



```
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 sklearn.naive_bayes import GaussianNB
from datetime import datetime
```

## ▼ Stack Overflow: Tag Prediction

### 1. Business Problem

#### 1.1 Description

##### Description

Stack Overflow is the largest, most trusted online community for developers to learn, share their programming know

Stack Overflow is something which every programmer use one way or another. Each month, over 50 million developers knowledge, and build their careers. It features questions and answers on a wide range of topics in computer program users to ask and answer questions, and, through membership and active participation, to vote questions and answers a fashion similar to a wiki or Digg. As of April 2014 Stack Overflow has over 4,000,000 registered users, and it exceeds Based on the type of tags assigned to questions, the top eight most discussed topics on the site are: Java, JavaScript,

##### Problem Statement

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.

**Train.csv** contains 4 columns: Id,Title,Body,Tags.

**Test.csv** contains the same columns but without the Tags, which you are to predict.

**Size of Train.csv** - 6.75GB

**Size of Test.csv** - 2GB

**Number of rows in Train.csv** = 6034195

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

## Data Field Explanation

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

**Id** - Unique identifier for each question

**Title** - The question's title

**Body** - The body of the question

**Tags** - The tags associated with the question in a space-separated format (all lowercase, ampersands '&')

### 2.1.2 Example Data point

**Title:** Implementing Boundary Value Analysis of Software Testing in a C++ program?

**Body :**

```
#include<
iostream>\n
#include<
stdlib.h>\n\n
using namespace std;\n\n
int main()\n
{\n
    int n,a[n],x,c,u[n],m[n],e[n][4];\n
    cout<<"Enter the number of variables";\n
    cin>>n;\n
    cout<<"Enter the Lower, and Upper Limits of the variables";\n
    for(int y=1; y<n+1; y++)\n
    {\n
        cin>>m[y];\n
        cin>>u[y];\n
    }\n
    for(x=1; x<n+1; x++)\n
    {\n
        a[x] = (m[x] + u[x])/2;\n
    }\n
    c=(n*4)-4;\n
    for(int a1=1; a1<n+1; a1++)\n
    {\n
        e[a1][0] = m[a1];\n
        e[a1][1] = m[a1]+1;\n
        e[a1][2] = u[a1]-1;\n
    }
}
```

```

e[a1][3] = u[a1];\n
}\n
for(int i=1; i<n+1; i++)\n
{\n
    for(int l=1; l<=i; l++)\n
    {\n
        if(l!=1)\n
        {\n
            cout<<a[l]<<"\\t";\n
        }\n
    }\n
    for(int j=0; j<4; j++)\n
    {\n
        cout<<e[i][j];\n
        for(int k=0; k<n-(i+1); k++)\n
        {\n
            cout<<a[k]<<"\\t";\n
        }\n
        cout<<"\\n";\n
    }\n
}\n\n
system("PAUSE");\n
return 0;\n
}\n

```

\n\n

<p>The answer should come in the form of a table like</p>\n\n

<pre><code>

|     |     |       |
|-----|-----|-------|
| 1   | 50  | 50\n  |
| 2   | 50  | 50\n  |
| 99  | 50  | 50\n  |
| 100 | 50  | 50\n  |
| 50  | 1   | 50\n  |
| 50  | 2   | 50\n  |
| 50  | 99  | 50\n  |
| 50  | 100 | 50\n  |
| 50  | 50  | 1\n   |
| 50  | 50  | 2\n   |
| 50  | 50  | 99\n  |
| 50  | 50  | 100\n |

</code></pre>\n\n

<p>if the no of inputs is 3 and their ranges are\n

1,100\n

1,100\n

1,100\n

(could be varied too)</p>\n\n

<p>The output is not coming, can anyone correct the code or tell me what's wrong?<

Tags : 'c++ c'

## 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 of as a many-to-many relationship between features and classes. Labels are not mutually exclusive, such as topics that are relevant for a document. A question on Stackoverflow might be about management, programming, and software development 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 value at 1 and recall at 0. The relative contribution of precision and recall to the F1 score are equal. The formula for microaveraging over a k-class binary problem 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 measure than accuracy especially if you 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.metrics.f1\\_score.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html)

**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

```
#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, i
        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

```
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 generate train.db f
```

 Number of rows in the database :  
6034196  
Time taken to count the number of rows : 0:00:04.957837

### 3.1.3 Checking for duplicates

```
#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
    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 generate train.db file")
```

 Time taken to run this cell : 0:19:47.557987

```
df_no_dup.head()
# we can observe that there are duplicates
```



	Title	Body
0	Implementing Boundary Value Analysis of S...	<pre><code>#include<iostream>\n#include<...</code></pre>
1	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical...
2	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical... c#
3	java.lang.NoClassDefFoundError: javax/serv...	<p>I followed the guide in <a href="http://sta...
4	java.sql.SQLException:[Microsoft][ODBC Dri...	<p>I use the following code</p>\n<pre><code>...

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

number of duplicate questions : 1827881 ( 30.292038906260256 % )

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

1	2656284
2	1272336
3	277575
4	90
5	25
6	5
Name:	cnt_dup, dtype: int64

```
df_no_dup = df_no_dup.dropna(how='any',axis=0)
```

```
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 : 0:00:06.819139

	Title	Body
0	Implementing Boundary Value Analysis of S...	<pre><code>#include<iostream>\n#include<...</code></pre>
1	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical...
2	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical... c#
3	java.lang.NoClassDefFoundError: javax/serv...	<p>I followed the guide in <a href="http://sta...
4	java.sql.SQLException:[Microsoft][ODBC Dri...	<p>I use the following code</p>\n<pre><code>...

```
# distribution of number of tags per question
df_no_dup.tag_count.value_counts()
```



```

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

```

```

#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)

#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 generate train.db")

```

 Time taken to run this cell : 0:00:32.373284

## 3.2 Analysis of Tags

### 3.2.1 Total number of unique tags

```

# 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'])

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

```

 Number of data points : 4206314  
Number of unique tags : 42048

`'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',

### 3.2.3 Number of times a tag appeared

```
# 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))
```

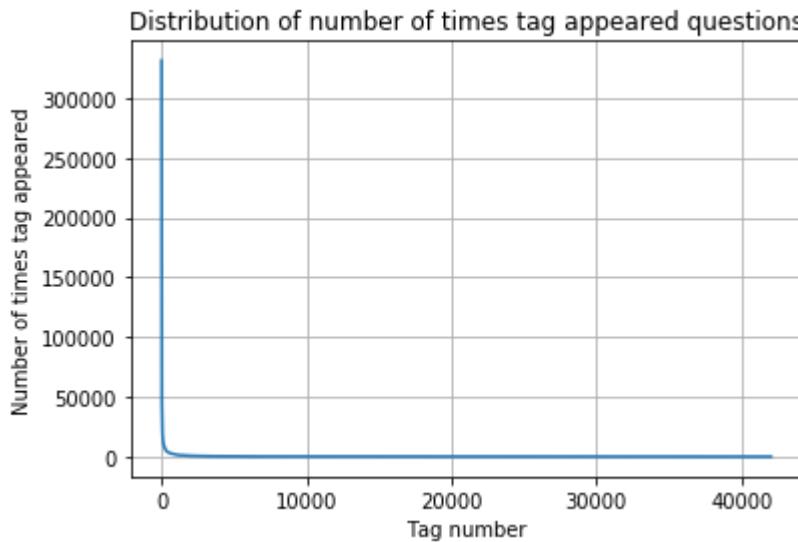
```
#Saving this dictionary to csv files.
if not os.path.isfile('tag_counts_dict_dtm.csv'):
    with open('tag_counts_dict_dtm.csv', 'w') as csv_file:
        writer = csv.writer(csv_file)
        for key, value in result.items():
            writer.writerow([key, value])
tag_df = pd.read_csv("tag_counts_dict_dtm.csv",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

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

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

```
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()
```

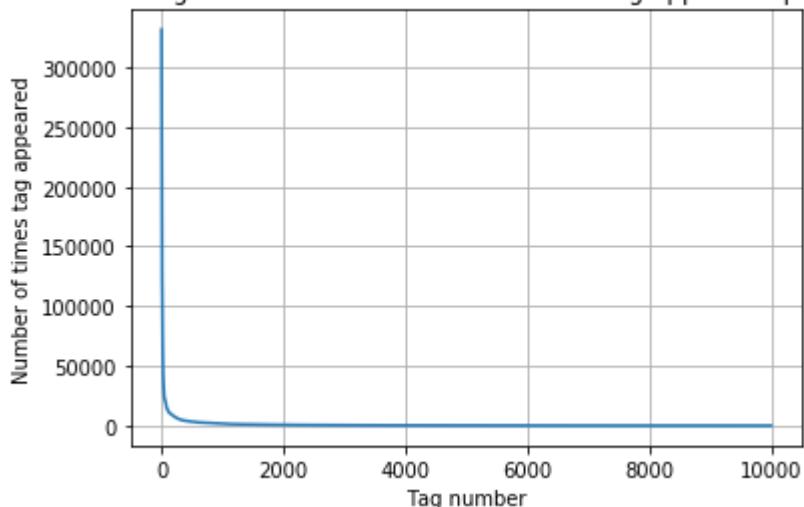




