

Sentiment Analysis using US-Airlines dataset

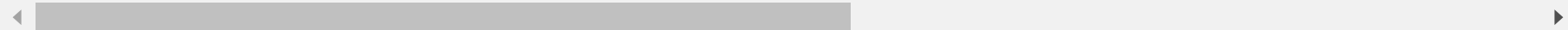
```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import string,re
        4 import nltk
        5 #nltk.download('all')
        6 import matplotlib.pyplot as plt
        7 import seaborn as sns
        8 %matplotlib inline
```

Downloading Dataset

```
In [2]: 1 data_source_url = "https://raw.githubusercontent.com/satyajeetkrjha/kaggle-Twitter-US-Airline-Sentiment-/master/T
2 airline_tweets = pd.read_csv(data_source_url)
3 airline_tweets.head()
```

Out[2]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN



```
In [3]: 1 airline_tweets.shape
```

Out[3]: (14640, 15)

Converting target labels into numerical form for ease

```
In [3]: 1 airline_tweets['sentiment_class'] =
2         airline_tweets['airline_sentiment'].map({'negative' : 0,
3                                                  'neutral' : 3,
4                                                  'positive' : 4})
5         airline_tweets.head()
```

Out[3]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN

```
In [4]: 1 airline_tweets.shape
```

Out[4]: (14640, 16)

In [5]: 1 airline_tweets.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14640 entries, 0 to 14639
Data columns (total 16 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            10522 non-null  float64
5   airline                              14640 non-null  object
6   airline_sentiment_gold               40 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                          1019 non-null   object
12  tweet_created                        14640 non-null  object
13  tweet_location                       9907 non-null   object
14  user_timezone                        9820 non-null   object
15  sentiment_class                      14640 non-null  int64
dtypes: float64(2), int64(3), object(11)
memory usage: 1.8+ MB
```

```
In [6]: 1 airline_tweets.isnull().sum()
```

```
Out[6]: tweet_id                0
airline_sentiment              0
airline_sentiment_confidence   0
negativereason                5462
negativereason_confidence     4118
airline                       0
airline_sentiment_gold        14600
name                          0
negativereason_gold          14608
retweet_count                  0
text                           0
tweet_coord                   13621
tweet_created                  0
tweet_location                4733
user_timezone                 4820
sentiment_class                0
dtype: int64
```

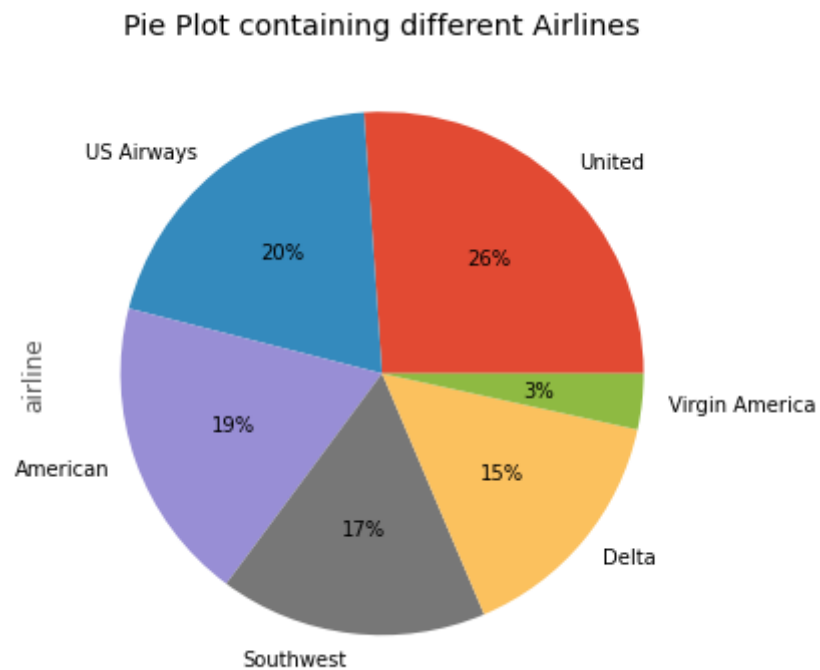
Pre-defining fix Fig size

