

```
In [1]: import nltk

In [2]: nltk.download('stopwords')

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

True

Out[2]:

In [4]: import numpy as np
import pandas as pd
import string
import nltk
from nltk.corpus import stopwords
from nltk import PorterStemmer
import string
import re
import matplotlib.pyplot as plt
import seaborn as sns

In [6]: df=pd.read_csv('C:\Users\pedur\OneDrive\Tweets.csv')
df

Out[6]:
   tweet_id  airline_sentiment  airline_sentiment_confidence  negativereason  negativereason_confidence  airline  airline_sentiment_gold  name  negativereason_g
0  57030613677760513          neutral                    1.0000             NaN              NaN              NaN              Virgin America      Na
1  570301130886122368          positive                    0.3486             NaN              0.0000              Virgin America      Na
2  570301083672813571          neutral                    0.6837             NaN              NaN              Virgin America      Na
3  570301031407624196          negative                    1.0000             Bad Flight              0.7033              Virgin America      Na
4  570300817074462722          negative                    1.0000             Can't Tell              1.0000              Virgin America      Na
...      ...
14635  569587686496825344          positive                    0.3487             NaN              0.0000             American      Na
14636  569587371693355008          negative                    1.0000             Customer Service Issue              1.0000             American      Na
14637  569587242672396336          neutral                    1.0000             NaN              NaN              American      Na
14638  569587186687634433          negative                    1.0000             Customer Service Issue              0.6659             American      Na
14639  569587140490866689          neutral                    0.6771             NaN              0.0000             American      Na

14640 rows x 15 columns

In [7]: df.head(10)

Out[7]:
   tweet_id  airline_sentiment  airline_sentiment_confidence  negativereason  negativereason_confidence  airline  airline_sentiment_gold  name  negativereason_g  retwe
0  57030613677760513          neutral                    1.0000             NaN              NaN              Virgin America      Na
1  570301130886122368          positive                    0.3486             NaN              0.0000              Virgin America      Na
2  570301083672813571          neutral                    0.6837             NaN              NaN              Virgin America      Na
3  570301031407624196          negative                    1.0000             Bad Flight              0.7033              Virgin America      Na
4  570300817074462722          negative                    1.0000             Can't Tell              1.0000              Virgin America      Na
5  570300707074181121          negative                    1.0000             Can't Tell              0.6842              Virgin America      Na
6  57030616901320704          positive                    0.6745             NaN              0.0000              Virgin America      Na
7  57030024855349120          neutral                    0.6340             NaN              NaN              Virgin America      Na
8  570299953286942721          positive                    0.6559             NaN              NaN              Virgin America      Na
9  570295459631263746          positive                    1.0000             NaN              NaN              Virgin America      Na

In [8]: df.columns

Out[8]: Index(['tweet_id', 'airline_sentiment', 'airline_sentiment_confidence', 'negativereason', 'negativereason_confidence', 'airline', 'airline_sentiment_gold', 'name', 'negativereason_gold', 'retweet_count', 'text', 'tweet_coord', 'tweet_created', 'tweet_location', 'user_timezone'],
      dtype='object')

In [9]: df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14640 entries, 0 to 14639
Data columns (total 15 columns):
#   Column                Non-Null Count  Dtype
---  ---
0  tweet_id              14640 non-null  int64
1  airline_sentiment     14640 non-null  object
2  airline_sentiment_confidence  14640 non-null  float64
3  negativereason        9178 non-null  object
4  negativereason_confidence  15522 non-null  float64
5  airline               14640 non-null  object
6  airline_sentiment_gold  48 non-null     object
7  name                 14640 non-null  object
8  negativereason_gold    32 non-null     object
9  retweet_count         14640 non-null  int64
10 text                14640 non-null  object
11 tweet_coord          1819 non-null   object
12 tweet_created         14640 non-null  object
13 tweet_location       9907 non-null   object
14 user_timezone        9820 non-null   object
dtypes: float64(2), int64(2), object(11)
memory usage: 1.7+ MB

In [10]: df.dtypes

Out[10]:
tweet_id              int64
airline_sentiment     object
airline_sentiment_confidence  float64
negativereason        object
negativereason_confidence  float64
airline               object
airline_sentiment_gold  object
name                 object
negativereason_gold    object
