## ⌄ Importing the Necessary Libraries

```
import re    #for regex
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt  #for various plots
import seaborn as sns
```

```
import nltk
import string
from nltk.corpus import stopwords

nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
True
```

```
!pip install bertopic
```

```
Requirement already satisfied: torch>=1.11.0 in /usr/local/lib/python3.10/dist-packages (from sentence-transformers>=0.4.1->bertopic) (2.3.0+cu121)
Requirement already satisfied: huggingface-hub>=0.15.1 in /usr/local/lib/python3.10/dist-packages (from sentence-transformers>=0.4.1->bertopic) (0.23.4)
Requirement already satisfied: Pillow in /usr/local/lib/python3.10/dist-packages (from sentence-transformers>=0.4.1->bertopic) (9.4.0)
Requirement already satisfied: numba>=0.51.2 in /usr/local/lib/python3.10/dist-packages (from umap-learn>=0.5.0->bertopic) (0.58.1)
Collecting pynndescent>=0.5 (from umap-learn>=0.5.0->bertopic)
  Downloading pynndescent-0.5.13-py3-none-any.whl (56 kB)
                                ──────────────── 56.9/56.9 kB 7.8 MB/s eta 0:00:00
Requirement already satisfied: filelock in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.15.1->sentence-transformers>=0.4.1->bertopic) (3.15.4
Requirement already satisfied: fsspec>=2023.5.0 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.15.1->sentence-transformers>=0.4.1->bertopic)
Requirement already satisfied: pyyaml>=5.1 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.15.1->sentence-transformers>=0.4.1->bertopic) (6.0
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.15.1->sentence-transformers>=0.4.1->bertopic) (2.31.0
Requirement already satisfied: typing-extensions>=3.7.4.3 in /usr/local/lib/python3.10/dist-packages (from huggingface-hub>=0.15.1->sentence-transformers>=0.4.1-
Requirement already satisfied: llvmlite<0.42,>=0.41.0dev0 in /usr/local/lib/python3.10/dist-packages (from numba>=0.51.2->umap-learn>=0.5.0->bertopic) (0.41.1)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-packages (from python-dateutil>=2.8.2->pandas>=1.1.5->bertopic) (1.16.0)
Requirement already satisfied: sympy in /usr/local/lib/python3.10/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (1.12.1)
Requirement already satisfied: networkx in /usr/local/lib/python3.10/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.3)
Requirement already satisfied: jinja2 in /usr/local/lib/python3.10/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (3.1.4)
Collecting nvidia-cuda-nvrtc-cu12==12.1.105 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cuda_nvrtc_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (23.7 MB)
Collecting nvidia-cuda-runtime-cu12==12.1.105 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cuda_runtime_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (823 kB)
Collecting nvidia-cuda-cupti-cu12==12.1.105 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cuda_cupti_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (14.1 MB)
Collecting nvidia-cudnn-cu12==8.9.2.26 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cudnn_cu12-8.9.2.26-py3-none-manylinux1_x86_64.whl (731.7 MB)
Collecting nvidia-cublas-cu12==12.1.3.1 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cublas_cu12-12.1.3.1-py3-none-manylinux1_x86_64.whl (410.6 MB)
Collecting nvidia-cufft-cu12==11.0.2.54 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cufft_cu12-11.0.2.54-py3-none-manylinux1_x86_64.whl (121.6 MB)
Collecting nvidia-curand-cu12==10.3.2.106 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_curand_cu12-10.3.2.106-py3-none-manylinux1_x86_64.whl (56.5 MB)
Collecting nvidia-cusolver-cu12==11.4.5.107 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cusolver_cu12-11.4.5.107-py3-none-manylinux1_x86_64.whl (124.2 MB)
Collecting nvidia-cusparse-cu12==12.1.0.106 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_cusparse_cu12-12.1.0.106-py3-none-manylinux1_x86_64.whl (196.0 MB)
Collecting nvidia-nccl-cu12==2.20.5 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_nccl_cu12-2.20.5-py3-none-manylinux2014_x86_64.whl (176.2 MB)
Collecting nvidia-nvtx-cu12==12.1.105 (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Using cached nvidia_nvtx_cu12-12.1.105-py3-none-manylinux1_x86_64.whl (99 kB)
Requirement already satisfied: triton==2.3.0 in /usr/local/lib/python3.10/dist-packages (from torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (2.3.0)
Collecting nvidia-nvjitlink-cu12 (from nvidia-cusolver-cu12==11.4.5.107->torch>=1.11.0->sentence-transformers>=0.4.1->bertopic)
  Downloading nvidia_nvjitlink_cu12-12.5.82-py3-none-manylinux2014_x86_64.whl (21.3 MB)
                                ──────────────── 21.3/21.3 MB 60.4 MB/s eta 0:00:00
Requirement already satisfied: regex!=2019.12.17 in /usr/local/lib/python3.10/dist-packages (from transformers<5.0.0,>=4.34.0->sentence-transformers>=0.4.1->bert
Requirement already satisfied: tokenizers<0.20,>=0.19 in /usr/local/lib/python3.10/dist-packages (from transformers<5.0.0,>=4.34.0->sentence-transformers>=0.4.1-
Requirement already satisfied: safetensors>=0.4.1 in /usr/local/lib/python3.10/dist-packages (from transformers<5.0.0,>=4.34.0->sentence-transformers>=0.4.1->ser
Requirement already satisfied: MarkupSafe>=2.0 in /usr/local/lib/python3.10/dist-packages (from jinja2->torch>=1.11.0->sentence-transformers>=0.4.1->bertopic) (2
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub>=0.15.1->sentence-transformers
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub>=0.15.1->sentence-transformers>=0.4.1->ber
Requirement already satisfied: urllib3<3,>=1.21.1 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub>=0.15.1->sentence-transformers>=0.4.
Requirement already satisfied: certifi>=2017.4.17 in /usr/local/lib/python3.10/dist-packages (from requests->huggingface-hub>=0.15.1->sentence-transformers>=0.4.
Requirement already satisfied: mpmath<1.4.0,>=1.1.0 in /usr/local/lib/python3.10/dist-packages (from sympy->torch>=1.11.0->sentence-transformers>=0.4.1->bertopic
Installing collected packages: nvidia-nvtx-cu12, nvidia-nvjitlink-cu12, nvidia-nccl-cu12, nvidia-curand-cu12, nvidia-cufft-cu12, nvidia-cuda-runtime-cu12, nvidia
  Attempting uninstall: cython
    Found existing installation: Cython 3.0.10
  Uninstalling Cython-3.0.10:
    Successfully uninstalled Cython-3.0.10
Successfully installed bertopic-0.16.2 cython-0.29.37 hdbscan-0.8.37 nvidia-cublas-cu12-12.1.3.1 nvidia-cuda-cupti-cu12-12.1.105 nvidia-cuda-nvrtc-cu12-12.1.105
```

```
from bertopic import BERTopic
from sklearn.feature_extraction import text
from sklearn.feature_extraction.text import ENGLISH_STOP_WORDS
from sklearn.feature_extraction.text import CountVectorizer, TfidfVectorizer
from bertopic.vectorizers import ClassTfidfTransformer
from sentence_transformers import SentenceTransformer
```

