

2211CS020524

V.Sai Charmika (SIGMA-ALML)

```
In [2]: import pandas as pd
import re
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import nltk
```

```
In [3]: nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')
```

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\charmi\AppData\Roaming\nltk_data...
[nltk_data] Unzipping tokenizers\punkt.zip.
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\charmi\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\charmi\AppData\Roaming\nltk_data...
```

Out[3]: True

```
In [15]: file_path = "Tweets.csv"
df = pd.read_csv(file_path)
```

```
In [16]: print("Original Dataset:")
print(df.head())
```

Original Dataset:

	tweet_id	airline_sentiment	airline_sentiment_confidence	\
0	570306133677760513	neutral	1.0000	
1	570301130888122368	positive	0.3486	
2	570301083672813571	neutral	0.6837	
3	570301031407624196	negative	1.0000	
4	570300817074462722	negative	1.0000	

	negativereason	negativereason_confidence	airline	\
0	NaN	NaN	Virgin America	
1	NaN	0.0000	Virgin America	
2	NaN	NaN	Virgin America	
3	Bad Flight	0.7033	Virgin America	
4	Can't Tell	1.0000	Virgin America	

	airline_sentiment_gold	name	negativereason_gold	retweet_count	\
0	NaN	cairdin	NaN	0	
1	NaN	jnardino	NaN	0	
2	NaN	yvonnalynn	NaN	0	
3	NaN	jnardino	NaN	0	
4	NaN	jnardino	NaN	0	

	text	tweet_coord	\
0	@VirginAmerica What @dhepburn said.	NaN	
1	@VirginAmerica plus you've added commercials t...	NaN	
2	@VirginAmerica I didn't today... Must mean I n...	NaN	
3	@VirginAmerica it's really aggressive to blast...	NaN	
4	@VirginAmerica and it's a really big bad thing...	NaN	

	tweet_created	tweet_location	user_timezone
0	2015-02-24 11:35:52 -0800	NaN	Eastern Time (US & Canada)
1	2015-02-24 11:15:59 -0800	NaN	Pacific Time (US & Canada)
2	2015-02-24 11:15:48 -0800	Lets Play	Central Time (US & Canada)
3	2015-02-24 11:15:36 -0800	NaN	Pacific Time (US & Canada)
4	2015-02-24 11:14:45 -0800	NaN	Pacific Time (US & Canada)

```
In [18]: df=df [['text', 'airline_sentiment']]
```

```
In [19]: print("\nSentiment Distribution:")
print(df ['airline_sentiment'].value_counts())
```

Sentiment Distribution:

```
negative    9178
neutral     3099
positive    2363
Name: airline_sentiment, dtype: int64
```

```
In [21]: stopwords_set = set(stopwords.words('english'))
         lemmatizer=WordNetLemmatizer()
```

```
In [28]: def preprocess_text(text):
         # Remove non-alphabetical characters
         text = re.sub(r' [^a-zA-Z]', '', text)
         #Convert text to lowercase
         text = text.lower()
         #Tokenize the text
         tokens = word_tokenize(text)
         #Remove stopwords and lemmatize tokens
         tokens = [lemmatizer.lemmatize (word) for word in tokens if word not in st
         return tokens
```

```
In [32]: def preprocess_text(text):
         # Example: simple preprocessing steps like lowercasing, removing punctuati
         text = text.lower() # Example step
         # Add other text processing steps here
         return text
```

```
In [33]: df['Processed_Text'] = df['text'].apply(preprocess_text)
```

```
In [34]: print("\nProcessed Data:")
         print(df [['text', 'Processed_Text', 'airline_sentiment']].head())
```

Processed Data:

	text \	Processed_Text	airline_sentiment
0	@VirginAmerica What @dhepburn said.	@virginamerica what @dhepburn said.	neutral
1	@VirginAmerica plus you've added commercials t...	@virginamerica plus you've added commercials t...	positive
2	@VirginAmerica I didn't today... Must mean I n...	@virginamerica i didn't today... must mean i n...	neutral
3	@VirginAmerica it's really aggressive to blast...	@virginamerica it's really aggressive to blast...	negative
4	@VirginAmerica and it's a really big bad thing...	@virginamerica and it's a really big bad thing...	negative

```
In [ ]:
```