# prabhudayala@gmail.com_27

November 3, 2019

```python
import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import sqlite3
import csv
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from wordcloud import WordCloud
import re
import os
from sqlalchemy import create_engine # database connection
import datetime as dt
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem.snowball import SnowballStemmer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import SGDClassifier
from sklearn import metrics
from sklearn.metrics import f1_score,precision_score,recall_score
from sklearn import svm
from sklearn.linear_model import LogisticRegression
from skmultilearn.adapt import mlknn
from skmultilearn.problem_transform import ClassifierChain
from skmultilearn.problem_transform import BinaryRelevance
from skmultilearn.problem_transform import LabelPowerset
from sklearn.naive_bayes import GaussianNB
from datetime import datetime
```

# 1 Stack Overflow: Tag Prediction

1. Business Problem

1.1 Description

1

Description

Stack Overflow is the largest, most trusted online community for developers to learn, share their programming knowledge, and build their careers. Stack Overflow is something which every programmer use one way or another. Each month, over 50 million developers come to Stack Overflow to learn, share their knowledge, and build their careers. It features questions and answers on a wide range of topics in computer programming. The website serves as a platform for users to ask and answer questions, and, through membership and active participation, to vote questions and answers up or down and edit questions and answers in a fashion similar to a wiki or Digg. As of April 2014 Stack Overflow has over 4,000,000 registered users, and it exceeded 10,000,000 questions in late August 2015. Based on the type of tags assigned to questions, the top eight most discussed topics on the site are: Java, JavaScript, C#, PHP, Android, jQuery, Python and HTML.

Problem Statemtent

Suggest the tags based on the content that was there in the question posted on Stackoverflow.

Source: https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/

1.2 Source / useful links

Data Source : https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data Youtube : https://youtu.be/nNDqbUhtIRg Research paper : https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tagging-1.pdf Research paper : https://dl.acm.org/citation.cfm?id=2660970&dl=ACM&coll=DL

1.3 Real World / Business Objectives and Constraints

1. Predict as many tags as possible with high precision and recall.
2. Incorrect tags could impact customer experience on StackOverflow.
3. No strict latency constraints.

2. Machine Learning problem

2.1 Data

2.1.1 Data Overview

Refer: https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data All of the data is in 2 files: Train and Test.

The questions are randomized and contains a mix of verbose text sites as well as sites related to math and programming. The number of questions from each site may vary, and no filtering has been performed on the questions (such as closed questions).

**Data Field Explaination**

Dataset contains 6,034,195 rows. The columns in the table are:

2.1.2 Example Data point

2.2 Mapping the real-world problem to a Machine Learning Problem

2.2.1 Type of Machine Learning Problem

It is a multi-label classification problem Multi-label Classification: Multilabel classification assigns to each sample a set of target labels. This can be thought as predicting properties of a datapoint that are not mutually exclusive, such as topics that are relevant for a document. A question on Stackoverflow might be about any of C, Pointers, FileIO and/or memory-management at the same time or none of these. **Credit**: http://scikit-learn.org/stable/modules/multiclass.html

2.2.2 Performance metric

Micro-Averaged F1-Score (Mean F Score) : The F1 score can be interpreted as a weighted average of the precision and recall, where an F1 score reaches its best value at 1 and worst score at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for the F1 score is:

F1 = 2 * (precision * recall) / (precision + recall)

In the multi-class and multi-label case, this is the weighted average of the F1 score of each class.

'Micro f1 score': Calculate metrics globally by counting the total true positives, false negatives and false positives. This is a better metric when we have class imbalance.

'Macro f1 score': Calculate metrics for each label, and find their unweighted mean. This does not take label imbalance into account.

https://www.kaggle.com/wiki/MeanFScore http://scikit-learn.org/stable/modules/generated/sklearn.me Hamming loss : The Hamming loss is the fraction of labels that are incorrectly predicted. https://www.kaggle.com/wiki/HammingLoss

3. Exploratory Data Analysis

3.1 Data Loading and Cleaning
3.1.1 Using Pandas with SQLite to Load the data

```python
[2]: #Creating db file from csv
     #Learn SQL: https://www.w3schools.com/sql/default.asp
     if not os.path.isfile('train.db'):
         start = datetime.now()
         disk_engine = create_engine('sqlite:///train.db')
         start = dt.datetime.now()
         chunksize = 180000
         j = 0
         index_start = 1
         for df in pd.read_csv('Train.csv', names=['Id', 'Title', 'Body', 'Tags'],
     ↪chunksize=chunksize, iterator=True, encoding='utf-8', ):
             df.index += index_start
             j+=1
             print('{} rows'.format(j*chunksize))
             df.to_sql('data', disk_engine, if_exists='append')
             index_start = df.index[-1] + 1
         print("Time taken to run this cell :", datetime.now() - start)
```

3.1.2 Counting the number of rows

```python
[3]: if os.path.isfile('train.db'):
         start = datetime.now()
         con = sqlite3.connect('train.db')
         num_rows = pd.read_sql_query("""SELECT count(*) FROM data""", con)
         #Always remember to close the database
         print("Number of rows in the database :","\n",num_rows['count(*)'].
     ↪values[0])
         con.close()
         print("Time taken to count the number of rows :", datetime.now() - start)
     else:
         print("Please download the train.db file from drive or run the above cell
     ↪to genarate train.db file")
```

```
Number of rows in the database :
 6034196
Time taken to count the number of rows : 0:00:04.984194
```

### 3.1.3 Checking for duplicates

```python
[4]: #Learn SQl: https://www.w3schools.com/sql/default.asp
     if os.path.isfile('train.db'):
         start = datetime.now()
         con = sqlite3.connect('train.db')
         df_no_dup = pd.read_sql_query('SELECT Title, Body, Tags, COUNT(*) as␣
      ↪cnt_dup FROM data GROUP BY Title, Body, Tags', con)
         con.close()
         print("Time taken to run this cell :", datetime.now() - start)
     else:
         print("Please download the train.db file from drive or run the first to␣
      ↪genarate train.db file")
```

```
Time taken to run this cell : 0:04:23.749168
```

```python
[5]: df_no_dup.head()
     # we can observe that there are duplicates
```

```
[5]:                                               Title  \
     0        Implementing Boundary Value Analysis of S...
     1             Dynamic Datagrid Binding in Silverlight?
     2             Dynamic Datagrid Binding in Silverlight?
     3        java.lang.NoClassDefFoundError: javax/serv...
     4        java.sql.SQLException:[Microsoft][ODBC Dri...

                                                Body  \
     0  <pre><code>#include&lt;iostream&gt;\n#include&...
     1  <p>I should do binding for datagrid dynamicall...
     2  <p>I should do binding for datagrid dynamicall...
     3  <p>I followed the guide in <a href="http://sta...
     4  <p>I use the following code</p>\n\n<pre><code>...

                                  Tags  cnt_dup
     0                          c++ c        1
     1          c# silverlight data-binding        1
     2  c# silverlight data-binding columns        1
     3                        jsp jstl        1
     4                        java jdbc        2
```

```python
[6]: df_no_dup.isna().any()
```

```
[6]: Title      False
     Body       False
     Tags        True
     cnt_dup    False
```

```
dtype: bool
```

```
[7]: df_no_dup.isnull().any()
```

```
[7]: Title      False
     Body       False
     Tags        True
     cnt_dup    False
     dtype: bool
```

```
[8]: df_no_dup.Tags.isna().sum()
```

```
[8]: 7
```

```
[9]: df_no_dup[df_no_dup.Tags.isna() ==True]
```

```
[9]:                                                 Title  \
     777547                         Do we really need NULL?
     962680    Find all values that are not null and not in a...
     1126558                              Handle NullObjects
     1256102                           How do Germans call null
     2430668   Page cannot be null. Please ensure that this o...
     3329908       What is the difference between NULL and "0"?
     3551595         a bit of difference between null and space


                                                 Body  Tags  cnt_dup
     777547    <blockquote>\n  <p><strong>Possible Duplicate:...  None        1
     962680    <p>I am running into a problem which results i...  None        1
     1126558   <p>I have done quite a bit of research on best...  None        1
     1256102   <p>In german null means 0, so how do they call...  None        1
     2430668   <p>I get this error when i remove dynamically ...  None        1
     3329908   <p>What is the difference from NULL and "0"?</...  None        1
     3551595   <p>I was just reading this quote</p>\n\n<block...  None        2
```

```
[10]: df_no_dup.shape
```

```
[10]: (4206315, 4)
```

```
[11]: df_no_dup.drop(df_no_dup[df_no_dup.Tags.isna() ==True].index, inplace=True)
```

```
[12]: df_no_dup.shape
```

```
[12]: (4206308, 4)
```

7 rows are deleted where there was no Tags

```
[13]: print("number of duplicate questions :", num_rows['count(*)'].values[0]-␣
      ↪df_no_dup.shape[0], "(",(1-((df_no_dup.shape[0])/(num_rows['count(*)'].
      ↪values[0])))*100,"% )")
```

number of duplicate questions : 1827888 ( 30.29215491177284 % )

```
[14]: # number of times each question appeared in our database
      df_no_dup.cnt_dup.value_counts()
```

```
[14]: 1    2656278
      2    1272335
      3     277575
      4         90
      5         25
      6          5
      Name: cnt_dup, dtype: int64
```

```
[15]: start = datetime.now()
      df_no_dup["tag_count"] = df_no_dup["Tags"].apply(lambda text: len(text.split("␣
        ↪")))
      # adding a new feature number of tags per question
      print("Time taken to run this cell :", datetime.now() - start)
      df_no_dup.head()
```

```
Time taken to run this cell : 0:00:02.062485
```

```
[15]:                                              Title  \
      0        Implementing Boundary Value Analysis of S...
      1            Dynamic Datagrid Binding in Silverlight?
      2            Dynamic Datagrid Binding in Silverlight?
      3       java.lang.NoClassDefFoundError: javax/serv...
      4       java.sql.SQLException:[Microsoft][ODBC Dri...

                                                    Body  \
      0  <pre><code>#include&lt;iostream&gt;\n#include&...
      1  <p>I should do binding for datagrid dynamicall...
      2  <p>I should do binding for datagrid dynamicall...
      3  <p>I followed the guide in <a href="http://sta...
      4  <p>I use the following code</p>\n\n<pre><code>...

                                    Tags  cnt_dup  tag_count
      0                          c++ c          1          2
      1          c# silverlight data-binding          1          3
      2  c# silverlight data-binding columns          1          4
      3                        jsp jstl          1          2
      4                       java jdbc          2          2
```

```
[16]: # distribution of number of tags per question
      df_no_dup.tag_count.value_counts()
```

```
[16]: 3    1206157
      2    1111706
      4     814996
      1     568291
      5     505158
      Name: tag_count, dtype: int64
```

```python
[17]: #Creating a new database with no duplicates
      if not os.path.isfile('train_no_dup.db'):
          disk_dup = create_engine("sqlite:///train_no_dup.db")
          no_dup = pd.DataFrame(df_no_dup, columns=['Title', 'Body', 'Tags'])
          no_dup.to_sql('no_dup_train',disk_dup)
```

```python
[18]: #This method seems more appropriate to work with this much data.
      #creating the connection with database file.
      if os.path.isfile('train_no_dup.db'):
          start = datetime.now()
          con = sqlite3.connect('train_no_dup.db')
          tag_data = pd.read_sql_query("""SELECT Tags FROM no_dup_train""", con)
          #Always remember to close the database
          con.close()

          # Let's now drop unwanted column.
          tag_data.drop(tag_data.index[0], inplace=True)
          #Printing first 5 columns from our data frame
          tag_data.head()
          print("Time taken to run this cell :", datetime.now() - start)
      else:
          print("Please download the train.db file from drive or run the above cells␣
       ↪to genarate train.db file")
```

```
Time taken to run this cell : 0:00:24.248834
```

3.2 Analysis of Tags
3.2.1 Total number of unique tags

```python
[19]: # Importing & Initializing the "CountVectorizer" object, which
      #is scikit-learn's bag of words tool.

      #by default 'split()' will tokenize each tag using space.
      vectorizer = CountVectorizer(tokenizer = lambda x: x.split())
      # fit_transform() does two functions: First, it fits the model
      # and learns the vocabulary; second, it transforms our training data
      # into feature vectors. The input to fit_transform should be a list of strings.
      tag_dtm = vectorizer.fit_transform(tag_data['Tags'])
```

```python
[20]: print("Number of data points :", tag_dtm.shape[0])
      print("Number of unique tags :", tag_dtm.shape[1])
```

```
Number of data points : 4206307
Number of unique tags : 42048
```

```python
[21]: #'get_feature_name()' gives us the vocabulary.
      tags = vectorizer.get_feature_names()
      #Lets look at the tags we have.
      print("Some of the tags we have :", tags[:10])
```

Some of the tags we have : ['.a', '.app', '.asp.net-mvc', '.aspxauth', '.bash-profile', '.class-file', '.cs-file', '.doc', '.drv', '.ds-store']

### 3.2.3 Number of times a tag appeared

```
[22]: # https://stackoverflow.com/questions/15115765/
      ↪how-to-access-sparse-matrix-elements
      #Lets now store the document term matrix in a dictionary.
      freqs = tag_dtm.sum(axis=0).A1
      result = dict(zip(tags, freqs))
```

```
[23]: #Saving this dictionary to csv files.
      if not os.path.isfile('tag_counts_dict_dtm.csv'):
          with open('tag_counts_dict_dtm.csv', 'w') as csv_file:
              writer = csv.writer(csv_file)
              for key, value in result.items():
                  writer.writerow([key, value])
      tag_df = pd.read_csv("tag_counts_dict_dtm.csv", names=['Tags', 'Counts'])
      tag_df.head()
```

```
[23]:           Tags  Counts
      0            .a      18
      1          .app      37
      2  .asp.net-mvc       1
      3     .aspxauth      21
      4  .bash-profile     138
```

```
[24]: tag_df_sorted = tag_df.sort_values(['Counts'], ascending=False)
      tag_counts = tag_df_sorted['Counts'].values
```

```
[25]: plt.plot(tag_counts)
      plt.title("Distribution of number of times tag appeared questions")
      plt.grid()
      plt.xlabel("Tag number")
      plt.ylabel("Number of times tag appeared")
      plt.show()
```

## Distribution of number of times tag appeared questions



[26]:
```python
plt.plot(tag_counts[0:10000])
plt.title('first 10k tags: Distribution of number of times tag appeared␣
 ↪questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:10000:25]), tag_counts[0:10000:25])
```

## first 10k tags: Distribution of number of times tag appeared questions



```
400  [331505   44829   22429   17728   13364   11162   10029    9148    8054    7151
      6466    5865    5370    4983    4526    4281    4144    3929    3750    3593
      3453    3299    3123    2989    2891    2738    2647    2527    2431    2331
      2259    2186    2097    2020    1959    1900    1828    1770    1723    1673
      1631    1574    1532    1479    1448    1406    1365    1328    1300    1266
      1245    1222    1197    1181    1158    1139    1121    1101    1076    1056
      1038    1023    1006     983     966     952     938     926     911     891
       882     869     856     841     830     816     804     789     779     770
       752     743     733     725     712     702     688     678     671     658
       650     643     634     627     616     607     598     589     583     577
       568     559     552     545     540     533     526     518     512     506
       500     495     490     485     480     477     469     465     457     450
       447     442     437     432     426     422     418     413     408     403
       398     393     388     385     381     378     374     370     367     365
       361     357     354     350     347     344     342     339     336     332
       330     326     323     319     315     312     309     307     304     301
       299     296     293     291     289     286     284     281     278     276
       275     272     270     268     265     262     260     258     256     254
       252     250     249     247     245     243     241     239     238     236
       234     233     232     230     228     226     224     222     220     219
       217     215     214     212     210     209     207     205     204     203
       201     200     199     198     196     194     193     192     191     189
       188     186     185     183     182     181     180     179     178     177
       175     174     172     171     170     169     168     167     166     165
       164     162     161     160     159     158     157     156     156     155
```

```
154    153    152    151    150    149    149    148    147    146
145    144    143    142    142    141    140    139    138    137
137    136    135    134    134    133    132    131    130    130
129    128    128    127    126    126    125    124    124    123
123    122    122    121    120    120    119    118    118    117
117    116    116    115    115    114    113    113    112    111
111    110    109    109    108    108    107    106    106    106
105    105    104    104    103    103    102    102    101    101
100    100     99     99     98     98     97     97     96     96
 95     95     94     94     93     93     93     92     92     91
 91     90     90     89     89     88     88     87     87     86
 86     86     85     85     84     84     83     83     83     82
 82     82     81     81     80     80     80     79     79     78
 78     78     78     77     77     76     76     76     75     75
 75     74     74     74     73     73     73     73     72     72]
```

[27]:
```python
plt.plot(tag_counts[0:1000])
plt.title('first 1k tags: Distribution of number of times tag appeared␣
 ↪questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



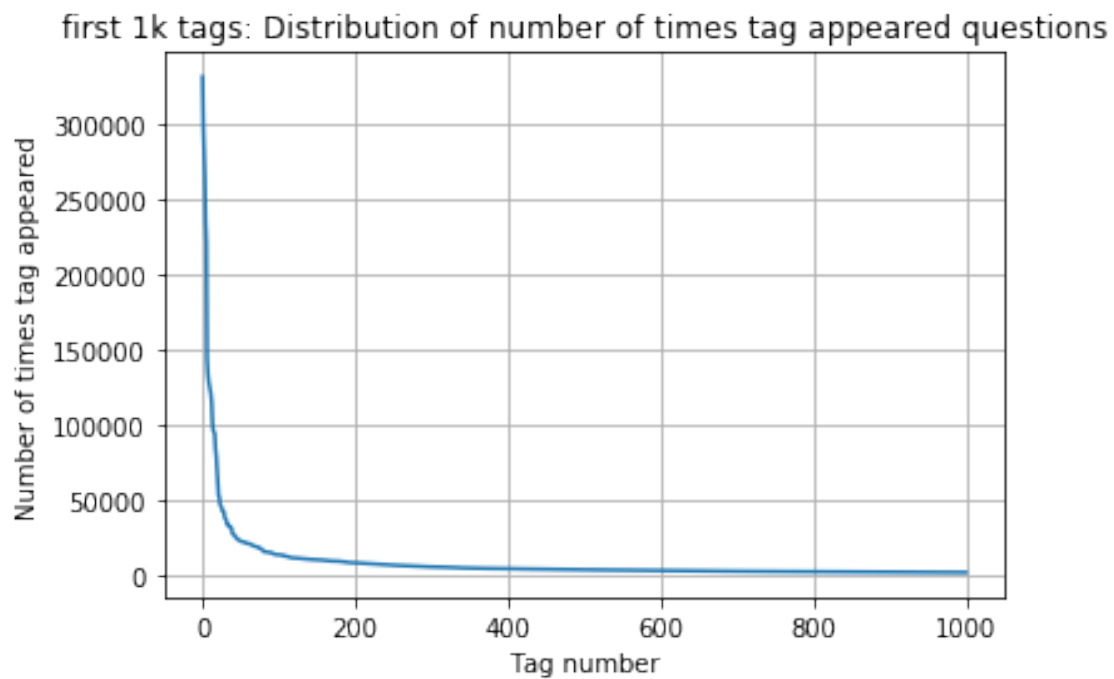first 1k tags: Distribution of number of times tag appeared questions

```
200 [331505 221533 122769  95160  62023  44829  37170  31897  26925  24537
  22429  21820  20957  19758  18905  17728  15533  15097  14884  13703
  13364  13157  12407  11658  11228  11162  10863  10600  10350  10224
  10029   9884   9719   9411   9252   9148   9040   8617   8361   8163
   8054   7867   7702   7564   7274   7151   7052   6847   6656   6553
   6466   6291   6183   6093   5971   5865   5760   5577   5490   5411
   5370   5283   5207   5107   5066   4983   4891   4785   4658   4549
   4526   4487   4429   4335   4310   4281   4239   4228   4195   4159
   4144   4088   4050   4002   3957   3929   3874   3849   3818   3797
   3750   3703   3685   3658   3615   3593   3564   3521   3505   3483
   3453   3427   3396   3363   3326   3299   3272   3232   3196   3168
   3123   3094   3073   3050   3012   2989   2984   2953   2934   2903
   2891   2844   2819   2784   2754   2738   2726   2708   2681   2669
   2647   2621   2604   2594   2556   2527   2510   2482   2460   2444
   2431   2409   2395   2380   2363   2331   2312   2297   2290   2281
   2259   2246   2222   2211   2198   2186   2162   2142   2132   2107
   2097   2078   2057   2045   2036   2020   2011   1994   1971   1965
   1959   1952   1940   1932   1912   1900   1879   1865   1855   1841
   1828   1821   1813   1801   1782   1770   1760   1747   1741   1734
   1723   1707   1697   1688   1683   1673   1665   1656   1646   1639]
```
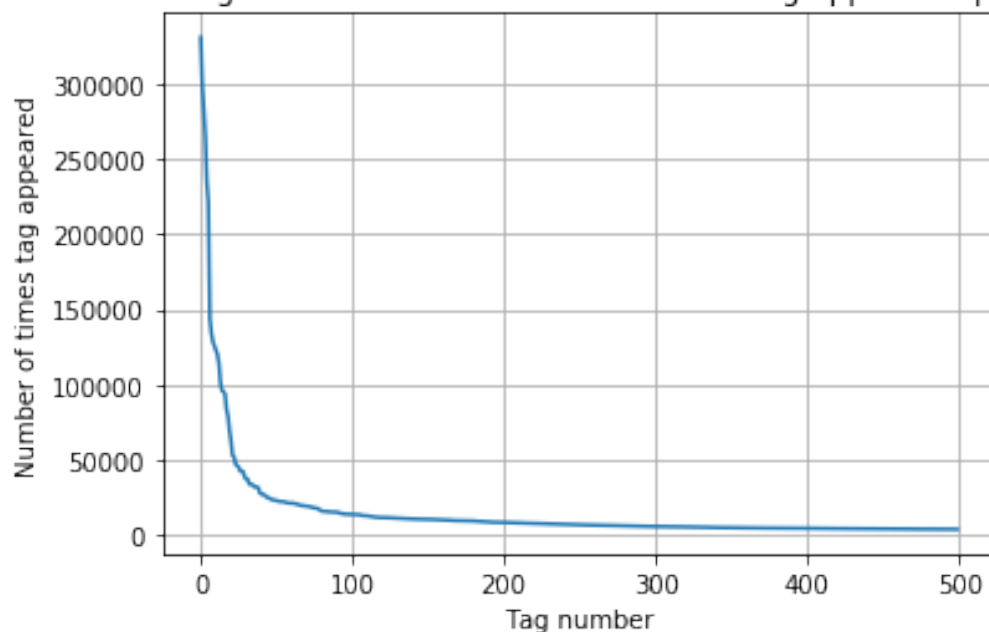
[28]:
```python
plt.plot(tag_counts[0:500])
plt.title('first 500 tags: Distribution of number of times tag appeared␣
 ↪questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:500:5]), tag_counts[0:500:5])
```

first 500 tags: Distribution of number of times tag appeared questions

```
100 [331505 221533 122769  95160  62023  44829  37170  31897  26925  24537
    22429  21820  20957  19758  18905  17728  15533  15097  14884  13703
    13364  13157  12407  11658  11228  11162  10863  10600  10350  10224
    10029   9884   9719   9411   9252   9148   9040   8617   8361   8163
     8054   7867   7702   7564   7274   7151   7052   6847   6656   6553
     6466   6291   6183   6093   5971   5865   5760   5577   5490   5411
     5370   5283   5207   5107   5066   4983   4891   4785   4658   4549
     4526   4487   4429   4335   4310   4281   4239   4228   4195   4159
     4144   4088   4050   4002   3957   3929   3874   3849   3818   3797
     3750   3703   3685   3658   3615   3593   3564   3521   3505   3483]
```