## ⌄ Loading the datasets

```
from google.colab import drive
drive.mount('/content/drive')
```

```
Mounted at /content/drive
```

```
test = pd.read_csv('/content/drive/MyDrive/Twitter Sentiment Analysis, EDA and Visualization/test.csv')
train = pd.read_csv('/content/drive/MyDrive/Twitter Sentiment Analysis, EDA and Visualization/train.csv')
ss = pd.read_csv('/content/drive/MyDrive/Twitter Sentiment Analysis, EDA and Visualization/sample_submission.csv')
```

## ⌄ Preprocessing the dataset

```
#combine the train and test dataset
df = [train, test]

df = pd.concat(df)

display(df.head(10))
```

| | textID | text | selected_text | sentiment |
|---|---|---|---|---|
| 0 | cb774db0d1 | I`d have responded, if I were going | I`d have responded, if I were going | neutral |
| 1 | 549e992a42 | Sooo SAD I will miss you here in San Diego!!! | Sooo SAD | negative |
| 2 | 088c60f138 | my boss is bullying me... | bullying me | negative |
| 3 | 9642c003ef | what interview! leave me alone | leave me alone | negative |
| 4 | 358bd9e861 | Sons of ****, why couldn`t they put them on t... | Sons of ****, | negative |
| 5 | 28b57f3990 | http://www.dothebouncy.com/smf - some shameles... | http://www.dothebouncy.com/smf - some shameles... | neutral |
| 6 | 6e0c6d75b1 | 2am feedings for the baby are fun when he is a... | fun | positive |
| 7 | 50e14c0bb8 | Soooo high | Soooo high | neutral |
| 8 | e050245fbd | Both of you | Both of you | neutral |
| 9 | fc2cbefa9d | Journey!? Wow... u just became cooler. hehe.... | Wow... u just became cooler. | positive |

```
print(df.shape)
```
```
(31015, 4)
```

```
df.isnull().sum()
```
```
textID             0
text               1
selected_text   3535
sentiment          0
dtype: int64
```

```
#drop the 1 row with null value in text column
df.dropna(inplace=True)
```