```
In [7]: 1 plt.rcParams["font.family"] = 'DejaVu Sans'
2 plt.style.use('ggplot')
3 plot_size = plt.rcParams["figure.figsize"]
4 # print(plot_size[0])
5 # print(plot_size[1])
6
7 plot_size[0] = 8
8 plot_size[1] = 6
9 plt.rcParams["figure.figsize"] = plot_size
```

EDA part

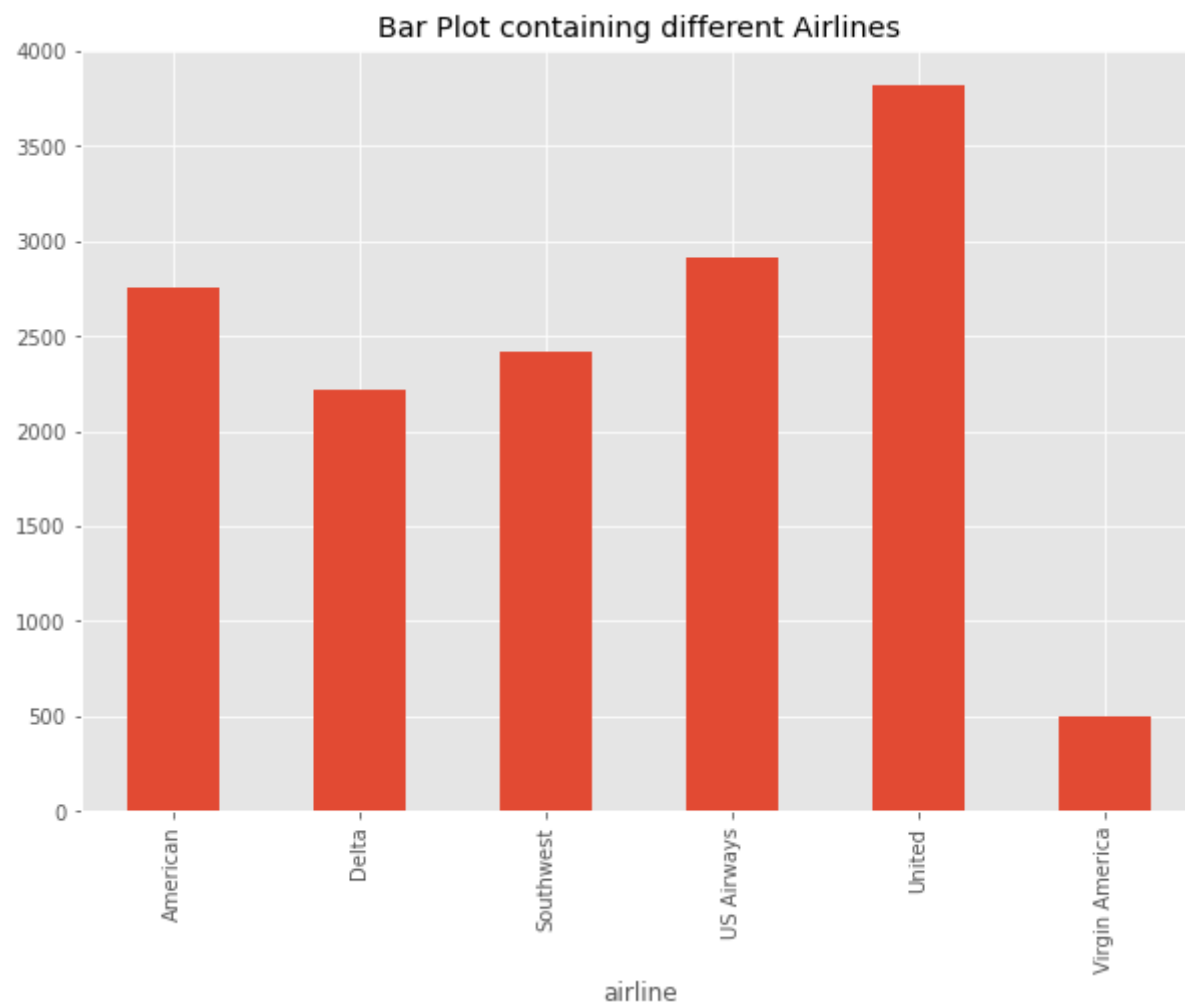
```
In [8]: 1 airline_tweets.airline.value_counts().plot(kind='pie',  
2         title='Pie Plot containing different Airlines',  
3         autopct='%1.0f%%')
```