```python
[29]: plt.plot(tag_counts[0:100], c='b')
      plt.scatter(x=list(range(0,100,5)), y=tag_counts[0:100:5], c='orange',
       ↪label="quantiles with 0.05 intervals")
      # quantiles with 0.25 difference
      plt.scatter(x=list(range(0,100,25)), y=tag_counts[0:100:25], c='m', label =
       ↪"quantiles with 0.25 intervals")

      for x,y in zip(list(range(0,100,25)), tag_counts[0:100:25]):
          plt.annotate(s="({} , {})".format(x,y), xy=(x,y), xytext=(x-0.05, y+500))

      plt.title('first 100 tags: Distribution of number of times tag appeared
       ↪questions')
      plt.grid()
      plt.xlabel("Tag number")
```

```
plt.ylabel("Number of times tag appeared")
plt.legend()
plt.show()
print(len(tag_counts[0:100:5]), tag_counts[0:100:5])
```

first 100 tags: Distribution of number of times tag appeared questions



```
20 [331505 221533 122769  95160  62023  44829  37170  31897  26925  24537
   22429  21820  20957  19758  18905  17728  15533  15097  14884  13703]
```

[30]:
```
# Store tags greater than 10K in one list
lst_tags_gt_10k = tag_df[tag_df.Counts>10000].Tags
#Print the length of the list
print ('{} Tags are used more than 10000 times'.format(len(lst_tags_gt_10k)))
# Store tags greater than 100K in one list
lst_tags_gt_100k = tag_df[tag_df.Counts>100000].Tags
#Print the length of the list.
print ('{} Tags are used more than 100000 times'.format(len(lst_tags_gt_100k)))
```

```
153 Tags are used more than 10000 times
14 Tags are used more than 100000 times
```

Observations: 1. There are total 153 tags which are used more than 10000 times. 2. 14 tags are used more than 100000 times. 3. Most frequent tag (i.e. c#) is used 331505 times. 4. Since some tags occur much more frequenctly than others, Micro-averaged F1-score is the appropriate metric for this probelm.

3.2.4 Tags Per Question

```
[31]: #Storing the count of tag in each question in list 'tag_count'
      tag_quest_count = tag_dtm.sum(axis=1).tolist()
      #Converting list of lists into single list, we will get [[3], [4], [2], [2],␣
      ↪[3]] and we are converting this to [3, 4, 2, 2, 3]
      tag_quest_count=[int(j) for i in tag_quest_count for j in i]
      print ('We have total {} datapoints.'.format(len(tag_quest_count)))

      print(tag_quest_count[:5])
```

```
We have total 4206307 datapoints.
[3, 4, 2, 2, 3]
```

```
[32]: print( "Maximum number of tags per question: %d"%max(tag_quest_count))
      print( "Minimum number of tags per question: %d"%min(tag_quest_count))
      print( "Avg. number of tags per question: %f"% ((sum(tag_quest_count)*1.0)/
      ↪len(tag_quest_count)))
```

```
Maximum number of tags per question: 5
Minimum number of tags per question: 1
Avg. number of tags per question: 2.899443
```

```
[33]: sns.countplot(tag_quest_count, palette='gist_rainbow')
      plt.title("Number of tags in the questions ")
      plt.xlabel("Number of Tags")
      plt.ylabel("Number of questions")
      plt.show()
```

Observations: 1. Maximum number of tags per question: 5 2. Minimum number of tags per question: 1 3. Avg. number of tags per question: 2.899 4. Most of the questions are having 2 or 3 tags

3.2.5 Most Frequent Tags

```python
# Ploting word cloud
start = datetime.now()

# Lets first convert the 'result' dictionary to 'list of tuples'
tup = dict(result.items())
#Initializing WordCloud using frequencies of tags.
wordcloud = WordCloud(    background_color='black',
                          width=1600,
                          height=800,
                    ).generate_from_frequencies(tup)

fig = plt.figure(figsize=(30,20))
plt.imshow(wordcloud)
plt.axis('off')
plt.tight_layout(pad=0)
fig.savefig("tag.png")
plt.show()
print("Time taken to run this cell :", datetime.now() - start)
```



```
Time taken to run this cell : 0:00:03.136614
```

Observations: A look at the word cloud shows that "c#", "java", "php", "asp.net", "javascript", "c++" are some of the most frequent tags.

16

### 3.2.6 The top 20 tags

```
[35]: i=np.arange(30)
      tag_df_sorted.head(30).plot(kind='bar')
      plt.title('Frequency of top 20 tags')
      plt.xticks(i, tag_df_sorted['Tags'])
      plt.xlabel('Tags')
      plt.ylabel('Counts')
      plt.show()
```



Observations: 1. Majority of the most frequent tags are programming language. 2. C# is the top most frequent programming language. 3. Android, IOS, Linux and windows are among the top most frequent operating systems.

3.3 Cleaning and preprocessing of Questions

3.3.1 Preprocessing

Sample 1M data points

Separate out code-snippets from Body

Remove Spcial characters from Question title and description (not in code)

Remove stop words (Except 'C')

Remove HTML Tags

Convert all the characters into small letters

17

Use SnowballStemmer to stem the words

```python
[36]: def striphtml(data):
          cleanr = re.compile('<.*?>')
          cleantext = re.sub(cleanr, ' ', str(data))
          return cleantext
      stop_words = set(stopwords.words('english'))
      stemmer = SnowballStemmer("english")
```

```python
[37]: #http://www.sqlitetutorial.net/sqlite-python/create-tables/
      def create_connection(db_file):
          """ create a database connection to the SQLite database
              specified by db_file
          :param db_file: database file
          :return: Connection object or None
          """
          try:
              conn = sqlite3.connect(db_file)
              return conn
          except Error as e:
              print(e)

          return None

      def create_table(conn, create_table_sql):
          """ create a table from the create_table_sql statement
          :param conn: Connection object
          :param create_table_sql: a CREATE TABLE statement
          :return:
          """
          try:
              c = conn.cursor()
              c.execute(create_table_sql)
          except Error as e:
              print(e)

      def checkTableExists(dbcon):
          cursr = dbcon.cursor()
          str = "select name from sqlite_master where type='table'"
          table_names = cursr.execute(str)
          print("Tables in the databse:")
          tables =table_names.fetchall()
          print(tables[0][0])
          return(len(tables))

      def create_database_table(database, query):
          conn = create_connection(database)
          if conn is not None:
              create_table(conn, query)
```

```
        checkTableExists(conn)
    else:
        print("Error! cannot create the database connection.")
    conn.close()

sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question␣
 ↪text NOT NULL, code text, tags text, words_pre integer, words_post integer,␣
 ↪is_code integer);"""
create_database_table("Processed.db", sql_create_table)
```

```
Tables in the databse:
QuestionsProcessed
```

[38]:
```
# http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/
 ↪select-random-row-from-a-sqlite-table
start = datetime.now()
read_db = 'train_no_dup.db'
write_db = 'Processed.db'
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        reader =conn_r.cursor()
        reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY␣
 ↪RANDOM() LIMIT 1000000;")

if os.path.isfile(write_db):
    conn_w = create_connection(write_db)
    if conn_w is not None:
        tables = checkTableExists(conn_w)
        writer =conn_w.cursor()
        if tables != 0:
            writer.execute("DELETE FROM QuestionsProcessed WHERE 1")
            print("Cleared All the rows")
print("Time taken to run this cell :", datetime.now() - start)
```

```
Tables in the databse:
QuestionsProcessed
Cleared All the rows
Time taken to run this cell : 0:09:38.947421
```

__ we create a new data base to store the sampled and preprocessed questions __

[39]:
```
#http://www.bernzilla.com/2008/05/13/
 ↪selecting-a-random-row-from-an-sqlite-table/

start = datetime.now()
```

```python
preprocessed_data_list=[]
reader.fetchone()
questions_with_code=0
len_pre=0
len_post=0
questions_proccesed = 0
for row in reader:

    is_code = 0

    title, question, tags = row[0], row[1], row[2]

    if '<code>' in question:
        questions_with_code+=1
        is_code = 1
    x = len(question)+len(title)
    len_pre+=x

    code = str(re.findall(r'<code>(.*?)</code>', question, flags=re.DOTALL))

    question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTILINE|re.
↪DOTALL)
    question=striphtml(question.encode('utf-8'))

    title=title.encode('utf-8')

    question=str(title)+" "+str(question)
    question=re.sub(r'[^A-Za-z]+',' ',question)
    words=word_tokenize(str(question.lower()))

    #Removing all single letter and and stopwords from question exceptt for the␣
↪letter 'c'
    question=' '.join(str(stemmer.stem(j)) for j in words if j not in␣
↪stop_words and (len(j)!=1 or j=='c'))

    len_post+=len(question)
    tup = (question,code,tags,x,len(question),is_code)
    questions_proccesed += 1
    writer.execute("insert into␣
↪QuestionsProcessed(question,code,tags,words_pre,words_post,is_code) values (?
↪,?,?,?,?,?)",tup)
    if (questions_proccesed%100000==0):
        print("number of questions completed=",questions_proccesed)

no_dup_avg_len_pre=(len_pre*1.0)/questions_proccesed
no_dup_avg_len_post=(len_post*1.0)/questions_proccesed
```

```
print( "Avg. length of questions(Title+Body) before processing:␣
  ↪%d"%no_dup_avg_len_pre)
print( "Avg. length of questions(Title+Body) after processing:␣
  ↪%d"%no_dup_avg_len_post)
print ("Percent of questions containing code: %d"%((questions_with_code*100.0)/
  ↪questions_proccesed))

print("Time taken to run this cell :", datetime.now() - start)
```

```
number of questions completed= 100000
number of questions completed= 200000
number of questions completed= 300000
number of questions completed= 400000
number of questions completed= 500000
number of questions completed= 600000
number of questions completed= 700000
number of questions completed= 800000
number of questions completed= 900000
Avg. length of questions(Title+Body) before processing: 1172
Avg. length of questions(Title+Body) after processing: 326
Percent of questions containing code: 57
Time taken to run this cell : 0:18:50.011295
```

[40]:
```
# dont forget to close the connections, or else you will end up with locks
conn_r.commit()
conn_w.commit()
conn_r.close()
conn_w.close()
```

[41]:
```
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        reader =conn_r.cursor()
        reader.execute("SELECT question From QuestionsProcessed LIMIT 10")
        print("Questions after preprocessed")
        print('='*100)
        reader.fetchone()
        for row in reader:
            print(row)
            print('-'*100)
conn_r.commit()
conn_r.close()
```

```
Questions after preprocessed
====================================================================================
====================
('core data imag desktop iphon built simpl mac data entri tool use iphon applic
```

recent ad thumbnail ad via imag well use simpl bind transform data type seem
work fine iphon applic howev show imag attribut null get imag appear follow
cellforrowatindexpath think either problem transform use default
nskeyedunarchivefromdata call thumbnail newbi help would great appreci',)
--------------------------------------------------------------------------------
--------------------
('navig call code visual studio want know use jump function call code return
code',)
--------------------------------------------------------------------------------
--------------------
('footer background extend bottom browser problem fix footer bottom browser
problem resolut chang window resiz footer content overlap content websit current
css footer div anybodi know fix thank updat need exact reason work web page work
cut past code blank page sinc page full content everyth import element includ
url trick work websit fill content let say line trick work updat ii websit dynam
content think use sort css sticki footer sometim websit line sometim pack
content that footer stick bottom webpag problem stick footer plenti content
websit problem without',)
--------------------------------------------------------------------------------
--------------------
('instal squeez suffer kernel bug product environ problem caus complet outag
degrad servic soft lockup like tri newer kernel howev squeez kernel org compil
sourc past realli necessari realli rather put admin manual track secur bug
backport seem longer kernel pita instal pull mountain depend linux base http
packag debian org squeez backport linux imag bpo pae depend could break abilti
boot back cut backout plan potenti lead long product outag refer onlin backport
newer kernel go thank',)
--------------------------------------------------------------------------------
--------------------
('mousemov effect anim found way make effect smoother enter contain idea list
icon number must relat center occur put way make softer posit without lose
mousemov plugin found web suggest would great http jsfiddl net rvalverd zw pj',)
--------------------------------------------------------------------------------
--------------------
('googl map caus font discrep blink text process elimin found googl map api side
effect web text problem appear effect safari platform idea screenshot',)
--------------------------------------------------------------------------------
--------------------
('access preload imag parent script use child script tri updat imag parent
window clickabl link child window preload imag parent window one javascript file
scriptss js nmi problem need access preload imag parent window childscript
scriptremot js thank js help nthe js scriptss js html parent window js child
window html child window',)
--------------------------------------------------------------------------------
--------------------
('entiti framework fluent map foreign key foreign object string key move edmx
map ef dbcontext fluent map want map string foreign key foreign object use
fluent api employe option offic would like officeid offic object employe class

read need abl save object object int key work fine tri sever string key get
result officeid field popul offic object come back null chekck sql profil data
queri offic object popul feedback ladislav map like onmodelcr assum subtleti
string key miss queri follow omit officeid field employe set map like offic
object popul need officeid field employe object',)
--------------------------------------------------------------------------------
--------------------
('spring autowir use object factori choos implement tri let piec runtim state
decid implement interfac use prefer sole autowir tri make object factori
interfac thet use dynam proxi use qualifi coerc autowir inject use factori
qualifi necessari factori implement respond interfac problem end annot everi
autowir refer qualifi realli want annot non factori implement someth like
notcandidateforautowiringbyinterfac fantasi annot even better make spring prefer
singl un qualifi bean inject un qualifi field may think along total wrong line
altern suggest welcom nanyon know make happen',)
--------------------------------------------------------------------------------
--------------------

```
[42]: #Taking 1 Million entries to a dataframe.
      write_db = 'Processed.db'
      if os.path.isfile(write_db):
          conn_r = create_connection(write_db)
          if conn_r is not None:
              preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM␣
      ↪QuestionsProcessed""", conn_r)
      conn_r.commit()
      conn_r.close()
```

```
[43]: preprocessed_data.head()
```

```
[43]:                                          question  \
      0  php opendir anoth server kinda new php got two...
      1  core data imag desktop iphon built simpl mac d...
      2  navig call code visual studio want know use ju...
      3  footer background extend bottom browser proble...
      4  instal squeez suffer kernel bug product enviro...

                                          tags
      0                    php directory opendir
      1  iphone core-data imageview nsmanagedobject
      2                            visual-studio
      3                                 html css
      4                  debian kernel backports
```

```
[44]: print("number of data points in sample :", preprocessed_data.shape[0])
      print("number of dimensions :", preprocessed_data.shape[1])
```

```
number of data points in sample : 999999
number of dimensions : 2
```

4. Machine Learning Models

4.5 Modeling with less data points (0.5M data points) and more weight to title and 500 tags only.

```
[45]: sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question␣
      ↪text NOT NULL, code text, tags text, words_pre integer, words_post integer,␣
      ↪is_code integer);"""
      create_database_table("Titlemoreweight.db", sql_create_table)
```

```
Tables in the databse:
QuestionsProcessed
```

```
[46]: # http://www.sqlitetutorial.net/sqlite-delete/
      # https://stackoverflow.com/questions/2279706/
       ↪select-random-row-from-a-sqlite-table

      read_db = 'train_no_dup.db'
      write_db = 'Titlemoreweight.db'
      train_datasize = 400000
      if os.path.isfile(read_db):
          conn_r = create_connection(read_db)
          if conn_r is not None:
              reader =conn_r.cursor()
              # for selecting first 0.5M rows
              reader.execute("SELECT Title, Body, Tags From no_dup_train LIMIT 500001;
       ↪")
              # for selecting random points
              #reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY␣
       ↪RANDOM() LIMIT 500001;")

      if os.path.isfile(write_db):
          conn_w = create_connection(write_db)
          if conn_w is not None:
              tables = checkTableExists(conn_w)
              writer =conn_w.cursor()
              if tables != 0:
                  writer.execute("DELETE FROM QuestionsProcessed WHERE 1")
                  print("Cleared All the rows")
```

```
Tables in the databse:
QuestionsProcessed
Cleared All the rows
```

4.5.1 Preprocessing of questions
Separate Code from Body
Remove Spcial characters from Question title and description (not in code)
Give more weightage to title : Add title three times to the question

24

Remove stop words (Except 'C')
Remove HTML Tags
Convert all the characters into small letters
Use SnowballStemmer to stem the words

```
[47]:  #http://www.bernzilla.com/2008/05/13/
        ↪selecting-a-random-row-from-an-sqlite-table/
       start = datetime.now()
       preprocessed_data_list=[]
       reader.fetchone()
       questions_with_code=0
       len_pre=0
       len_post=0
       questions_proccesed = 0
       for row in reader:

           is_code = 0

           title, question, tags = row[0], row[1], str(row[2])

           if '<code>' in question:
               questions_with_code+=1
               is_code = 1
           x = len(question)+len(title)
           len_pre+=x

           code = str(re.findall(r'<code>(.*?)</code>', question, flags=re.DOTALL))

           question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTILINE|re.
        ↪DOTALL)
           question=striphtml(question.encode('utf-8'))

           title=title.encode('utf-8')

           # adding title three time to the data to increase its weight
           # add tags string to the training data

           question=str(title)+" "+str(title)+" "+str(title)+" "+question

       #     if questions_proccesed<=train_datasize:
       #         question=str(title)+" "+str(title)+" "+str(title)+" "+question+"␣
        ↪"+str(tags)
       #     else:
       #         question=str(title)+" "+str(title)+" "+str(title)+" "+question

           question=re.sub(r'[^A-Za-z0-9#+.\-]+',' ',question)
           words=word_tokenize(str(question.lower()))
```

```
    #Removing all single letter and and stopwords from question exceptt for the
 ↪letter 'c'
    question=' '.join(str(stemmer.stem(j)) for j in words if j not in
 ↪stop_words and (len(j)!=1 or j=='c'))

    len_post+=len(question)
    tup = (question,code,tags,x,len(question),is_code)
    questions_proccesed += 1
    writer.execute("insert into
 ↪QuestionsProcessed(question,code,tags,words_pre,words_post,is_code) values (?
 ↪,?,?,?,?,?)",tup)
    if (questions_proccesed%100000==0):
        print("number of questions completed=",questions_proccesed)

no_dup_avg_len_pre=(len_pre*1.0)/questions_proccesed
no_dup_avg_len_post=(len_post*1.0)/questions_proccesed

print( "Avg. length of questions(Title+Body) before processing:
 ↪%d"%no_dup_avg_len_pre)
print( "Avg. length of questions(Title+Body) after processing:
 ↪%d"%no_dup_avg_len_post)
print ("Percent of questions containing code: %d"%((questions_with_code*100.0)/
 ↪questions_proccesed))

print("Time taken to run this cell :", datetime.now() - start)
```

```
number of questions completed= 100000
number of questions completed= 200000
number of questions completed= 300000
number of questions completed= 400000
number of questions completed= 500000
Avg. length of questions(Title+Body) before processing: 1239
Avg. length of questions(Title+Body) after processing: 424
Percent of questions containing code: 57
Time taken to run this cell : 0:12:36.056637
```

```
[48]: # never forget to close the conections or else we will end up with database
 ↪locks
conn_r.commit()
conn_w.commit()
conn_r.close()
conn_w.close()
```

__ Sample quesitons after preprocessing of data __

```
[49]: if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
```

```python
        reader =conn_r.cursor()
        reader.execute("SELECT question From QuestionsProcessed LIMIT 10")
        print("Questions after preprocessed")
        print('='*100)
        reader.fetchone()
        for row in reader:
            print(row)
            print('-'*100)
conn_r.commit()
conn_r.close()
```

Questions after preprocessed
====================================================================================================
('dynam datagrid bind silverlight dynam datagrid bind silverlight dynam datagrid
bind silverlight bind datagrid dynam code wrote code debug code block seem bind
correct grid come column form come grid column although necessari bind nthank
repli advance..',)
----------------------------------------------------------------------------------------------------
('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid
java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid
java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid follow
guid link instal jstl got follow error tri launch jsp page
java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid taglib
declar instal jstl 1.1 tomcat webapp tri project work also tri version 1.2 jstl
still messag caus solv',)
----------------------------------------------------------------------------------------------------
('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index
java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index
java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index use
follow code display caus solv',)
----------------------------------------------------------------------------------------------------
('better way updat feed fb php sdk better way updat feed fb php sdk better way
updat feed fb php sdk novic facebook api read mani tutori still confused.i find
post feed api method like correct second way use curl someth like way better',)
----------------------------------------------------------------------------------------------------
('btnadd click event open two window record ad btnadd click event open two
window record ad btnadd click event open two window record ad open window
search.aspx use code hav add button search.aspx nwhen insert record btnadd click
event open anoth window nafter insert record close window',)
----------------------------------------------------------------------------------------------------
('sql inject issu prevent correct form submiss php sql inject issu prevent

correct form submiss php sql inject issu prevent correct form submiss php check
everyth think make sure input field safe type sql inject good news safe bad news
one tag mess form submiss place even touch life figur exact html use templat
file forgiv okay entir php script get execut see data post none forum field post
problem use someth titl field none data get post current use print post see
submit noth work flawless statement though also mention script work flawless
local machin use host come across problem state list input test mess',)
--------------------------------------------------------------------------------
--------------------

('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl
subaddit lebesgu measur let lbrace rbrace sequenc set sigma -algebra mathcal
want show left bigcup right leq sum left right countabl addit measur defin set
sigma algebra mathcal think use monoton properti somewher proof start appreci
littl help nthank ad han answer make follow addit construct given han answer
clear bigcup bigcup cap emptyset neq left bigcup right left bigcup right sum
left right also construct subset monoton left right leq left right final would
sum leq sum result follow',)
--------------------------------------------------------------------------------
--------------------

('hql equival sql queri hql equival sql queri hql equival sql queri hql queri
replac name class properti name error occur hql error',)
--------------------------------------------------------------------------------
--------------------

('undefin symbol architectur i386 objc class skpsmtpmessag referenc error
undefin symbol architectur i386 objc class skpsmtpmessag referenc error undefin
symbol architectur i386 objc class skpsmtpmessag referenc error import framework
send email applic background import framework i.e skpsmtpmessag somebodi suggest
get error collect2 ld return exit status import framework correct sorc taken
framework follow mfmailcomposeviewcontrol question lock field updat answer drag
drop folder project click copi nthat',)
--------------------------------------------------------------------------------
--------------------

__ Saving Preprocessed data to a Database __

```
[50]: #Taking 0.5 Million entries to a dataframe.
      write_db = 'Titlemoreweight.db'
      if os.path.isfile(write_db):
          conn_r = create_connection(write_db)
          if conn_r is not None:
              preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM␣
       ↪QuestionsProcessed""", conn_r)
      conn_r.commit()
      conn_r.close()
```

```
[51]: preprocessed_data.head()
```

```
[51]:                                              question  \
      0  dynam datagrid bind silverlight dynam datagrid...
```

```
1   dynam datagrid bind silverlight dynam datagrid...
2   java.lang.noclassdeffounderror javax servlet j...
3   java.sql.sqlexcept microsoft odbc driver manag...
4   better way updat feed fb php sdk better way up...