**Dataset Analysis :**

1. textID : A unique identifier
2. text : The actual tweet made by the user
3. selected_text : The stripped down tweet that can be fed to a Machine Learning model to make assumptions.
4. sentiment : The sentiment of the tweet that a Maching learning model tries to predict.

```
# Sort the sentiments by ascending order
df = df.sort_values(by='sentiment')

#drop the selected_text column
df = df.drop('selected_text', axis=1)

#The reason to drop the column is to get more words from the text column to verify whether the sentiment predicted is correct or not
```

```
df.head(10)
```

| | textID | text | sentiment |
|---|---|---|---|
| 15563 | 2543065d78 | Is there a way I can sleep for the next 8 or 9... | negative |
| 6044 | ee267131b1 | ok... twitter I almost pass out because of you... | negative |
| 21221 | 5b4cf5d1c6 | watching The Biggest Loser on Hallmark. Never ... | negative |
| 6041 | 856e0029b7 | Greg Pritchard should have got threw to the fi... | negative |
| 21223 | 5c83af1147 | Gourmet pizza = BLEH. Pizza is SUPPOSED to be... | negative |
| 21227 | 8581262345 | There isn`t any right now. They need to make... | negative |
| 10884 | a435e058ae | srry can`t go paintballing tonight and there... | negative |
| 16125 | 3cbcb82071 | LOL too bad he`s taken!!!!!!! | negative |
| 21230 | 7416c5eee3 | hypnotyst .... hmmmm... i should beware.. | negative |
| 21231 | a24c1d14d7 | http://twitpic.com/67nxe - Yeah..I`m bored XD ... | negative |

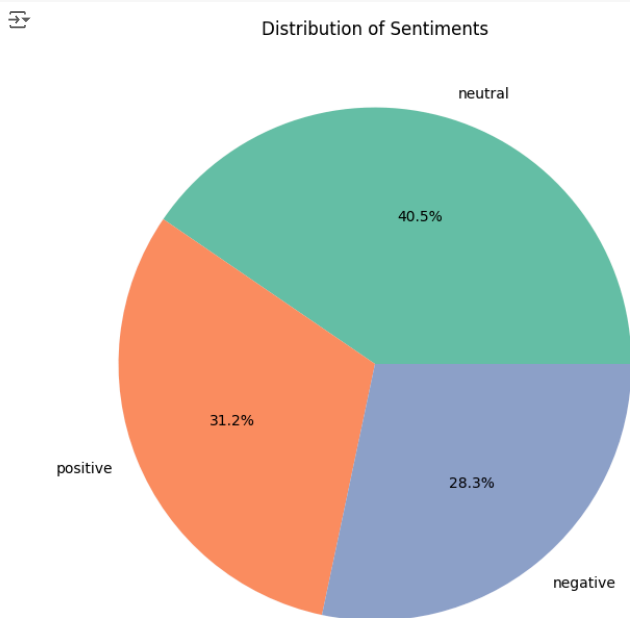## ⌄ Exploratory Data Analysis (EDA)

```
df.describe()
```

| | textID | text | sentiment |
|---|---|---|---|
| count | 27480 | 27480 | 27480 |
| unique | 27480 | 27429 | 3 |
| top | 2543065d78 | thanks | neutral |
| freq | 1 | 5 | 11117 |

```
check = df.groupby('sentiment').count()['text'].reset_index().sort_values(by='text',ascending=False)
check.style.background_gradient(cmap='Greens')
```

| | sentiment | text |
|---|---|---|
| 1 | neutral | 11117 |
| 2 | positive | 8582 |
| 0 | negative | 7781 |

```
# Count the occurrences of each sentiment
sentiment_counts = df['sentiment'].value_counts()

plt.figure(figsize=(8, 8))
plt.pie(sentiment_counts, labels=sentiment_counts.index, colors=sns.color_palette('Set2'), autopct='%1.1f%%')
plt.title('Distribution of Sentiments')
plt.show()
```

Distribution of Sentiments



Now let's go deeper into the sentiments, every sentiment has distinguised emotions.