```
plt.plot(tag_counts[0:10000])
plt.title('first 10k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:10000:25]), tag_counts[0:10000:25])
```



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

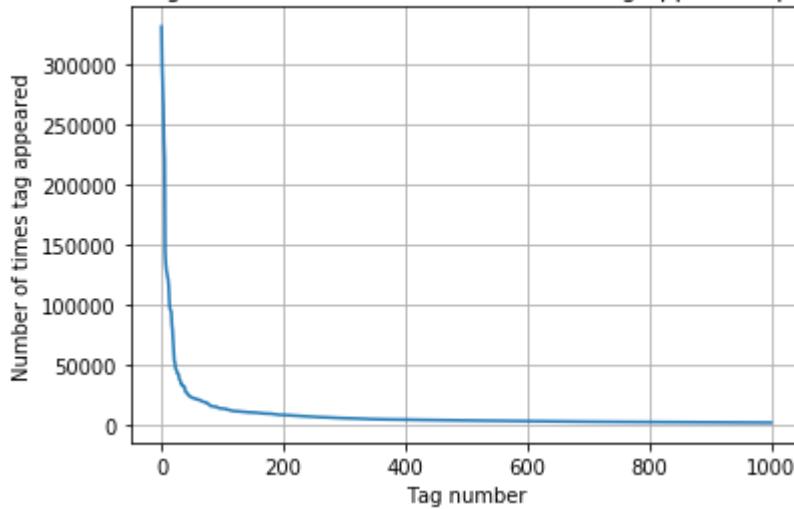


400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7151
6466	5865	5370	4983	4526	4281	4144	3929	3750	3593	
3453	3299	3123	2989	2891	2738	2647	2527	2431	2331	
2259	2186	2097	2020	1959	1900	1828	1770	1723	1673	
1631	1574	1532	1479	1448	1406	1365	1328	1300	1266	
1245	1222	1197	1181	1158	1139	1121	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	

```
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

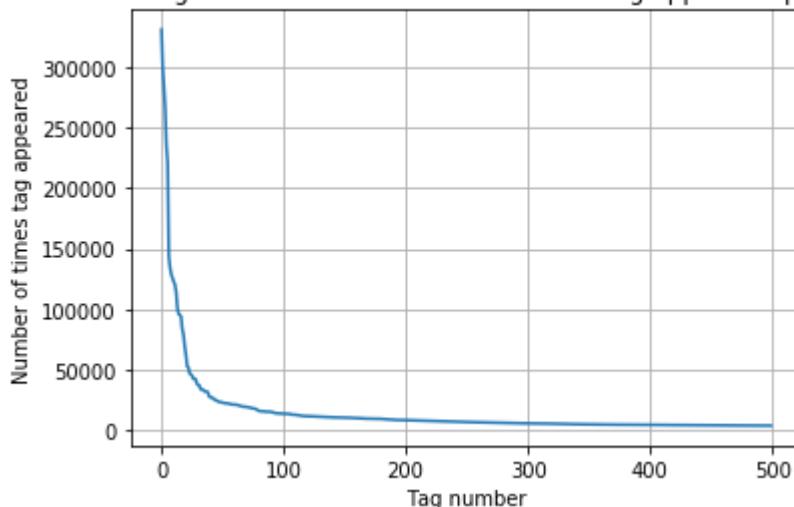


200	[331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429	21820 20957 19758 18905 17728 15533 15097 14884 13703
13364	13157 12407 11658 11228 11162 10863 10600 10350 10224
10029	9884 9719 9411 9252 9148 9040 8617 8361 8163
8054	7867 7702 7564 7274 7151 7052 6847 6656 6553
6466	6291 6183 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]

```
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]

```

```

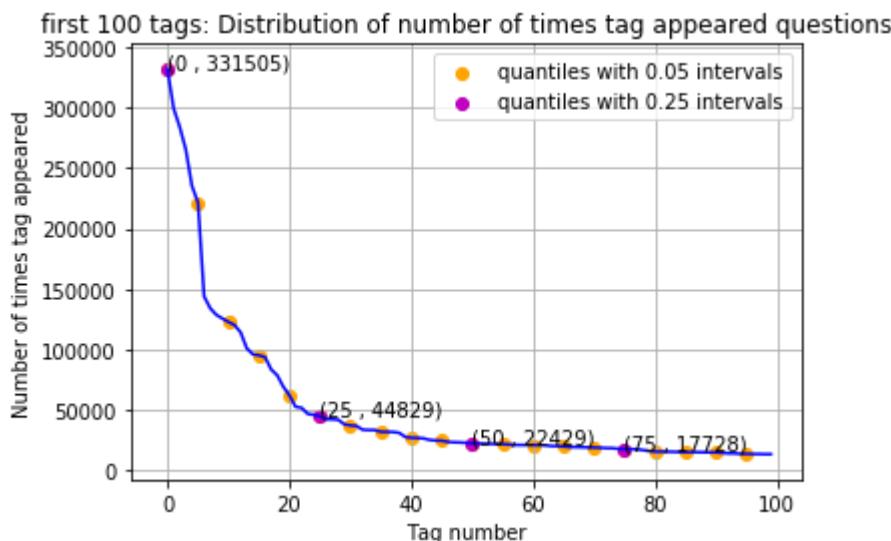
plt.plot(tag_counts[0:100], c='b')
plt.scatter(x=list(range(0,100,5)), y=tag_counts[0:100:5], c='orange', label="quantiles with 0.05 in
# quantiles with 0.25 difference
plt.scatter(x=list(range(0,100,25)), y=tag_counts[0:100:25], c='m', label = "quantiles with 0.25 int

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])

```





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

```
# 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 frequently than others, Micro-averaged F1-score is the appropriate metric

### 3.2.4 Tags Per Question

```
#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 convert
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])
```

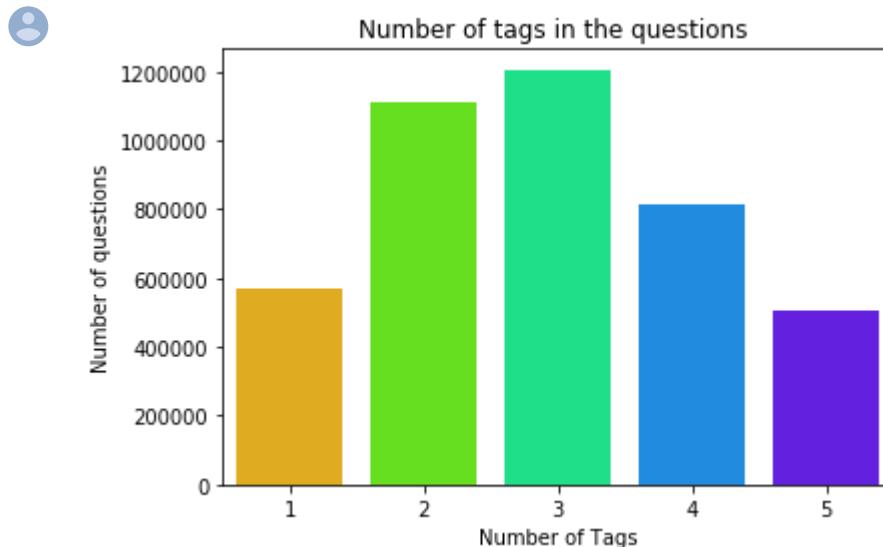
We have total 4206314 datapoints.  
[3, 4, 2, 2, 3]

```
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: 2.899440

```
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

```
# 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)
```







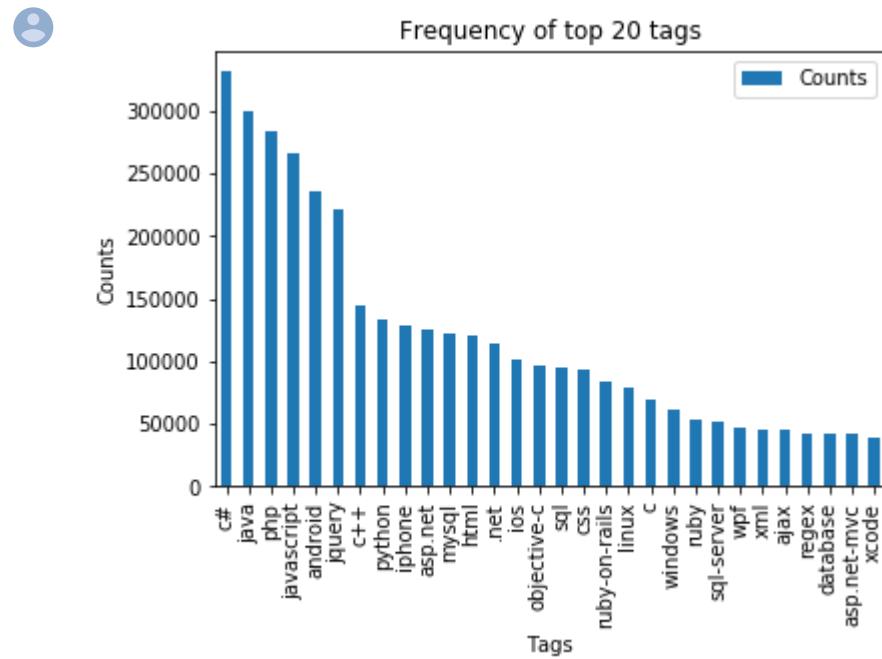
Time taken to run this cell : 0:00:08.993280

### Observations:

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

### 3.2.6 The top 20 tags