                                        tags
0             c# silverlight data-binding
1   c# silverlight data-binding columns
2                                  jsp jstl
3                                 java jdbc
4         facebook api facebook-php-sdk
```

```python
[52]: print("number of data points in sample :", preprocessed_data.shape[0])
      print("number of dimensions :", preprocessed_data.shape[1])
```

```
number of data points in sample : 500000
number of dimensions : 2
```

__ Converting string Tags to multilable output variables __

```python
[53]: vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
      multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

__ Selecting 500 Tags __

```python
[54]: def tags_to_choose(n):
          t = multilabel_y.sum(axis=0).tolist()[0]
          sorted_tags_i = sorted(range(len(t)), key=lambda i: t[i], reverse=True)
          multilabel_yn=multilabel_y[:,sorted_tags_i[:n]]
          return multilabel_yn

      def questions_explained_fn(n):
          multilabel_yn = tags_to_choose(n)
          x= multilabel_yn.sum(axis=1)
          return (np.count_nonzero(x==0))
```

```python
[55]: questions_explained = []
      total_tags=multilabel_y.shape[1]
      total_qs=preprocessed_data.shape[0]
      for i in range(500, total_tags, 100):
          questions_explained.append(np.round(((total_qs-questions_explained_fn(i))/
       →total_qs)*100,3))
```

```python
[56]: fig, ax = plt.subplots()
      ax.plot(questions_explained)
      xlabel = list(500+np.array(range(-50,450,50))*50)
      ax.set_xticklabels(xlabel)
      plt.xlabel("Number of tags")
      plt.ylabel("Number Questions coverd partially")
      plt.grid()
      plt.show()
```

```
# you can choose any number of tags based on your computing power, minimun is␣
 ↪500(it covers 90% of the tags)
print("with ",5500,"tags we are covering ",questions_explained[50],"% of␣
 ↪questions")
print("with ",500,"tags we are covering ",questions_explained[0],"% of␣
 ↪questions")
```



```
with   5500 tags we are covering   99.157 % of questions
with   500 tags we are covering   90.956 % of questions
```

[57]:
```
# we will be taking 500 tags
multilabel_yx = tags_to_choose(500)
print("number of questions that are not covered :",␣
 ↪questions_explained_fn(500),"out of ", total_qs)
```

```
number of questions that are not covered : 45221 out of   500000
```

[58]:
```
x_train=preprocessed_data.head(train_datasize)
x_test=preprocessed_data.tail(preprocessed_data.shape[0] - 400000)

y_train = multilabel_yx[0:train_datasize,:]
y_test = multilabel_yx[train_datasize:preprocessed_data.shape[0],:]
```

```
[59]: print("Number of data points in train data :", y_train.shape)
      print("Number of data points in test data :", y_test.shape)
```

```
Number of data points in train data : (400000, 500)
Number of data points in test data : (100000, 500)
```

### 4.5.2 Featurizing data with TfIdf vectorizer

```
[60]: start = datetime.now()
      vectorizer = TfidfVectorizer(min_df=0.00009, max_features=200000,␣
        ↪smooth_idf=True, norm="l2", \
                                     tokenizer = lambda x: x.split(),␣
        ↪sublinear_tf=False, ngram_range=(1,3))
      x_train_multilabel = vectorizer.fit_transform(x_train['question'])
      x_test_multilabel = vectorizer.transform(x_test['question'])
      print("Time taken to run this cell :", datetime.now() - start)
```

```
Time taken to run this cell : 0:03:07.249178
```

```
[61]: start = datetime.now()
      vectorizer_BOW = CountVectorizer(max_features=200000,ngram_range=(1,4))
      x_train_multilabel_BOW = vectorizer_BOW.fit_transform(x_train['question'])
      x_test_multilabel_BOW = vectorizer_BOW.transform(x_test['question'])
      print("Time taken to run this cell :", datetime.now() - start)
```

```
Time taken to run this cell : 0:05:34.109731
```

```
[62]: print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.
        ↪shape)
      print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

```
Dimensions of train data X: (400000, 94927) Y : (400000, 500)
Dimensions of test data X: (100000, 94927) Y: (100000, 500)
```

```
[63]: print("Dimensions of train data X:",x_train_multilabel_BOW.shape, "Y :",y_train.
        ↪shape)
      print("Dimensions of test data X:",x_test_multilabel_BOW.shape,"Y:",y_test.
        ↪shape)
```

```
Dimensions of train data X: (400000, 200000) Y : (400000, 500)
Dimensions of test data X: (100000, 200000) Y: (100000, 500)
```

### 4.5.3 Applying Logistic Regression with OneVsRest Classifier

```
[65]: start = datetime.now()
      classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001,␣
        ↪penalty='l1'), n_jobs=1)
```

```
classifier.fit(x_train_multilabel_BOW, y_train)
predictions = classifier.predict (x_test_multilabel_BOW)


print("Accuracy :",metrics.accuracy_score(y_test, predictions))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions))


precision = precision_score(y_test, predictions, average='micro')
recall = recall_score(y_test, predictions, average='micro')
f1 = f1_score(y_test, predictions, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,␣
 ↪recall, f1))

precision = precision_score(y_test, predictions, average='macro')
recall = recall_score(y_test, predictions, average='macro')
f1 = f1_score(y_test, predictions, average='macro')

print("Macro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,␣
 ↪recall, f1))

print (metrics.classification_report(y_test, predictions))
print("Time taken to run this cell :", datetime.now() - start)
```

```
Accuracy : 0.0995
Hamming loss  0.00566152
Micro-average quality numbers
Precision: 0.2997, Recall: 0.4702, F1-measure: 0.3661
Macro-average quality numbers
Precision: 0.2200, Recall: 0.4098, F1-measure: 0.2785
           precision    recall  f1-score   support

        0       0.78      0.78      0.78      5519
        1       0.33      0.51      0.40      8190
        2       0.54      0.48      0.51      6529
        3       0.36      0.52      0.43      3231
        4       0.58      0.51      0.54      6430
        5       0.45      0.47      0.46      2879
        6       0.62      0.58      0.60      5086
        7       0.64      0.64      0.64      4533
        8       0.23      0.23      0.23      3000
        9       0.57      0.64      0.60      2765
       10       0.32      0.30      0.31      3051
       11       0.46      0.47      0.46      3009
```

| | | | | |
|----|------|------|------|------|
| 12 | 0.40 | 0.42 | 0.41 | 2630 |
| 13 | 0.26 | 0.37 | 0.30 | 1426 |
| 14 | 0.67 | 0.69 | 0.68 | 2548 |
| 15 | 0.38 | 0.36 | 0.37 | 2371 |
| 16 | 0.28 | 0.38 | 0.32 | 873 |
| 17 | 0.59 | 0.70 | 0.64 | 2151 |
| 18 | 0.30 | 0.34 | 0.32 | 2204 |
| 19 | 0.29 | 0.50 | 0.37 | 831 |
| 20 | 0.54 | 0.56 | 0.55 | 1860 |
| 21 | 0.20 | 0.23 | 0.22 | 2023 |
| 22 | 0.29 | 0.36 | 0.32 | 1513 |
| 23 | 0.55 | 0.67 | 0.61 | 1207 |
| 24 | 0.28 | 0.39 | 0.32 | 506 |
| 25 | 0.28 | 0.45 | 0.34 | 425 |
| 26 | 0.36 | 0.49 | 0.42 | 793 |
| 27 | 0.37 | 0.47 | 0.41 | 1291 |
| 28 | 0.42 | 0.51 | 0.46 | 1208 |
| 29 | 0.12 | 0.22 | 0.15 | 406 |
| 30 | 0.22 | 0.33 | 0.26 | 504 |
| 31 | 0.13 | 0.20 | 0.16 | 732 |
| 32 | 0.23 | 0.41 | 0.30 | 441 |
| 33 | 0.34 | 0.42 | 0.38 | 1645 |
| 34 | 0.28 | 0.35 | 0.31 | 1058 |
| 35 | 0.49 | 0.64 | 0.56 | 946 |
| 36 | 0.25 | 0.39 | 0.30 | 644 |
| 37 | 0.27 | 0.75 | 0.40 | 136 |
| 38 | 0.30 | 0.51 | 0.38 | 570 |
| 39 | 0.33 | 0.40 | 0.36 | 766 |
| 40 | 0.37 | 0.45 | 0.41 | 1132 |
| 41 | 0.11 | 0.29 | 0.16 | 174 |
| 42 | 0.31 | 0.60 | 0.41 | 210 |
| 43 | 0.36 | 0.52 | 0.43 | 433 |
| 44 | 0.36 | 0.54 | 0.43 | 626 |
| 45 | 0.33 | 0.44 | 0.38 | 852 |
| 46 | 0.37 | 0.55 | 0.44 | 534 |
| 47 | 0.16 | 0.32 | 0.21 | 350 |
| 48 | 0.38 | 0.57 | 0.46 | 496 |
| 49 | 0.56 | 0.68 | 0.61 | 785 |
| 50 | 0.13 | 0.23 | 0.16 | 475 |
| 51 | 0.12 | 0.25 | 0.16 | 305 |
| 52 | 0.10 | 0.17 | 0.13 | 251 |
| 53 | 0.34 | 0.46 | 0.39 | 914 |
| 54 | 0.22 | 0.29 | 0.25 | 728 |
| 55 | 0.05 | 0.09 | 0.06 | 258 |
| 56 | 0.22 | 0.36 | 0.27 | 821 |
| 57 | 0.17 | 0.28 | 0.21 | 541 |
| 58 | 0.29 | 0.43 | 0.35 | 748 |
| 59 | 0.62 | 0.80 | 0.70 | 724 |

| | | | | |
|---|---|---|---|---|
| 60 | 0.17 | 0.26 | 0.20 | 660 |
| 61 | 0.16 | 0.32 | 0.21 | 235 |
| 62 | 0.62 | 0.80 | 0.70 | 718 |
| 63 | 0.53 | 0.72 | 0.61 | 468 |
| 64 | 0.21 | 0.50 | 0.29 | 191 |
| 65 | 0.14 | 0.25 | 0.18 | 429 |
| 66 | 0.13 | 0.23 | 0.17 | 415 |
| 67 | 0.28 | 0.58 | 0.38 | 274 |
| 68 | 0.40 | 0.59 | 0.48 | 510 |
| 69 | 0.32 | 0.51 | 0.40 | 466 |
| 70 | 0.13 | 0.24 | 0.17 | 305 |
| 71 | 0.12 | 0.30 | 0.17 | 247 |
| 72 | 0.38 | 0.58 | 0.46 | 401 |
| 73 | 0.35 | 0.86 | 0.50 | 86 |
| 74 | 0.21 | 0.49 | 0.30 | 120 |
| 75 | 0.26 | 0.71 | 0.38 | 129 |
| 76 | 0.09 | 0.12 | 0.11 | 473 |
| 77 | 0.10 | 0.36 | 0.16 | 143 |
| 78 | 0.39 | 0.60 | 0.47 | 347 |
| 79 | 0.23 | 0.36 | 0.28 | 479 |
| 80 | 0.21 | 0.49 | 0.29 | 279 |
| 81 | 0.26 | 0.43 | 0.32 | 461 |
| 82 | 0.08 | 0.17 | 0.11 | 298 |
| 83 | 0.36 | 0.63 | 0.46 | 396 |
| 84 | 0.21 | 0.50 | 0.30 | 184 |
| 85 | 0.27 | 0.38 | 0.31 | 573 |
| 86 | 0.12 | 0.19 | 0.15 | 325 |
| 87 | 0.18 | 0.42 | 0.25 | 273 |
| 88 | 0.10 | 0.29 | 0.15 | 135 |
| 89 | 0.12 | 0.27 | 0.16 | 232 |
| 90 | 0.31 | 0.52 | 0.39 | 409 |
| 91 | 0.25 | 0.40 | 0.31 | 420 |
| 92 | 0.41 | 0.62 | 0.49 | 408 |
| 93 | 0.23 | 0.57 | 0.33 | 241 |
| 94 | 0.07 | 0.16 | 0.10 | 211 |
| 95 | 0.15 | 0.25 | 0.19 | 277 |
| 96 | 0.15 | 0.22 | 0.18 | 410 |
| 97 | 0.48 | 0.63 | 0.55 | 501 |
| 98 | 0.27 | 0.68 | 0.38 | 136 |
| 99 | 0.22 | 0.47 | 0.30 | 239 |
| 100 | 0.14 | 0.28 | 0.19 | 324 |
| 101 | 0.53 | 0.80 | 0.63 | 277 |
| 102 | 0.64 | 0.82 | 0.72 | 613 |
| 103 | 0.13 | 0.31 | 0.18 | 157 |
| 104 | 0.12 | 0.22 | 0.15 | 295 |
| 105 | 0.33 | 0.53 | 0.41 | 334 |
| 106 | 0.30 | 0.43 | 0.35 | 335 |
| 107 | 0.36 | 0.59 | 0.44 | 389 |

| | | | | |
|---|---|---|---|---|
| 108 | 0.18 | 0.39 | 0.25 | 251 |
| 109 | 0.32 | 0.50 | 0.39 | 317 |
| 110 | 0.06 | 0.18 | 0.09 | 187 |
| 111 | 0.10 | 0.29 | 0.15 | 140 |
| 112 | 0.24 | 0.60 | 0.34 | 154 |
| 113 | 0.23 | 0.35 | 0.27 | 332 |
| 114 | 0.20 | 0.39 | 0.27 | 323 |
| 115 | 0.20 | 0.42 | 0.27 | 344 |
| 116 | 0.37 | 0.57 | 0.44 | 370 |
| 117 | 0.22 | 0.36 | 0.28 | 313 |
| 118 | 0.61 | 0.77 | 0.68 | 874 |
| 119 | 0.15 | 0.33 | 0.21 | 293 |
| 120 | 0.04 | 0.14 | 0.07 | 200 |
| 121 | 0.42 | 0.65 | 0.51 | 463 |
| 122 | 0.11 | 0.25 | 0.15 | 119 |
| 123 | 0.01 | 0.02 | 0.01 | 256 |
| 124 | 0.41 | 0.77 | 0.54 | 195 |
| 125 | 0.12 | 0.35 | 0.18 | 138 |
| 126 | 0.39 | 0.62 | 0.48 | 376 |
| 127 | 0.03 | 0.09 | 0.04 | 122 |
| 128 | 0.07 | 0.13 | 0.09 | 252 |
| 129 | 0.20 | 0.42 | 0.27 | 144 |
| 130 | 0.09 | 0.24 | 0.13 | 150 |
| 131 | 0.06 | 0.11 | 0.08 | 210 |
| 132 | 0.20 | 0.37 | 0.26 | 361 |
| 133 | 0.63 | 0.69 | 0.66 | 453 |
| 134 | 0.34 | 0.78 | 0.48 | 124 |
| 135 | 0.05 | 0.16 | 0.08 | 91 |
| 136 | 0.13 | 0.44 | 0.20 | 128 |
| 137 | 0.22 | 0.47 | 0.30 | 218 |
| 138 | 0.12 | 0.26 | 0.17 | 243 |
| 139 | 0.12 | 0.31 | 0.17 | 149 |
| 140 | 0.40 | 0.56 | 0.47 | 318 |
| 141 | 0.07 | 0.20 | 0.10 | 159 |
| 142 | 0.33 | 0.51 | 0.40 | 274 |
| 143 | 0.56 | 0.85 | 0.67 | 362 |
| 144 | 0.11 | 0.36 | 0.17 | 118 |
| 145 | 0.18 | 0.49 | 0.26 | 164 |
| 146 | 0.26 | 0.46 | 0.33 | 461 |
| 147 | 0.26 | 0.53 | 0.35 | 159 |
| 148 | 0.12 | 0.25 | 0.16 | 166 |
| 149 | 0.55 | 0.61 | 0.58 | 346 |
| 150 | 0.15 | 0.28 | 0.19 | 350 |
| 151 | 0.32 | 0.82 | 0.46 | 55 |
| 152 | 0.42 | 0.58 | 0.49 | 387 |
| 153 | 0.19 | 0.33 | 0.24 | 150 |
| 154 | 0.10 | 0.19 | 0.13 | 281 |
| 155 | 0.10 | 0.26 | 0.14 | 202 |

| | | | | |
|-----|------|------|------|-----|
| 156 | 0.30 | 0.69 | 0.41 | 130 |
| 157 | 0.12 | 0.19 | 0.15 | 245 |
| 158 | 0.44 | 0.76 | 0.56 | 177 |
| 159 | 0.15 | 0.52 | 0.24 | 130 |
| 160 | 0.16 | 0.31 | 0.22 | 336 |
| 161 | 0.40 | 0.75 | 0.52 | 220 |
| 162 | 0.12 | 0.31 | 0.17 | 229 |
| 163 | 0.45 | 0.61 | 0.52 | 316 |
| 164 | 0.29 | 0.55 | 0.38 | 283 |
| 165 | 0.20 | 0.40 | 0.27 | 197 |
| 166 | 0.24 | 0.57 | 0.34 | 101 |
| 167 | 0.13 | 0.26 | 0.17 | 231 |
| 168 | 0.21 | 0.40 | 0.27 | 370 |
| 169 | 0.19 | 0.33 | 0.24 | 258 |
| 170 | 0.07 | 0.29 | 0.12 | 101 |
| 171 | 0.09 | 0.30 | 0.14 | 89 |
| 172 | 0.18 | 0.42 | 0.25 | 193 |
| 173 | 0.25 | 0.43 | 0.32 | 309 |
| 174 | 0.09 | 0.28 | 0.14 | 172 |
| 175 | 0.32 | 0.78 | 0.45 | 95 |
| 176 | 0.54 | 0.71 | 0.61 | 346 |
| 177 | 0.36 | 0.62 | 0.46 | 322 |
| 178 | 0.27 | 0.56 | 0.36 | 232 |
| 179 | 0.08 | 0.18 | 0.11 | 125 |
| 180 | 0.19 | 0.46 | 0.27 | 145 |
| 181 | 0.05 | 0.25 | 0.08 | 77 |
| 182 | 0.11 | 0.26 | 0.15 | 182 |
| 183 | 0.28 | 0.47 | 0.35 | 257 |
| 184 | 0.09 | 0.20 | 0.12 | 216 |
| 185 | 0.14 | 0.30 | 0.19 | 242 |
| 186 | 0.16 | 0.34 | 0.22 | 165 |
| 187 | 0.34 | 0.62 | 0.44 | 263 |
| 188 | 0.11 | 0.25 | 0.15 | 174 |
| 189 | 0.32 | 0.44 | 0.37 | 136 |
| 190 | 0.46 | 0.70 | 0.56 | 202 |
| 191 | 0.10 | 0.27 | 0.15 | 134 |
| 192 | 0.26 | 0.51 | 0.34 | 230 |
| 193 | 0.09 | 0.30 | 0.13 | 90 |
| 194 | 0.27 | 0.57 | 0.37 | 185 |
| 195 | 0.06 | 0.15 | 0.09 | 156 |
| 196 | 0.07 | 0.19 | 0.10 | 160 |
| 197 | 0.13 | 0.24 | 0.17 | 266 |
| 198 | 0.15 | 0.23 | 0.18 | 284 |
| 199 | 0.07 | 0.16 | 0.09 | 145 |
| 200 | 0.50 | 0.83 | 0.62 | 212 |
| 201 | 0.24 | 0.42 | 0.30 | 317 |
| 202 | 0.47 | 0.67 | 0.55 | 427 |
| 203 | 0.13 | 0.29 | 0.18 | 232 |

| | | | | |
|---|---|---|---|---|
| 204 | 0.17 | 0.37 | 0.24 | 217 |
| 205 | 0.39 | 0.56 | 0.46 | 527 |
| 206 | 0.06 | 0.15 | 0.08 | 124 |
| 207 | 0.20 | 0.46 | 0.28 | 103 |
| 208 | 0.32 | 0.56 | 0.41 | 287 |
| 209 | 0.08 | 0.19 | 0.11 | 193 |
| 210 | 0.22 | 0.45 | 0.30 | 220 |
| 211 | 0.14 | 0.35 | 0.20 | 140 |
| 212 | 0.07 | 0.21 | 0.11 | 161 |
| 213 | 0.20 | 0.58 | 0.29 | 72 |
| 214 | 0.44 | 0.55 | 0.49 | 396 |
| 215 | 0.18 | 0.53 | 0.27 | 134 |
| 216 | 0.20 | 0.33 | 0.25 | 400 |
| 217 | 0.12 | 0.41 | 0.19 | 75 |
| 218 | 0.59 | 0.81 | 0.68 | 219 |
| 219 | 0.22 | 0.47 | 0.30 | 210 |
| 220 | 0.54 | 0.77 | 0.63 | 298 |
| 221 | 0.61 | 0.78 | 0.69 | 266 |
| 222 | 0.39 | 0.59 | 0.47 | 290 |
| 223 | 0.04 | 0.13 | 0.07 | 128 |
| 224 | 0.21 | 0.52 | 0.30 | 159 |
| 225 | 0.14 | 0.44 | 0.22 | 164 |
| 226 | 0.21 | 0.46 | 0.29 | 144 |
| 227 | 0.34 | 0.57 | 0.42 | 276 |
| 228 | 0.06 | 0.15 | 0.09 | 235 |
| 229 | 0.04 | 0.10 | 0.06 | 216 |
| 230 | 0.10 | 0.29 | 0.15 | 228 |
| 231 | 0.21 | 0.59 | 0.31 | 64 |
| 232 | 0.04 | 0.13 | 0.06 | 103 |
| 233 | 0.27 | 0.53 | 0.36 | 216 |
| 234 | 0.14 | 0.33 | 0.19 | 116 |
| 235 | 0.17 | 0.51 | 0.26 | 77 |
| 236 | 0.44 | 0.70 | 0.54 | 67 |
| 237 | 0.10 | 0.23 | 0.14 | 218 |
| 238 | 0.09 | 0.28 | 0.13 | 139 |
| 239 | 0.04 | 0.10 | 0.06 | 94 |
| 240 | 0.11 | 0.32 | 0.16 | 77 |
| 241 | 0.05 | 0.15 | 0.08 | 167 |
| 242 | 0.26 | 0.52 | 0.35 | 86 |
| 243 | 0.08 | 0.33 | 0.13 | 58 |
| 244 | 0.28 | 0.43 | 0.34 | 269 |
| 245 | 0.07 | 0.18 | 0.10 | 112 |
| 246 | 0.66 | 0.84 | 0.74 | 255 |
| 247 | 0.08 | 0.34 | 0.13 | 58 |
| 248 | 0.02 | 0.12 | 0.04 | 81 |
| 249 | 0.07 | 0.15 | 0.09 | 131 |
| 250 | 0.11 | 0.37 | 0.17 | 93 |
| 251 | 0.18 | 0.47 | 0.26 | 154 |

| | | | | |
|-----|------|------|------|-----|
| 252 | 0.04 | 0.12 | 0.06 | 129 |
| 253 | 0.13 | 0.43 | 0.20 | 83 |
| 254 | 0.12 | 0.23 | 0.16 | 191 |
| 255 | 0.10 | 0.16 | 0.12 | 219 |
| 256 | 0.05 | 0.16 | 0.08 | 130 |
| 257 | 0.11 | 0.35 | 0.17 | 93 |
| 258 | 0.32 | 0.58 | 0.41 | 217 |
| 259 | 0.09 | 0.27 | 0.14 | 141 |
| 260 | 0.21 | 0.43 | 0.28 | 143 |
| 261 | 0.14 | 0.25 | 0.18 | 219 |
| 262 | 0.18 | 0.52 | 0.27 | 107 |
| 263 | 0.24 | 0.39 | 0.30 | 236 |
| 264 | 0.10 | 0.33 | 0.15 | 119 |
| 265 | 0.13 | 0.46 | 0.21 | 72 |
| 266 | 0.07 | 0.23 | 0.10 | 70 |
| 267 | 0.13 | 0.34 | 0.19 | 107 |
| 268 | 0.24 | 0.56 | 0.34 | 169 |
| 269 | 0.16 | 0.36 | 0.22 | 129 |
| 270 | 0.36 | 0.67 | 0.47 | 159 |
| 271 | 0.31 | 0.64 | 0.42 | 190 |
| 272 | 0.19 | 0.41 | 0.26 | 248 |
| 273 | 0.61 | 0.78 | 0.69 | 264 |
| 274 | 0.48 | 0.81 | 0.60 | 105 |
| 275 | 0.07 | 0.24 | 0.11 | 104 |
| 276 | 0.04 | 0.12 | 0.06 | 115 |
| 277 | 0.36 | 0.71 | 0.47 | 170 |
| 278 | 0.33 | 0.59 | 0.42 | 145 |
| 279 | 0.52 | 0.80 | 0.63 | 230 |
| 280 | 0.15 | 0.40 | 0.21 | 80 |
| 281 | 0.41 | 0.64 | 0.50 | 217 |
| 282 | 0.31 | 0.63 | 0.42 | 175 |
| 283 | 0.15 | 0.32 | 0.21 | 269 |
| 284 | 0.19 | 0.51 | 0.28 | 74 |
| 285 | 0.31 | 0.64 | 0.42 | 206 |
| 286 | 0.50 | 0.74 | 0.60 | 227 |
| 287 | 0.20 | 0.62 | 0.30 | 130 |
| 288 | 0.06 | 0.16 | 0.09 | 129 |
| 289 | 0.04 | 0.21 | 0.07 | 80 |
| 290 | 0.06 | 0.23 | 0.10 | 99 |
| 291 | 0.26 | 0.48 | 0.33 | 208 |
| 292 | 0.03 | 0.18 | 0.05 | 67 |
| 293 | 0.26 | 0.61 | 0.37 | 109 |
| 294 | 0.13 | 0.39 | 0.19 | 140 |
| 295 | 0.12 | 0.26 | 0.17 | 241 |
| 296 | 0.07 | 0.24 | 0.11 | 72 |
| 297 | 0.09 | 0.30 | 0.14 | 107 |
| 298 | 0.31 | 0.57 | 0.40 | 61 |
| 299 | 0.36 | 0.66 | 0.46 | 77 |

| | | | | |
|---|---|---|---|---|
| 300 | 0.07 | 0.20 | 0.10 | 111 |
| 301 | 0.01 | 0.02 | 0.01 | 126 |
| 302 | 0.05 | 0.15 | 0.07 | 73 |
| 303 | 0.26 | 0.45 | 0.33 | 176 |
| 304 | 0.77 | 0.86 | 0.81 | 230 |
| 305 | 0.54 | 0.83 | 0.65 | 156 |
| 306 | 0.23 | 0.45 | 0.30 | 146 |
| 307 | 0.08 | 0.22 | 0.12 | 98 |
| 308 | 0.01 | 0.05 | 0.02 | 78 |
| 309 | 0.08 | 0.23 | 0.12 | 94 |
| 310 | 0.21 | 0.45 | 0.29 | 162 |
| 311 | 0.27 | 0.65 | 0.38 | 116 |
| 312 | 0.12 | 0.40 | 0.18 | 57 |
| 313 | 0.07 | 0.14 | 0.09 | 65 |
| 314 | 0.20 | 0.41 | 0.27 | 138 |
| 315 | 0.21 | 0.44 | 0.29 | 195 |
| 316 | 0.11 | 0.41 | 0.17 | 69 |
| 317 | 0.11 | 0.29 | 0.16 | 134 |
| 318 | 0.22 | 0.47 | 0.30 | 148 |
| 319 | 0.44 | 0.60 | 0.51 | 161 |
| 320 | 0.07 | 0.25 | 0.11 | 104 |
| 321 | 0.32 | 0.60 | 0.42 | 156 |
| 322 | 0.20 | 0.48 | 0.28 | 134 |
| 323 | 0.33 | 0.52 | 0.40 | 232 |
| 324 | 0.07 | 0.21 | 0.11 | 92 |
| 325 | 0.18 | 0.39 | 0.25 | 197 |
| 326 | 0.07 | 0.21 | 0.11 | 126 |
| 327 | 0.02 | 0.05 | 0.03 | 115 |
| 328 | 0.64 | 0.74 | 0.69 | 198 |
| 329 | 0.16 | 0.42 | 0.23 | 125 |
| 330 | 0.17 | 0.37 | 0.23 | 81 |
| 331 | 0.08 | 0.17 | 0.11 | 94 |
| 332 | 0.10 | 0.25 | 0.14 | 56 |
| 333 | 0.08 | 0.18 | 0.11 | 260 |
| 334 | 0.08 | 0.32 | 0.13 | 60 |
| 335 | 0.14 | 0.29 | 0.18 | 110 |
| 336 | 0.21 | 0.54 | 0.30 | 71 |
| 337 | 0.04 | 0.14 | 0.06 | 66 |
| 338 | 0.20 | 0.48 | 0.28 | 150 |
| 339 | 0.01 | 0.06 | 0.02 | 54 |
| 340 | 0.43 | 0.63 | 0.51 | 195 |
| 341 | 0.38 | 0.62 | 0.47 | 79 |
| 342 | 0.10 | 0.45 | 0.16 | 38 |
| 343 | 0.13 | 0.44 | 0.20 | 43 |
| 344 | 0.21 | 0.34 | 0.26 | 68 |
| 345 | 0.24 | 0.53 | 0.33 | 73 |
| 346 | 0.06 | 0.19 | 0.09 | 116 |
| 347 | 0.23 | 0.61 | 0.33 | 111 |

| | | | | |
|---|---|---|---|---|
| 348 | 0.04 | 0.19 | 0.07 | 63 |
| 349 | 0.37 | 0.73 | 0.49 | 104 |
| 350 | 0.19 | 0.48 | 0.27 | 44 |
| 351 | 0.10 | 0.28 | 0.15 | 40 |
| 352 | 0.36 | 0.62 | 0.46 | 136 |
| 353 | 0.09 | 0.30 | 0.14 | 54 |
| 354 | 0.09 | 0.24 | 0.13 | 134 |
| 355 | 0.21 | 0.53 | 0.30 | 120 |
| 356 | 0.26 | 0.44 | 0.33 | 228 |
| 357 | 0.36 | 0.54 | 0.43 | 269 |
| 358 | 0.19 | 0.46 | 0.27 | 80 |
| 359 | 0.34 | 0.69 | 0.45 | 140 |
| 360 | 0.11 | 0.26 | 0.15 | 125 |
| 361 | 0.56 | 0.80 | 0.66 | 169 |
| 362 | 0.04 | 0.16 | 0.07 | 56 |
| 363 | 0.53 | 0.80 | 0.63 | 154 |
| 364 | 0.10 | 0.28 | 0.15 | 58 |
| 365 | 0.14 | 0.32 | 0.20 | 71 |
| 366 | 0.61 | 0.74 | 0.67 | 54 |
| 367 | 0.05 | 0.16 | 0.07 | 116 |
| 368 | 0.05 | 0.19 | 0.08 | 54 |
| 369 | 0.05 | 0.18 | 0.08 | 71 |
| 370 | 0.02 | 0.10 | 0.03 | 61 |
| 371 | 0.06 | 0.21 | 0.09 | 71 |
| 372 | 0.22 | 0.56 | 0.31 | 52 |
| 373 | 0.40 | 0.60 | 0.48 | 150 |
| 374 | 0.09 | 0.32 | 0.14 | 93 |
| 375 | 0.06 | 0.21 | 0.09 | 67 |
| 376 | 0.02 | 0.05 | 0.02 | 76 |
| 377 | 0.17 | 0.37 | 0.23 | 106 |
| 378 | 0.02 | 0.06 | 0.03 | 86 |
| 379 | 0.01 | 0.14 | 0.02 | 14 |
| 380 | 0.37 | 0.62 | 0.46 | 122 |
| 381 | 0.03 | 0.10 | 0.05 | 104 |
| 382 | 0.05 | 0.23 | 0.08 | 66 |
| 383 | 0.21 | 0.45 | 0.29 | 110 |
| 384 | 0.04 | 0.09 | 0.06 | 155 |
| 385 | 0.09 | 0.38 | 0.14 | 50 |
| 386 | 0.10 | 0.30 | 0.15 | 64 |
| 387 | 0.12 | 0.24 | 0.16 | 93 |
| 388 | 0.19 | 0.50 | 0.27 | 102 |
| 389 | 0.05 | 0.13 | 0.07 | 108 |
| 390 | 0.67 | 0.77 | 0.72 | 178 |
| 391 | 0.13 | 0.27 | 0.18 | 115 |
| 392 | 0.22 | 0.57 | 0.32 | 42 |
| 393 | 0.02 | 0.04 | 0.03 | 134 |
| 394 | 0.07 | 0.18 | 0.10 | 112 |
| 395 | 0.18 | 0.40 | 0.25 | 176 |

| | | | | |
|-----|------|------|------|-----|
| 396 | 0.10 | 0.30 | 0.15 | 125 |
| 397 | 0.40 | 0.58 | 0.48 | 224 |
| 398 | 0.26 | 0.70 | 0.37 | 63 |
| 399 | 0.01 | 0.05 | 0.02 | 59 |
| 400 | 0.15 | 0.46 | 0.23 | 63 |
| 401 | 0.11 | 0.43 | 0.17 | 98 |
| 402 | 0.18 | 0.35 | 0.24 | 162 |
| 403 | 0.09 | 0.24 | 0.13 | 83 |
| 404 | 0.31 | 0.89 | 0.47 | 19 |
| 405 | 0.07 | 0.28 | 0.11 | 92 |
| 406 | 0.08 | 0.51 | 0.13 | 41 |
| 407 | 0.27 | 0.56 | 0.37 | 43 |
| 408 | 0.25 | 0.50 | 0.33 | 160 |
| 409 | 0.08 | 0.22 | 0.12 | 50 |
| 410 | 0.00 | 0.05 | 0.01 | 19 |
| 411 | 0.12 | 0.23 | 0.15 | 175 |
| 412 | 0.05 | 0.15 | 0.08 | 72 |
| 413 | 0.06 | 0.17 | 0.09 | 95 |
| 414 | 0.10 | 0.29 | 0.15 | 97 |
| 415 | 0.05 | 0.19 | 0.08 | 48 |
| 416 | 0.19 | 0.40 | 0.25 | 83 |
| 417 | 0.04 | 0.12 | 0.06 | 40 |
| 418 | 0.11 | 0.34 | 0.16 | 91 |
| 419 | 0.22 | 0.47 | 0.30 | 90 |
| 420 | 0.07 | 0.30 | 0.12 | 37 |
| 421 | 0.08 | 0.24 | 0.12 | 66 |
| 422 | 0.12 | 0.45 | 0.20 | 73 |
| 423 | 0.10 | 0.32 | 0.15 | 56 |
| 424 | 0.48 | 0.91 | 0.62 | 33 |
| 425 | 0.02 | 0.05 | 0.03 | 76 |
| 426 | 0.04 | 0.14 | 0.06 | 81 |
| 427 | 0.61 | 0.79 | 0.69 | 150 |
| 428 | 0.28 | 0.79 | 0.41 | 29 |
| 429 | 0.95 | 0.95 | 0.95 | 389 |
| 430 | 0.31 | 0.49 | 0.38 | 167 |
| 431 | 0.10 | 0.20 | 0.13 | 123 |
| 432 | 0.12 | 0.44 | 0.18 | 39 |
| 433 | 0.18 | 0.52 | 0.27 | 82 |
| 434 | 0.36 | 0.73 | 0.48 | 66 |
| 435 | 0.18 | 0.41 | 0.25 | 93 |
| 436 | 0.22 | 0.52 | 0.31 | 87 |
| 437 | 0.06 | 0.15 | 0.08 | 86 |
| 438 | 0.36 | 0.62 | 0.46 | 104 |
| 439 | 0.08 | 0.26 | 0.13 | 100 |
| 440 | 0.06 | 0.13 | 0.09 | 141 |
| 441 | 0.23 | 0.55 | 0.33 | 110 |
| 442 | 0.09 | 0.23 | 0.12 | 123 |
| 443 | 0.15 | 0.41 | 0.22 | 71 |

| | | | |
|---|---|---|---|
| 444 | 0.13 | 0.27 | 0.18 | 109 |
| 445 | 0.09 | 0.42 | 0.15 | 48 |
| 446 | 0.15 | 0.49 | 0.23 | 76 |
| 447 | 0.05 | 0.29 | 0.08 | 38 |
| 448 | 0.27 | 0.67 | 0.38 | 81 |
| 449 | 0.24 | 0.39 | 0.30 | 132 |
| 450 | 0.20 | 0.44 | 0.27 | 81 |
| 451 | 0.19 | 0.43 | 0.26 | 76 |
| 452 | 0.05 | 0.11 | 0.07 | 44 |
| 453 | 0.04 | 0.14 | 0.06 | 44 |
| 454 | 0.18 | 0.54 | 0.27 | 70 |
| 455 | 0.09 | 0.32 | 0.14 | 155 |
| 456 | 0.10 | 0.30 | 0.15 | 43 |
| 457 | 0.14 | 0.47 | 0.22 | 72 |
| 458 | 0.05 | 0.24 | 0.09 | 62 |
| 459 | 0.10 | 0.33 | 0.16 | 69 |
| 460 | 0.03 | 0.08 | 0.05 | 119 |
| 461 | 0.33 | 0.43 | 0.38 | 79 |
| 462 | 0.09 | 0.32 | 0.14 | 47 |
| 463 | 0.17 | 0.39 | 0.23 | 104 |
| 464 | 0.22 | 0.44 | 0.30 | 106 |
| 465 | 0.09 | 0.30 | 0.13 | 64 |
| 466 | 0.26 | 0.45 | 0.33 | 173 |
| 467 | 0.23 | 0.45 | 0.30 | 107 |
| 468 | 0.17 | 0.41 | 0.24 | 126 |
| 469 | 0.02 | 0.04 | 0.03 | 114 |
| 470 | 0.67 | 0.87 | 0.76 | 140 |
| 471 | 0.25 | 0.49 | 0.33 | 79 |
| 472 | 0.26 | 0.44 | 0.33 | 143 |
| 473 | 0.34 | 0.56 | 0.42 | 158 |
| 474 | 0.11 | 0.20 | 0.14 | 138 |
| 475 | 0.04 | 0.17 | 0.07 | 59 |
| 476 | 0.18 | 0.41 | 0.25 | 88 |
| 477 | 0.45 | 0.70 | 0.55 | 176 |
| 478 | 0.38 | 0.83 | 0.52 | 24 |
| 479 | 0.07 | 0.18 | 0.10 | 92 |
| 480 | 0.30 | 0.62 | 0.41 | 100 |
| 481 | 0.22 | 0.48 | 0.30 | 103 |
| 482 | 0.07 | 0.28 | 0.11 | 74 |
| 483 | 0.33 | 0.70 | 0.45 | 105 |
| 484 | 0.07 | 0.20 | 0.10 | 83 |
| 485 | 0.05 | 0.13 | 0.08 | 82 |
| 486 | 0.07 | 0.28 | 0.11 | 71 |
| 487 | 0.15 | 0.31 | 0.20 | 120 |
| 488 | 0.07 | 0.14 | 0.09 | 105 |
| 489 | 0.21 | 0.47 | 0.29 | 87 |
| 490 | 0.46 | 0.81 | 0.58 | 32 |
| 491 | 0.02 | 0.07 | 0.03 | 69 |