```
Out[8]: <AxesSubplot:title={'center': 'Pie Plot containing different Airlines'}, ylabel='airline'>
```



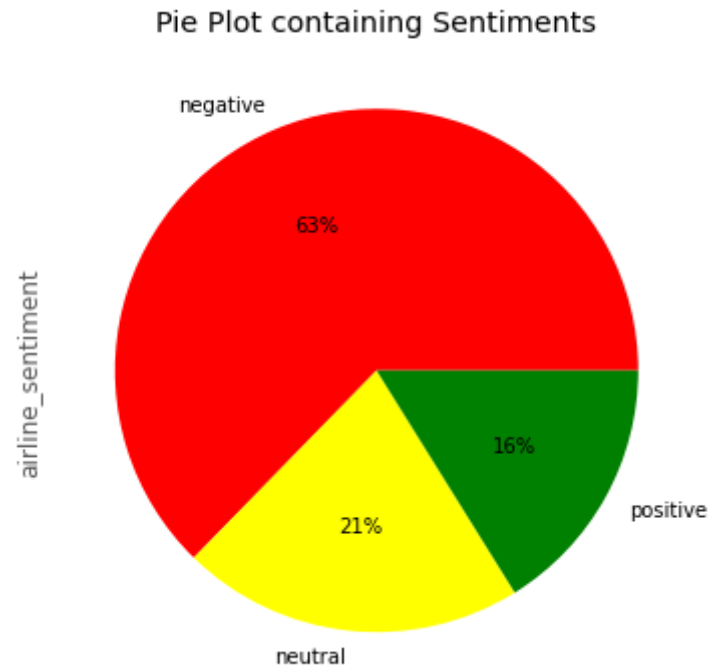
```
In [9]: 1 airline_tweets.groupby('airline')['airline'].count().plot(kind='bar',  
2 title='Bar Plot containing different Airlines',  
3 figsize=(10,7),grid=True)
```

```
Out[9]: <AxesSubplot:title={'center': 'Bar Plot containing different Airlines'}, xlabel='airline'>
```



```
In [10]: 1 airline_tweets.airline_sentiment.value_counts().plot(kind='pie',  
2         title='Pie Plot containing Sentiments',  
3         autopct='%1.0f%%',  
4         colors=["red", "yellow", "green"])
```

```
Out[10]: <AxesSubplot:title={'center':'Pie Plot containing Sentiments'}, ylabel='airline_sentiment'>
```



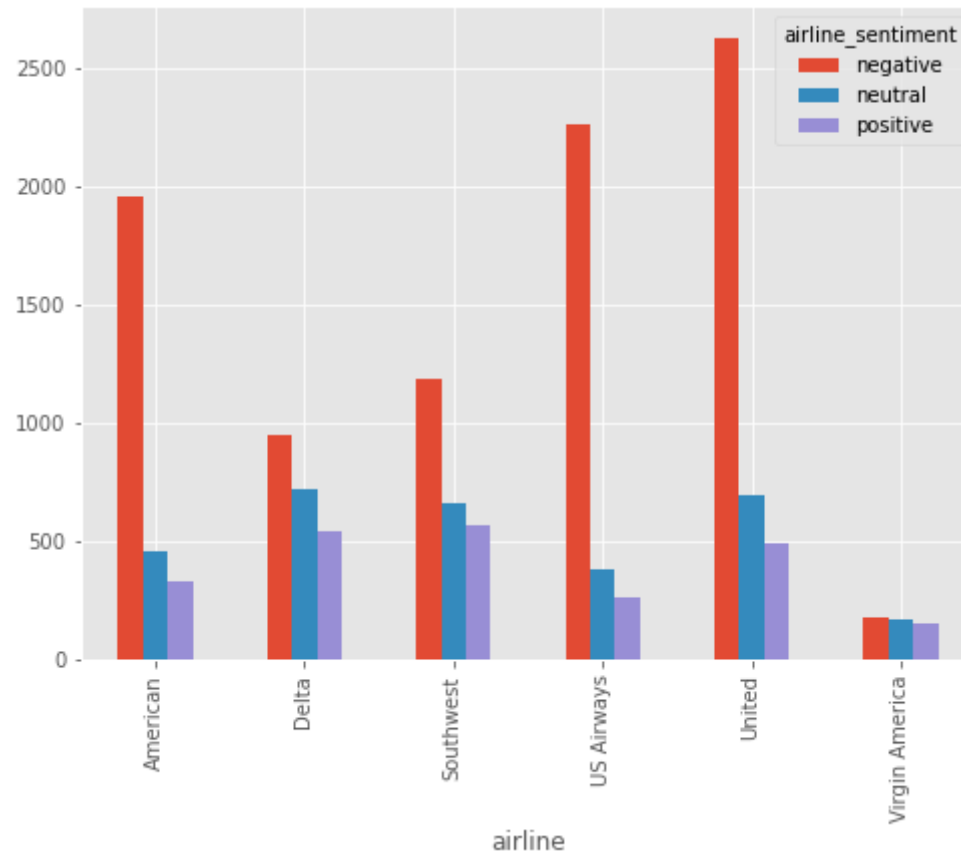