```
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

1. Sample 1M data points
2. Separate out code-snippets from Body
3. Remove Special characters from Question title and description (not in code)
4. Remove stop words (Except 'C')
5. Remove HTML Tags
6. Convert all the characters into small letters
7. Use SnowballStemmer to stem the words

```
import nltk
nltk.download('stopwords')
```

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

 [nltk\_data] Error loading stopwords: <urlopen error [WinError 10060] A connection attempt failed because the connected party did not properly respond after a period of time, or established connection failed because connected host has failed to respond>

```
#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 database:")
```

```

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 te
create_database_table("Processed.db", sql_create_table)

```

 Tables in the database:  
QuestionsProcessed

```

# 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 100000;")

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 database:  
QuestionsProcessed  
Cleared All the rows  
Time taken to run this cell : 0:02:32.869286

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

```

#http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/
import nltk
nltk.download('punkt')

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], 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[0]=='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)
if (questions_proccesed%10000==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)

```



```

[nltk_data] Error loading punkt: <urlopen error [WinError 10060] A
[nltk_data]     connection attempt failed because the connected party
[nltk_data]     did not properly respond after a period of time, or
[nltk_data]     established connection failed because connected host
[nltk_data]     has failed to respond>
number of questions completed= 10000
number of questions completed= 20000
number of questions completed= 30000
number of questions completed= 40000
number of questions completed= 50000
number of questions completed= 60000
number of questions completed= 70000
number of questions completed= 80000
number of questions completed= 90000
Avg. length of questions(Title+Body) before processing: 1179
Avg. length of questions(Title+Body) after processing: 327
Percent of questions containing code: 57
Time taken to run this cell : 0:04:36.820751

```

```

# 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()

```

```
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

```
===== ('queri work empti field access wrong queri run ms access databas work fine field empti ----- ('man section access specif section man page',) ----- ('unrecogn option profil tri run command part attempt produc mp play video tag get messa ----- ('font emac cento open emac ssh connect get error font replac empti squar anyon know cou ----- ('set css properti develop vodafone mwp work widget vodafone mobil widget platform tri cha ----- ('creat temporari url temporari pdf generat session data websit creat pdf generat user a ----- ('place caption insid imag recent seen follow imag question look way creat nice caption ----- ('direct open mmf cash drawer rj port aunt bought cash drawer auction recent task figur ----- ('fullcalendar flot resiz conflict success integr flot line graph instanc fullcalendar s
```

```
#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""", c
conn_r.commit()
conn_r.close()
```

```
preprocessed_data.head()
```



```
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 : 99999  
number of dimensions : 2

## 4. Machine Learning Models

### 4.1 Converting tags for multilabel problems

X	y1	y2	y3	y4
x1	0	1	1	0
x1	1	0	0	0
x1	0	1	0	0

```
# binary='true' will give a binary vectorizer
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

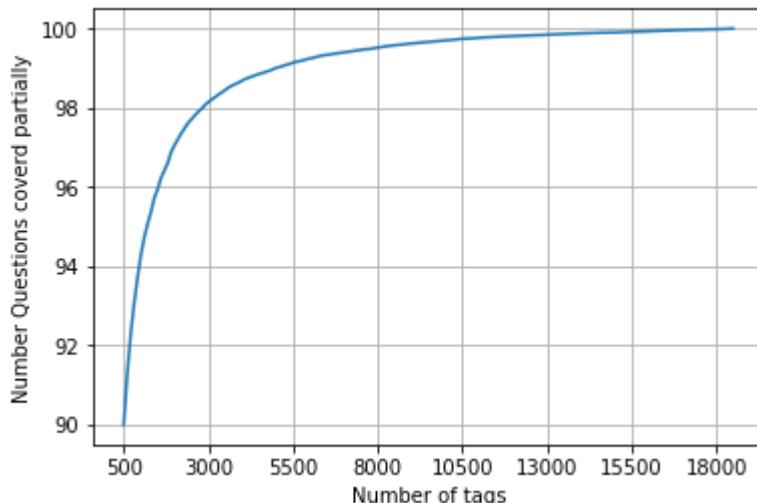
— We will sample the number of tags instead considering all of them (due to limitation of computing power) —

```
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))
```

```
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))
```

```
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 50(it covers 90% of th
print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



with 500 tags we are covering 89.988 % of questions

```
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 : 10012 out of 99999

```
print("Number of tags in sample :", multilabel_y.shape[1])
print("number of tags taken :", multilabel_yx.shape[1],"(,(multilabel_yx.shape[1]/multilabel_y.shap
```

Number of tags in sample : 18534  
number of tags taken : 500 ( 2.6977446854429696 %)

— We consider top 2.6% tags which covers 90% of the questions —

## 4.2 Split the data into test and train (80:20)

```
total_size=preprocessed_data.shape[0]
train_size=int(0.80*total_size)

x_train=preprocessed_data.head(train_size)
x_test=preprocessed_data.tail(total_size - train_size)

y_train = multilabel_yx[0:train_size,:]
y_test = multilabel_yx[train_size:total_size,:]

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 : (79999, 500)  
Number of data points in test data : (20000, 500)

## 4.3 Featurizing data

```

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 : 0:01:02.609616

```

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: (79999, 89596) Y : (79999, 500)  
Dimensions of test data X: (20000, 89596) Y: (20000, 500)

```

# https://www.analyticsvidhya.com/blog/2017/08/introduction-to-multi-label-classification/
#https://stats.stackexchange.com/questions/117796/scikit-multi-label-classification
# classifier = LabelPowerset(GaussianNB())
"""

from skmultilearn.adapt import MLkNN
classifier = MLkNN(k=21)

# train
classifier.fit(x_train_multilabel, y_train)

# predict
predictions = classifier.predict(x_test_multilabel)
print(accuracy_score(y_test,predictions))
print(metrics.f1_score(y_test, predictions, average = 'macro'))
print(metrics.f1_score(y_test, predictions, average = 'micro'))
print(metrics.hamming_loss(y_test,predictions))

"""
# we are getting memory error because the multilearn package
# is trying to convert the data into dense matrix
# -----
#MemoryError                                         Traceback (most recent call last)
#<ipython-input-170-f0e7c7f3e0be> in <module>()
#----> classifier.fit(x_train_multilabel, y_train)

```

 "\nfrom skmultilearn.adapt import MLkNN\nclassifier = MLkNN(k=21)\n\n# train\nclassifier

## 4.4 Applying Logistic Regression with OneVsRest Classifier

```

# this will be taking so much time try not to run it, download the lr_with_equal_weight.pkl file and
# This takes about 6-7 hours to run.
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