Such as : anger, fear, anticipation, trust, surprise, sadness, joy, and disgust

## ⌄ Cleaning the text

```
#Make all the characters lower case
df['text'] = df['text'].str.lower()
```

```
#Remove multple spaces
df['text'] = df['text'].map(lambda x: re.sub("\s{2,6}", " ", x))
```

```
df['text']
```

```
15563    is there a way i can sleep for the next 8 or 9...
6044     ok... twitter i almost pass out because of you...
21221    watching the biggest loser on hallmark. never ...
6041     greg pritchard should have got threw to the fi...
21223     gourmet pizza = bleh. pizza is supposed to be...
                               ...
6253      he needs to go back to his scotty. that is wh...
20996    radio:active never gets old and never will thi...
10999                     is maxin and relaxin... ahhh
21020    just wanted to say that i <3 ur music(both th...
6759     ooohhh well you could always borrow and burn ...
Name: text, Length: 27480, dtype: object
```

```python
#Delete Url's in the post
df['text'] = df['text'].map(lambda x: re.sub('http[s]?:/\/\[^\s]*', ' ',x))
```

```python
df['text']
```

```
15563     is there a way i can sleep for the next 8 or 9...
 6044     ok... twitter i almost pass out because of you...
21221     watching the biggest loser on hallmark. never ...
 6041     greg pritchard should have got threw to the fi...
21223      gourmet pizza = bleh. pizza is supposed to be...
                              ...
 6253       he needs to go back to his scotty. that is wh...
20996     radio:active never gets old and never will thi...
10999                 is maxin and relaxin... ahhh
21020     just wanted to say that i <3 ur music(both th...
 6759      ooohhh well you could always borrow and burn ...
Name: text, Length: 27480, dtype: object
```

```python
df_text = df['text']
df_text
```

```
15563     is there a way i can sleep for the next 8 or 9...
 6044     ok... twitter i almost pass out because of you...
21221     watching the biggest loser on hallmark. never ...
 6041     greg pritchard should have got threw to the fi...
21223      gourmet pizza = bleh. pizza is supposed to be...
                              ...
 6253       he needs to go back to his scotty. that is wh...
20996     radio:active never gets old and never will thi...
10999                 is maxin and relaxin... ahhh
21020     just wanted to say that i <3 ur music(both th...
 6759      ooohhh well you could always borrow and burn ...
Name: text, Length: 27480, dtype: object
```

```python
stop_words = set(stopwords.words('english'))
punctuation = set(string.punctuation)

def remove_stopwords(text):
    words = text.split()
    filtered_words = [word for word in words if word.lower() not in stop_words and word not in punctuation]
    return ' '.join(filtered_words)

df['text'] = df_text.apply(remove_stopwords)
```

```python
default_stop_words = set(TfidfVectorizer(stop_words="english").get_stop_words())
all_stop_words = list(default_stop_words.union(stop_words))
```

```python
vectorizer_model =  TfidfVectorizer(stop_words= list(stop_words),
                                    ngram_range=(2,3), sublinear_tf=True)
sentence_model = SentenceTransformer("paraphrase-MiniLM-L6-v2")
topic_model = BERTopic(vectorizer_model = vectorizer_model)
topics, probs = topic_model.fit_transform(df['text'])
```

```
modules.json: 100%                              229/229 [00:00<00:00, 6.54kB/s]

config_sentence_transformers.json: 100%              122/122 [00:00<00:00, 3.64kB/s]

README.md: 100%                           3.73k/3.73k [00:00<00:00, 151kB/s]

sentence_bert_config.json: 100%                  53.0/53.0 [00:00<00:00, 2.64kB/s]

config.json: 100%                         629/629 [00:00<00:00, 20.2kB/s]

model.safetensors: 100%                      90.9M/90.9M [00:01<00:00, 62.0MB/s]

tokenizer_config.json: 100%                    314/314 [00:00<00:00, 16.6kB/s]

vocab.txt: 100%                       232k/232k [00:00<00:00, 1.34MB/s]

tokenizer.json: 100%                      466k/466k [00:00<00:00, 906kB/s]

special_tokens_map.json: 100%                  112/112 [00:00<00:00, 6.03kB/s]

1_Pooling/config.json: 100%                    190/190 [00:00<00:00, 7.02kB/s]

modules.json: 100%                         349/349 [00:00<00:00, 10.4kB/s]

config_sentence_transformers.json: 100%              116/116 [00:00<00:00, 3.14kB/s]

README.md: 100%                           10.7k/10.7k [00:00<00:00, 295kB/s]

sentence_bert_config.json: 100%                  53.0/53.0 [00:00<00:00, 1.66kB/s]

config.json: 100%                         612/612 [00:00<00:00, 34.4kB/s]

model.safetensors: 100%                      90.9M/90.9M [00:00<00:00, 195MB/s]

tokenizer_config.json: 100%                    350/350 [00:00<00:00, 12.3kB/s]

vocab.txt: 100%                       232k/232k [00:00<00:00, 1.36MB/s]

tokenizer.json: 100%                      466k/466k [00:00<00:00, 1.36MB/s]

special_tokens_map.json: 100%                  112/112 [00:00<00:00, 6.22kB/s]

1_Pooling/config.json: 100%                    190/190 [00:00<00:00, 8.64kB/s]

/usr/local/lib/python3.10/dist-packages/joblib/externals/loky/backend/fork_exec.py:38: RuntimeWarning: os.fork() was called. os.fork() is incompatible with multi
  pid = os.fork()
```

```
from transformers import pipeline

classifier = pipeline("text-classification", model = "j-hartmann/emotion-english-distilroberta-base", top_k = 8) #using top_k to get the top 8 sentiment score
sentiment = classifier('df_f')
sentiment
```