```
In [11]: 1 airline_tweets.groupby('airline_sentiment')['airline_sentiment'].count().plot(kind='bar',  
2 title='Target class',  
3 figsize=(10,7),  
4 grid=True)
```

```
Out[11]: <AxesSubplot:title={'center':'Target class'}, xlabel='airline_sentiment'>
```



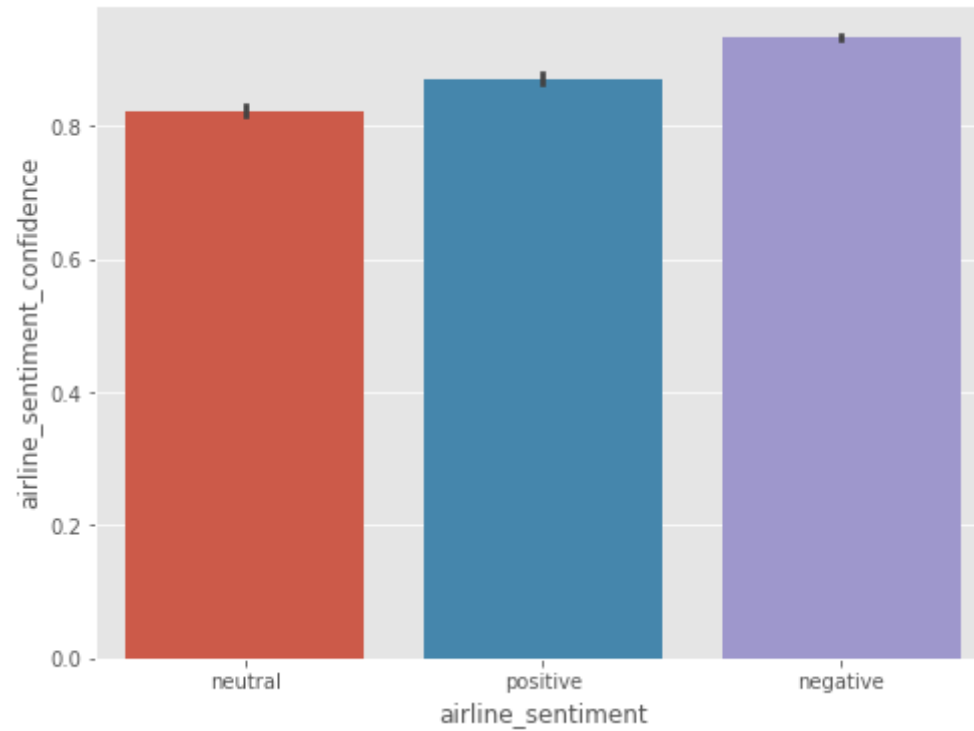
```
In [12]: 1 airline_sentiment = airline_tweets.groupby(['airline', 'airline_sentiment']).airline_sentiment.count().unstack()  
        2 airline_sentiment.plot(kind='bar')
```

Out[12]: <AxesSubplot:xlabel='airline'>



```
In [13]: 1 sns.barplot(x='airline_sentiment',  
2               y='airline_sentiment_confidence' ,  
3               data=airline_tweets)
```

```
Out[13]: <AxesSubplot:xlabel='airline_sentiment', ylabel='airline_sentiment_confidence'>
```



Data Pre-Processing

Removing Punctuation, Numbers, Special Characters and Short Words

```
In [14]: 1 def remove_pattern(text,pattern):
2
3         r = re.findall(pattern,text)
4
5         for i in r:
6             text = re.sub(i,"",text)
7
8         return text
```

```
In [15]: 1 airline_tweets['text'] = np.vectorize(remove_pattern)(airline_tweets['text'], "@[\w]*")
2
3         airline_tweets.head()
```

Out[15]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN

^[a-zA-Z] means any a-z or A-Z at the start of a line

[^a-zA-Z] means any character that IS NOT a-z OR A-Z

the first means "match all strings that start with a letter", the second means "match all strings that contain a non-letter". The caret ("^") is used in two different ways, one to signal the start of the text, one to negate a character match inside square brackets.