print("accuracy :",metrics.accuracy_score(y_test,predictions))
print("macro f1 score :",metrics.f1_score(y_test, predictions, average = 'macro'))
print("micro f1 score :",metrics.f1_score(y_test, predictions, average = 'micro'))
print("hamming loss :",metrics.hamming_loss(y_test,predictions))
print("Precision recall report :\n",metrics.classification_report(y_test, predictions))

```



accuracy : 0.21725  
macro f1 score : 0.3098526917188083  
micro f1 scoore : 0.4505976244894901  
hamming loss : 0.0029326  
Precision recall report :

	precision	recall	f1-score	support
0	0.62	0.26	0.37	1553
1	0.78	0.48	0.59	1402
2	0.82	0.53	0.64	1337
3	0.72	0.42	0.53	1291
4	0.95	0.75	0.84	1064
5	0.87	0.66	0.75	1095
6	0.72	0.34	0.46	675
7	0.70	0.43	0.53	604
8	0.90	0.57	0.70	650
9	0.75	0.42	0.54	579
10	0.84	0.63	0.72	586
11	0.50	0.20	0.29	573
12	0.61	0.13	0.21	521
13	0.55	0.24	0.34	458
14	0.54	0.22	0.31	447
15	0.56	0.31	0.40	446
16	0.77	0.56	0.65	455
17	0.78	0.48	0.59	396
18	0.62	0.25	0.35	347
19	0.63	0.15	0.24	372
20	0.37	0.05	0.09	320
21	0.54	0.24	0.34	246
22	0.71	0.42	0.53	244
23	0.86	0.64	0.73	210
24	0.63	0.41	0.49	211
25	0.66	0.44	0.53	206
26	0.90	0.65	0.75	209
27	0.60	0.33	0.43	212
28	0.34	0.06	0.11	202
29	0.66	0.24	0.35	179
30	0.46	0.26	0.33	184
31	0.95	0.80	0.87	173
32	0.57	0.14	0.23	166
33	0.75	0.31	0.44	162
34	0.76	0.61	0.68	145
35	0.72	0.61	0.66	146
36	0.52	0.22	0.30	157
37	0.36	0.13	0.19	144
38	0.76	0.48	0.58	149
39	0.41	0.11	0.18	153
40	0.70	0.46	0.55	136
41	0.65	0.35	0.46	142
42	0.42	0.07	0.13	148
43	0.38	0.09	0.14	139
44	0.70	0.30	0.42	123
45	0.31	0.12	0.17	102
46	0.14	0.01	0.02	122
47	0.50	0.09	0.15	99
48	0.31	0.10	0.15	111
49	0.86	0.80	0.83	97

50	0.24	0.04	0.06	107
51	0.35	0.15	0.21	104
52	0.56	0.30	0.39	96
53	0.69	0.37	0.48	107
54	0.88	0.21	0.33	102
55	0.64	0.17	0.27	107
56	0.79	0.50	0.61	96
57	0.85	0.38	0.52	116
58	0.33	0.10	0.16	97
59	0.41	0.11	0.17	114
60	0.29	0.06	0.10	101
61	0.28	0.07	0.11	99
62	0.44	0.13	0.20	106
63	0.47	0.15	0.23	101
64	0.71	0.17	0.27	103
65	0.91	0.52	0.66	100
66	0.80	0.45	0.57	105
67	0.87	0.40	0.55	97
68	0.91	0.81	0.86	95
69	0.77	0.26	0.38	90
70	0.66	0.29	0.41	92
71	0.00	0.00	0.00	80
72	0.88	0.43	0.57	101
73	0.67	0.51	0.58	78
74	0.48	0.34	0.40	76
75	0.83	0.54	0.66	90
76	0.40	0.16	0.23	74
77	0.74	0.39	0.51	67
78	0.12	0.01	0.02	85
79	0.50	0.38	0.43	77
80	0.89	0.69	0.77	70
81	0.86	0.56	0.68	75
82	0.79	0.39	0.53	66
83	0.52	0.25	0.33	69
84	0.50	0.02	0.05	82
85	0.89	0.58	0.70	71
86	0.79	0.38	0.51	71
87	0.32	0.07	0.12	80
88	0.45	0.13	0.21	75
89	0.79	0.43	0.55	61
90	0.60	0.47	0.53	68
91	0.60	0.21	0.31	72
92	0.64	0.31	0.42	74
93	0.42	0.20	0.27	69
94	0.97	0.57	0.72	63
95	0.09	0.01	0.02	75
96	0.72	0.51	0.60	70
97	0.22	0.09	0.13	55
98	0.25	0.01	0.03	76
99	0.77	0.14	0.24	71
100	0.58	0.43	0.50	67
101	0.33	0.04	0.07	55
102	0.92	0.55	0.69	64
103	0.88	0.61	0.72	61
104	0.00	0.00	0.00	70
105	0.43	0.20	0.27	66
106	0.68	0.45	0.54	62
107	0.65	0.22	0.33	50

108	0.44	0.11	0.17	64
109	0.67	0.52	0.58	58
110	0.97	0.61	0.75	56
111	0.18	0.04	0.07	67
112	0.33	0.11	0.17	62
113	0.29	0.08	0.12	53
114	0.25	0.17	0.20	47
115	0.71	0.07	0.14	67
116	0.96	0.70	0.81	61
117	0.32	0.11	0.16	55
118	0.66	0.35	0.46	54
119	0.47	0.40	0.43	45
120	0.44	0.21	0.28	58
121	0.44	0.22	0.29	55
122	0.46	0.12	0.18	52
123	0.40	0.15	0.22	53
124	0.53	0.20	0.29	46
125	0.29	0.04	0.07	48
126	0.61	0.27	0.37	41
127	0.93	0.79	0.86	68
128	0.50	0.30	0.37	47
129	0.98	0.78	0.87	58
130	0.67	0.44	0.53	32
131	0.00	0.00	0.00	49
132	0.85	0.60	0.71	48
133	0.40	0.11	0.17	54
134	0.28	0.19	0.23	37
135	0.95	0.72	0.82	57
136	0.37	0.17	0.24	40
137	0.38	0.11	0.17	53
138	0.27	0.06	0.10	47
139	0.44	0.22	0.30	49
140	0.12	0.03	0.05	35
141	0.19	0.06	0.09	52
142	0.68	0.30	0.41	44
143	0.17	0.02	0.04	50
144	0.00	0.00	0.00	44
145	0.54	0.40	0.46	48
146	0.68	0.53	0.59	51
147	0.53	0.22	0.31	37
148	0.94	0.59	0.73	54
149	0.60	0.46	0.52	39
150	0.67	0.24	0.35	50
151	0.89	0.62	0.73	53
152	0.97	0.63	0.77	49
153	0.50	0.15	0.23	60
154	0.50	0.11	0.17	57
155	0.61	0.41	0.49	49
156	0.50	0.11	0.18	45
157	0.61	0.35	0.45	48
158	0.33	0.08	0.13	48
159	0.29	0.04	0.07	47
160	0.46	0.13	0.20	46
161	0.62	0.12	0.20	41
162	0.50	0.14	0.22	43
163	0.76	0.55	0.64	40
164	0.86	0.59	0.70	41
165	0.20	0.02	0.12	10

100	0.20	0.00	0.12	47
166	0.70	0.12	0.21	58
167	0.53	0.21	0.30	48
168	0.92	0.45	0.60	49
169	0.48	0.26	0.33	43
170	0.29	0.08	0.12	51
171	0.17	0.02	0.04	41
172	0.23	0.08	0.12	38
173	0.40	0.08	0.13	52
174	0.60	0.14	0.22	44
175	0.65	0.42	0.52	40
176	0.43	0.07	0.13	40
177	0.71	0.14	0.23	37
178	0.41	0.12	0.19	58
179	0.00	0.00	0.00	36
180	0.38	0.07	0.11	45
181	0.81	0.37	0.51	46
182	0.60	0.41	0.49	44
183	0.53	0.17	0.25	48
184	0.33	0.04	0.08	47
185	1.00	0.63	0.77	38
186	0.80	0.58	0.67	48
187	0.62	0.32	0.42	47
188	0.90	0.49	0.63	39
189	0.66	0.56	0.60	34
190	0.40	0.08	0.14	49
191	0.73	0.30	0.42	37
192	0.31	0.09	0.14	43
193	0.84	0.49	0.62	43
194	0.95	0.45	0.61	42
195	0.00	0.00	0.00	45
196	0.47	0.17	0.25	47
197	0.50	0.16	0.25	49
198	0.29	0.06	0.10	34
199	0.00	0.00	0.00	34
200	0.09	0.03	0.05	33
201	0.33	0.07	0.12	42
202	0.29	0.05	0.09	37
203	0.71	0.13	0.22	39
204	0.83	0.35	0.49	43
205	0.93	0.68	0.79	38
206	0.73	0.22	0.33	37
207	0.64	0.42	0.51	43
208	0.64	0.21	0.32	42
209	0.29	0.06	0.09	36
210	1.00	0.02	0.05	43
211	0.25	0.03	0.05	39
212	0.75	0.14	0.24	42
213	0.70	0.37	0.48	38
214	0.70	0.19	0.30	36
215	0.88	0.41	0.56	37
216	0.67	0.05	0.09	40
217	0.25	0.03	0.05	33
218	0.25	0.10	0.15	39
219	0.14	0.03	0.05	31
220	0.43	0.20	0.27	30
221	0.75	0.32	0.45	28
222	0.75	0.42	0.54	36

223	0.00	0.00	0.00	36
224	0.11	0.02	0.04	46
225	0.79	0.59	0.68	39
226	0.38	0.29	0.33	28
227	0.65	0.32	0.43	34
228	0.00	0.00	0.00	28
229	0.33	0.13	0.19	38
230	0.53	0.28	0.36	36
231	0.95	0.59	0.73	34
232	0.00	0.00	0.00	44
233	0.96	0.76	0.85	33
234	0.57	0.40	0.47	30
235	0.85	0.62	0.72	37
236	0.00	0.00	0.00	32
237	0.25	0.15	0.19	26
238	0.10	0.03	0.04	39
239	0.00	0.00	0.00	23
240	0.92	0.71	0.80	31