| | | |
|---|---|---|
| config.json: 100% | 1.00k/1.00k [00:00<00:00, 23.8kB/s] | |
| pytorch_model.bin: 100% | 329M/329M [00:06<00:00, 43.2MB/s] | |
| tokenizer_config.json: 100% | 294/294 [00:00<00:00, 16.3kB/s] | |
| vocab.json: 100% | 798k/798k [00:00<00:00, 1.15MB/s] | |
| merges.txt: 100% | 456k/456k [00:00<00:00, 872kB/s] | |
| tokenizer.json: 100% | 1.36M/1.36M [00:00<00:00, 1.94MB/s] | |
| special_tokens_map.json: 100% | 239/239 [00:00<00:00, 6.77kB/s] | |

```
[[{'label': 'neutral', 'score': 0.7798877358436584},
  {'label': 'surprise', 'score': 0.07403440773487091},
  {'label': 'sadness', 'score': 0.054175395518541336},
  {'label': 'anger', 'score': 0.04028287157416344},
  {'label': 'joy', 'score': 0.020953591912984848},
  {'label': 'disgust', 'score': 0.019315825775265694},
  {'label': 'fear', 'score': 0.011350187472999096}]]
```

```
topic_model.get_topic_info()
```

| | Topic | Count | Name | Representation | Representative_Docs |
|---|---|---|---|---|---|
| 0 | -1 | 10901 | -1_low low_im crying_low low low_getting ready... | [low low, im crying, low low low, getting read... | [wonders i`m 1 dat n church sumtime nt knowng ... |
| 1 | 0 | 930 | 0_day happy mothers_mothers day happy_mothers ... | [day happy mothers, mothers day happy, mothers... | [happy mothers` day moms!, happy mothers day m... |
| 2 | 1 | 926 | 1_still awake_bed night_sleep night_get sleep | [still awake, bed night, sleep night, get slee... | [ha think got like two hours sleep last night ... |
| 3 | 2 | 491 | 2_welcome twitter_im twitter_new twitter_saw t... | [welcome twitter, im twitter, new twitter, saw... | [certain tweets write spot messages person, gr... |
| 4 | 3 | 405 | 3_new album_lost voice_listening music_listeni... | [new album, lost voice, listening music, liste... | [saw fiddler topol! girls looooved it! next mo... |
| ... | ... | ... | ... | ... | ... |
| 320 | 319 | 10 | 319_nxt wari fun_us scorpian_naisee bad see_ni... | [nxt wari fun, us scorpian, naisee bad see, ni... | [naisee. bad see lens flares arond listening i... |
| 321 | 320 | 10 | 320_luck finals_good luck finals_luck finals e... | [luck finals, good luck finals, luck finals ev... | [could barely sleep last night, ugh...anyways ... |
| 322 | 321 | 10 | 321_jerrys loved place_9am ok cause_day open_c... | [jerrys loved place, 9am ok cause, day open, c... | [ahhhh! cant find anything way much open, haha... |
| 323 | 322 | 10 | 322_engineer making tracks_believed created mi... | [engineer making tracks, believed created mill... | [whole time ton things u would believe u would... |
| 324 | 323 | 10 | 323_grass evening fun_one south_mowing grass_l... | [grass evening fun, one south, mowing grass, l... | [sunny morning big k, lawns mow 2 mile run att... |