```
In [16]: 1 airline_tweets['text'] = airline_tweets['text'].str.replace("[^a-zA-Z#]", " ")
          2
          3 airline_tweets.head()
```

C:\Users\Eric\AppData\Roaming\Python\Python37\site-packages\ipykernel_launcher.py:1: FutureWarning: The default value of regex will change from True to False in a future version.
 """Entry point for launching an IPython kernel.

Out[16]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN



```
In [17]: 1 airline_tweets['text'] = airline_tweets['text'].str.lower()
2
3 airline_tweets.head()
```

Out[17]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN

```
In [18]: 1 airline_tweets['text'] = airline_tweets['text'].apply(lambda x: ' '.join([w for w in x.split() if len(w)>3]))
2
3 airline_tweets.head()
```

Out[18]:

	tweet_id	airline_sentiment	airline_sentiment_confidence	negativereason	negativereason_confidence	airline	airline_sentiment_gold
0	570306133677760513	neutral	1.0000	NaN	NaN	Virgin America	NaN
1	570301130888122368	positive	0.3486	NaN	0.0000	Virgin America	NaN
2	570301083672813571	neutral	0.6837	NaN	NaN	Virgin America	NaN
3	570301031407624196	negative	1.0000	Bad Flight	0.7033	Virgin America	NaN
4	570300817074462722	negative	1.0000	Can't Tell	1.0000	Virgin America	NaN

Data Visualisation using WordCloud

```
In [19]: 1 from wordcloud import WordCloud, ImageColorGenerator
2 from PIL import Image
3 import urllib
4 import requests
```

Generating WordCloud for tweets with postive labels only

```
In [20]: 1 all_positive_words = ' '.join(text for text in airline_tweets['text'][airline_tweets['airline_sentiment'] == 'pos']
          2
          3 all_positive_words
```

```
Out[20]: 'plus added commercials experience tacky nearly every time this worm away well didn amazing arrived hour early goo
d pretty graphics much better than minimal iconography this such great deal already thinking about trip haven even
gone trip flying your #fabulous #seductive skies again take #stress away from travel http ahlxhhiyn thanks excite
d first cross country flight heard nothing great things about virgin america daystogo flying know what would amazi
ngly awesome please want with only love this graphic http grrwaaa love hipster innovation feel good brand this gre
at news america could start flights hawaii year http moodlighting only best experience ever cool calming #moodlitm
onday done done best airline around hands down view downtown angeles hollywood sign beyond that rain mountains htt
p ibtr #elevategold good reason rock this just blew mind julie andrews though very impressive know need spotify st
at #guiltypleasures lady gaga amazing love three really beat classics congrats winning award best deals from airli
ne http iljaebv worried been great ride plane with great crew airlines should like this awesome flew yall morning
correct bill applied more then once member #inflight crew team interested #flightattendant #dreampath amazing cust
omer service again raeann best #customerservice #virginamerica #flying getaway deals through from lots cool cities
http tzzjhuibch #cheapflights #farecompare getaway deals through from lots cool cities http #cheapflights #farecom
pare have great week come back #phl already need take this horrible cold #pleasecomeback http glxfwp best airline
have flown easy change your reservation helpful representatives comfortable flying experience again another kicked
butt naelah represents your team beautifully thank your beautiful front design down right cool still book ticket y
our back secure love team running gate tonight waited delayed flight they kept things entertaining thanks your out
standing crew moved mountains home francisco tonight have absolute best team customer service ever every time
delighted thank completely awesome experience last month nonstop thanks such awesome flight depart time #vabeatsjb
1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 
```

```
In [21]: 1 Mask = np.array(Image.open(requests.get
2                                     ('http://clipart-library.com/image_gallery2/Twitter-PNG-Image.png', stream=True).raw))
3
4 image_colors = ImageColorGenerator(Mask)
5
6 wc = WordCloud(background_color='black',
7                 height=1500,
8                 width=4000, mask=Mask).generate(all_positive_words)
```



```
In [23]: 1 all_negative_words = ' '.join(text for text in airline_tweets['text'][airline_tweets['airline_sentiment'] == 'neg']
2
3 all_negative_words
```