```
               492        0.04        0.12        0.06          49
               493        0.04        0.14        0.06         117
               494        0.15        0.33        0.20          61
               495        0.83        0.94        0.88         344
               496        0.16        0.33        0.22          52
               497        0.16        0.39        0.23         137
               498        0.18        0.41        0.25          98
               499        0.08        0.30        0.13          79

    micro avg        0.30        0.47        0.37      173812
    macro avg        0.22        0.41        0.28      173812
 weighted avg        0.36        0.47        0.40      173812
  samples avg        0.35        0.44        0.35      173812

Time taken to run this cell : 4:04:01.787154
```

[67]:
```python
from sklearn.externals import joblib
joblib.dump(classifier, 'lr_with_more_title_weight.pkl')
```

```
C:\Users\user\Anaconda3\lib\site-
packages\sklearn\externals\joblib\__init__.py:15: DeprecationWarning:
sklearn.externals.joblib is deprecated in 0.21 and will be removed in 0.23.
Please import this functionality directly from joblib, which can be installed
with: pip install joblib. If this warning is raised when loading pickled models,
you may need to re-serialize those models with scikit-learn 0.21+.
  warnings.warn(msg, category=DeprecationWarning)
```

[67]: ['lr_with_more_title_weight.pkl']

[68]:
```python
start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.00001,
 ↪penalty='l1'))
classifier.fit(x_train_multilabel_BOW, y_train)
predictions = classifier.predict (x_test_multilabel_BOW)


print("Accuracy :",metrics.accuracy_score(y_test, predictions))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions))


precision = precision_score(y_test, predictions, average='micro')
recall = recall_score(y_test, predictions, average='micro')
f1 = f1_score(y_test, predictions, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,
 ↪recall, f1))
```

```python
precision = precision_score(y_test, predictions, average='macro')
recall = recall_score(y_test, predictions, average='macro')
f1 = f1_score(y_test, predictions, average='macro')

print("Macro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,
 →recall, f1))

print (metrics.classification_report(y_test, predictions))
print("Time taken to run this cell :", datetime.now() - start)
```

```
Accuracy : 0.10059
Hamming loss  0.00562846
Micro-average quality numbers
Precision: 0.3017, Recall: 0.4709, F1-measure: 0.3678
Macro-average quality numbers
Precision: 0.2192, Recall: 0.4080, F1-measure: 0.2775
           precision    recall  f1-score   support