241	0.33	0.04	0.07	26
242	0.62	0.33	0.43	30
243	0.22	0.08	0.11	26
244	0.40	0.16	0.23	38
245	0.75	0.33	0.46	27
246	0.72	0.48	0.58	27
247	1.00	0.11	0.20	37
248	0.62	0.19	0.29	27
249	0.76	0.43	0.55	37
250	0.78	0.58	0.67	24
251	0.55	0.43	0.48	37
252	0.00	0.00	0.00	33
253	0.78	0.64	0.70	22
254	0.50	0.21	0.29	34
255	0.25	0.03	0.05	38
256	0.00	0.00	0.00	37
257	0.67	0.06	0.12	31
258	0.81	0.45	0.58	29
259	0.50	0.06	0.11	33
260	0.67	0.13	0.22	30
261	0.33	0.06	0.11	31
262	0.62	0.43	0.51	23
263	0.85	0.34	0.49	32
264	0.71	0.38	0.50	26
265	0.00	0.00	0.00	28
266	0.48	0.36	0.41	28
267	0.59	0.48	0.53	21
268	0.57	0.14	0.22	29
269	0.75	0.28	0.41	32
270	0.80	0.39	0.52	31
271	0.00	0.00	0.00	30
272	0.61	0.39	0.48	28
273	0.14	0.03	0.05	31
274	0.50	0.03	0.06	34
275	0.69	0.27	0.39	33
276	0.95	0.74	0.83	27
277	0.43	0.12	0.19	24
278	0.82	0.50	0.62	28
279	0.00	0.00	0.00	25
280	0.50	0.04	0.08	24

281	0.47	0.27	0.34	30
282	0.25	0.05	0.08	20
283	0.33	0.04	0.08	23
284	0.65	0.52	0.58	29
285	0.43	0.33	0.38	18
286	0.50	0.04	0.08	24
287	0.75	0.12	0.21	25
288	0.29	0.09	0.14	22
289	0.50	0.12	0.19	25
290	0.53	0.33	0.41	24
291	0.73	0.39	0.51	28
292	0.12	0.04	0.06	28
293	0.00	0.00	0.00	26
294	0.64	0.31	0.42	29
295	0.44	0.13	0.20	31
296	0.92	0.52	0.67	23
297	0.50	0.03	0.06	34
298	0.00	0.00	0.00	32
299	0.40	0.08	0.13	25
300	0.27	0.10	0.14	31
301	0.00	0.00	0.00	24
302	1.00	0.84	0.92	32
303	0.70	0.28	0.40	25
304	0.87	0.65	0.74	20
305	0.38	0.14	0.20	22
306	1.00	0.39	0.56	23
307	0.67	0.09	0.15	23
308	0.00	0.00	0.00	24
309	0.00	0.00	0.00	20
310	0.00	0.00	0.00	19
311	0.20	0.04	0.07	25
312	0.73	0.36	0.48	22
313	0.00	0.00	0.00	30
314	0.40	0.22	0.29	18
315	0.67	0.35	0.46	23
316	0.56	0.19	0.29	26
317	0.41	0.28	0.33	25
318	0.00	0.00	0.00	31
319	0.45	0.19	0.26	27
320	0.25	0.04	0.07	23
321	0.67	0.10	0.17	20
322	1.00	0.19	0.33	36
323	0.88	0.60	0.71	25
324	0.33	0.04	0.07	24
325	0.20	0.05	0.08	21
326	0.75	0.22	0.34	27
327	0.88	0.35	0.50	20
328	0.56	0.19	0.28	27
329	0.78	0.29	0.42	24
330	0.45	0.24	0.31	21
331	0.12	0.04	0.06	25
332	0.70	0.28	0.40	25
333	0.40	0.15	0.22	26
334	0.75	0.47	0.58	32
335	0.33	0.08	0.13	25
336	0.00	0.00	0.00	17
337	0.78	0.33	0.47	21
338	0.10	0.17	0.21	21

330	0.40	0.17	0.24	24
339	0.43	0.32	0.36	19
340	0.00	0.00	0.00	18
341	1.00	0.04	0.08	25
342	0.00	0.00	0.00	21
343	0.00	0.00	0.00	20
344	0.60	0.33	0.43	9
345	0.09	0.06	0.07	17
346	0.82	0.47	0.60	19
347	0.00	0.00	0.00	18
348	0.00	0.00	0.00	29
349	0.50	0.14	0.22	21
350	0.29	0.10	0.15	20
351	0.36	0.22	0.28	18
352	0.50	0.17	0.25	30
353	0.50	0.25	0.33	20
354	0.60	0.14	0.23	21
355	0.00	0.00	0.00	19
356	0.33	0.11	0.17	18
357	0.93	0.50	0.65	26
358	0.00	0.00	0.00	26
359	0.33	0.05	0.09	19
360	0.57	0.21	0.31	19
361	0.00	0.00	0.00	20
362	0.75	0.21	0.33	14
363	0.71	0.23	0.34	22
364	1.00	0.47	0.64	15
365	0.00	0.00	0.00	16
366	0.00	0.00	0.00	14
367	0.36	0.19	0.24	27
368	0.69	0.47	0.56	19
369	0.00	0.00	0.00	19
370	0.40	0.15	0.22	26
371	0.86	0.43	0.57	14
372	0.40	0.11	0.17	19
373	0.30	0.19	0.23	16
374	0.56	0.60	0.58	15
375	0.50	0.18	0.26	17
376	0.88	0.68	0.77	22
377	0.50	0.11	0.18	18
378	0.44	0.22	0.30	18
379	0.67	0.09	0.16	22
380	0.33	0.09	0.14	22
381	0.67	0.11	0.18	19
382	0.00	0.00	0.00	18
383	1.00	0.61	0.76	28
384	0.90	0.47	0.62	19
385	0.90	0.50	0.64	18
386	0.00	0.00	0.00	21
387	0.29	0.10	0.15	20
388	0.00	0.00	0.00	23
389	0.71	0.22	0.33	23
390	0.33	0.03	0.06	32
391	0.75	0.21	0.33	14
392	0.85	0.41	0.55	27
393	0.00	0.00	0.00	14
394	0.33	0.07	0.12	28
395	0.67	0.31	0.42	13

396	0.67	0.10	0.17	20
397	0.50	0.42	0.45	12
398	1.00	0.20	0.33	20
399	0.33	0.05	0.09	19
400	0.00	0.00	0.00	20
401	0.83	0.38	0.53	26
402	0.00	0.00	0.00	25
403	0.25	0.04	0.07	24
404	0.50	0.12	0.19	17
405	0.00	0.00	0.00	18
406	0.70	0.41	0.52	17
407	0.00	0.00	0.00	22
408	0.67	0.47	0.55	17
409	0.75	0.64	0.69	14
410	0.00	0.00	0.00	21
411	0.00	0.00	0.00	23
412	0.73	0.38	0.50	21
413	0.50	0.05	0.09	21
414	0.88	0.47	0.61	15
415	0.00	0.00	0.00	15
416	1.00	0.74	0.85	19
417	0.00	0.00	0.00	13
418	0.50	0.09	0.15	22
419	0.25	0.06	0.09	18
420	0.00	0.00	0.00	13
421	0.25	0.05	0.08	21
422	0.91	0.48	0.62	21
423	0.33	0.11	0.17	18
424	0.75	0.14	0.23	22
425	0.60	0.17	0.26	18
426	1.00	0.17	0.30	23
427	0.00	0.00	0.00	19
428	0.14	0.04	0.06	24
429	0.33	0.08	0.13	24
430	1.00	0.69	0.82	13
431	0.50	0.04	0.07	25
432	0.50	0.12	0.19	17
433	0.14	0.08	0.11	12
434	1.00	0.41	0.58	22
435	0.83	0.29	0.43	17
436	0.00	0.00	0.00	14
437	1.00	0.37	0.54	19
438	0.00	0.00	0.00	19
439	0.33	0.24	0.28	17
440	0.44	0.21	0.29	19
441	0.50	0.10	0.16	21
442	0.20	0.06	0.09	17
443	0.50	0.05	0.10	19
444	0.40	0.18	0.25	11
445	0.44	0.21	0.29	19
446	0.00	0.00	0.00	22
447	0.83	0.45	0.59	11
448	0.00	0.00	0.00	14
449	0.38	0.16	0.22	19
450	0.00	0.00	0.00	16
451	0.75	0.35	0.48	17
452	0.00	0.00	0.00	15
453	0.64	0.41	0.50	17

454	0.89	0.38	0.53	21
455	0.00	0.00	0.00	14
456	0.50	0.25	0.33	16
457	0.91	0.62	0.74	16
458	0.43	0.23	0.30	13
459	0.25	0.06	0.09	18
460	0.88	0.44	0.58	16
461	0.81	0.65	0.72	20
462	0.73	0.53	0.62	15
463	0.00	0.00	0.00	25
464	0.33	0.15	0.21	13
465	0.00	0.00	0.00	17
466	0.18	0.12	0.14	17
467	1.00	0.26	0.42	19
468	0.58	0.33	0.42	21
469	1.00	0.37	0.54	19
470	0.88	0.47	0.61	15
471	0.79	0.65	0.71	17
472	0.92	0.57	0.71	21
473	0.78	0.50	0.61	14
474	0.00	0.00	0.00	14
475	0.00	0.00	0.00	12
476	0.00	0.00	0.00	16
477	0.67	0.62	0.64	13
478	0.38	0.21	0.27	14
479	0.43	0.20	0.27	15
480	0.33	0.12	0.18	16
481	0.00	0.00	0.00	15
482	0.17	0.05	0.08	19
483	0.71	0.26	0.38	19
484	0.50	0.07	0.12	14
485	0.00	0.00	0.00	15
486	0.50	0.27	0.35	15
487	0.00	0.00	0.00	11
488	0.00	0.00	0.00	22
489	0.80	0.21	0.33	19
490	0.80	0.22	0.35	18
491	0.40	0.10	0.15	21
492	0.36	0.25	0.30	16
493	0.86	0.35	0.50	17
494	0.00	0.00	0.00	13
495	0.00	0.00	0.00	15
496	0.71	0.33	0.45	15
497	0.50	0.09	0.15	22
498	0.67	0.13	0.22	15
499	0.57	0.27	0.36	15
micro avg	0.70	0.33	0.45	36148
macro avg	0.51	0.24	0.31	36148
weighted avg	0.62	0.33	0.42	36148
samples avg	0.43	0.32	0.35	36148