```
Out[23]: 'really aggressive blast obnoxious entertainment your guests faces they have little recourse really thing about se
riously would flight seats that didn't have this playing really only thing about flying schedule still flew from las
t week couldn't fully seat large gentleman either side help your first fares over three times more than other carrie
rs when seats available select guys messed seating reserved seating with friends guys gave seat away want free int
ernet status match program applied been three weeks called emailed with response what happened vegan food options
least site know able anything next #fail amazing that cold from vents #noair #worstflightever #roasted #sfotobos j
ust bked cool birthday trip with elevate cause entered middle name during flight booking problems help left expens
ive headphones flight today seat answering number awaiting return phone call just would prefer your online self se
rvice option your chat support working your site http gtdwpk first time flyer next week excited having hard time g
etting flights added elevate account help excited about your deal been trying book since last week page never load
s called weeks about adding flights from elevate they still haven't shown help heyyyy guyyys been trying through ho
ur someone call please virgin hold minutes there earlier flights from tonight earlier than everything fine until l
ost your airline awesome your loft needs step game dirty tables floors http vrfhjht what going with customer servi
ce there anyway speak human asap thank supp traveler like have customer service like #neverflyvirginforbusiness be
st whenever begrudgingly other airline delayed late flight have interesting flying with after this will cancelled
flight next four flights planned #neverflyvirginforbusiness disappointing experience which will shared with every
business traveler meet #neverflyvirgin having trouble adding this flight wife booked elevate account help http hqo
ks site down when will back like interesting video just disappointed cancelled flightled flight when other flights
went saturday just landed hour after should be here your late flight check business travel friendly #nomorevirgi
```

```
In [24]: 1 Mask = np.array(Image.open(requests.get('http://clipart-library.com/image_gallery2/Twitter-PNG-Image.png', stream
2
3 image_colors = ImageColorGenerator(Mask)
4
5 wc = WordCloud(background_color='black',
6                 height=1500, width=4000, mask=Mask).generate(all_negative_words)
```



```
In [26]: 1 from nltk.corpus import stopwords
2 from sklearn.feature_extraction.text import TfidfVectorizer
3 from sklearn.feature_extraction.text import CountVectorizer
```

max_df is used for removing terms that appear too frequently, also known as "corpus-specific stop words". For example:

max_df = 0.50 means "ignore terms that appear in more than 50% of the documents". max_df = 25 means "ignore terms that appear in more than 25 documents". The default max_df is 1.0, which means "ignore terms that appear in more than 100% of the documents". Thus, the default setting does not ignore any terms.

min_df is used for removing terms that appear too infrequently. For example:

min_df = 0.01 means "ignore terms that appear in less than 1% of the documents". min_df = 5 means "ignore terms that appear in less than 5 documents". The default min_df is 1, which means "ignore terms that appear in less than 1 document".

```
In [27]: 1 vectorizer = CountVectorizer(max_df=0.90, min_df=2,
2                               max_features=2500, stop_words = 'english')
3
4 airline_vector = vectorizer.fit_transform(airline_tweets['text'])
5
6 airline_vector_df = pd.DataFrame(airline_vector.todense())
7
8 airline_vector_df.shape
```

Out[27]: (14640, 2500)

```
In [28]: 1 tfidf = TfidfVectorizer(max_df=0.90, min_df=2,
2                               max_features=2500, stop_words='english')
3
4 airline_tfidf = tfidf.fit_transform(airline_tweets['text'])
5
6 airline_tfidf_df = pd.DataFrame(airline_tfidf.todense())
7
8 airline_tfidf_df.shape
```

Out[28]: (14640, 2500)

Importing necessary packages

```
In [29]: 1 from sklearn.linear_model import LogisticRegression
2 from sklearn.ensemble import RandomForestClassifier
3 from sklearn.svm import LinearSVC
4 from sklearn.svm import SVC
5 from sklearn.model_selection import train_test_split
6 from sklearn.metrics import f1_score
7 from sklearn.metrics import classification_report, confusion_matrix, accuracy_score
```

Splitting our dataset into Training and Validation Set for both TFIDF and CountVectorizer

```
In [30]: 1 X_train,X_test,Y_train,Y_test = train_test_split(airline_vector,airline_tweets['sentiment_class'],test_size=0.3,
2
3 X_train.shape,Y_train.shape,X_test.shape,Y_test.shape
```

```
Out[30]: ((10248, 2500), (10248,), (4392, 2500), (4392,))
```

```
In [31]: 1 x_train,x_test,y_train,y_test = train_test_split(airline_tfidf,airline_tweets['sentiment_class'],test_size=0.3, r
2
3 x_train.shape,y_train.shape,x_test.shape,y_test.shape
```

```
Out[31]: ((10248, 2500), (10248,), (4392, 2500), (4392,))
```

LogisticRegression

```
In [32]: 1 Log_Reg = LogisticRegression(random_state=0,solver='lbfgs')
```

```
In [33]: 1 Log_Reg.fit(X_train,Y_train)
         2 Log_Reg.fit(x_train,y_train)
```