        0       0.78      0.79      0.79      5519
        1       0.34      0.51      0.41      8190
        2       0.55      0.49      0.52      6529
        3       0.36      0.52      0.43      3231
        4       0.57      0.51      0.54      6430
        5       0.46      0.48      0.47      2879
        6       0.62      0.60      0.61      5086
        7       0.64      0.65      0.64      4533
        8       0.22      0.22      0.22      3000
        9       0.60      0.62      0.61      2765
       10       0.32      0.31      0.32      3051
       11       0.48      0.49      0.48      3009
       12       0.41      0.41      0.41      2630
       13       0.25      0.36      0.30      1426
       14       0.64      0.70      0.67      2548
       15       0.37      0.36      0.37      2371
       16       0.29      0.34      0.31       873
       17       0.59      0.69      0.64      2151
       18       0.32      0.36      0.34      2204
       19       0.29      0.45      0.35       831
       20       0.54      0.57      0.56      1860
       21       0.20      0.25      0.22      2023
       22       0.30      0.34      0.32      1513
       23       0.57      0.65      0.60      1207
       24       0.24      0.39      0.30       506
       25       0.25      0.46      0.32       425
       26       0.39      0.53      0.45       793
       27       0.38      0.47      0.42      1291
```

| | | | | |
|---|---|---|---|---|
| 28 | 0.44 | 0.52 | 0.48 | 1208 |
| 29 | 0.12 | 0.23 | 0.16 | 406 |
| 30 | 0.23 | 0.32 | 0.27 | 504 |
| 31 | 0.14 | 0.20 | 0.16 | 732 |
| 32 | 0.20 | 0.38 | 0.26 | 441 |
| 33 | 0.35 | 0.41 | 0.37 | 1645 |
| 34 | 0.30 | 0.36 | 0.33 | 1058 |
| 35 | 0.47 | 0.62 | 0.54 | 946 |
| 36 | 0.25 | 0.43 | 0.32 | 644 |
| 37 | 0.30 | 0.75 | 0.43 | 136 |
| 38 | 0.35 | 0.52 | 0.42 | 570 |
| 39 | 0.33 | 0.39 | 0.36 | 766 |
| 40 | 0.33 | 0.48 | 0.39 | 1132 |
| 41 | 0.13 | 0.31 | 0.18 | 174 |
| 42 | 0.34 | 0.62 | 0.44 | 210 |
| 43 | 0.33 | 0.53 | 0.41 | 433 |
| 44 | 0.36 | 0.52 | 0.42 | 626 |
| 45 | 0.31 | 0.46 | 0.37 | 852 |
| 46 | 0.37 | 0.53 | 0.44 | 534 |
| 47 | 0.17 | 0.32 | 0.22 | 350 |
| 48 | 0.37 | 0.56 | 0.45 | 496 |
| 49 | 0.58 | 0.72 | 0.64 | 785 |
| 50 | 0.13 | 0.24 | 0.17 | 475 |
| 51 | 0.10 | 0.21 | 0.13 | 305 |
| 52 | 0.09 | 0.17 | 0.11 | 251 |
| 53 | 0.35 | 0.48 | 0.40 | 914 |
| 54 | 0.24 | 0.30 | 0.27 | 728 |
| 55 | 0.06 | 0.10 | 0.08 | 258 |
| 56 | 0.24 | 0.36 | 0.29 | 821 |
| 57 | 0.15 | 0.23 | 0.18 | 541 |
| 58 | 0.35 | 0.43 | 0.38 | 748 |
| 59 | 0.63 | 0.79 | 0.70 | 724 |
| 60 | 0.19 | 0.29 | 0.23 | 660 |
| 61 | 0.17 | 0.32 | 0.22 | 235 |
| 62 | 0.59 | 0.78 | 0.67 | 718 |
| 63 | 0.47 | 0.74 | 0.58 | 468 |
| 64 | 0.17 | 0.47 | 0.25 | 191 |
| 65 | 0.13 | 0.21 | 0.16 | 429 |
| 66 | 0.13 | 0.23 | 0.17 | 415 |
| 67 | 0.35 | 0.58 | 0.43 | 274 |
| 68 | 0.42 | 0.62 | 0.50 | 510 |
| 69 | 0.32 | 0.54 | 0.40 | 466 |
| 70 | 0.13 | 0.24 | 0.16 | 305 |
| 71 | 0.13 | 0.28 | 0.18 | 247 |
| 72 | 0.39 | 0.59 | 0.47 | 401 |
| 73 | 0.36 | 0.84 | 0.51 | 86 |
| 74 | 0.20 | 0.58 | 0.30 | 120 |
| 75 | 0.36 | 0.72 | 0.48 | 129 |

| | | | | |
|---|---|---|---|---|
| 76 | 0.10 | 0.12 | 0.11 | 473 |
| 77 | 0.11 | 0.39 | 0.17 | 143 |
| 78 | 0.36 | 0.59 | 0.45 | 347 |
| 79 | 0.22 | 0.35 | 0.27 | 479 |
| 80 | 0.23 | 0.51 | 0.32 | 279 |
| 81 | 0.28 | 0.44 | 0.34 | 461 |
| 82 | 0.10 | 0.21 | 0.13 | 298 |
| 83 | 0.36 | 0.62 | 0.46 | 396 |
| 84 | 0.21 | 0.48 | 0.29 | 184 |
| 85 | 0.25 | 0.40 | 0.31 | 573 |
| 86 | 0.13 | 0.18 | 0.15 | 325 |
| 87 | 0.22 | 0.47 | 0.30 | 273 |
| 88 | 0.12 | 0.34 | 0.18 | 135 |
| 89 | 0.13 | 0.29 | 0.18 | 232 |
| 90 | 0.29 | 0.48 | 0.36 | 409 |
| 91 | 0.24 | 0.40 | 0.30 | 420 |
| 92 | 0.41 | 0.61 | 0.49 | 408 |
| 93 | 0.26 | 0.57 | 0.36 | 241 |
| 94 | 0.07 | 0.15 | 0.09 | 211 |
| 95 | 0.15 | 0.25 | 0.19 | 277 |
| 96 | 0.11 | 0.17 | 0.13 | 410 |
| 97 | 0.48 | 0.62 | 0.54 | 501 |
| 98 | 0.29 | 0.71 | 0.41 | 136 |
| 99 | 0.24 | 0.48 | 0.32 | 239 |
| 100 | 0.14 | 0.27 | 0.19 | 324 |
| 101 | 0.49 | 0.78 | 0.60 | 277 |
| 102 | 0.63 | 0.80 | 0.71 | 613 |
| 103 | 0.12 | 0.31 | 0.17 | 157 |
| 104 | 0.10 | 0.22 | 0.14 | 295 |
| 105 | 0.33 | 0.50 | 0.40 | 334 |
| 106 | 0.31 | 0.44 | 0.36 | 335 |
| 107 | 0.33 | 0.59 | 0.42 | 389 |
| 108 | 0.21 | 0.42 | 0.28 | 251 |
| 109 | 0.27 | 0.50 | 0.35 | 317 |
| 110 | 0.07 | 0.21 | 0.10 | 187 |
| 111 | 0.09 | 0.25 | 0.13 | 140 |
| 112 | 0.22 | 0.60 | 0.33 | 154 |
| 113 | 0.23 | 0.36 | 0.28 | 332 |
| 114 | 0.21 | 0.39 | 0.28 | 323 |
| 115 | 0.18 | 0.37 | 0.24 | 344 |
| 116 | 0.38 | 0.56 | 0.45 | 370 |
| 117 | 0.24 | 0.38 | 0.29 | 313 |
| 118 | 0.61 | 0.77 | 0.68 | 874 |
| 119 | 0.17 | 0.32 | 0.22 | 293 |
| 120 | 0.06 | 0.15 | 0.08 | 200 |
| 121 | 0.38 | 0.59 | 0.46 | 463 |
| 122 | 0.12 | 0.32 | 0.17 | 119 |
| 123 | 0.01 | 0.03 | 0.02 | 256 |

| | | | | |
|---|---|---|---|---|
| 124 | 0.42 | 0.73 | 0.53 | 195 |
| 125 | 0.08 | 0.28 | 0.13 | 138 |
| 126 | 0.42 | 0.61 | 0.50 | 376 |
| 127 | 0.05 | 0.12 | 0.07 | 122 |
| 128 | 0.06 | 0.12 | 0.08 | 252 |
| 129 | 0.24 | 0.40 | 0.30 | 144 |
| 130 | 0.09 | 0.22 | 0.13 | 150 |
| 131 | 0.06 | 0.13 | 0.08 | 210 |
| 132 | 0.24 | 0.37 | 0.29 | 361 |
| 133 | 0.60 | 0.70 | 0.64 | 453 |
| 134 | 0.32 | 0.76 | 0.45 | 124 |
| 135 | 0.04 | 0.15 | 0.07 | 91 |
| 136 | 0.11 | 0.39 | 0.18 | 128 |
| 137 | 0.20 | 0.46 | 0.28 | 218 |
| 138 | 0.12 | 0.25 | 0.17 | 243 |
| 139 | 0.12 | 0.32 | 0.17 | 149 |
| 140 | 0.37 | 0.54 | 0.44 | 318 |
| 141 | 0.06 | 0.16 | 0.09 | 159 |
| 142 | 0.31 | 0.47 | 0.37 | 274 |
| 143 | 0.60 | 0.87 | 0.71 | 362 |
| 144 | 0.11 | 0.34 | 0.17 | 118 |
| 145 | 0.19 | 0.52 | 0.28 | 164 |
| 146 | 0.24 | 0.45 | 0.32 | 461 |
| 147 | 0.29 | 0.57 | 0.38 | 159 |
| 148 | 0.12 | 0.25 | 0.16 | 166 |
| 149 | 0.48 | 0.64 | 0.54 | 346 |
| 150 | 0.19 | 0.30 | 0.23 | 350 |
| 151 | 0.27 | 0.82 | 0.41 | 55 |
| 152 | 0.39 | 0.59 | 0.47 | 387 |
| 153 | 0.16 | 0.27 | 0.20 | 150 |
| 154 | 0.13 | 0.21 | 0.16 | 281 |
| 155 | 0.11 | 0.30 | 0.16 | 202 |
| 156 | 0.31 | 0.69 | 0.42 | 130 |
| 157 | 0.15 | 0.21 | 0.18 | 245 |
| 158 | 0.43 | 0.77 | 0.55 | 177 |
| 159 | 0.15 | 0.48 | 0.22 | 130 |
| 160 | 0.17 | 0.30 | 0.22 | 336 |
| 161 | 0.37 | 0.70 | 0.48 | 220 |
| 162 | 0.07 | 0.19 | 0.11 | 229 |
| 163 | 0.45 | 0.60 | 0.51 | 316 |
| 164 | 0.28 | 0.54 | 0.37 | 283 |
| 165 | 0.19 | 0.41 | 0.26 | 197 |
| 166 | 0.23 | 0.55 | 0.33 | 101 |
| 167 | 0.17 | 0.31 | 0.22 | 231 |
| 168 | 0.21 | 0.41 | 0.28 | 370 |
| 169 | 0.20 | 0.34 | 0.25 | 258 |
| 170 | 0.04 | 0.15 | 0.07 | 101 |
| 171 | 0.09 | 0.31 | 0.14 | 89 |

| | | | | |
|---|---|---|---|---|
| 172 | 0.17 | 0.40 | 0.24 | 193 |
| 173 | 0.25 | 0.43 | 0.32 | 309 |
| 174 | 0.10 | 0.25 | 0.14 | 172 |
| 175 | 0.30 | 0.77 | 0.43 | 95 |
| 176 | 0.57 | 0.73 | 0.64 | 346 |
| 177 | 0.39 | 0.62 | 0.48 | 322 |
| 178 | 0.30 | 0.60 | 0.40 | 232 |
| 179 | 0.07 | 0.19 | 0.11 | 125 |
| 180 | 0.20 | 0.52 | 0.28 | 145 |
| 181 | 0.05 | 0.27 | 0.08 | 77 |
| 182 | 0.07 | 0.18 | 0.10 | 182 |
| 183 | 0.22 | 0.42 | 0.29 | 257 |
| 184 | 0.08 | 0.18 | 0.11 | 216 |
| 185 | 0.11 | 0.26 | 0.16 | 242 |
| 186 | 0.14 | 0.31 | 0.19 | 165 |
| 187 | 0.37 | 0.61 | 0.46 | 263 |
| 188 | 0.10 | 0.20 | 0.14 | 174 |
| 189 | 0.37 | 0.48 | 0.41 | 136 |
| 190 | 0.44 | 0.72 | 0.55 | 202 |
| 191 | 0.08 | 0.24 | 0.12 | 134 |
| 192 | 0.25 | 0.48 | 0.33 | 230 |
| 193 | 0.07 | 0.23 | 0.11 | 90 |
| 194 | 0.27 | 0.57 | 0.37 | 185 |
| 195 | 0.05 | 0.13 | 0.07 | 156 |
| 196 | 0.06 | 0.17 | 0.09 | 160 |
| 197 | 0.11 | 0.19 | 0.14 | 266 |
| 198 | 0.14 | 0.27 | 0.18 | 284 |
| 199 | 0.08 | 0.19 | 0.11 | 145 |
| 200 | 0.51 | 0.82 | 0.63 | 212 |
| 201 | 0.20 | 0.40 | 0.27 | 317 |
| 202 | 0.45 | 0.66 | 0.54 | 427 |
| 203 | 0.12 | 0.29 | 0.17 | 232 |
| 204 | 0.22 | 0.44 | 0.29 | 217 |
| 205 | 0.38 | 0.56 | 0.45 | 527 |
| 206 | 0.06 | 0.18 | 0.09 | 124 |
| 207 | 0.20 | 0.45 | 0.28 | 103 |
| 208 | 0.31 | 0.53 | 0.39 | 287 |
| 209 | 0.09 | 0.19 | 0.12 | 193 |
| 210 | 0.26 | 0.50 | 0.34 | 220 |
| 211 | 0.13 | 0.33 | 0.19 | 140 |
| 212 | 0.08 | 0.19 | 0.11 | 161 |
| 213 | 0.21 | 0.60 | 0.32 | 72 |
| 214 | 0.39 | 0.52 | 0.45 | 396 |
| 215 | 0.16 | 0.51 | 0.25 | 134 |
| 216 | 0.23 | 0.32 | 0.27 | 400 |
| 217 | 0.14 | 0.45 | 0.21 | 75 |
| 218 | 0.60 | 0.81 | 0.69 | 219 |
| 219 | 0.23 | 0.48 | 0.31 | 210 |

| | | | | |
|---|---|---|---|---|
| 220 | 0.48 | 0.77 | 0.59 | 298 |
| 221 | 0.53 | 0.78 | 0.63 | 266 |
| 222 | 0.35 | 0.55 | 0.43 | 290 |
| 223 | 0.03 | 0.11 | 0.05 | 128 |
| 224 | 0.21 | 0.48 | 0.30 | 159 |
| 225 | 0.13 | 0.43 | 0.20 | 164 |
| 226 | 0.17 | 0.44 | 0.25 | 144 |
| 227 | 0.36 | 0.59 | 0.45 | 276 |
| 228 | 0.04 | 0.08 | 0.05 | 235 |
| 229 | 0.06 | 0.16 | 0.09 | 216 |
| 230 | 0.10 | 0.27 | 0.15 | 228 |
| 231 | 0.21 | 0.56 | 0.31 | 64 |
| 232 | 0.07 | 0.23 | 0.11 | 103 |
| 233 | 0.27 | 0.49 | 0.35 | 216 |
| 234 | 0.12 | 0.29 | 0.17 | 116 |
| 235 | 0.16 | 0.48 | 0.24 | 77 |
| 236 | 0.41 | 0.78 | 0.54 | 67 |
| 237 | 0.11 | 0.24 | 0.15 | 218 |
| 238 | 0.08 | 0.22 | 0.11 | 139 |
| 239 | 0.01 | 0.03 | 0.02 | 94 |
| 240 | 0.12 | 0.42 | 0.19 | 77 |
| 241 | 0.04 | 0.12 | 0.06 | 167 |
| 242 | 0.22 | 0.50 | 0.31 | 86 |
| 243 | 0.06 | 0.29 | 0.10 | 58 |
| 244 | 0.31 | 0.48 | 0.38 | 269 |
| 245 | 0.08 | 0.19 | 0.11 | 112 |
| 246 | 0.68 | 0.86 | 0.76 | 255 |
| 247 | 0.06 | 0.24 | 0.10 | 58 |
| 248 | 0.03 | 0.16 | 0.05 | 81 |
| 249 | 0.04 | 0.14 | 0.07 | 131 |
| 250 | 0.10 | 0.32 | 0.16 | 93 |
| 251 | 0.17 | 0.42 | 0.24 | 154 |
| 252 | 0.03 | 0.09 | 0.05 | 129 |
| 253 | 0.12 | 0.39 | 0.19 | 83 |
| 254 | 0.11 | 0.20 | 0.14 | 191 |
| 255 | 0.09 | 0.16 | 0.11 | 219 |
| 256 | 0.04 | 0.08 | 0.05 | 130 |
| 257 | 0.14 | 0.34 | 0.20 | 93 |
| 258 | 0.37 | 0.65 | 0.47 | 217 |
| 259 | 0.10 | 0.30 | 0.15 | 141 |
| 260 | 0.19 | 0.39 | 0.25 | 143 |
| 261 | 0.17 | 0.32 | 0.23 | 219 |
| 262 | 0.16 | 0.48 | 0.24 | 107 |
| 263 | 0.23 | 0.39 | 0.29 | 236 |
| 264 | 0.12 | 0.38 | 0.18 | 119 |
| 265 | 0.13 | 0.38 | 0.19 | 72 |
| 266 | 0.07 | 0.26 | 0.11 | 70 |
| 267 | 0.13 | 0.32 | 0.19 | 107 |

| | | | | |
|------|------|------|------|-----|
| 268 | 0.27 | 0.57 | 0.37 | 169 |
| 269 | 0.18 | 0.39 | 0.24 | 129 |
| 270 | 0.32 | 0.64 | 0.43 | 159 |
| 271 | 0.28 | 0.63 | 0.39 | 190 |
| 272 | 0.19 | 0.35 | 0.25 | 248 |
| 273 | 0.58 | 0.82 | 0.68 | 264 |
| 274 | 0.43 | 0.78 | 0.56 | 105 |
| 275 | 0.08 | 0.26 | 0.12 | 104 |
| 276 | 0.05 | 0.14 | 0.07 | 115 |
| 277 | 0.37 | 0.65 | 0.47 | 170 |
| 278 | 0.28 | 0.59 | 0.38 | 145 |
| 279 | 0.52 | 0.79 | 0.63 | 230 |
| 280 | 0.16 | 0.45 | 0.23 | 80 |
| 281 | 0.44 | 0.66 | 0.53 | 217 |
| 282 | 0.39 | 0.64 | 0.48 | 175 |
| 283 | 0.14 | 0.25 | 0.18 | 269 |
| 284 | 0.14 | 0.38 | 0.21 | 74 |
| 285 | 0.32 | 0.67 | 0.43 | 206 |
| 286 | 0.44 | 0.70 | 0.55 | 227 |
| 287 | 0.20 | 0.58 | 0.30 | 130 |
| 288 | 0.05 | 0.14 | 0.08 | 129 |
| 289 | 0.05 | 0.25 | 0.08 | 80 |
| 290 | 0.07 | 0.26 | 0.11 | 99 |
| 291 | 0.24 | 0.51 | 0.32 | 208 |
| 292 | 0.03 | 0.15 | 0.04 | 67 |
| 293 | 0.27 | 0.63 | 0.38 | 109 |
| 294 | 0.16 | 0.39 | 0.23 | 140 |
| 295 | 0.13 | 0.28 | 0.18 | 241 |
| 296 | 0.08 | 0.29 | 0.13 | 72 |
| 297 | 0.09 | 0.26 | 0.13 | 107 |
| 298 | 0.24 | 0.56 | 0.33 | 61 |
| 299 | 0.32 | 0.64 | 0.43 | 77 |
| 300 | 0.07 | 0.24 | 0.11 | 111 |
| 301 | 0.00 | 0.01 | 0.00 | 126 |
| 302 | 0.05 | 0.15 | 0.07 | 73 |
| 303 | 0.27 | 0.55 | 0.37 | 176 |
| 304 | 0.69 | 0.86 | 0.77 | 230 |
| 305 | 0.52 | 0.79 | 0.63 | 156 |
| 306 | 0.21 | 0.45 | 0.29 | 146 |
| 307 | 0.06 | 0.15 | 0.09 | 98 |
| 308 | 0.02 | 0.08 | 0.03 | 78 |
| 309 | 0.06 | 0.20 | 0.10 | 94 |
| 310 | 0.28 | 0.55 | 0.37 | 162 |
| 311 | 0.31 | 0.57 | 0.40 | 116 |
| 312 | 0.15 | 0.42 | 0.22 | 57 |
| 313 | 0.05 | 0.17 | 0.08 | 65 |
| 314 | 0.18 | 0.42 | 0.25 | 138 |
| 315 | 0.27 | 0.37 | 0.31 | 195 |

| 316 | 0.13 | 0.33 | 0.19 | 69 |
|-----|------|------|------|-----|
| 317 | 0.09 | 0.31 | 0.13 | 134 |
| 318 | 0.22 | 0.43 | 0.30 | 148 |
| 319 | 0.39 | 0.61 | 0.47 | 161 |
| 320 | 0.11 | 0.38 | 0.16 | 104 |
| 321 | 0.34 | 0.65 | 0.45 | 156 |
| 322 | 0.21 | 0.41 | 0.28 | 134 |
| 323 | 0.30 | 0.51 | 0.38 | 232 |
| 324 | 0.10 | 0.28 | 0.15 | 92 |
| 325 | 0.17 | 0.34 | 0.23 | 197 |
| 326 | 0.05 | 0.17 | 0.08 | 126 |
| 327 | 0.02 | 0.05 | 0.03 | 115 |
| 328 | 0.58 | 0.79 | 0.67 | 198 |
| 329 | 0.19 | 0.43 | 0.27 | 125 |
| 330 | 0.11 | 0.27 | 0.16 | 81 |
| 331 | 0.11 | 0.27 | 0.16 | 94 |
| 332 | 0.09 | 0.27 | 0.13 | 56 |
| 333 | 0.06 | 0.13 | 0.08 | 260 |
| 334 | 0.08 | 0.27 | 0.12 | 60 |
| 335 | 0.10 | 0.25 | 0.14 | 110 |
| 336 | 0.20 | 0.56 | 0.30 | 71 |
| 337 | 0.04 | 0.17 | 0.06 | 66 |
| 338 | 0.18 | 0.45 | 0.26 | 150 |
| 339 | 0.02 | 0.09 | 0.03 | 54 |
| 340 | 0.42 | 0.66 | 0.51 | 195 |
| 341 | 0.31 | 0.59 | 0.41 | 79 |
| 342 | 0.10 | 0.39 | 0.17 | 38 |
| 343 | 0.11 | 0.40 | 0.17 | 43 |
| 344 | 0.21 | 0.44 | 0.28 | 68 |
| 345 | 0.29 | 0.52 | 0.37 | 73 |
| 346 | 0.06 | 0.21 | 0.09 | 116 |
| 347 | 0.24 | 0.63 | 0.35 | 111 |
| 348 | 0.03 | 0.16 | 0.06 | 63 |
| 349 | 0.43 | 0.77 | 0.55 | 104 |
| 350 | 0.22 | 0.66 | 0.33 | 44 |
| 351 | 0.11 | 0.35 | 0.17 | 40 |
| 352 | 0.45 | 0.65 | 0.53 | 136 |
| 353 | 0.08 | 0.28 | 0.13 | 54 |
| 354 | 0.08 | 0.22 | 0.11 | 134 |
| 355 | 0.23 | 0.55 | 0.32 | 120 |
| 356 | 0.29 | 0.48 | 0.36 | 228 |
| 357 | 0.32 | 0.48 | 0.38 | 269 |
| 358 | 0.21 | 0.46 | 0.29 | 80 |
| 359 | 0.37 | 0.73 | 0.49 | 140 |
| 360 | 0.13 | 0.28 | 0.17 | 125 |
| 361 | 0.68 | 0.80 | 0.73 | 169 |
| 362 | 0.05 | 0.20 | 0.08 | 56 |
| 363 | 0.63 | 0.79 | 0.70 | 154 |

| | | | | |
|---|---|---|---|---|
| 364 | 0.12 | 0.34 | 0.18 | 58 |
| 365 | 0.09 | 0.30 | 0.14 | 71 |
| 366 | 0.47 | 0.80 | 0.59 | 54 |
| 367 | 0.05 | 0.16 | 0.07 | 116 |
| 368 | 0.08 | 0.22 | 0.12 | 54 |
| 369 | 0.03 | 0.11 | 0.04 | 71 |
| 370 | 0.03 | 0.11 | 0.04 | 61 |
| 371 | 0.05 | 0.18 | 0.08 | 71 |
| 372 | 0.16 | 0.54 | 0.25 | 52 |
| 373 | 0.35 | 0.60 | 0.44 | 150 |
| 374 | 0.13 | 0.40 | 0.19 | 93 |
| 375 | 0.04 | 0.15 | 0.06 | 67 |
| 376 | 0.02 | 0.05 | 0.03 | 76 |
| 377 | 0.18 | 0.41 | 0.25 | 106 |
| 378 | 0.03 | 0.09 | 0.05 | 86 |
| 379 | 0.01 | 0.14 | 0.02 | 14 |
| 380 | 0.36 | 0.66 | 0.47 | 122 |
| 381 | 0.04 | 0.12 | 0.06 | 104 |
| 382 | 0.05 | 0.21 | 0.08 | 66 |
| 383 | 0.18 | 0.38 | 0.24 | 110 |
| 384 | 0.04 | 0.11 | 0.06 | 155 |
| 385 | 0.13 | 0.56 | 0.21 | 50 |
| 386 | 0.07 | 0.19 | 0.10 | 64 |
| 387 | 0.09 | 0.24 | 0.13 | 93 |
| 388 | 0.16 | 0.41 | 0.23 | 102 |
| 389 | 0.05 | 0.13 | 0.07 | 108 |
| 390 | 0.63 | 0.77 | 0.70 | 178 |
| 391 | 0.15 | 0.39 | 0.21 | 115 |
| 392 | 0.23 | 0.50 | 0.31 | 42 |
| 393 | 0.02 | 0.04 | 0.03 | 134 |
| 394 | 0.07 | 0.20 | 0.10 | 112 |
| 395 | 0.15 | 0.45 | 0.23 | 176 |
| 396 | 0.10 | 0.26 | 0.14 | 125 |
| 397 | 0.40 | 0.53 | 0.45 | 224 |
| 398 | 0.25 | 0.73 | 0.37 | 63 |
| 399 | 0.01 | 0.03 | 0.01 | 59 |
| 400 | 0.17 | 0.54 | 0.26 | 63 |
| 401 | 0.11 | 0.33 | 0.17 | 98 |
| 402 | 0.15 | 0.40 | 0.22 | 162 |
| 403 | 0.08 | 0.34 | 0.13 | 83 |
| 404 | 0.25 | 0.89 | 0.40 | 19 |
| 405 | 0.08 | 0.26 | 0.12 | 92 |
| 406 | 0.09 | 0.51 | 0.15 | 41 |
| 407 | 0.30 | 0.51 | 0.38 | 43 |
| 408 | 0.28 | 0.50 | 0.36 | 160 |
| 409 | 0.09 | 0.22 | 0.12 | 50 |
| 410 | 0.02 | 0.21 | 0.04 | 19 |
| 411 | 0.13 | 0.27 | 0.18 | 175 |

| | | | | |
|---|---|---|---|---|
| 412 | 0.05 | 0.14 | 0.07 | 72 |
| 413 | 0.07 | 0.20 | 0.11 | 95 |
| 414 | 0.10 | 0.27 | 0.15 | 97 |
| 415 | 0.07 | 0.21 | 0.10 | 48 |
| 416 | 0.19 | 0.42 | 0.26 | 83 |
| 417 | 0.06 | 0.20 | 0.09 | 40 |
| 418 | 0.11 | 0.30 | 0.16 | 91 |
| 419 | 0.23 | 0.52 | 0.32 | 90 |
| 420 | 0.04 | 0.22 | 0.07 | 37 |
| 421 | 0.07 | 0.21 | 0.11 | 66 |
| 422 | 0.12 | 0.41 | 0.19 | 73 |
| 423 | 0.10 | 0.27 | 0.14 | 56 |
| 424 | 0.40 | 0.94 | 0.56 | 33 |
| 425 | 0.04 | 0.09 | 0.05 | 76 |
| 426 | 0.02 | 0.07 | 0.04 | 81 |
| 427 | 0.57 | 0.75 | 0.65 | 150 |
| 428 | 0.45 | 0.72 | 0.55 | 29 |
| 429 | 0.91 | 0.92 | 0.91 | 389 |
| 430 | 0.25 | 0.51 | 0.34 | 167 |
| 431 | 0.10 | 0.20 | 0.14 | 123 |
| 432 | 0.10 | 0.41 | 0.16 | 39 |
| 433 | 0.14 | 0.35 | 0.20 | 82 |
| 434 | 0.48 | 0.73 | 0.58 | 66 |
| 435 | 0.24 | 0.55 | 0.34 | 93 |
| 436 | 0.23 | 0.62 | 0.34 | 87 |
| 437 | 0.06 | 0.19 | 0.09 | 86 |
| 438 | 0.34 | 0.64 | 0.44 | 104 |
| 439 | 0.13 | 0.32 | 0.19 | 100 |
| 440 | 0.05 | 0.12 | 0.07 | 141 |
| 441 | 0.22 | 0.52 | 0.31 | 110 |
| 442 | 0.11 | 0.23 | 0.15 | 123 |
| 443 | 0.14 | 0.30 | 0.19 | 71 |
| 444 | 0.12 | 0.24 | 0.16 | 109 |
| 445 | 0.10 | 0.46 | 0.17 | 48 |
| 446 | 0.19 | 0.45 | 0.26 | 76 |
| 447 | 0.08 | 0.37 | 0.13 | 38 |
| 448 | 0.29 | 0.63 | 0.40 | 81 |
| 449 | 0.21 | 0.34 | 0.26 | 132 |
| 450 | 0.17 | 0.41 | 0.24 | 81 |
| 451 | 0.16 | 0.42 | 0.23 | 76 |
| 452 | 0.11 | 0.27 | 0.15 | 44 |
| 453 | 0.02 | 0.07 | 0.03 | 44 |
| 454 | 0.17 | 0.56 | 0.26 | 70 |
| 455 | 0.12 | 0.35 | 0.18 | 155 |
| 456 | 0.10 | 0.53 | 0.17 | 43 |
| 457 | 0.14 | 0.51 | 0.22 | 72 |
| 458 | 0.04 | 0.19 | 0.07 | 62 |
| 459 | 0.07 | 0.32 | 0.11 | 69 |