```
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'precision', 'predicted', average, warn_for)
```

```
from sklearn.externals import joblib
joblib.dump(classifier, 'lr_with_equal_weight.pkl')
```

 C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\ex  
  warnings.warn(msg, category=DeprecationWarning)  
['lr\_with\_equal\_weight.pkl']

## 4.5 Modeling with less data points (0.1M data points) and more weight to title

```
sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question text NOT NULL, code text  
create_database_table("Titlemoreweight1.db", sql_create_table)
```

 Tables in the database:  
QuestionsProcessed

```
# 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 = 'Titlemoreweight1.db'
train_datasize = 80000
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 100001;")
        # 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 database:  
QuestionsProcessed  
Cleared All the rows

### 4.5.1 Preprocessing of questions

1. Separate Code from Body
2. Remove Special characters from Question title and description (not in code)
3. Give more weightage to title : Add title three times to the question

- <li> Remove stop words (Except 'C') </li>
- <li> Remove HTML Tags </li>

<li> Convert all the characters into small letters </li>

```
#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

    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")
    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

Avg. length of questions>Title+Body) before processing: 1232

Avg. length of questions>Title+Body) after processing: 441

Percent of questions containing code: 57

Time taken to run this cell : 0:07:24.804862

```
# never forget to close the connections or else we will end up with database locks
conn_r.commit()
conn_w.commit()
conn_r.close()
conn_w.close()
```

## \_\_ Sample questions after preprocessing of data \_\_

```
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

```
=====
('dynam datagrid bind silverlight dynam datagrid bind silverlight dynam datagrid bind si
-----
('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid java.lang.nocl
-----
('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index java.sql.sqlex
-----
('better way updat feed fb php sdk better way updat feed fb php sdk better way updat fee
-----
('btnadd click event open two window record ad btnadd click event open two window record
-----
('sql inject issu prevent correct form submiss php sql inject issu prevent correct form
-----
('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl subaddit le
-----
('hql equival sql queri hql equival sql queri hql equival sql queri hql queri replac nam
-----
('undefin symbol architectur i386 objc class skpsmtpmessag referenc error undefin symbol
```

## \_\_ Saving Preprocessed data to a Database \_\_

```
#Taking 0.1 Million entries to a dataframe.
write_db = 'Titlemoreweight1.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""", c
conn_r.commit()
conn_r.close()
```

```
preprocessed_data.head()
```

	question	tags
0	dynam datagrid bind silverlight dynam datagrid...	c# silverlight data-binding
1	dynam datagrid bind silverlight dynam datagrid...	c# silverlight data-binding columns
2	java.lang.noclassdeffounderror javax servlet j...	jsp jstl
3	java.sql.SQLException microsoft odbc driver manag...	java jdbc
4	better way updat feed fb php sdk better way up...	facebook api facebook-php-sdk

```
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 : 100000  
 number of dimensions : 2

\_\_ Converting string Tags to multilable output variables \_\_

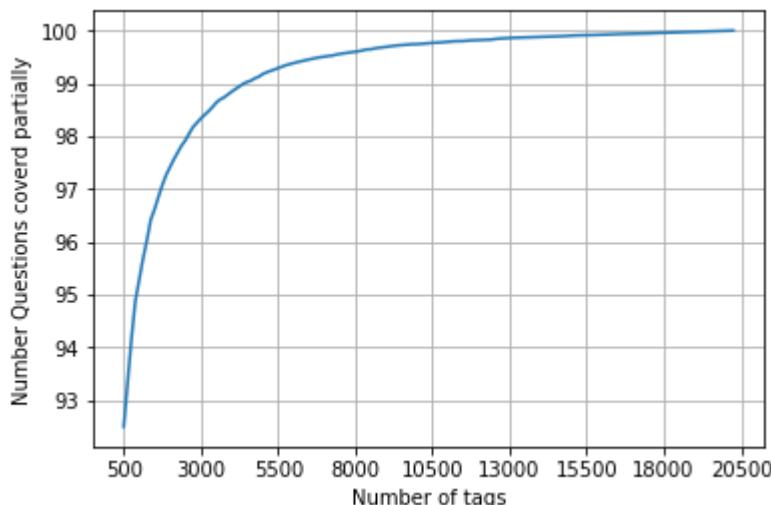
```
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

\_\_ Selecting 500 Tags \_\_

```
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))
```

```
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 covered partially")
plt.grid()
plt.show()
# you can choose any number of tags based on your computing power, minimum is 500(it covers 90% of t
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 99.481 % of questions  
 with 500 tags we are covering 92.5 % of questions

```
# 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 : 7500 out of 100000

```
x_train=preprocessed_data.head(train_datasize)
x_test=preprocessed_data.tail(preprocessed_data.shape[0] - 80000)

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

```
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 : (80000, 500)  
 Number of data points in test data : (20000, 500)

## 4.5.2 Featurizing data with TfIdf vectorizer

```
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 : 0:01:27.552167

```
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: (80000, 100247) Y : (80000, 500)  
Dimensions of test data X: (20000, 100247) Y: (20000, 500)