```
C:\Users\Eric\anaconda3\envs\new_env\lib\site-packages\sklearn\linear_model\_logistic.py:765: ConvergenceWarning: lb
fgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)

```
C:\Users\Eric\anaconda3\envs\new_env\lib\site-packages\sklearn\linear_model\_logistic.py:765: ConvergenceWarning: lb
fgs failed to converge (status=1):
```

```
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
```

Increase the number of iterations (max_iter) or scale the data as shown in:

<https://scikit-learn.org/stable/modules/preprocessing.html> (<https://scikit-learn.org/stable/modules/preprocessing.html>)

Please also refer to the documentation for alternative solver options:

https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression (https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression)

extra_warning_msg=_LOGISTIC_SOLVER_CONVERGENCE_MSG)

```
Out[33]: LogisticRegression(random_state=0)
```

```
In [34]: 1 log_predict_vector = Log_Reg.predict(X_test)
         2 log_predict_tfidf = Log_Reg.predict(x_test)
```

```
1 print("##### Scores for CountVectorizer #####", "\n")
2 print("_____")
3 print(confusion_matrix(Y_test, log_predict_vector))
4 print(classification_report(Y_test, log_predict_vector))
5 print(accuracy_score(Y_test, log_predict_vector))
6 print("_____")
```

[[2427 163 141]					
[442 381 138]					
[131 52 517]]					
	precision	recall	f1-score	support	
	0	0.81	0.89	0.85	2731
	3	0.64	0.40	0.49	961
	4	0.65	0.74	0.69	700
accuracy				0.76	4392
macro avg		0.70	0.67	0.68	4392
weighted avg		0.75	0.76	0.74	4392
0.7570582877959927					

```
1 print("##### Scores for TFIDF #####", "\n")
2 print("_____")
3 print(confusion_matrix(y_test, log_predict_tfidf))
4 print(classification_report(y_test, log_predict_tfidf))
5 print(accuracy_score(y_test, log_predict_tfidf))
6 print("_____")
```

[[2559	134	38]			
[513	395	53]			
[199	93	408]]			
	precision		recall	f1-score	support
	0	0.78	0.94	0.85	2731
	3	0.64	0.41	0.50	961
	4	0.82	0.58	0.68	700
accuracy				0.77	4392
macro avg	0.75	0.64	0.68		4392
weighted avg	0.76	0.77	0.75		4392
0.7654826958105647					

RandomForestClassifier


```
In [ ]: 1 text_classifier = RandomForestClassifier(n_estimators=50, random_state=0)
        2 text_classifier.fit(X_train, Y_train)
        3 text_classifier.fit(x_train, y_train)
```

```
Out[87]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                                criterion='gini', max_depth=None, max_features='auto',
                                max_leaf_nodes=None, max_samples=None,
                                min_impurity_decrease=0.0, min_impurity_split=None,
                                min_samples_leaf=1, min_samples_split=2,
                                min_weight_fraction_leaf=0.0, n_estimators=50,
                                n_jobs=None, oob_score=False, random_state=0, verbose=0,
                                warm_start=False)
```

```
In [ ]: 1 ran_pred_vector = text_classifier.predict(X_test)
        2 ran_pred_tfidf = text_classifier.predict(x_test)
```

```
1 print("##### Scores for CountVectorizer #####", "\n")
2 print("_____")
3 print(confusion_matrix(Y_test, ran_pred_vector))
4 print(classification_report(Y_test, ran_pred_vector))
5 print(accuracy_score(Y_test, ran_pred_vector))
6 print("_____")
```

[[1916 667 148]					
[219 623 119]					
[81 167 452]]					
	precision	recall	f1-score	support	
0	0.86	0.70	0.77	2731	
3	0.43	0.65	0.52	961	
4	0.63	0.65	0.64	700	
accuracy			0.68	4392	
macro avg	0.64	0.67	0.64	4392	
weighted avg	0.73	0.68	0.70	4392	
0.6810109289617486					

```
1 print("##### Scores for TFIDF #####", "\n")
2 print("_____")
3 print(confusion_matrix(y_test, ran_pred_tfidf))
4 print(classification_report(y_test, ran_pred_tfidf))
5 print(accuracy_score(y_test, ran_pred_tfidf))
6 print("_____")
```