```
460        0.04      0.08      0.05        119
461        0.32      0.51      0.39         79
462        0.07      0.26      0.11         47
463        0.14      0.46      0.21        104
464        0.21      0.35      0.26        106
465        0.12      0.33      0.18         64
466        0.29      0.50      0.36        173
467        0.23      0.54      0.32        107
468        0.19      0.39      0.26        126
469        0.05      0.08      0.06        114
470        0.72      0.85      0.78        140
471        0.24      0.53      0.33         79
472        0.22      0.39      0.28        143
473        0.30      0.46      0.36        158
474        0.06      0.12      0.08        138
475        0.02      0.10      0.04         59
476        0.23      0.45      0.30         88
477        0.45      0.70      0.55        176
478        0.28      0.67      0.40         24
479        0.08      0.18      0.11         92
480        0.39      0.62      0.48        100
481        0.24      0.45      0.31        103
482        0.08      0.24      0.12         74
483        0.40      0.70      0.51        105
484        0.05      0.12      0.07         83
485        0.03      0.13      0.05         82
486        0.08      0.28      0.12         71
487        0.14      0.30      0.19        120
488        0.07      0.18      0.10        105
489        0.20      0.45      0.28         87
490        0.34      0.84      0.48         32
491        0.03      0.12      0.05         69
492        0.03      0.10      0.05         49
493        0.03      0.09      0.04        117
494        0.12      0.33      0.17         61
495        0.81      0.89      0.85        344
496        0.13      0.29      0.18         52
497        0.20      0.32      0.24        137
498        0.14      0.31      0.20         98
499        0.12      0.29      0.17         79

   micro avg       0.30      0.47      0.37     173812
   macro avg       0.22      0.41      0.28     173812
weighted avg       0.36      0.47      0.40     173812
 samples avg       0.36      0.44      0.35     173812

Time taken to run this cell : 2:20:10.610383
```

```
[70]: joblib.dump(classifier, 'SVM_with_more_title_weight.pkl')
```

```
[70]: ['SVM_with_more_title_weight.pkl']
```

```
[75]: #This cell took around 24 hours to run.
      from sklearn.model_selection import GridSearchCV
      classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1',C=1.0))
      parameters = {'estimator__C':[0.0001,0.001,0.01,0.1,1,10,100]}
      gridSCV = GridSearchCV(estimator = classifier_2, param_grid = parameters, cv=3,␣
       ↪scoring='f1_micro')
      gridSCV.fit(x_train_multilabel, y_train)
      print(gridSCV.best_params_)
      print(gridSCV.best_score_)
```

```
{'estimator__C': 1}
0.5254786609395317
```

```
[77]: start = datetime.now()
      classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1',C=1.0))
      classifier_2.fit(x_train_multilabel, y_train)
      predictions_2 = classifier_2.predict(x_test_multilabel)
      print("Accuracy :",metrics.accuracy_score(y_test, predictions_2))
      print("Hamming loss ",metrics.hamming_loss(y_test,predictions_2))


      precision = precision_score(y_test, predictions_2, average='micro')
      recall = recall_score(y_test, predictions_2, average='micro')
      f1 = f1_score(y_test, predictions_2, average='micro')

      print("Micro-average quality numbers")
      print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,␣
       ↪recall, f1))


      precision = precision_score(y_test, predictions_2, average='macro')
      recall = recall_score(y_test, predictions_2, average='macro')
      f1 = f1_score(y_test, predictions_2, average='macro')

      print("Macro-average quality numbers")
      print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,␣
       ↪recall, f1))

      print (metrics.classification_report(y_test, predictions_2))
      print("Time taken to run this cell :", datetime.now() - start)
```

```
Accuracy : 0.25113
Hamming loss  0.00270284
Micro-average quality numbers
Precision: 0.7172, Recall: 0.3673, F1-measure: 0.4858
```

```
Macro-average quality numbers
Precision: 0.5570, Recall: 0.2951, F1-measure: 0.3710
        precision   recall   f1-score   support

    0       0.94     0.72      0.82       5519
    1       0.70     0.34      0.45       8190
    2       0.80     0.42      0.55       6529
    3       0.82     0.49      0.61       3231
    4       0.80     0.44      0.57       6430
    5       0.82     0.38      0.52       2879
    6       0.86     0.53      0.66       5086
    7       0.87     0.58      0.70       4533
    8       0.60     0.13      0.22       3000
    9       0.82     0.57      0.67       2765
   10       0.60     0.20      0.30       3051
   11       0.68     0.38      0.49       3009
   12       0.62     0.29      0.40       2630
   13       0.73     0.30      0.43       1426
   14       0.89     0.57      0.70       2548
   15       0.65     0.23      0.34       2371
   16       0.65     0.25      0.37        873
   17       0.89     0.63      0.74       2151
   18       0.60     0.25      0.35       2204
   19       0.71     0.41      0.52        831
   20       0.76     0.47      0.58       1860
   21       0.29     0.09      0.14       2023
   22       0.52     0.24      0.33       1513
   23       0.89     0.55      0.68       1207
   24       0.56     0.28      0.38        506
   25       0.69     0.34      0.45        425
   26       0.65     0.43      0.52        793
   27       0.62     0.38      0.47       1291
   28       0.74     0.39      0.51       1208
   29       0.46     0.10      0.17        406
   30       0.76     0.21      0.33        504
   31       0.26     0.08      0.12        732
   32       0.60     0.29      0.39        441
   33       0.60     0.27      0.38       1645
   34       0.69     0.26      0.38       1058
   35       0.83     0.58      0.68        946
   36       0.65     0.24      0.35        644
   37       0.98     0.65      0.78        136
   38       0.62     0.38      0.47        570
   39       0.84     0.31      0.45        766
   40       0.59     0.35      0.44       1132
   41       0.47     0.18      0.26        174
   42       0.76     0.49      0.59        210
   43       0.75     0.42      0.54        433
```

| | | | | |
|---|---|---|---|---|
| 44 | 0.66 | 0.52 | 0.58 | 626 |
| 45 | 0.71 | 0.36 | 0.47 | 852 |
| 46 | 0.77 | 0.45 | 0.57 | 534 |
| 47 | 0.37 | 0.15 | 0.22 | 350 |
| 48 | 0.75 | 0.52 | 0.62 | 496 |
| 49 | 0.78 | 0.64 | 0.71 | 785 |
| 50 | 0.21 | 0.06 | 0.09 | 475 |
| 51 | 0.37 | 0.13 | 0.19 | 305 |
| 52 | 0.42 | 0.03 | 0.06 | 251 |
| 53 | 0.66 | 0.40 | 0.50 | 914 |
| 54 | 0.50 | 0.18 | 0.26 | 728 |
| 55 | 0.47 | 0.03 | 0.05 | 258 |
| 56 | 0.45 | 0.24 | 0.31 | 821 |
| 57 | 0.46 | 0.10 | 0.17 | 541 |
| 58 | 0.76 | 0.31 | 0.45 | 748 |
| 59 | 0.94 | 0.66 | 0.77 | 724 |
| 60 | 0.35 | 0.10 | 0.15 | 660 |
| 61 | 0.78 | 0.20 | 0.31 | 235 |
| 62 | 0.92 | 0.74 | 0.82 | 718 |
| 63 | 0.83 | 0.69 | 0.75 | 468 |
| 64 | 0.55 | 0.36 | 0.43 | 191 |
| 65 | 0.33 | 0.11 | 0.17 | 429 |
| 66 | 0.29 | 0.06 | 0.10 | 415 |
| 67 | 0.74 | 0.50 | 0.59 | 274 |
| 68 | 0.82 | 0.53 | 0.64 | 510 |
| 69 | 0.67 | 0.45 | 0.54 | 466 |
| 70 | 0.30 | 0.09 | 0.13 | 305 |
| 71 | 0.49 | 0.17 | 0.25 | 247 |
| 72 | 0.78 | 0.53 | 0.64 | 401 |
| 73 | 0.99 | 0.77 | 0.86 | 86 |
| 74 | 0.72 | 0.42 | 0.53 | 120 |
| 75 | 0.92 | 0.67 | 0.78 | 129 |
| 76 | 0.47 | 0.02 | 0.04 | 473 |
| 77 | 0.40 | 0.29 | 0.33 | 143 |
| 78 | 0.79 | 0.49 | 0.60 | 347 |
| 79 | 0.69 | 0.25 | 0.36 | 479 |
| 80 | 0.56 | 0.34 | 0.43 | 279 |
| 81 | 0.70 | 0.23 | 0.34 | 461 |
| 82 | 0.34 | 0.04 | 0.07 | 298 |
| 83 | 0.78 | 0.50 | 0.61 | 396 |
| 84 | 0.55 | 0.29 | 0.38 | 184 |
| 85 | 0.61 | 0.24 | 0.35 | 573 |
| 86 | 0.50 | 0.07 | 0.12 | 325 |
| 87 | 0.51 | 0.29 | 0.37 | 273 |
| 88 | 0.49 | 0.21 | 0.30 | 135 |
| 89 | 0.36 | 0.11 | 0.17 | 232 |
| 90 | 0.56 | 0.34 | 0.43 | 409 |
| 91 | 0.61 | 0.27 | 0.37 | 420 |

| | | | | |
|-----|------|------|------|-----|
| 92  | 0.78 | 0.57 | 0.66 | 408 |
| 93  | 0.66 | 0.44 | 0.53 | 241 |
| 94  | 0.30 | 0.04 | 0.07 | 211 |
| 95  | 0.37 | 0.10 | 0.15 | 277 |
| 96  | 0.28 | 0.04 | 0.07 | 410 |
| 97  | 0.86 | 0.43 | 0.57 | 501 |
| 98  | 0.75 | 0.63 | 0.69 | 136 |
| 99  | 0.54 | 0.34 | 0.42 | 239 |
| 100 | 0.57 | 0.15 | 0.24 | 324 |
| 101 | 0.91 | 0.68 | 0.78 | 277 |
| 102 | 0.91 | 0.75 | 0.82 | 613 |
| 103 | 0.47 | 0.17 | 0.25 | 157 |
| 104 | 0.22 | 0.06 | 0.10 | 295 |
| 105 | 0.75 | 0.43 | 0.55 | 334 |
| 106 | 0.88 | 0.28 | 0.43 | 335 |
| 107 | 0.75 | 0.54 | 0.63 | 389 |
| 108 | 0.58 | 0.27 | 0.37 | 251 |
| 109 | 0.58 | 0.45 | 0.51 | 317 |
| 110 | 0.68 | 0.10 | 0.18 | 187 |
| 111 | 0.73 | 0.11 | 0.20 | 140 |
| 112 | 0.67 | 0.43 | 0.52 | 154 |
| 113 | 0.58 | 0.20 | 0.29 | 332 |
| 114 | 0.46 | 0.27 | 0.34 | 323 |
| 115 | 0.47 | 0.26 | 0.33 | 344 |
| 116 | 0.75 | 0.55 | 0.63 | 370 |
| 117 | 0.58 | 0.24 | 0.34 | 313 |
| 118 | 0.78 | 0.73 | 0.75 | 874 |
| 119 | 0.45 | 0.21 | 0.29 | 293 |
| 120 | 0.11 | 0.01 | 0.01 | 200 |
| 121 | 0.77 | 0.51 | 0.61 | 463 |
| 122 | 0.32 | 0.10 | 0.15 | 119 |
| 123 | 0.67 | 0.02 | 0.03 | 256 |
| 124 | 0.91 | 0.70 | 0.79 | 195 |
| 125 | 0.44 | 0.14 | 0.21 | 138 |
| 126 | 0.81 | 0.54 | 0.65 | 376 |
| 127 | 0.27 | 0.03 | 0.06 | 122 |
| 128 | 0.20 | 0.04 | 0.07 | 252 |
| 129 | 0.48 | 0.22 | 0.30 | 144 |
| 130 | 0.42 | 0.11 | 0.18 | 150 |
| 131 | 0.33 | 0.03 | 0.06 | 210 |
| 132 | 0.65 | 0.28 | 0.39 | 361 |
| 133 | 0.92 | 0.59 | 0.72 | 453 |
| 134 | 0.89 | 0.77 | 0.82 | 124 |
| 135 | 0.31 | 0.05 | 0.09 | 91  |
| 136 | 0.69 | 0.28 | 0.40 | 128 |
| 137 | 0.55 | 0.38 | 0.45 | 218 |
| 138 | 0.67 | 0.18 | 0.28 | 243 |
| 139 | 0.45 | 0.18 | 0.26 | 149 |

| | | | | |
|---|---|---|---|---|
| 140 | 0.77 | 0.46 | 0.58 | 318 |
| 141 | 0.32 | 0.10 | 0.15 | 159 |
| 142 | 0.63 | 0.38 | 0.47 | 274 |
| 143 | 0.85 | 0.79 | 0.82 | 362 |
| 144 | 0.54 | 0.21 | 0.30 | 118 |
| 145 | 0.63 | 0.39 | 0.48 | 164 |
| 146 | 0.54 | 0.31 | 0.39 | 461 |
| 147 | 0.68 | 0.45 | 0.54 | 159 |
| 148 | 0.30 | 0.12 | 0.17 | 166 |
| 149 | 0.97 | 0.55 | 0.70 | 346 |
| 150 | 0.64 | 0.13 | 0.21 | 350 |
| 151 | 0.93 | 0.67 | 0.78 | 55 |
| 152 | 0.78 | 0.52 | 0.63 | 387 |
| 153 | 0.51 | 0.17 | 0.25 | 150 |
| 154 | 0.58 | 0.12 | 0.21 | 281 |
| 155 | 0.25 | 0.06 | 0.10 | 202 |
| 156 | 0.81 | 0.67 | 0.73 | 130 |
| 157 | 0.28 | 0.06 | 0.10 | 245 |
| 158 | 0.93 | 0.63 | 0.75 | 177 |
| 159 | 0.53 | 0.34 | 0.41 | 130 |
| 160 | 0.48 | 0.18 | 0.26 | 336 |
| 161 | 0.90 | 0.65 | 0.75 | 220 |
| 162 | 0.28 | 0.06 | 0.09 | 229 |
| 163 | 0.87 | 0.44 | 0.58 | 316 |
| 164 | 0.78 | 0.44 | 0.56 | 283 |
| 165 | 0.60 | 0.34 | 0.44 | 197 |
| 166 | 0.65 | 0.43 | 0.51 | 101 |
| 167 | 0.45 | 0.18 | 0.26 | 231 |
| 168 | 0.56 | 0.27 | 0.36 | 370 |
| 169 | 0.40 | 0.21 | 0.27 | 258 |
| 170 | 0.36 | 0.08 | 0.13 | 101 |
| 171 | 0.38 | 0.24 | 0.29 | 89 |
| 172 | 0.53 | 0.36 | 0.43 | 193 |
| 173 | 0.47 | 0.26 | 0.33 | 309 |
| 174 | 0.62 | 0.14 | 0.23 | 172 |
| 175 | 0.92 | 0.73 | 0.81 | 95 |
| 176 | 0.93 | 0.62 | 0.74 | 346 |
| 177 | 0.86 | 0.57 | 0.69 | 322 |
| 178 | 0.65 | 0.51 | 0.57 | 232 |
| 179 | 0.20 | 0.04 | 0.07 | 125 |
| 180 | 0.65 | 0.33 | 0.44 | 145 |
| 181 | 0.44 | 0.10 | 0.17 | 77 |
| 182 | 0.26 | 0.06 | 0.10 | 182 |
| 183 | 0.60 | 0.32 | 0.41 | 257 |
| 184 | 0.21 | 0.03 | 0.05 | 216 |
| 185 | 0.35 | 0.09 | 0.14 | 242 |
| 186 | 0.43 | 0.18 | 0.25 | 165 |
| 187 | 0.75 | 0.59 | 0.66 | 263 |

| | | | | |
|---|---|---|---|---|
| 188 | 0.39 | 0.12 | 0.18 | 174 |
| 189 | 0.75 | 0.40 | 0.53 | 136 |
| 190 | 0.89 | 0.55 | 0.68 | 202 |
| 191 | 0.44 | 0.16 | 0.24 | 134 |
| 192 | 0.68 | 0.40 | 0.51 | 230 |
| 193 | 0.44 | 0.18 | 0.25 | 90 |
| 194 | 0.57 | 0.48 | 0.52 | 185 |
| 195 | 0.26 | 0.05 | 0.09 | 156 |
| 196 | 0.33 | 0.07 | 0.11 | 160 |
| 197 | 0.49 | 0.10 | 0.16 | 266 |
| 198 | 0.47 | 0.13 | 0.20 | 284 |
| 199 | 0.32 | 0.04 | 0.07 | 145 |
| 200 | 0.93 | 0.74 | 0.82 | 212 |
| 201 | 0.65 | 0.26 | 0.37 | 317 |
| 202 | 0.78 | 0.59 | 0.67 | 427 |
| 203 | 0.36 | 0.11 | 0.17 | 232 |
| 204 | 0.51 | 0.29 | 0.37 | 217 |
| 205 | 0.50 | 0.46 | 0.48 | 527 |
| 206 | 0.24 | 0.03 | 0.06 | 124 |
| 207 | 0.50 | 0.17 | 0.26 | 103 |
| 208 | 0.85 | 0.53 | 0.65 | 287 |
| 209 | 0.33 | 0.11 | 0.16 | 193 |
| 210 | 0.75 | 0.38 | 0.50 | 220 |
| 211 | 0.71 | 0.21 | 0.32 | 140 |
| 212 | 0.12 | 0.02 | 0.03 | 161 |
| 213 | 0.63 | 0.43 | 0.51 | 72 |
| 214 | 0.64 | 0.45 | 0.53 | 396 |
| 215 | 0.87 | 0.34 | 0.49 | 134 |
| 216 | 0.61 | 0.17 | 0.27 | 400 |
| 217 | 0.51 | 0.24 | 0.33 | 75 |
| 218 | 0.96 | 0.76 | 0.85 | 219 |
| 219 | 0.77 | 0.42 | 0.54 | 210 |
| 220 | 0.88 | 0.64 | 0.74 | 298 |
| 221 | 0.96 | 0.70 | 0.81 | 266 |
| 222 | 0.76 | 0.45 | 0.57 | 290 |
| 223 | 0.11 | 0.01 | 0.01 | 128 |
| 224 | 0.78 | 0.45 | 0.57 | 159 |
| 225 | 0.55 | 0.29 | 0.38 | 164 |
| 226 | 0.58 | 0.31 | 0.41 | 144 |
| 227 | 0.56 | 0.29 | 0.38 | 276 |
| 228 | 0.19 | 0.03 | 0.05 | 235 |
| 229 | 0.33 | 0.03 | 0.06 | 216 |
| 230 | 0.40 | 0.17 | 0.23 | 228 |
| 231 | 0.70 | 0.48 | 0.57 | 64 |
| 232 | 0.48 | 0.10 | 0.16 | 103 |
| 233 | 0.72 | 0.35 | 0.47 | 216 |
| 234 | 0.72 | 0.11 | 0.19 | 116 |
| 235 | 0.54 | 0.36 | 0.43 | 77 |

| | | | | |
|-----|------|------|------|-----|
| 236 | 0.90 | 0.67 | 0.77 | 67 |
| 237 | 0.57 | 0.12 | 0.20 | 218 |
| 238 | 0.40 | 0.14 | 0.20 | 139 |
| 239 | 0.00 | 0.00 | 0.00 | 94 |
| 240 | 0.54 | 0.34 | 0.42 | 77 |
| 241 | 0.47 | 0.08 | 0.14 | 167 |
| 242 | 0.78 | 0.37 | 0.50 | 86 |
| 243 | 0.40 | 0.10 | 0.16 | 58 |
| 244 | 0.62 | 0.27 | 0.38 | 269 |
| 245 | 0.16 | 0.04 | 0.07 | 112 |
| 246 | 0.95 | 0.76 | 0.84 | 255 |
| 247 | 0.44 | 0.24 | 0.31 | 58 |
| 248 | 0.44 | 0.05 | 0.09 | 81 |
| 249 | 0.23 | 0.02 | 0.04 | 131 |
| 250 | 0.43 | 0.24 | 0.31 | 93 |
| 251 | 0.61 | 0.29 | 0.39 | 154 |
| 252 | 0.36 | 0.04 | 0.07 | 129 |
| 253 | 0.69 | 0.40 | 0.50 | 83 |
| 254 | 0.34 | 0.08 | 0.13 | 191 |
| 255 | 0.15 | 0.03 | 0.05 | 219 |
| 256 | 0.32 | 0.05 | 0.09 | 130 |
| 257 | 0.48 | 0.26 | 0.34 | 93 |
| 258 | 0.65 | 0.48 | 0.55 | 217 |
| 259 | 0.41 | 0.13 | 0.20 | 141 |
| 260 | 0.86 | 0.17 | 0.29 | 143 |
| 261 | 0.62 | 0.17 | 0.27 | 219 |
| 262 | 0.55 | 0.27 | 0.36 | 107 |
| 263 | 0.41 | 0.27 | 0.32 | 236 |
| 264 | 0.32 | 0.22 | 0.26 | 119 |
| 265 | 0.57 | 0.24 | 0.33 | 72 |
| 266 | 0.00 | 0.00 | 0.00 | 70 |
| 267 | 0.36 | 0.14 | 0.20 | 107 |
| 268 | 0.67 | 0.44 | 0.53 | 169 |
| 269 | 0.32 | 0.14 | 0.19 | 129 |
| 270 | 0.74 | 0.53 | 0.62 | 159 |
| 271 | 0.88 | 0.48 | 0.62 | 190 |
| 272 | 0.61 | 0.27 | 0.37 | 248 |
| 273 | 0.90 | 0.75 | 0.82 | 264 |
| 274 | 0.90 | 0.68 | 0.77 | 105 |
| 275 | 0.52 | 0.12 | 0.20 | 104 |
| 276 | 0.08 | 0.01 | 0.02 | 115 |
| 277 | 0.83 | 0.63 | 0.72 | 170 |
| 278 | 0.74 | 0.41 | 0.52 | 145 |
| 279 | 0.90 | 0.70 | 0.78 | 230 |
| 280 | 0.58 | 0.42 | 0.49 | 80 |
| 281 | 0.66 | 0.54 | 0.59 | 217 |
| 282 | 0.75 | 0.50 | 0.60 | 175 |
| 283 | 0.33 | 0.13 | 0.18 | 269 |

| | | | | |
|---|---|---|---|---|
| 284 | 0.65 | 0.32 | 0.43 | 74 |
| 285 | 0.82 | 0.49 | 0.61 | 206 |
| 286 | 0.89 | 0.66 | 0.75 | 227 |
| 287 | 0.84 | 0.41 | 0.55 | 130 |
| 288 | 0.32 | 0.07 | 0.11 | 129 |
| 289 | 0.57 | 0.05 | 0.09 | 80 |
| 290 | 0.21 | 0.09 | 0.13 | 99 |
| 291 | 0.76 | 0.35 | 0.48 | 208 |
| 292 | 0.42 | 0.07 | 0.13 | 67 |
| 293 | 0.84 | 0.48 | 0.61 | 109 |
| 294 | 0.46 | 0.26 | 0.34 | 140 |
| 295 | 0.24 | 0.12 | 0.16 | 241 |
| 296 | 0.31 | 0.12 | 0.18 | 72 |
| 297 | 0.44 | 0.11 | 0.18 | 107 |
| 298 | 0.77 | 0.49 | 0.60 | 61 |
| 299 | 0.89 | 0.51 | 0.64 | 77 |
| 300 | 0.21 | 0.08 | 0.12 | 111 |
| 301 | 0.00 | 0.00 | 0.00 | 126 |
| 302 | 0.25 | 0.01 | 0.03 | 73 |
| 303 | 0.57 | 0.43 | 0.49 | 176 |
| 304 | 0.91 | 0.79 | 0.85 | 230 |
| 305 | 0.92 | 0.72 | 0.81 | 156 |
| 306 | 0.50 | 0.37 | 0.43 | 146 |
| 307 | 0.34 | 0.11 | 0.17 | 98 |
| 308 | 0.00 | 0.00 | 0.00 | 78 |
| 309 | 0.80 | 0.13 | 0.22 | 94 |
| 310 | 0.74 | 0.41 | 0.53 | 162 |
| 311 | 0.79 | 0.51 | 0.62 | 116 |
| 312 | 0.52 | 0.28 | 0.36 | 57 |
| 313 | 0.83 | 0.08 | 0.14 | 65 |
| 314 | 0.52 | 0.36 | 0.42 | 138 |
| 315 | 0.54 | 0.22 | 0.31 | 195 |
| 316 | 0.56 | 0.35 | 0.43 | 69 |
| 317 | 0.29 | 0.13 | 0.18 | 134 |
| 318 | 0.56 | 0.39 | 0.46 | 148 |
| 319 | 0.84 | 0.50 | 0.63 | 161 |
| 320 | 0.24 | 0.19 | 0.21 | 104 |
| 321 | 0.82 | 0.61 | 0.70 | 156 |
| 322 | 0.60 | 0.37 | 0.46 | 134 |
| 323 | 0.58 | 0.44 | 0.50 | 232 |
| 324 | 0.34 | 0.15 | 0.21 | 92 |
| 325 | 0.41 | 0.24 | 0.31 | 197 |
| 326 | 0.14 | 0.03 | 0.05 | 126 |
| 327 | 0.20 | 0.03 | 0.05 | 115 |
| 328 | 0.99 | 0.70 | 0.82 | 198 |
| 329 | 0.59 | 0.32 | 0.41 | 125 |
| 330 | 0.73 | 0.20 | 0.31 | 81 |
| 331 | 0.45 | 0.10 | 0.16 | 94 |

| | | | |
|---|---|---|---|
| 332 | 0.54 | 0.12 | 0.20 | 56 |
| 333 | 0.19 | 0.05 | 0.08 | 260 |
| 334 | 0.42 | 0.13 | 0.20 | 60 |
| 335 | 0.35 | 0.08 | 0.13 | 110 |
| 336 | 0.62 | 0.49 | 0.55 | 71 |
| 337 | 0.18 | 0.05 | 0.07 | 66 |
| 338 | 0.47 | 0.36 | 0.41 | 150 |
| 339 | 0.00 | 0.00 | 0.00 | 54 |
| 340 | 0.84 | 0.57 | 0.68 | 195 |
| 341 | 0.91 | 0.52 | 0.66 | 79 |
| 342 | 0.38 | 0.26 | 0.31 | 38 |
| 343 | 0.62 | 0.42 | 0.50 | 43 |
| 344 | 0.56 | 0.29 | 0.38 | 68 |
| 345 | 0.62 | 0.33 | 0.43 | 73 |
| 346 | 0.14 | 0.03 | 0.04 | 116 |
| 347 | 0.86 | 0.43 | 0.57 | 111 |
| 348 | 0.33 | 0.11 | 0.17 | 63 |
| 349 | 0.84 | 0.65 | 0.74 | 104 |
| 350 | 0.62 | 0.48 | 0.54 | 44 |
| 351 | 0.57 | 0.30 | 0.39 | 40 |
| 352 | 0.93 | 0.57 | 0.70 | 136 |
| 353 | 0.38 | 0.15 | 0.21 | 54 |
| 354 | 0.39 | 0.09 | 0.15 | 134 |
| 355 | 0.64 | 0.35 | 0.45 | 120 |
| 356 | 0.54 | 0.30 | 0.38 | 228 |
| 357 | 0.66 | 0.36 | 0.47 | 269 |
| 358 | 0.62 | 0.38 | 0.47 | 80 |
| 359 | 0.84 | 0.59 | 0.69 | 140 |
| 360 | 0.39 | 0.18 | 0.24 | 125 |
| 361 | 0.90 | 0.71 | 0.79 | 169 |
| 362 | 0.14 | 0.05 | 0.08 | 56 |
| 363 | 0.92 | 0.73 | 0.82 | 154 |
| 364 | 0.46 | 0.10 | 0.17 | 58 |
| 365 | 0.22 | 0.08 | 0.12 | 71 |
| 366 | 1.00 | 0.69 | 0.81 | 54 |
| 367 | 0.31 | 0.07 | 0.11 | 116 |
| 368 | 0.38 | 0.06 | 0.10 | 54 |
| 369 | 0.33 | 0.03 | 0.05 | 71 |
| 370 | 0.00 | 0.00 | 0.00 | 61 |
| 371 | 0.40 | 0.08 | 0.14 | 71 |
| 372 | 0.72 | 0.44 | 0.55 | 52 |
| 373 | 0.78 | 0.41 | 0.54 | 150 |
| 374 | 0.41 | 0.14 | 0.21 | 93 |
| 375 | 0.20 | 0.04 | 0.07 | 67 |
| 376 | 0.00 | 0.00 | 0.00 | 76 |
| 377 | 0.58 | 0.28 | 0.38 | 106 |
| 378 | 0.25 | 0.02 | 0.04 | 86 |
| 379 | 0.50 | 0.14 | 0.22 | 14 |

| | | | | |
|---|---|---|---|---|
| 380 | 0.93 | 0.52 | 0.67 | 122 |
| 381 | 0.23 | 0.07 | 0.10 | 104 |
| 382 | 0.46 | 0.20 | 0.28 | 66 |
| 383 | 0.54 | 0.35 | 0.42 | 110 |
| 384 | 0.14 | 0.01 | 0.01 | 155 |
| 385 | 0.69 | 0.22 | 0.33 | 50 |
| 386 | 0.20 | 0.06 | 0.10 | 64 |
| 387 | 0.32 | 0.08 | 0.12 | 93 |
| 388 | 0.53 | 0.24 | 0.33 | 102 |
| 389 | 0.07 | 0.01 | 0.02 | 108 |
| 390 | 0.96 | 0.68 | 0.80 | 178 |
| 391 | 0.49 | 0.17 | 0.26 | 115 |
| 392 | 0.81 | 0.40 | 0.54 | 42 |
| 393 | 0.00 | 0.00 | 0.00 | 134 |
| 394 | 0.22 | 0.04 | 0.06 | 112 |
| 395 | 0.54 | 0.27 | 0.36 | 176 |
| 396 | 0.47 | 0.13 | 0.20 | 125 |
| 397 | 0.74 | 0.37 | 0.49 | 224 |
| 398 | 0.84 | 0.67 | 0.74 | 63 |
| 399 | 0.30 | 0.05 | 0.09 | 59 |
| 400 | 0.51 | 0.32 | 0.39 | 63 |
| 401 | 0.50 | 0.24 | 0.33 | 98 |
| 402 | 0.51 | 0.19 | 0.27 | 162 |
| 403 | 0.38 | 0.14 | 0.21 | 83 |
| 404 | 0.76 | 0.84 | 0.80 | 19 |
| 405 | 0.34 | 0.11 | 0.17 | 92 |
| 406 | 0.69 | 0.22 | 0.33 | 41 |
| 407 | 0.64 | 0.37 | 0.47 | 43 |
| 408 | 0.80 | 0.46 | 0.58 | 160 |
| 409 | 0.20 | 0.12 | 0.15 | 50 |
| 410 | 0.00 | 0.00 | 0.00 | 19 |
| 411 | 0.36 | 0.11 | 0.17 | 175 |
| 412 | 0.28 | 0.07 | 0.11 | 72 |
| 413 | 0.38 | 0.05 | 0.09 | 95 |
| 414 | 0.12 | 0.02 | 0.04 | 97 |
| 415 | 0.33 | 0.10 | 0.16 | 48 |
| 416 | 0.53 | 0.35 | 0.42 | 83 |
| 417 | 0.43 | 0.07 | 0.13 | 40 |
| 418 | 0.48 | 0.16 | 0.25 | 91 |
| 419 | 0.53 | 0.37 | 0.43 | 90 |
| 420 | 0.38 | 0.27 | 0.32 | 37 |
| 421 | 0.04 | 0.02 | 0.02 | 66 |
| 422 | 0.69 | 0.45 | 0.55 | 73 |
| 423 | 0.48 | 0.25 | 0.33 | 56 |
| 424 | 0.94 | 0.88 | 0.91 | 33 |
| 425 | 0.00 | 0.00 | 0.00 | 76 |
| 426 | 0.27 | 0.05 | 0.08 | 81 |
| 427 | 0.98 | 0.73 | 0.84 | 150 |

| | | | |
|---|---|---|---|
| 428 | 0.95 | 0.69 | 0.80 | 29 |
| 429 | 0.99 | 0.93 | 0.96 | 389 |
| 430 | 0.63 | 0.40 | 0.49 | 167 |
| 431 | 0.57 | 0.11 | 0.18 | 123 |
| 432 | 0.52 | 0.31 | 0.39 | 39 |
| 433 | 0.33 | 0.21 | 0.25 | 82 |
| 434 | 1.00 | 0.70 | 0.82 | 66 |
| 435 | 0.55 | 0.38 | 0.45 | 93 |
| 436 | 0.56 | 0.37 | 0.44 | 87 |
| 437 | 0.10 | 0.02 | 0.04 | 86 |
| 438 | 0.72 | 0.53 | 0.61 | 104 |
| 439 | 0.54 | 0.13 | 0.21 | 100 |
| 440 | 0.38 | 0.04 | 0.06 | 141 |
| 441 | 0.43 | 0.33 | 0.37 | 110 |
| 442 | 0.37 | 0.15 | 0.22 | 123 |
| 443 | 0.57 | 0.18 | 0.28 | 71 |
| 444 | 0.32 | 0.06 | 0.11 | 109 |
| 445 | 0.45 | 0.31 | 0.37 | 48 |
| 446 | 0.47 | 0.29 | 0.36 | 76 |
| 447 | 0.39 | 0.18 | 0.25 | 38 |
| 448 | 0.67 | 0.54 | 0.60 | 81 |
| 449 | 0.67 | 0.26 | 0.37 | 132 |
| 450 | 0.42 | 0.27 | 0.33 | 81 |
| 451 | 0.89 | 0.32 | 0.47 | 76 |
| 452 | 0.00 | 0.00 | 0.00 | 44 |
| 453 | 0.00 | 0.00 | 0.00 | 44 |
| 454 | 0.84 | 0.51 | 0.64 | 70 |
| 455 | 0.39 | 0.18 | 0.25 | 155 |
| 456 | 0.50 | 0.21 | 0.30 | 43 |
| 457 | 0.54 | 0.28 | 0.37 | 72 |
| 458 | 0.35 | 0.13 | 0.19 | 62 |
| 459 | 0.63 | 0.25 | 0.35 | 69 |
| 460 | 0.00 | 0.00 | 0.00 | 119 |
| 461 | 0.71 | 0.19 | 0.30 | 79 |
| 462 | 0.61 | 0.23 | 0.34 | 47 |
| 463 | 0.39 | 0.14 | 0.21 | 104 |
| 464 | 0.70 | 0.42 | 0.52 | 106 |
| 465 | 0.64 | 0.22 | 0.33 | 64 |
| 466 | 0.55 | 0.35 | 0.43 | 173 |
| 467 | 0.78 | 0.42 | 0.55 | 107 |
| 468 | 0.56 | 0.26 | 0.36 | 126 |
| 469 | 0.20 | 0.01 | 0.02 | 114 |
| 470 | 0.93 | 0.81 | 0.87 | 140 |
| 471 | 0.85 | 0.42 | 0.56 | 79 |
| 472 | 0.40 | 0.35 | 0.37 | 143 |
| 473 | 0.67 | 0.37 | 0.47 | 158 |
| 474 | 0.48 | 0.10 | 0.17 | 138 |
| 475 | 0.00 | 0.00 | 0.00 | 59 |