325 rows × 5 columns

```
topic_grams = []
num_topics = topic_model.get_topic_info().shape[0]
for k in range(num_topics):
    cur_top = topic_model.get_topic(k)
    if cur_top:
        cur_d = {'topic number': k}
        for j in range(10):
            cur_d[f'topic ngram {j+1}'] = cur_top[j][0]
        topic_grams.append(cur_d)
topics_df = pd.DataFrame(topic_grams)
```

```
topics_df
```

| | topic number | topic ngram 1 | topic ngram 2 | topic ngram 3 | topic ngram 4 | topic ngram 5 | topic ngram 6 | topic ngram 7 | topic ngram 8 | topic ngram 9 | topic ngram 10 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | day happy mothers | mothers day happy | mothers day mothers | day moms | mothers day moms | day mothers | moms happy | mom happy | moms day | moms happy mothers |
| 1 | 1 | still awake | bed night | sleep night | get sleep | cant sleep | sleep slept | night everyone | im going bed | time bed | sleep still |
| 2 | 2 | welcome twitter | im twitter | new twitter | saw tweet | twitter lol | twitter account | twitter im | got twitter | twitter love | tweet later |
| 3 | 3 | new album | lost voice | listening music | listening new | love music | im listening | new song | song listening | good song | fav song |
| 4 | 4 | stuck traffic | parking lot | traffic jam | new car | public transport | speeding ticket | car drive | car late | inspection sticker | tha bus |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 319 | 319 | nxt wari fun | us scorpian | naisee bad see | nite naisee bad | hi hopin | jass warn | nxt wari | warn b4 | avin gud | bad see lens |
| 320 | 320 | luck finals | good luck finals | luck finals everyone | second week june | im starting finals | bummer yo im | finals soonso guess | yo im starting | guess ill see | starting finals soonso |
| 321 | 321 | jerrys loved place | 9am ok cause | day open | cause one orange | ok cause one | today hmv | hmv opens | opens half | half ourbetter | ourbetter go |
| 322 | 322 | engineer making tracks | believed created million | created million people | howser got | idea captains | log truly | truly enlightening | enlightening pollard | pollard denial | denial truth |
| 323 | 323 | grass evening fun | one south | mowing grass | lawn getting | run shops | little mowing | kids back | definitely grass cutting | grass cutting cole | cutting cole committed |

324 rows × 11 columns

```
topic_model.get_topic(2)
```

```
[('welcome twitter', 0.002884449609169769),
 ('im twitter', 0.0026547569132210846),
 ('new twitter', 0.002579003386770751),
 ('saw tweet', 0.0025069856155680095),
 ('twitter lol', 0.0023235970202756958),
 ('twitter account', 0.0022786288049295067),
 ('twitter im', 0.0021816030966903483),
 ('got twitter', 0.0021513971716379536),
 ('twitter love', 0.0021049904562209306),
 ('tweet later', 0.0020979755897523964)]
```

```python
from transformers import pipeline
classifier = pipeline("zero-shot-classification", model="MoritzLaurer/DeBERTa-v3-base-mnli-fever-anli")
```