#### 4.5.3 OneVsRestClassifier with Logistic Regression (SGDClassifier with log-loss )

```
start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict (x_test_multilabel)

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 : 0.1882  
 Hamming loss 0.0030721  
 Micro-average quality numbers  
 Precision: 0.7115, Recall: 0.3031, F1-measure: 0.4251  
 Macro-average quality numbers  
 Precision: 0.5257, Recall: 0.2370, F1-measure: 0.3051

	precision	recall	f1-score	support
0	0.83	0.37	0.51	820
1	0.71	0.18	0.29	1931
2	0.61	0.11	0.19	544
3	0.65	0.23	0.33	222
4	0.82	0.46	0.59	1311
5	0.89	0.48	0.62	1014
6	0.78	0.37	0.50	1374
7	0.89	0.56	0.69	702
8	0.93	0.59	0.72	1424
9	0.73	0.46	0.57	1037
10	0.79	0.50	0.61	797
11	0.74	0.37	0.50	156
12	0.73	0.31	0.43	36
13	0.76	0.38	0.51	610
14	0.46	0.18	0.26	405
15	0.81	0.17	0.29	144
16	0.64	0.23	0.34	425
17	0.69	0.22	0.34	485
18	0.82	0.55	0.66	269
19	0.90	0.56	0.69	518
20	0.53	0.21	0.30	529
21	0.83	0.57	0.68	294
22	0.85	0.40	0.55	520
23	0.68	0.30	0.41	246
24	0.67	0.28	0.39	312
25	0.55	0.27	0.37	314
26	0.68	0.23	0.34	190
27	0.27	0.06	0.10	342
28	0.58	0.19	0.28	96
29	0.43	0.09	0.15	32
30	0.56	0.25	0.35	747
31	0.83	0.36	0.50	14
32	0.66	0.59	0.62	166
33	0.61	0.36	0.45	171
34	0.78	0.24	0.37	256
35	0.84	0.57	0.68	199
36	0.25	0.02	0.03	60
37	0.37	0.21	0.26	203
38	0.68	0.45	0.54	201
39	0.58	0.20	0.30	208
40	0.60	0.23	0.33	13
41	0.62	0.12	0.20	154
42	0.41	0.30	0.35	69
43	0.39	0.06	0.10	426
44	0.53	0.30	0.38	77
45	0.54	0.29	0.38	223
46	0.72	0.22	0.33	144
47	0.87	0.41	0.56	245
48	0.60	0.16	0.26	91

49	0.77	0.29	0.42	157
50	0.91	0.68	0.78	132
51	0.91	0.71	0.79	41
52	0.63	0.42	0.50	124
53	0.27	0.17	0.21	96
54	0.26	0.08	0.12	128
55	0.70	0.35	0.46	46
56	0.69	0.58	0.63	151
57	0.00	0.00	0.00	80
58	0.32	0.12	0.18	65
59	0.52	0.16	0.25	182
60	0.95	0.65	0.77	148
61	0.34	0.06	0.10	196
62	0.46	0.21	0.29	58
63	0.75	0.21	0.33	43
64	0.76	0.27	0.40	197
65	0.66	0.40	0.50	82
66	0.68	0.34	0.45	50
67	0.68	0.50	0.57	105
68	0.24	0.08	0.12	98
69	0.27	0.03	0.06	238
70	0.67	0.06	0.11	35
71	0.61	0.31	0.41	54
72	0.22	0.08	0.12	25
73	0.46	0.21	0.29	29
74	0.50	0.03	0.06	29
75	0.59	0.25	0.35	40
76	0.81	0.56	0.66	105
77	0.59	0.36	0.44	28
78	0.33	0.03	0.06	202
79	0.53	0.46	0.49	37
80	0.75	0.20	0.32	15
81	0.41	0.33	0.37	52
82	0.40	0.20	0.27	50
83	0.14	0.02	0.03	56
84	0.71	0.54	0.61	54
85	0.52	0.44	0.48	34
86	0.50	0.17	0.25	30
87	0.53	0.34	0.42	29
88	0.86	0.75	0.80	24
89	0.90	0.80	0.85	117
90	0.18	0.05	0.07	66
91	0.42	0.12	0.18	68
92	0.87	0.30	0.44	67
93	0.47	0.29	0.36	28
94	0.67	0.24	0.35	17
95	0.89	0.47	0.62	51
96	0.76	0.36	0.49	53
97	0.50	0.02	0.03	61
98	0.00	0.00	0.00	79
99	0.75	0.33	0.46	18
100	0.00	0.00	0.00	11
101	0.71	0.46	0.56	207
102	0.00	0.00	0.00	6
103	0.50	0.03	0.06	30
104	0.50	0.02	0.04	54
105	0.81	0.44	0.57	39
106	0.36	0.11	0.17	70

107	1.00	0.14	0.25	14
108	0.67	0.18	0.29	66
109	0.56	0.18	0.27	50
110	0.80	0.14	0.24	87
111	0.43	0.47	0.45	51
112	0.00	0.00	0.00	291
113	0.97	0.78	0.86	49
114	0.00	0.00	0.00	110
115	0.00	0.00	0.00	28
116	0.00	0.00	0.00	5
117	0.42	0.09	0.15	56
118	0.80	0.45	0.57	125
119	0.88	0.32	0.47	44
120	0.89	0.19	0.31	42
121	0.65	0.20	0.31	55
122	0.85	0.50	0.63	68
123	0.00	0.00	0.00	82
124	0.00	0.00	0.00	0
125	1.00	0.71	0.83	7
126	0.29	0.11	0.16	18
127	0.60	0.19	0.29	31
128	0.86	0.46	0.60	13
129	0.72	0.52	0.60	50
130	0.11	0.01	0.02	91
131	0.79	0.63	0.70	35
132	0.29	0.08	0.12	26
133	0.33	0.03	0.06	32
134	0.71	0.43	0.54	35
135	0.91	0.54	0.68	37
136	0.00	0.00	0.00	55
137	0.26	0.27	0.27	41
138	0.43	0.20	0.27	15
139	0.41	0.11	0.17	99
140	0.93	0.50	0.65	86
141	0.63	0.23	0.33	53
142	0.31	0.11	0.16	36
143	0.57	0.48	0.52	66
144	0.71	0.42	0.53	64
145	0.33	0.04	0.07	25
146	0.00	0.00	0.00	125
147	0.21	0.20	0.21	15
148	0.74	0.48	0.58	48
149	0.36	0.25	0.29	65
150	0.00	0.00	0.00	11
151	0.36	0.27	0.31	15
152	0.36	0.15	0.22	52
153	0.55	0.33	0.41	18
154	1.00	0.19	0.32	16
155	0.80	0.20	0.32	20
156	0.52	0.12	0.20	121
157	0.53	0.30	0.38	107
158	0.50	0.13	0.21	15
159	0.78	0.49	0.60	105
160	0.50	0.25	0.33	69
161	0.71	0.21	0.33	56
162	0.00	0.00	0.00	47
163	0.38	0.02	0.05	121
164	0.50	0.21	0.22	11

104	0.50	0.44	0.33	41
165	0.00	0.00	0.00	229
166	0.88	0.14	0.25	98
167	0.60	0.18	0.28	33
168	0.69	0.25	0.37	44
169	0.62	0.47	0.53	45
170	0.88	0.41	0.56	51
171	0.00	0.00	0.00	18
172	0.67	0.50	0.57	48
173	0.25	0.17	0.20	12
174	0.35	0.11	0.17	62
175	0.71	0.57	0.63	44
176	0.95	0.70	0.81	30
177	0.61	0.37	0.46	30
178	0.00	0.00	0.00	0
179	1.00	1.00	1.00	1
180	0.59	0.25	0.35	40
181	0.33	0.05	0.08	44
182	0.00	0.00	0.00	2
183	0.65	0.40	0.50	75
184	1.00	0.25	0.40	4
185	0.80	0.31	0.45	64
186	0.50	0.33	0.40	12
187	1.00	0.60	0.75	55
188	0.80	0.61	0.69	64
189	0.44	0.11	0.18	96
190	0.00	0.00	0.00	22
191	0.90	0.24	0.38	76
192	0.64	0.40	0.49	45
193	0.71	0.36	0.48	14
194	0.75	0.48	0.59	50
195	0.80	0.20	0.32	20
196	0.84	0.60	0.70	35
197	0.69	0.23	0.35	94
198	1.00	0.07	0.13	14
199	0.00	0.00	0.00	25
200	1.00	0.04	0.07	54
201	0.60	0.14	0.22	22
202	0.75	0.14	0.24	43
203	0.00	0.00	0.00	43
204	0.97	0.55	0.70	62
205	0.00	0.00	0.00	3
206	0.33	0.05	0.08	43
207	0.00	0.00	0.00	7
208	0.33	0.12	0.18	8
209	0.21	0.07	0.11	42
210	0.36	0.40	0.38	10
211	0.41	0.17	0.25	40
212	0.80	0.35	0.48	23
213	0.00	0.00	0.00	6
214	0.72	0.38	0.50	47
215	0.46	0.10	0.16	62
216	0.66	0.27	0.39	77
217	0.67	0.18	0.29	22
218	0.00	0.00	0.00	3
219	0.00	0.00	0.00	28
220	0.75	0.04	0.07	81
221	0.67	0.06	0.12	31

222	1.00	0.03	0.06	34
223	1.00	0.38	0.55	60
224	1.00	0.20	0.33	10
225	0.86	0.60	0.71	10
226	0.74	0.62	0.67	92
227	0.80	0.31	0.44	13
228	0.50	0.08	0.13	13
229	0.89	0.77	0.82	43
230	0.45	0.14	0.22	35
231	0.00	0.00	0.00	4
232	0.42	0.25	0.31	20
233	0.60	0.02	0.04	145
234	0.91	0.53	0.67	55
235	0.00	0.00	0.00	2
236	0.44	0.19	0.26	37
237	0.82	0.47	0.60	90
238	1.00	0.03	0.07	58
239	1.00	0.25	0.40	20
240	0.97	0.52	0.68	61
241	0.82	0.64	0.72	42
242	0.55	0.77	0.64	30
243	0.87	0.50	0.63	66
244	0.56	0.24	0.33	42
245	0.10	0.03	0.05	31
246	1.00	0.33	0.50	6
247	0.75	0.17	0.27	18
248	0.82	0.53	0.64	51
249	0.73	0.47	0.57	17
250	0.53	0.36	0.43	22
251	0.78	0.35	0.48	52
252	0.50	0.03	0.06	29
253	0.08	0.04	0.05	28
254	0.00	0.00	0.00	10
255	0.25	0.20	0.22	5
256	0.17	0.33	0.22	3
257	0.80	0.20	0.31	41
258	0.57	0.13	0.22	30
259	1.00	0.33	0.50	3
260	0.00	0.00	0.00	38
261	0.00	0.00	0.00	1
262	0.70	0.37	0.48	19
263	0.00	0.00	0.00	14
264	1.00	0.03	0.05	37
265	0.14	0.11	0.12	9
266	0.29	0.09	0.14	45
267	0.55	0.48	0.52	33
268	0.75	0.56	0.64	16
269	0.68	0.49	0.57	35
270	0.60	0.27	0.37	11
271	0.00	0.00	0.00	30
272	0.50	0.38	0.43	8
273	0.11	0.05	0.07	21
274	1.00	0.01	0.02	123
275	0.63	0.28	0.39	67
276	0.89	0.80	0.84	20
277	0.00	0.00	0.00	14
278	1.00	0.05	0.10	19
279	1.00	0.50	0.67	12

280	0.00	0.00	0.00	15
281	0.92	0.65	0.76	17
282	1.00	0.66	0.79	41
283	0.67	0.13	0.22	15
284	0.63	0.26	0.37	74
285	0.56	0.13	0.21	38
286	0.29	0.12	0.17	16
287	0.40	0.07	0.11	30
288	0.94	0.57	0.71	28
289	0.00	0.00	0.00	21
290	0.81	0.51	0.63	41
291	0.38	0.25	0.30	12
292	0.75	0.12	0.21	24
293	0.60	0.45	0.51	20
294	0.00	0.00	0.00	23
295	0.00	0.00	0.00	29
296	0.50	0.04	0.07	28
297	0.53	0.19	0.28	42
298	0.00	0.00	0.00	53
299	0.00	0.00	0.00	36
300	0.50	0.12	0.20	41
301	0.75	0.41	0.53	37
302	0.83	0.38	0.53	26
303	0.60	0.27	0.37	11
304	0.36	0.16	0.22	31
305	0.50	0.18	0.26	17
306	1.00	0.11	0.20	9
307	1.00	0.17	0.29	6
308	0.00	0.00	0.00	34
309	0.72	0.42	0.53	43
310	0.17	0.03	0.06	30
311	0.40	0.12	0.18	50
312	0.00	0.00	0.00	24
313	0.00	0.00	0.00	42
314	0.67	0.18	0.29	22
315	1.00	0.02	0.03	58
316	0.00	0.00	0.00	10
317	0.42	0.18	0.25	57
318	1.00	0.40	0.57	10
319	0.00	0.00	0.00	11
320	0.00	0.00	0.00	11
321	0.40	0.25	0.31	8
322	0.80	0.36	0.50	22
323	0.94	0.61	0.74	28
324	0.70	0.46	0.55	50
325	0.67	0.11	0.19	18
326	1.00	0.03	0.06	33
327	0.29	0.12	0.17	17
328	0.67	0.07	0.12	29
329	1.00	0.14	0.25	7
330	0.62	0.50	0.56	10
331	0.12	0.04	0.06	25
332	0.67	1.00	0.80	2
333	0.80	0.36	0.50	11
334	0.00	0.00	0.00	24
335	1.00	0.20	0.33	5
336	0.67	0.06	0.11	33

327	0.15	0.20	0.24	28
338	0.96	0.52	0.68	42
339	0.50	0.08	0.13	26
340	0.48	0.39	0.43	36
341	1.00	0.46	0.63	13
342	0.80	0.36	0.50	11
343	1.00	0.40	0.57	10
344	0.35	0.38	0.36	21
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	6
347	0.50	0.08	0.14	12
348	0.40	0.15	0.22	13
349	0.67	0.17	0.27	24
350	0.67	0.30	0.41	27
351	0.50	0.09	0.16	43
352	0.00	0.00	0.00	30
353	0.55	0.27	0.36	22
354	1.00	0.03	0.06	31
355	0.83	0.50	0.62	10
356	0.57	0.20	0.30	20
357	0.86	0.60	0.71	20
358	0.44	0.29	0.35	28
359	0.62	0.48	0.54	21
360	0.00	0.00	0.00	25
361	0.58	0.31	0.41	35
362	0.83	0.56	0.67	36
363	0.33	0.18	0.23	17
364	1.00	0.23	0.38	13
365	0.00	0.00	0.00	21
366	0.00	0.00	0.00	18
367	0.00	0.00	0.00	97
368	0.69	0.38	0.49	29
369	1.00	0.58	0.74	12
370	0.75	0.23	0.35	13
371	0.33	0.17	0.22	18
372	0.00	0.00	0.00	6
373	0.50	0.33	0.40	6
374	0.25	0.03	0.06	30
375	0.50	0.15	0.23	27
376	0.33	0.04	0.06	28
377	0.00	0.00	0.00	2
378	0.50	0.25	0.33	4
379	0.00	0.00	0.00	19
380	0.20	0.20	0.20	5
381	1.00	0.33	0.50	18
382	0.78	0.32	0.45	22
383	0.00	0.00	0.00	16
384	0.38	0.23	0.29	13
385	0.33	0.17	0.22	18
386	0.90	0.82	0.86	11
387	0.41	0.25	0.31	88
388	1.00	0.15	0.27	13
389	0.00	0.00	0.00	6
390	0.00	0.00	0.00	6
391	1.00	0.57	0.72	51
392	0.00	0.00	0.00	13
393	0.63	0.32	0.43	37
394	0.00	0.00	0.00	6

395	1.00	0.11	0.20	9
396	0.33	0.08	0.12	13
397	1.00	0.50	0.67	6
398	0.55	0.38	0.45	29
399	0.96	0.70	0.81	33
400	0.50	0.03	0.06	31
401	0.67	0.04	0.08	50
402	0.92	0.61	0.73	18
403	0.50	0.14	0.22	7
404	0.68	0.58	0.62	26
405	1.00	0.71	0.83	56
406	1.00	0.25	0.40	4
407	0.17	0.06	0.09	17
408	1.00	0.27	0.43	11
409	0.00	0.00	0.00	18
410	0.60	0.30	0.40	10
411	0.27	0.07	0.11	45
412	0.83	0.25	0.38	20
413	0.50	0.04	0.07	25
414	0.50	0.05	0.09	20
415	0.00	0.00	0.00	6
416	1.00	0.04	0.07	26
417	1.00	0.10	0.18	10
418	0.00	0.00	0.00	18
419	1.00	0.17	0.29	6
420	0.50	0.47	0.48	17
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	6
423	0.00	0.00	0.00	12
424	1.00	0.25	0.40	4
425	1.00	0.09	0.17	11
426	0.00	0.00	0.00	11
427	0.83	0.62	0.71	8
428	0.50	0.12	0.19	26
429	0.55	0.40	0.46	40
430	0.00	0.00	0.00	2
431	0.00	0.00	0.00	35
432	1.00	0.20	0.33	15
433	0.00	0.00	0.00	18
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.25	0.04	0.06	28
437	0.50	0.18	0.27	33
438	0.83	0.50	0.62	20
439	0.00	0.00	0.00	36
440	0.67	0.11	0.19	18
441	0.57	0.44	0.50	18
442	0.82	0.56	0.67	16
443	0.10	0.05	0.06	22
444	0.00	0.00	0.00	6
445	0.90	0.43	0.58	21
446	0.89	0.52	0.66	46
447	0.00	0.00	0.00	69
448	0.00	0.00	0.00	7
449	0.00	0.00	0.00	3
450	0.00	0.00	0.00	52
451	0.00	0.00	0.00	16
452	1.00	0.59	0.74	17

453	0.00	0.00	0.00	13
454	0.80	0.36	0.50	11
455	0.00	0.00	0.00	12
456	0.50	0.17	0.25	6
457	0.17	0.06	0.08	18
458	0.00	0.00	0.00	15
459	0.93	0.46	0.62	28
460	0.00	0.00	0.00	18
461	0.80	0.40	0.53	10
462	0.40	0.08	0.14	24
463	1.00	0.11	0.20	18
464	0.95	0.46	0.62	39
465	0.36	0.36	0.36	11
466	0.12	0.03	0.05	35
467	0.12	0.05	0.07	21
468	0.67	0.11	0.19	37
469	1.00	0.20	0.33	5
470	0.50	0.12	0.20	8
471	0.80	0.32	0.46	37
472	0.17	0.02	0.04	47
473	0.44	0.29	0.35	14
474	1.00	0.61	0.76	23
475	0.58	0.56	0.57	66
476	0.00	0.00	0.00	3
477	0.55	0.32	0.40	19
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	23
480	0.00	0.00	0.00	60
481	0.33	0.08	0.12	26
482	1.00	0.25	0.40	4
483	0.50	0.12	0.20	8
484	0.89	0.35	0.50	23
485	0.54	0.39	0.45	18
486	0.50	0.25	0.33	12
487	0.78	0.24	0.37	29
488	1.00	1.00	1.00	1
489	0.50	0.17	0.25	6
490	0.40	0.29	0.33	7
491	0.00	0.00	0.00	3
492	0.33	0.20	0.25	10
493	0.42	0.26	0.32	19
494	0.00	0.00	0.00	7
495	1.00	0.50	0.67	8
496	0.47	0.39	0.42	18
497	0.00	0.00	0.00	72
498	0.00	0.00	0.00	8
499	0.31	0.12	0.18	32
micro avg	0.71	0.30	0.43	37472
macro avg	0.53	0.24	0.31	37472
weighted avg	0.64	0.30	0.40	37472
samples avg	0.41	0.30	0.32	37472

Time taken to run this cell : 0:05:27.578063