[[2433	215	83]				
[452	403	106]				
[193	87	420]]				
		precision	recall	f1-score	support	
	0	0.79	0.89	0.84	2731	
	3	0.57	0.42	0.48	961	
	4	0.69	0.60	0.64	700	
	accuracy			0.74	4392	
	macro avg	0.68	0.64	0.65	4392	
	weighted avg	0.73	0.74	0.73	4392	
						0.7413479052823315

```
1 clf_1 = LinearSVC()
```

```
1 clf_1.fit(X_train, Y_train)
2 clf_1.fit(x_train, y_train)
```

```
Out[92]: LinearSVC(C=1.0, class_weight=None, dual=True, fit_intercept=True,
                  intercept_scaling=1, loss='squared_hinge', max_iter=1000,
                  multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,
                  verbose=0)
```

```
In [ ]: 1 svc_pred_vector = clf_1.predict(X_test)
        2 svc_pred_tfidf = clf_1.predict(x_test)
```

```
In [ ]: 1 print("##### Scores for CountVectorizer #####", "\n")
2 print("_____")
3 print(confusion_matrix(Y_test, svc_pred_vector))
4 print(classification_report(Y_test, svc_pred_vector))
5 print(accuracy_score(Y_test, svc_pred_vector))
6 print("_____")
```

```
##### Scores for CountVectorizer #####
```

[[2441 151 139]					
[400 403 158]					
[111 63 526]]					
		precision	recall	f1-score	support
	0	0.83	0.89	0.86	2731
	3	0.65	0.42	0.51	961
	4	0.64	0.75	0.69	700
accuracy				0.77	4392
macro avg		0.71	0.69	0.69	4392
weighted avg		0.76	0.77	0.76	4392

0.767304189435337

```
In [ ]: 1 print("##### Scores for TFIDF #####", "\n")
2 print("_____")
3 print(confusion_matrix(y_test, svc_pred_tfidf))
4 print(classification_report(y_test, svc_pred_tfidf))
5 print(accuracy_score(y_test, svc_pred_tfidf))
6 print("_____")
```

Scores for TFIDF

```
[[2439  221   71]
 [ 384  493   84]
 [ 143  108  449]]
      precision    recall  f1-score   support

      0       0.82      0.89      0.86      2731
      3       0.60      0.51      0.55       961
      4       0.74      0.64      0.69      700

 accuracy          0.77      4392
 macro avg       0.72      0.68      0.70      4392
 weighted avg    0.76      0.77      0.76      4392

0.769808743169399
```

SVC

```
In [ ]: 1 clf_2 = SVC(kernel = 'linear')
```

```
In [ ]: 1 clf_2.fit(X_train, Y_train)
2 clf_2.fit(x_train, y_train)
```

Out[97]: SVC(C=1.0, break_ties=False, cache_size=200, class_weight=None, coef0=0.0, decision_function_shape='ovr', degree=3, gamma='scale', kernel='linear', max_iter=-1, probability=False, random_state=None, shrinking=True, tol=0.001, verbose=False)

```
In [ ]: 1 svm_pred_vector = clf_2.predict(X_test)
        2 svm_pred_tfidf = clf_2.predict(x_test)
```

```
In [ ]: 1 print("##### Scores for CountVectorizer #####", "\n")
2 print("_____")
3 print(confusion_matrix(Y_test, svm_pred_vector))
4 print(classification_report(Y_test, svm_pred_vector))
5 print(accuracy_score(Y_test, svm_pred_vector))
6 print("_____")
```

```
##### Scores for CountVectorizer #####
```

[[2395 157 179]					
[462 346 153]					
[110 49 541]]					
	precision	recall	f1-score	support	
0	0.81	0.88	0.84	2731	
3	0.63	0.36	0.46	961	
4	0.62	0.77	0.69	700	
accuracy			0.75	4392	
macro avg	0.68	0.67	0.66	4392	
weighted avg	0.74	0.75	0.73	4392	

0.7472677595628415

```
1 print("##### Scores for TFIDF #####", "\n")
2 print("_____")
3 print(confusion_matrix(y_test, svm_pred_tfidf))
4 print(classification_report(y_test, svm_pred_tfidf))
5 print(accuracy_score(y_test, svm_pred_tfidf))
6 print("_____")
```

[[2523 149 59]					
[517 384 60]					
[188 73 439]]					
	precision	recall	f1-score	support	
0	0.78	0.92	0.85	2731	
3	0.63	0.40	0.49	961	
4	0.79	0.63	0.70	700	
accuracy			0.76	4392	
macro avg	0.73	0.65	0.68	4392	
weighted avg	0.75	0.76	0.75	4392	
0.76183970856102					