

```
In [1]: 1 import pandas as pd
        2 dataset = pd.read_csv('hate_speech.csv')
        3 dataset.head()
```

Out[1]:

	id	label	tweet
0	1	0	@user when a father is dysfunctional and is s...
1	2	0	@user @user thanks for #lyft credit i can't us...
2	3	0	bihday your majesty
3	4	0	#model i love u take with u all the time in ...
4	5	0	factsguide: society now #motivation

```
In [2]: 1 dataset.label.value_counts()
```

Out[2]: 0 3000
1 2242
Name: label, dtype: int64

```
In [3]: 1 for index, tweet in enumerate(dataset["tweet"] [10:15]):
        2     print(index+1,"-", tweet)
```

1 - â #ireland consumer price index (mom) climbed from previous 0.2% to 0.5% in may #blog #silver #gold #forex
2 - we are so selfish. #orlando #standwithorlando #pulseshooting #orlandos hooting #biggerproblems #selfish #heabreaking #values #love #
3 - i get to see my daddy today!! #80days #gettingfed
4 - ouch...junior is angryð#got7 #junior #yugyoem #omg
5 - i am thankful for having a paner. #thankful #positive

```
In [5]: 1 import re
        2 def clean_text(text):
        3     text = re.sub(r'^a-zA-Z\'', '', text)
        4     text = re.sub(r'^\x00-\x7F+', '', text)
        5     text = text.lower()
        6     return text
```

```
In [6]: 1 dataset['clean_text'] = dataset.tweet.apply(lambda x: clean_text(x))
```

```
In [7]: 1 from nltk.corpus import stopwords
        2 len(stopwords.words('english'))
```

Out[7]: 179

```
In [8]: 1 import nltk
        2 nltk.download('stopwords')
```

[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\nihar\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!

Out[8]: True

```
In [21]: 1 def gen_freq(text):
2         word_list = []
3         for tw_words in text.split():
4             word_list.extend(tw_words)
5         word_freq = pd.Series(word_list).value_counts()
6         stop = stopwords.words('english')
7         word_freq = word_freq.drop(stop, errors='ignore')
8         return word_freq
```

```
In [22]: 1 def any_neg(words):
2         for word in words:
3             if word in ['n', 'no', 'non', 'not'] or re.search(r"\wn't", word):
4                 return 1
5             else:
6                 return 0
```

```
In [23]: 1 def any_rare(words, rare_100):
2         for word in words:
3             if word in rare_100:
4                 return 1
5             else:
6                 return 0
```

```
In [24]: 1 def is_question(words):
2         for word in words:
3             if word in ["when", "what", "how", "why", "who"]:
4                 return 1
5             else:
6                 return 0
```

```
In [25]: 1 word_freq=gen_freq(dataset.clean_text.str)
2 rare_100=word_freq[-100:]
3 dataset['word_count']=dataset.clean_text.str.split().apply(lambda x:len(x))
4 dataset['any_neg']=dataset.clean_text.str.split().apply(lambda x:any_neg(x))
5 dataset['is_question']=dataset.clean_text.str.split().apply(lambda x:is_question(x))
6 dataset['any_rare']=dataset.clean_text.str.split().apply(lambda x:any_rare(x))
7 dataset['char_count']=dataset.clean_text.apply(lambda x:len(x))
```

```
In [ ]: 1
```