```
476       0.63       0.33       0.43        88
477       0.83       0.65       0.73       176
478       0.95       0.79       0.86        24
479       0.22       0.04       0.07        92
480       0.79       0.50       0.61       100
481       0.51       0.28       0.36       103
482       0.40       0.22       0.28        74
483       0.78       0.63       0.69       105
484       0.20       0.02       0.04        83
485       0.20       0.02       0.04        82
486       0.48       0.15       0.23        71
487       0.45       0.21       0.29       120
488       0.50       0.06       0.10       105
489       0.73       0.37       0.49        87
490       1.00       0.81       0.90        32
491       0.33       0.03       0.05        69
492       0.33       0.02       0.04        49
493       0.11       0.02       0.03       117
494       0.52       0.23       0.32        61
495       0.95       0.79       0.87       344
496       0.32       0.13       0.19        52
497       0.59       0.28       0.38       137
498       0.31       0.10       0.15        98
499       0.48       0.20       0.29        79

   micro avg       0.72       0.37       0.49    173812
   macro avg       0.56       0.30       0.37    173812
weighted avg       0.67       0.37       0.46    173812
 samples avg       0.46       0.35       0.37    173812

Time taken to run this cell : 0:57:40.281828
```

[78]:
```python
start = datetime.now()
classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1',C=1.0))
classifier_2.fit(x_train_multilabel_BOW, y_train)
predictions_2 = classifier_2.predict(x_test_multilabel_BOW)
print("Accuracy :",metrics.accuracy_score(y_test, predictions_2))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions_2))


precision = precision_score(y_test, predictions_2, average='micro')
recall = recall_score(y_test, predictions_2, average='micro')
f1 = f1_score(y_test, predictions_2, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,
 →recall, f1))
```

```
precision = precision_score(y_test, predictions_2, average='macro')
recall = recall_score(y_test, predictions_2, average='macro')
f1 = f1_score(y_test, predictions_2, average='macro')

print("Macro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision,
 →recall, f1))

print (metrics.classification_report(y_test, predictions_2))
print("Time taken to run this cell :", datetime.now() - start)
```