| | |
|---|---|
| config.json: 100% | 1.09k/1.09k [00:00<00:00, 34.4kB/s] |
| model.safetensors: 100% | 369M/369M [00:29<00:00, 22.6MB/s] |
| tokenizer_config.json: 100% | 1.28k/1.28k [00:00<00:00, 65.8kB/s] |
| spm.model: 100% | 2.46M/2.46M [00:01<00:00, 2.39MB/s] |
| tokenizer.json: 100% | 8.66M/8.66M [00:00<00:00, 22.8MB/s] |
| added_tokens.json: 100% | 23.0/23.0 [00:00<00:00, 1.53kB/s] |
| special_tokens_map.json: 100% | 286/286 [00:00<00:00, 14.7kB/s] |

```python
text_labels = ["anticipation", "anger", "fear", "joy", "trust", "surprise","sadness"]
```

```python
sample_text = " I don't know if I should be excited or worried right now. I mean, I'm thrilled about the possibilities and opportunities ahead, but at the same time,
```

```python
classifier(sample_text, text_labels, multi_label = False)
```

```
Asking to truncate to max_length but no maximum length is provided and the model has no predefined maximum length. Default to no truncation.
{'sequence': " I don't know if I should be excited or worried right now. I mean, I'm thrilled about the possibilities and opportunities ahead, but at the same
time, there's this gnawing sense of uncertainty and fear of the unknown. It's like I'm standing on the edge of a cliff, eager to leap into the adventure, yet
hesitant because I can't see what's below.",
 'labels': ['fear',
  'anticipation',
  'surprise',
  'joy',
  'sadness',
  'trust',
  'anger'],
 'scores': [0.5346909165382385,
  0.4214157462120056,
  0.03750550001859665,
  0.0020592096261680126,
  0.002021969761699438,
  0.001491756527684629,
  0.0008147756452672184]}
```

```python
def predict_sentiment(df, text_column, text_labels):
    if text_column not in df.columns:
        raise ValueError(f"The DataFrame does not contain the column: {text_column}")

    results = []
    for index, row in df.iterrows():
        sequence_to_classify = row[text_column]
        result = classifier(sequence_to_classify, text_labels, multi_label=False)
        results.append({
            text_column: sequence_to_classify,
            'sentiment': result['labels'][0],
            'score': result['scores'][0]
        })

    result_df = pd.DataFrame(results)
    result_df = df.merge(result_df, left_on=text_column, right_on='text', how='outer')
    return result_df
```

```python
results_df = predict_sentiment(df.head(10), text_column="text", text_labels=text_labels)
print(results_df.head(10))
```

```
      textID                                               text sentiment_x  \
0  2543065d78  way sleep next 8 9 days? way wake up, she`ll r...    negative
1  ee267131b1          ok... twitter almost pass you!! **** :`(    negative
2  5b4cf5d1c6  watching biggest loser hallmark. never fails m...    negative
3  856e0029b7  greg pritchard got threw final britains got ta...    negative
4  5c83af1147  gourmet pizza bleh. pizza supposed greasy filt...    negative
5  8581262345            isn`t right now. need make more. sorry.    negative
6  a435e058ae     srry can`t go paintballing tonight good movies    negative
7  3cbcb82071                     lol bad he`s taken!!!!!!!       negative
8  7416c5eee3                 hypnotyst .... hmmmm.. beware..      negative
9  a24c1d14d7  http://twitpic.com/67nxe yeah..i`m bored xd pi...    negative

    sentiment_y     score
0  anticipation  0.930641
1  anticipation  0.785657
2       sadness  0.872735
3  anticipation  0.300687
4       sadness  0.674485
5       sadness  0.812764
6          fear  0.448848
7       sadness  0.299860
8      surprise  0.383578
9       sadness  0.546874
```

```
output_file_path = '/content/drive/MyDrive/Twitter Sentiment Analysis, EDA and Visualization/results_df.csv'  # Ensure this path matches the directory of your
results_df.to_csv(output_file_path, index=False)
```

## Conclusion:

The Twitter Sentiment Analysis project using BERTopic provided valuable insights into the sentiments expressed in tweets. By leveraging **advanced NLP techniques and the BERTopic model (Zero-Shot-classification)**, we successfully classified tweets into distinct emotions, including neutral, surprise, sadness, anger, joy, disgust, and fear. The project demonstrated the effectiveness of using BERTopic for topic modeling and sentiment analysis.

Through thorough data preprocessing, we ensured the quality and accuracy of the analysis. The cleaning process involved removing URLs, stop words, and punctuation, which significantly improved the model's performance. The exploratory data analysis (EDA) revealed the distribution of sentiments across the dataset, highlighting the prevalence of neutral, positive, and negative sentiments.

The implementation of the BERTopic model allowed for the extraction of topics from the tweets, providing a deeper understanding of the underlying themes and emotions. The model's ability to cluster similar tweets and identify representative topics using **Ngram** was crucial in gaining insights into the public's sentiments on various issues.

The use of the **SentenceTransformer** for sentence embeddings and the **TfidfVectorizer** for vectorization played a pivotal role in enhancing the model's accuracy. These techniques enabled the model to capture the semantic meaning of the tweets, resulting in more accurate sentiment predictions.