retweet_count         int64
text                 object
tweet_coord           object
tweet_created          object
tweet_location         object
user_timezone          object
dtype: object

In [11]: df.shape

Out[11]: (14640, 15)

In [12]: df.isnull().sum()

Out[12]:
tweet_id              0
airline_sentiment     0
airline_sentiment_confidence  0
negativereason        5462
negativereason_confidence  4128
airline               0
airline_sentiment_gold  14608
name                 0
negativereason_gold    14608
retweet_count         0
text                 0
tweet_coord           13621
tweet_created          0
tweet_location         4733
user_timezone          4820
dtype: int64

In [13]: df['airline_sentiment_gold'].unique()

Out[13]: array([nan, 'negative', 'neutral', 'positive'], dtype=object)

In [14]: df['negativereason'].unique()

Out[14]: array([nan, 'Bad Flight', 'Can't Tell', 'Late Flight', 'Customer Service Issue', 'Flight Booking Problems', 'Lost Luggage', 'Flight Attendant Complaints', 'Cancelled Flight', 'Damaged Luggage', 'longlines'], dtype=object)

In [15]: df.drop('negativereason',axis=1,inplace=True)

In [16]: cols=['airline_sentiment_gold','tweet_coord','negativereason_confidence']
df.drop(cols,axis=1,inplace=True)

In [17]: df

Out[17]:
   tweet_id  airline_sentiment  airline_sentiment_confidence  airline  name  negativereason_gold  retweet_count  text  tweet_created  tweet_location  u
0  57030613677760513          neutral                    1.0000             Virgin America      cairdin              NaN              0  @VirginAmerica What @dhepburn said. 2015-02-24 11:35:52 -0800      Na
1  570301130886122368          positive                    0.3486             Virgin America      jwardino              NaN              0  @VirginAmerica plus you've added commercials L... 2015-02-24 11:15:59 -0800      Na
2  570301083672813571          neutral                    0.6837             Virgin America      yonnalynn              NaN              0  @VirginAmerica I didn't today... Must mean I n... 2015-02-24 11:15:48 -0800      Lets Play
3  570301031407624196          negative                    1.0000             Virgin America      jwardino              NaN              0  @VirginAmerica it's really aggressive to blast... 2015-02-24 11:35:36 -0800      Na
4  570300817074462722          negative                    1.0000             Virgin America      jwardino              NaN              0  @VirginAmerica and it's a really big bad thing... 2015-02-24 11:44:45 -0800      Na
...      ...
14635  569587686496825344          positive                    0.3487             American      KristenReenders              NaN              0  @AmericanAir thank you we got on a different l... 2015-02-22 12:01:01 -0800      Na
14636  569587371693355008          negative                    1.0000             American      itsropes              NaN              0  @AmericanAir leaving over 20 minutes Late Flig... 2015-02-22 11:59:46 -0800      Texas
14637  569587242672396336          neutral                    1.0000             American      sanyabun              NaN              0  @AmericanAir Please bring American Airlines to... 2015-02-22 11:59:15 -0800      Nigeria,Iagos
14638  569587186687634433          negative                    1.0000             American      SraJackson              NaN              0  @AmericanAir you have my money, you change my ... 2015-02-22 11:59:02 -0800      New Jersey
14639  569587140490866689          neutral                    0.6771             American      davidtdwu              NaN              0  @AmericanAir we have 8 ppl so we need 2 know h... 2015-02-22 11:58:51 -0800      dallas, TX

14640 rows x 11 columns

In [18]: df.isnull().sum()

Out[18]:
tweet_id              0
airline_sentiment     0
airline_sentiment_confidence  0
airline               0
name                 0
negativereason_gold    14608
retweet_count         0
text                 0
tweet_created          0
tweet_location         4733
user_timezone          4820
dtype: int64

In [19]: df.drop('negativereason_gold',axis=1,inplace=True)

In [20]: cols=['tweet_location','user_timezone']
df.drop(cols,axis=1,inplace=True)

In [21]: df.isnull().sum()

Out[21]:
tweet_id              0
airline_sentiment     0
airline_sentiment_confidence  0
airline               0
name                 0
retweet_count         0
text                 0
tweet_created          0
dtype: int64

In [22]: display(df.shape)
display(df.info())

(14640, 8)
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14640 entries, 0 to 14639
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  ---
0  tweet_id              14640 non-null  int64
1  airline_sentiment     14640 non-null  object
2  airline_sentiment_confidence  14640 non-null  float64
3  airline               14640 non-null  object
4  name                 14640 non-null  object
5  retweet_count         14640 non-null  int64
6  text                 14640 non-null  object
7  tweet_created         14640 non-null  object
dtypes: float64(1), int64(2), object(5)
memory usage: 915.1+ KB
None

In [23]: df = df[['airline_sentiment','text']]
df

Out[23]:
   airline_sentiment  text
0      neutral      @VirginAmerica What @dhepburn said.
1      positive      @VirginAmerica plus you've added commercials L...
2      neutral      @VirginAmerica I didn't today... Must mean I n...
3      negative      @VirginAmerica it's really aggressive to blast...
4      negative      @VirginAmerica and it's a really big bad thing...
...      ...
14635  positive      @AmericanAir thank you we got on a different L...
14636  negative      @AmericanAir leaving over 20 minutes Late Flig...
14637      neutral      @AmericanAir Please bring American Airlines to...
14638  negative      @AmericanAir you have my money, you change my ...
14639      neutral      @AmericanAir we have 8 ppl so we need 2 know h...

14640 rows x 2 columns

In [24]: sns.countplot(data=df,x='airline_sentiment')
plt.title('Graph-1-Airline Sentiment Distribution')

Out[24]:
Text(0.5, 1.0, 'Graph-1-Airline Sentiment Distribution')

Graph-1-Airline Sentiment Distribution
Count
8000
6000
4000
2000
0
neutral positive negative
airline_sentiment

In [25]: df['count_word'] = df['text'].apply(lambda x : len(x.split(' ')))
sns.histplot(data = df , x='count_word',kde=True)
plt.title('Graph-2-Number de Word Distribution without any Cleaning Task')
plt.show()

C:\Users\pedur\AppData\Local\Temp\ipykernel_15784\4846151189.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df['count_word'] = df['text'].apply(lambda x : len(x.split(' ')))

Graph-2-Number de Word Distribution without any Cleaning Task
Count
1000
800
600
400
200
0
5 10 15 20 25 30 35
count_word

In [26]: sns.histplot(data = df , x='count_word',hue='airline_sentiment',alpha=0.6,kde=True)
plt.title('Graph-3-Number de Word Distribution without any Cleaning Task')
plt.show()

Graph-3-Number de Word Distribution without any Cleaning Task
Count
800
700
600
500
400
300
200
100
0
5 10 15 20 25 30 35
count_word
airline_sentiment
neutral positive negative

In [27]: sns.boxplot(data = df , y='count_word',x='airline_sentiment')
plt.title('Graph-3(1)=Boxplot Number of word Across Tweets Categories')
plt.show()

Graph-3(1)=Boxplot Number of Word Across Tweets Categories
Count
35
30
25
20
15
10
5
0
neutral positive negative
airline_sentiment

In [28]: df.loc[np.logical_or(df['count_word']>35,df['count_word']<=5),:]

Out[28]:
   airline_sentiment  text  count_word
0      neutral      @VirginAmerica What @dhepburn said. 4
14      positive      @VirginAmerica Thanks! 2
18      positive      I ♥️ Flying @VirginAmerica. 🍷 5
46      neutral      @VirginAmerica DREAM http://t.co/A2gRIaQ2H... 5
58      neutral      @VirginAmerica @ladygaga @cameunderwood - Ca... 5
...      ...
14312  positive      @AmericanAir awesome! Thx 3
14314  negative      @AmericanAir yes, and rebooked incorrectly. 5
14443      neutral      @AmericanAir hi how are you 5
14600      neutral      http://t.co/EIw2yYbFu roberts&amp;s=1 @Americ... 3
14630      positive      @AmericanAir Thanks! He is. 4

817 rows x 3 columns

In [34]: df.loc[np.logical_or(df['count_word']>35,df['count_word']<=5),:]

Out[34]:
   airline_sentiment  text  count_word
0      neutral      @VirginAmerica What @dhepburn said. 4
14      positive      I ♥️ Flying @VirginAmerica. 🍷 5
46      neutral      @VirginAmerica DREAM http://t.co/A2gRIaQ2H... 5
58      neutral      @VirginAmerica @ladygaga @cameunderwood - Ca... 5
...      ...
14312  positive      @AmericanAir awesome! Thx 3
14314  negative      @AmericanAir yes, and rebooked incorrectly. 5
14443      neutral      @AmericanAir hi how are you 5
14600      neutral      http://t.co/EIw2yYbFu roberts&amp;s=1 @Americ... 3
14630      positive      @AmericanAir Thanks! He is. 4

817 rows x 3 columns

In [36]: import re
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from nltk.tokenize import word_tokenize

In [37]: # punctuation Removal
def remove_punctuation(text):
    return re.sub(r'[\W\s]','',text)

#stopword removal
def remove_stopwords(text):
    stop_words = set(stopwords.words('english'))
    tokens = word_tokenize(text)
    filter_tokens = [word for word in tokens if word.lower() not in stop_words]
    return " ".join(filter_tokens)

#remove numeric
def remove_numeric(text):
    return re.sub(r'\d+','',text)

#stemming
def apply_stemming(text):
    stemmer = PorterStemmer()
    tokens = word_tokenize(text)
    stemmed_tokens = [stemmer.stem(word) for word in tokens]
    return " ".join(stemmed_tokens)

def remove_mentions(text):
    return re.sub(r'@w+', '',text)

In [38]: import nltk
nltk.download('punkt')
from nltk.tokenize import word_tokenize
from nltk.stem import PorterStemmer
def apply_stemming(text):
    stemmer = PorterStemmer()
    tokens = word_tokenize(text)
    stemmed_tokens = [stemmer.stem(word) for word in tokens]
    return " ".join(stemmed_tokens)

input_text = "walking throw the street, a passenger walked toward me, talking,about a walked chicken on the streets"
stemmed_text = apply_stemming(input_text)
print(stemmed_text)

walk throw the street , a passeng walk toward me , talking,about a walk chicken on the street

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

In [39]: apply_stemming('walking throw the street , a passenger walked toward me talking about a walked chicken on the streets')

In [40]: 'walk throw the street , a passeng walk toward me talk about a walk chicken on the street'

In [41]: def text_preprocessing(text):
    sentence = remove_mentions(text)
    sentence = remove_punctuation(sentence)
    sentence = remove_stopwords(sentence)
    sentence = remove_numeric(sentence)
    sentence = apply_stemming(sentence)
    return sentence

In [41]: text_preprocessing('walking throw the street , a passenger walked toward me,talking about a walked chicken on the streets')

Out[41]: 'walk throw the street passeng walk toward mealk walk chicken street'

In [42]: df.loc['new_text'] = df['text'].apply(lambda x : text_preprocessing(x))

C:\Users\pedur\AppData\Local\Temp\ipykernel_15784\3410990210.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df.loc['new_text'] = df['text'].apply(lambda x : text_preprocessing(x))

In [43]: df.loc['new_count_word'] = df['new_text'].apply(lambda x : len(x.split(' ')))
sns.histplot(data = df , x='new_count_word',kde=True)
plt.title('Graph-5-Number of word Distribution after cleaning Task')
plt.show()

C:\Users\pedur\AppData\Local\Temp\ipykernel_15784\1052658891.py:1: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame.
Try using .loc[row_indexer,col_indexer] = value instead

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
df.loc['new_count_word'] = df['new_text'].apply(lambda x : len(x.split(' ')))

Graph-5-Number of Word Distribution after Cleaning Task
Count
1400
1200
1000
800
600
400
200
0
0.0 2.5 5.0 7.5 10.0 12.5 15.0 17.5 20.0
new_count_word

In [44]: sns.countplot(data=df,x='airline_sentiment')
plt.title('Graph-1(a)-Airline Sentiment Distribution-after cleaning the data')
plt.show()

Graph-1(a)-Airline Sentiment Distribution-after cleaning the data
Count
8000
6000
4000
2000
0
neutral positive negative
airline_sentiment

In [ ]:

In [ ]:
```