```
Accuracy : 0.20973
Hamming loss  0.00316104
Micro-average quality numbers
Precision: 0.5611, Recall: 0.4161, F1-measure: 0.4779
Macro-average quality numbers
Precision: 0.4587, Recall: 0.3409, F1-measure: 0.3870
```

|    | precision | recall | f1-score | support |
|----|-----------|--------|----------|---------|
| 0  | 0.91      | 0.74   | 0.82     | 5519    |
| 1  | 0.40      | 0.51   | 0.45     | 8190    |
| 2  | 0.66      | 0.45   | 0.53     | 6529    |
| 3  | 0.48      | 0.50   | 0.49     | 3231    |
| 4  | 0.67      | 0.49   | 0.57     | 6430    |
| 5  | 0.65      | 0.43   | 0.52     | 2879    |
| 6  | 0.76      | 0.56   | 0.64     | 5086    |
| 7  | 0.77      | 0.62   | 0.69     | 4533    |
| 8  | 0.34      | 0.19   | 0.24     | 3000    |
| 9  | 0.72      | 0.58   | 0.64     | 2765    |
| 10 | 0.43      | 0.29   | 0.35     | 3051    |
| 11 | 0.61      | 0.44   | 0.51     | 3009    |
| 12 | 0.50      | 0.35   | 0.41     | 2630    |
| 13 | 0.37      | 0.31   | 0.34     | 1426    |
| 14 | 0.82      | 0.63   | 0.71     | 2548    |
| 15 | 0.48      | 0.28   | 0.35     | 2371    |
| 16 | 0.52      | 0.28   | 0.37     | 873     |
| 17 | 0.80      | 0.67   | 0.73     | 2151    |
| 18 | 0.43      | 0.29   | 0.34     | 2204    |
| 19 | 0.57      | 0.44   | 0.50     | 831     |
| 20 | 0.71      | 0.52   | 0.60     | 1860    |
| 21 | 0.28      | 0.18   | 0.22     | 2023    |
| 22 | 0.40      | 0.29   | 0.34     | 1513    |
| 23 | 0.79      | 0.60   | 0.68     | 1207    |
| 24 | 0.45      | 0.35   | 0.39     | 506     |
| 25 | 0.55      | 0.35   | 0.43     | 425     |
| 26 | 0.59      | 0.46   | 0.51     | 793     |

| | | | | |
|---|---|---|---|---|
| 27 | 0.53 | 0.40 | 0.45 | 1291 |
| 28 | 0.62 | 0.43 | 0.51 | 1208 |
| 29 | 0.28 | 0.17 | 0.22 | 406 |
| 30 | 0.48 | 0.25 | 0.33 | 504 |
| 31 | 0.23 | 0.14 | 0.18 | 732 |
| 32 | 0.49 | 0.35 | 0.41 | 441 |
| 33 | 0.52 | 0.33 | 0.40 | 1645 |
| 34 | 0.48 | 0.29 | 0.36 | 1058 |
| 35 | 0.77 | 0.60 | 0.67 | 946 |
| 36 | 0.47 | 0.33 | 0.39 | 644 |
| 37 | 0.95 | 0.74 | 0.83 | 136 |
| 38 | 0.53 | 0.41 | 0.46 | 570 |
| 39 | 0.63 | 0.34 | 0.44 | 766 |
| 40 | 0.55 | 0.41 | 0.47 | 1132 |
| 41 | 0.34 | 0.25 | 0.29 | 174 |
| 42 | 0.68 | 0.55 | 0.61 | 210 |
| 43 | 0.66 | 0.45 | 0.54 | 433 |
| 44 | 0.57 | 0.47 | 0.52 | 626 |
| 45 | 0.57 | 0.40 | 0.47 | 852 |
| 46 | 0.66 | 0.46 | 0.54 | 534 |
| 47 | 0.30 | 0.26 | 0.28 | 350 |
| 48 | 0.67 | 0.53 | 0.59 | 496 |
| 49 | 0.77 | 0.62 | 0.68 | 785 |
| 50 | 0.20 | 0.11 | 0.14 | 475 |
| 51 | 0.28 | 0.22 | 0.24 | 305 |
| 52 | 0.24 | 0.10 | 0.14 | 251 |
| 53 | 0.56 | 0.42 | 0.48 | 914 |
| 54 | 0.40 | 0.24 | 0.30 | 728 |
| 55 | 0.18 | 0.10 | 0.13 | 258 |
| 56 | 0.36 | 0.29 | 0.32 | 821 |
| 57 | 0.34 | 0.18 | 0.24 | 541 |
| 58 | 0.60 | 0.35 | 0.44 | 748 |
| 59 | 0.90 | 0.75 | 0.82 | 724 |
| 60 | 0.32 | 0.19 | 0.24 | 660 |
| 61 | 0.52 | 0.29 | 0.37 | 235 |
| 62 | 0.87 | 0.75 | 0.81 | 718 |
| 63 | 0.77 | 0.69 | 0.73 | 468 |
| 64 | 0.47 | 0.36 | 0.40 | 191 |
| 65 | 0.30 | 0.17 | 0.22 | 429 |
| 66 | 0.23 | 0.15 | 0.18 | 415 |
| 67 | 0.68 | 0.54 | 0.61 | 274 |
| 68 | 0.75 | 0.57 | 0.65 | 510 |
| 69 | 0.60 | 0.49 | 0.54 | 466 |
| 70 | 0.25 | 0.15 | 0.19 | 305 |
| 71 | 0.33 | 0.22 | 0.26 | 247 |
| 72 | 0.66 | 0.53 | 0.59 | 401 |
| 73 | 0.87 | 0.83 | 0.85 | 86 |
| 74 | 0.57 | 0.46 | 0.51 | 120 |

| | | | | |
|-----|------|------|------|-----|
| 75 | 0.87 | 0.71 | 0.78 | 129 |
| 76 | 0.19 | 0.07 | 0.10 | 473 |
| 77 | 0.40 | 0.36 | 0.38 | 143 |
| 78 | 0.71 | 0.47 | 0.56 | 347 |
| 79 | 0.51 | 0.28 | 0.36 | 479 |
| 80 | 0.46 | 0.44 | 0.45 | 279 |
| 81 | 0.57 | 0.31 | 0.40 | 461 |
| 82 | 0.14 | 0.06 | 0.09 | 298 |
| 83 | 0.68 | 0.54 | 0.60 | 396 |
| 84 | 0.37 | 0.41 | 0.39 | 184 |
| 85 | 0.48 | 0.31 | 0.37 | 573 |
| 86 | 0.27 | 0.12 | 0.17 | 325 |
| 87 | 0.49 | 0.37 | 0.42 | 273 |
| 88 | 0.39 | 0.27 | 0.32 | 135 |
| 89 | 0.29 | 0.19 | 0.23 | 232 |
| 90 | 0.48 | 0.40 | 0.44 | 409 |
| 91 | 0.49 | 0.30 | 0.38 | 420 |
| 92 | 0.70 | 0.60 | 0.64 | 408 |
| 93 | 0.55 | 0.52 | 0.54 | 241 |
| 94 | 0.21 | 0.08 | 0.12 | 211 |
| 95 | 0.30 | 0.18 | 0.22 | 277 |
| 96 | 0.22 | 0.12 | 0.16 | 410 |
| 97 | 0.78 | 0.51 | 0.62 | 501 |
| 98 | 0.71 | 0.66 | 0.68 | 136 |
| 99 | 0.42 | 0.34 | 0.38 | 239 |
| 100 | 0.33 | 0.21 | 0.25 | 324 |
| 101 | 0.87 | 0.75 | 0.81 | 277 |
| 102 | 0.88 | 0.78 | 0.83 | 613 |
| 103 | 0.34 | 0.19 | 0.24 | 157 |
| 104 | 0.26 | 0.16 | 0.20 | 295 |
| 105 | 0.67 | 0.46 | 0.55 | 334 |
| 106 | 0.74 | 0.36 | 0.48 | 335 |
| 107 | 0.69 | 0.59 | 0.64 | 389 |
| 108 | 0.55 | 0.34 | 0.42 | 251 |
| 109 | 0.53 | 0.47 | 0.50 | 317 |
| 110 | 0.33 | 0.11 | 0.16 | 187 |
| 111 | 0.39 | 0.14 | 0.21 | 140 |
| 112 | 0.62 | 0.50 | 0.55 | 154 |
| 113 | 0.46 | 0.27 | 0.34 | 332 |
| 114 | 0.42 | 0.32 | 0.36 | 323 |
| 115 | 0.39 | 0.29 | 0.33 | 344 |
| 116 | 0.68 | 0.57 | 0.62 | 370 |
| 117 | 0.45 | 0.29 | 0.35 | 313 |
| 118 | 0.75 | 0.74 | 0.74 | 874 |
| 119 | 0.34 | 0.29 | 0.32 | 293 |
| 120 | 0.13 | 0.10 | 0.11 | 200 |
| 121 | 0.70 | 0.55 | 0.62 | 463 |
| 122 | 0.28 | 0.14 | 0.19 | 119 |

| 123 | 0.03 | 0.01 | 0.01 | 256 |
|-----|------|------|------|-----|
| 124 | 0.85 | 0.70 | 0.77 | 195 |
| 125 | 0.33 | 0.20 | 0.25 | 138 |
| 126 | 0.67 | 0.59 | 0.63 | 376 |
| 127 | 0.14 | 0.05 | 0.07 | 122 |
| 128 | 0.15 | 0.07 | 0.09 | 252 |
| 129 | 0.37 | 0.29 | 0.33 | 144 |
| 130 | 0.32 | 0.18 | 0.23 | 150 |
| 131 | 0.18 | 0.07 | 0.10 | 210 |
| 132 | 0.50 | 0.34 | 0.40 | 361 |
| 133 | 0.87 | 0.67 | 0.75 | 453 |
| 134 | 0.78 | 0.77 | 0.77 | 124 |
| 135 | 0.18 | 0.14 | 0.16 | 91 |
| 136 | 0.51 | 0.39 | 0.44 | 128 |
| 137 | 0.47 | 0.42 | 0.45 | 218 |
| 138 | 0.29 | 0.17 | 0.22 | 243 |
| 139 | 0.32 | 0.21 | 0.26 | 149 |
| 140 | 0.71 | 0.54 | 0.62 | 318 |
| 141 | 0.19 | 0.13 | 0.15 | 159 |
| 142 | 0.57 | 0.38 | 0.46 | 274 |
| 143 | 0.82 | 0.83 | 0.83 | 362 |
| 144 | 0.43 | 0.25 | 0.32 | 118 |
| 145 | 0.52 | 0.43 | 0.47 | 164 |
| 146 | 0.52 | 0.39 | 0.45 | 461 |
| 147 | 0.63 | 0.41 | 0.50 | 159 |
| 148 | 0.28 | 0.17 | 0.22 | 166 |
| 149 | 0.89 | 0.58 | 0.70 | 346 |
| 150 | 0.51 | 0.22 | 0.30 | 350 |
| 151 | 0.86 | 0.78 | 0.82 | 55 |
| 152 | 0.71 | 0.52 | 0.60 | 387 |
| 153 | 0.41 | 0.40 | 0.40 | 150 |
| 154 | 0.34 | 0.14 | 0.20 | 281 |
| 155 | 0.26 | 0.17 | 0.21 | 202 |
| 156 | 0.73 | 0.65 | 0.69 | 130 |
| 157 | 0.29 | 0.13 | 0.18 | 245 |
| 158 | 0.88 | 0.71 | 0.79 | 177 |
| 159 | 0.43 | 0.41 | 0.42 | 130 |
| 160 | 0.40 | 0.21 | 0.28 | 336 |
| 161 | 0.80 | 0.67 | 0.73 | 220 |
| 162 | 0.20 | 0.13 | 0.16 | 229 |
| 163 | 0.85 | 0.55 | 0.67 | 316 |
| 164 | 0.65 | 0.48 | 0.55 | 283 |
| 165 | 0.49 | 0.34 | 0.40 | 197 |
| 166 | 0.62 | 0.57 | 0.59 | 101 |
| 167 | 0.37 | 0.25 | 0.29 | 231 |
| 168 | 0.45 | 0.35 | 0.39 | 370 |
| 169 | 0.38 | 0.25 | 0.30 | 258 |
| 170 | 0.28 | 0.18 | 0.22 | 101 |

| | | | |
|---|---|---|---|
| 171 | 0.33 | 0.25 | 0.28 | 89 |
| 172 | 0.48 | 0.33 | 0.39 | 193 |
| 173 | 0.47 | 0.33 | 0.39 | 309 |
| 174 | 0.30 | 0.16 | 0.21 | 172 |
| 175 | 0.84 | 0.77 | 0.80 | 95 |
| 176 | 0.86 | 0.64 | 0.73 | 346 |
| 177 | 0.73 | 0.54 | 0.62 | 322 |
| 178 | 0.55 | 0.45 | 0.50 | 232 |
| 179 | 0.17 | 0.10 | 0.13 | 125 |
| 180 | 0.47 | 0.40 | 0.43 | 145 |
| 181 | 0.32 | 0.21 | 0.25 | 77 |
| 182 | 0.18 | 0.10 | 0.13 | 182 |
| 183 | 0.51 | 0.38 | 0.43 | 257 |
| 184 | 0.19 | 0.10 | 0.13 | 216 |
| 185 | 0.31 | 0.19 | 0.24 | 242 |
| 186 | 0.36 | 0.23 | 0.28 | 165 |
| 187 | 0.66 | 0.56 | 0.61 | 263 |
| 188 | 0.27 | 0.14 | 0.18 | 174 |
| 189 | 0.67 | 0.45 | 0.54 | 136 |
| 190 | 0.83 | 0.59 | 0.69 | 202 |
| 191 | 0.39 | 0.22 | 0.28 | 134 |
| 192 | 0.56 | 0.45 | 0.50 | 230 |
| 193 | 0.25 | 0.20 | 0.22 | 90 |
| 194 | 0.57 | 0.52 | 0.54 | 185 |
| 195 | 0.23 | 0.13 | 0.17 | 156 |
| 196 | 0.14 | 0.08 | 0.10 | 160 |
| 197 | 0.20 | 0.10 | 0.13 | 266 |
| 198 | 0.30 | 0.14 | 0.19 | 284 |
| 199 | 0.19 | 0.09 | 0.12 | 145 |
| 200 | 0.86 | 0.78 | 0.82 | 212 |
| 201 | 0.54 | 0.33 | 0.41 | 317 |
| 202 | 0.70 | 0.62 | 0.65 | 427 |
| 203 | 0.19 | 0.12 | 0.15 | 232 |
| 204 | 0.41 | 0.30 | 0.35 | 217 |
| 205 | 0.49 | 0.46 | 0.47 | 527 |
| 206 | 0.11 | 0.05 | 0.07 | 124 |
| 207 | 0.37 | 0.29 | 0.33 | 103 |
| 208 | 0.76 | 0.53 | 0.63 | 287 |
| 209 | 0.24 | 0.14 | 0.18 | 193 |
| 210 | 0.57 | 0.42 | 0.48 | 220 |
| 211 | 0.56 | 0.24 | 0.34 | 140 |
| 212 | 0.20 | 0.14 | 0.16 | 161 |
| 213 | 0.49 | 0.56 | 0.52 | 72 |
| 214 | 0.62 | 0.45 | 0.52 | 396 |
| 215 | 0.59 | 0.38 | 0.46 | 134 |
| 216 | 0.47 | 0.26 | 0.34 | 400 |
| 217 | 0.36 | 0.24 | 0.29 | 75 |
| 218 | 0.95 | 0.76 | 0.84 | 219 |

| 219 | 0.61 | 0.42 | 0.50 | 210 |
|-----|------|------|------|-----|
| 220 | 0.89 | 0.69 | 0.78 | 298 |
| 221 | 0.90 | 0.73 | 0.81 | 266 |
| 222 | 0.66 | 0.46 | 0.54 | 290 |
| 223 | 0.13 | 0.05 | 0.08 | 128 |
| 224 | 0.65 | 0.45 | 0.54 | 159 |
| 225 | 0.41 | 0.31 | 0.35 | 164 |
| 226 | 0.47 | 0.36 | 0.41 | 144 |
| 227 | 0.55 | 0.38 | 0.45 | 276 |
| 228 | 0.11 | 0.05 | 0.07 | 235 |
| 229 | 0.14 | 0.05 | 0.07 | 216 |
| 230 | 0.34 | 0.19 | 0.24 | 228 |
| 231 | 0.66 | 0.52 | 0.58 | 64 |
| 232 | 0.20 | 0.16 | 0.17 | 103 |
| 233 | 0.60 | 0.38 | 0.46 | 216 |
| 234 | 0.44 | 0.21 | 0.28 | 116 |
| 235 | 0.46 | 0.34 | 0.39 | 77 |
| 236 | 0.83 | 0.78 | 0.80 | 67 |
| 237 | 0.30 | 0.18 | 0.23 | 218 |
| 238 | 0.24 | 0.16 | 0.19 | 139 |
| 239 | 0.27 | 0.06 | 0.10 | 94 |
| 240 | 0.43 | 0.27 | 0.33 | 77 |
| 241 | 0.15 | 0.07 | 0.09 | 167 |
| 242 | 0.66 | 0.47 | 0.54 | 86 |
| 243 | 0.28 | 0.24 | 0.26 | 58 |
| 244 | 0.53 | 0.38 | 0.44 | 269 |
| 245 | 0.14 | 0.09 | 0.11 | 112 |
| 246 | 0.92 | 0.81 | 0.86 | 255 |
| 247 | 0.23 | 0.22 | 0.23 | 58 |
| 248 | 0.21 | 0.09 | 0.12 | 81 |
| 249 | 0.06 | 0.02 | 0.03 | 131 |
| 250 | 0.34 | 0.26 | 0.29 | 93 |
| 251 | 0.50 | 0.32 | 0.39 | 154 |
| 252 | 0.13 | 0.05 | 0.08 | 129 |
| 253 | 0.50 | 0.34 | 0.40 | 83 |
| 254 | 0.17 | 0.09 | 0.12 | 191 |
| 255 | 0.10 | 0.05 | 0.07 | 219 |
| 256 | 0.16 | 0.09 | 0.12 | 130 |
| 257 | 0.45 | 0.30 | 0.36 | 93 |
| 258 | 0.63 | 0.52 | 0.57 | 217 |
| 259 | 0.26 | 0.22 | 0.24 | 141 |
| 260 | 0.45 | 0.27 | 0.34 | 143 |
| 261 | 0.45 | 0.18 | 0.26 | 219 |
| 262 | 0.52 | 0.39 | 0.45 | 107 |
| 263 | 0.43 | 0.28 | 0.34 | 236 |
| 264 | 0.25 | 0.19 | 0.22 | 119 |
| 265 | 0.44 | 0.35 | 0.39 | 72 |
| 266 | 0.20 | 0.09 | 0.12 | 70 |

| 267 | 0.38 | 0.23 | 0.29 | 107 |
| 268 | 0.55 | 0.48 | 0.51 | 169 |
| 269 | 0.30 | 0.16 | 0.21 | 129 |
| 270 | 0.71 | 0.53 | 0.61 | 159 |
| 271 | 0.76 | 0.53 | 0.63 | 190 |
| 272 | 0.46 | 0.32 | 0.38 | 248 |
| 273 | 0.86 | 0.77 | 0.81 | 264 |
| 274 | 0.85 | 0.69 | 0.76 | 105 |
| 275 | 0.18 | 0.11 | 0.13 | 104 |
| 276 | 0.04 | 0.02 | 0.02 | 115 |
| 277 | 0.79 | 0.64 | 0.71 | 170 |
| 278 | 0.69 | 0.54 | 0.60 | 145 |
| 279 | 0.84 | 0.76 | 0.80 | 230 |
| 280 | 0.53 | 0.45 | 0.49 | 80 |
| 281 | 0.63 | 0.51 | 0.56 | 217 |
| 282 | 0.67 | 0.53 | 0.59 | 175 |
| 283 | 0.28 | 0.18 | 0.22 | 269 |
| 284 | 0.52 | 0.35 | 0.42 | 74 |
| 285 | 0.73 | 0.56 | 0.63 | 206 |
| 286 | 0.83 | 0.70 | 0.76 | 227 |
| 287 | 0.70 | 0.45 | 0.55 | 130 |
| 288 | 0.25 | 0.10 | 0.14 | 129 |
| 289 | 0.22 | 0.12 | 0.16 | 80 |
| 290 | 0.19 | 0.15 | 0.17 | 99 |
| 291 | 0.65 | 0.40 | 0.50 | 208 |
| 292 | 0.31 | 0.13 | 0.19 | 67 |
| 293 | 0.76 | 0.57 | 0.65 | 109 |
| 294 | 0.34 | 0.31 | 0.33 | 140 |
| 295 | 0.22 | 0.14 | 0.17 | 241 |
| 296 | 0.25 | 0.15 | 0.19 | 72 |
| 297 | 0.27 | 0.16 | 0.20 | 107 |
| 298 | 0.63 | 0.61 | 0.62 | 61 |
| 299 | 0.78 | 0.55 | 0.64 | 77 |
| 300 | 0.16 | 0.14 | 0.15 | 111 |
| 301 | 0.06 | 0.02 | 0.02 | 126 |
| 302 | 0.12 | 0.08 | 0.10 | 73 |
| 303 | 0.58 | 0.44 | 0.50 | 176 |
| 304 | 0.90 | 0.84 | 0.87 | 230 |
| 305 | 0.86 | 0.71 | 0.77 | 156 |
| 306 | 0.40 | 0.36 | 0.38 | 146 |
| 307 | 0.30 | 0.17 | 0.22 | 98 |
| 308 | 0.09 | 0.03 | 0.04 | 78 |
| 309 | 0.40 | 0.18 | 0.25 | 94 |
| 310 | 0.55 | 0.40 | 0.46 | 162 |
| 311 | 0.69 | 0.53 | 0.60 | 116 |
| 312 | 0.48 | 0.35 | 0.40 | 57 |
| 313 | 0.25 | 0.06 | 0.10 | 65 |
| 314 | 0.42 | 0.38 | 0.40 | 138 |

| | | | | |
|---|---|---|---|---|
| 315 | 0.49 | 0.34 | 0.40 | 195 |
| 316 | 0.37 | 0.30 | 0.33 | 69 |
| 317 | 0.26 | 0.21 | 0.23 | 134 |
| 318 | 0.50 | 0.38 | 0.43 | 148 |
| 319 | 0.78 | 0.57 | 0.66 | 161 |
| 320 | 0.21 | 0.23 | 0.22 | 104 |
| 321 | 0.71 | 0.62 | 0.66 | 156 |
| 322 | 0.55 | 0.48 | 0.51 | 134 |
| 323 | 0.54 | 0.47 | 0.51 | 232 |
| 324 | 0.20 | 0.16 | 0.18 | 92 |
| 325 | 0.39 | 0.25 | 0.31 | 197 |
| 326 | 0.07 | 0.06 | 0.06 | 126 |
| 327 | 0.10 | 0.03 | 0.05 | 115 |
| 328 | 0.94 | 0.77 | 0.85 | 198 |
| 329 | 0.44 | 0.33 | 0.38 | 125 |
| 330 | 0.51 | 0.27 | 0.35 | 81 |
| 331 | 0.33 | 0.15 | 0.20 | 94 |
| 332 | 0.32 | 0.21 | 0.26 | 56 |
| 333 | 0.18 | 0.10 | 0.13 | 260 |
| 334 | 0.32 | 0.20 | 0.24 | 60 |
| 335 | 0.26 | 0.13 | 0.17 | 110 |
| 336 | 0.62 | 0.54 | 0.58 | 71 |
| 337 | 0.10 | 0.06 | 0.08 | 66 |
| 338 | 0.39 | 0.41 | 0.40 | 150 |
| 339 | 0.14 | 0.06 | 0.08 | 54 |
| 340 | 0.75 | 0.58 | 0.66 | 195 |
| 341 | 0.65 | 0.54 | 0.59 | 79 |
| 342 | 0.39 | 0.55 | 0.46 | 38 |
| 343 | 0.63 | 0.44 | 0.52 | 43 |
| 344 | 0.44 | 0.29 | 0.35 | 68 |
| 345 | 0.62 | 0.44 | 0.51 | 73 |
| 346 | 0.11 | 0.06 | 0.08 | 116 |
| 347 | 0.76 | 0.57 | 0.65 | 111 |
| 348 | 0.28 | 0.21 | 0.24 | 63 |
| 349 | 0.82 | 0.73 | 0.77 | 104 |
| 350 | 0.61 | 0.52 | 0.56 | 44 |
| 351 | 0.29 | 0.25 | 0.27 | 40 |
| 352 | 0.82 | 0.66 | 0.73 | 136 |
| 353 | 0.30 | 0.15 | 0.20 | 54 |
| 354 | 0.30 | 0.11 | 0.16 | 134 |
| 355 | 0.56 | 0.47 | 0.51 | 120 |
| 356 | 0.45 | 0.31 | 0.37 | 228 |
| 357 | 0.55 | 0.40 | 0.46 | 269 |
| 358 | 0.64 | 0.44 | 0.52 | 80 |
| 359 | 0.84 | 0.70 | 0.76 | 140 |
| 360 | 0.31 | 0.22 | 0.25 | 125 |
| 361 | 0.85 | 0.75 | 0.80 | 169 |
| 362 | 0.16 | 0.12 | 0.14 | 56 |

| 363 | 0.82 | 0.75 | 0.78 | 154 |
| 364 | 0.31 | 0.29 | 0.30 | 58 |
| 365 | 0.28 | 0.15 | 0.20 | 71 |
| 366 | 0.86 | 0.67 | 0.75 | 54 |
| 367 | 0.14 | 0.09 | 0.11 | 116 |
| 368 | 0.23 | 0.15 | 0.18 | 54 |
| 369 | 0.09 | 0.06 | 0.07 | 71 |
| 370 | 0.29 | 0.11 | 0.16 | 61 |
| 371 | 0.20 | 0.08 | 0.12 | 71 |
| 372 | 0.52 | 0.48 | 0.50 | 52 |
| 373 | 0.65 | 0.54 | 0.59 | 150 |
| 374 | 0.26 | 0.18 | 0.22 | 93 |
| 375 | 0.17 | 0.10 | 0.13 | 67 |
| 376 | 0.00 | 0.00 | 0.00 | 76 |
| 377 | 0.41 | 0.33 | 0.37 | 106 |
| 378 | 0.11 | 0.05 | 0.07 | 86 |
| 379 | 0.00 | 0.00 | 0.00 | 14 |
| 380 | 0.87 | 0.64 | 0.74 | 122 |
| 381 | 0.09 | 0.07 | 0.08 | 104 |
| 382 | 0.20 | 0.14 | 0.16 | 66 |
| 383 | 0.47 | 0.43 | 0.45 | 110 |
| 384 | 0.20 | 0.06 | 0.09 | 155 |
| 385 | 0.45 | 0.40 | 0.43 | 50 |
| 386 | 0.14 | 0.08 | 0.10 | 64 |
| 387 | 0.19 | 0.10 | 0.13 | 93 |
| 388 | 0.43 | 0.34 | 0.38 | 102 |
| 389 | 0.09 | 0.03 | 0.04 | 108 |
| 390 | 0.90 | 0.72 | 0.80 | 178 |
| 391 | 0.38 | 0.23 | 0.28 | 115 |
| 392 | 0.69 | 0.48 | 0.56 | 42 |
| 393 | 0.04 | 0.01 | 0.01 | 134 |
| 394 | 0.28 | 0.12 | 0.17 | 112 |
| 395 | 0.37 | 0.28 | 0.32 | 176 |
| 396 | 0.33 | 0.18 | 0.24 | 125 |
| 397 | 0.66 | 0.50 | 0.56 | 224 |
| 398 | 0.83 | 0.68 | 0.75 | 63 |
| 399 | 0.12 | 0.07 | 0.09 | 59 |
| 400 | 0.48 | 0.51 | 0.49 | 63 |
| 401 | 0.43 | 0.30 | 0.35 | 98 |
| 402 | 0.44 | 0.23 | 0.30 | 162 |
| 403 | 0.37 | 0.24 | 0.29 | 83 |
| 404 | 0.64 | 0.84 | 0.73 | 19 |
| 405 | 0.13 | 0.11 | 0.12 | 92 |
| 406 | 0.68 | 0.46 | 0.55 | 41 |
| 407 | 0.60 | 0.42 | 0.49 | 43 |
| 408 | 0.66 | 0.48 | 0.56 | 160 |
| 409 | 0.19 | 0.12 | 0.15 | 50 |
| 410 | 0.10 | 0.05 | 0.07 | 19 |

| | | | | |
|---|---|---|---|---|
| 411 | 0.24 | 0.19 | 0.21 | 175 |
| 412 | 0.24 | 0.17 | 0.20 | 72 |
| 413 | 0.25 | 0.13 | 0.17 | 95 |
| 414 | 0.26 | 0.12 | 0.17 | 97 |
| 415 | 0.15 | 0.10 | 0.12 | 48 |
| 416 | 0.49 | 0.40 | 0.44 | 83 |
| 417 | 0.24 | 0.20 | 0.22 | 40 |
| 418 | 0.28 | 0.14 | 0.19 | 91 |
| 419 | 0.50 | 0.43 | 0.46 | 90 |
| 420 | 0.34 | 0.30 | 0.32 | 37 |
| 421 | 0.15 | 0.11 | 0.12 | 66 |
| 422 | 0.49 | 0.40 | 0.44 | 73 |
| 423 | 0.45 | 0.32 | 0.38 | 56 |
| 424 | 0.93 | 0.85 | 0.89 | 33 |
| 425 | 0.14 | 0.07 | 0.09 | 76 |
| 426 | 0.06 | 0.02 | 0.03 | 81 |
| 427 | 0.95 | 0.75 | 0.84 | 150 |
| 428 | 0.85 | 0.76 | 0.80 | 29 |
| 429 | 0.99 | 0.98 | 0.99 | 389 |
| 430 | 0.57 | 0.43 | 0.49 | 167 |
| 431 | 0.41 | 0.13 | 0.20 | 123 |
| 432 | 0.27 | 0.18 | 0.22 | 39 |
| 433 | 0.28 | 0.26 | 0.27 | 82 |
| 434 | 0.92 | 0.68 | 0.78 | 66 |
| 435 | 0.52 | 0.44 | 0.48 | 93 |
| 436 | 0.51 | 0.44 | 0.47 | 87 |
| 437 | 0.21 | 0.12 | 0.15 | 86 |
| 438 | 0.60 | 0.48 | 0.53 | 104 |
| 439 | 0.43 | 0.21 | 0.28 | 100 |
| 440 | 0.16 | 0.07 | 0.10 | 141 |
| 441 | 0.44 | 0.46 | 0.45 | 110 |
| 442 | 0.28 | 0.23 | 0.25 | 123 |
| 443 | 0.39 | 0.25 | 0.31 | 71 |
| 444 | 0.22 | 0.12 | 0.15 | 109 |
| 445 | 0.43 | 0.33 | 0.38 | 48 |
| 446 | 0.43 | 0.36 | 0.39 | 76 |
| 447 | 0.19 | 0.21 | 0.20 | 38 |
| 448 | 0.57 | 0.57 | 0.57 | 81 |
| 449 | 0.47 | 0.27 | 0.34 | 132 |
| 450 | 0.42 | 0.37 | 0.39 | 81 |
| 451 | 0.66 | 0.38 | 0.48 | 76 |
| 452 | 0.04 | 0.02 | 0.03 | 44 |
| 453 | 0.10 | 0.02 | 0.04 | 44 |
| 454 | 0.58 | 0.56 | 0.57 | 70 |
| 455 | 0.32 | 0.25 | 0.28 | 155 |
| 456 | 0.44 | 0.33 | 0.37 | 43 |
| 457 | 0.40 | 0.32 | 0.35 | 72 |
| 458 | 0.27 | 0.18 | 0.21 | 62 |

```
459        0.40      0.28      0.33         69
460        0.07      0.03      0.04        119
461        0.71      0.43      0.54         79
462        0.32      0.23      0.27         47
463        0.30      0.22      0.26        104
464        0.64      0.44      0.53        106
465        0.42      0.27      0.33         64
466        0.48      0.32      0.39        173
467        0.58      0.46      0.51        107
468        0.48      0.31      0.38        126
469        0.10      0.04      0.05        114
470        0.93      0.81      0.87        140
471        0.67      0.49      0.57         79
472        0.40      0.42      0.41        143
473        0.62      0.38      0.47        158
474        0.20      0.08      0.11        138
475        0.20      0.14      0.16         59
476        0.62      0.43      0.51         88
477        0.79      0.69      0.73        176
478        0.90      0.75      0.82         24
479        0.26      0.16      0.20         92
480        0.68      0.58      0.63        100
481        0.44      0.39      0.41        103
482        0.30      0.16      0.21         74
483        0.71      0.67      0.69        105
484        0.25      0.10      0.14         83
485        0.04      0.02      0.03         82
486        0.29      0.21      0.24         71
487        0.37      0.21      0.27        120
488        0.28      0.09      0.13        105
489        0.54      0.31      0.39         87
490        0.93      0.84      0.89         32
491        0.07      0.03      0.04         69
492        0.12      0.04      0.06         49
493        0.09      0.05      0.07        117
494        0.43      0.31      0.36         61
495        0.94      0.86      0.90        344
496        0.20      0.13      0.16         52
497        0.46      0.37      0.41        137
498        0.36      0.17      0.23         98
499        0.35      0.22      0.27         79

   micro avg       0.56      0.42      0.48     173812
   macro avg       0.46      0.34      0.39     173812
weighted avg       0.55      0.42      0.47     173812
 samples avg       0.44      0.39      0.39     173812

Time taken to run this cell : 6:18:52.228275
```

```
[79]: from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Model", "TFIDF/BOW", "micro average"]
x.add_row(["SGD Linear Regression", "BOW", 0.37])
x.add_row(["SGD SVM", "BOW", 0.37])
x.add_row(["Linear Regression", "TFIDF", 0.49])
x.add_row(["Linear Regression", "BOW", 0.48])

x.border=True
print(x)
```

```
+-----------------------+-----------+---------------+
|         Model         | TFIDF/BOW | micro average |
+-----------------------+-----------+---------------+
| SGD Linear Regression |    BOW    |      0.37     |
|        SGD SVM        |    BOW    |      0.37     |
|   Linear Regression   |   TFIDF   |      0.49     |
|   Linear Regression   |    BOW    |      0.48     |
+-----------------------+-----------+---------------+
```

[ ]: