deadly', 'reports', 'including', 'members', 'community', 'hindus', 'religious', 'recent', 'attacks', 'hindu', 'bangladesh']

THE TOP 15 WORDS FOR TOPIC #6

['hindus', 'violence', 'iskcon', 'amp', 'businesses', 'recent', 'condemn', 'violent', 'celebrations', 'temples', 'attacks', 'durga', 'puja', 'hindu', 'bangladesh']

```
In [602]: 1 dtm.shape
```

Out[602]: (20000, 4494)

```
In [603]: 1 topic_results = LDA.transform(dtm)  
2 topic_results.shape
```

Out[603]: (20000, 7)

```
In [604]: 1 topic_results[0:10]
```

```
Out[604]: array([[0.01102874, 0.01099436, 0.01102568, 0.01100621, 0.01100623,
                  0.9339299 , 0.01100889],
                 [0.0238539 , 0.02399315, 0.02388642, 0.02401353, 0.02385474,
                  0.8565171 , 0.02388116],
                 [0.01190975, 0.01190838, 0.0119149 , 0.01190897, 0.01191327,
                  0.92852337, 0.01192136],
                 [0.01788423, 0.89270404, 0.01788151, 0.01788782, 0.01787331,
                  0.01788509, 0.017884 ],
                 [0.14285714, 0.14285714, 0.14285714, 0.14285714, 0.14285714,
                  0.14285714, 0.14285714],
                 [0.82809848, 0.02880152, 0.02859439, 0.02860095, 0.02860074,
                  0.02860636, 0.02869756],
                 [0.01190975, 0.01190838, 0.0119149 , 0.01190897, 0.01191327,
                  0.92852337, 0.01192136],
                 [0.92844933, 0.01190946, 0.01194992, 0.01191212, 0.01193029,
                  0.01192721, 0.01192166],
                 [0.9141727 , 0.01430631, 0.01431343, 0.01430167, 0.01429274,
                  0.01431449, 0.01429866],
                 [0.14285714, 0.14285714, 0.14285714, 0.14285714, 0.14285714,
                  0.14285714, 0.14285714]])
```

In [605]:

```
1 df.head()
```

Out[605]:

	text	created	id	statusSource	screenName	retweetCount	isRetwe
0	BIG At the request of the Bangladesh Govt. witter deletes Bangladesh Hindu Unity Council's twitter ha	2021-10-20 01:52:15	1450640978625236992	Twitter for iPhone	arvind291	292	Tr
1	Hundreds protest in Bangladesh over religious violence	2021-10-20 01:52:14	1450640973856317445	Twitter Web App	gojharan	280	Tr
2	We are appalled by recent reports of deadly attacks on the Hindu community in Bangladesh. All, including members of religious	2021-10-20 01:52:13	1450640971083907073	Twitter for Android	amodbhardwaj	898	Tr
3	Now, protests in USA against the violence on emples & Hindus in Bangladesh.	2021-10-20 01:52:12	1450640963878096901	Twitter for Android	VamsiKandula2	4260	Tr
4	in	2021-10-20 01:52:11	1450640962468737024	Twitter for Android	Karan_Hu_Mei	2012	Tr

In []:

```
1
```

In []:

```
1
```

In []:

```
1
```


In []:

1