```
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'precision', 'predicted', average, warn_for)
```

```

    recall , true , average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me

```

```
joblib.dump(classifier, 'lr_with_more_title_weight.pkl')
```

 ['lr\_with\_more\_title\_weight.pkl']

## Applying Logistic Regression with OneVsRest Classifier

```

start = datetime.now()
classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1'), n_jobs=-1)
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)

```



```
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\joblib\ext
    "timeout or by a memory leak.", UserWarning
```

Accuracy : 0.1878

Hamming loss 0.0030753

Micro-average quality numbers

Precision: 0.7072, Recall: 0.3060, F1-measure: 0.4272

Macro-average quality numbers

Precision: 0.5250, Recall: 0.2486, F1-measure: 0.3164

	precision	recall	f1-score	support
0	0.85	0.35	0.50	820
1	0.72	0.17	0.27	1931
2	0.60	0.12	0.20	544
3	0.64	0.21	0.32	222
4	0.82	0.46	0.59	1311
5	0.89	0.47	0.62	1014
6	0.77	0.37	0.50	1374
7	0.90	0.56	0.69	702
8	0.93	0.58	0.72	1424
9	0.72	0.47	0.57	1037
10	0.79	0.50	0.61	797
11	0.74	0.37	0.49	156
12	0.79	0.31	0.44	36
13	0.76	0.39	0.51	610
14	0.46	0.18	0.26	405
15	0.81	0.17	0.29	144
16	0.65	0.21	0.32	425
17	0.66	0.22	0.33	485
18	0.81	0.54	0.65	269
19	0.90	0.55	0.68	518
20	0.55	0.21	0.31	529
21	0.83	0.57	0.68	294
22	0.85	0.40	0.54	520
23	0.66	0.30	0.41	246
24	0.69	0.29	0.41	312
25	0.56	0.25	0.35	314
26	0.71	0.25	0.37	190
27	0.25	0.06	0.09	342
28	0.62	0.21	0.31	96
29	0.62	0.16	0.25	32
30	0.57	0.27	0.37	747
31	0.83	0.36	0.50	14
32	0.66	0.58	0.62	166
33	0.63	0.36	0.46	171
34	0.73	0.26	0.39	256
35	0.83	0.56	0.67	199
36	0.33	0.02	0.03	60
37	0.37	0.18	0.24	203
38	0.68	0.45	0.54	201
39	0.53	0.23	0.32	208
40	0.60	0.23	0.33	13
41	0.53	0.12	0.19	154
42	0.43	0.29	0.35	69
43	0.39	0.08	0.13	426
44	0.55	0.30	0.39	77
45	0.52	0.27	0.36	223
46	0.77	0.21	0.33	144

47	0.88	0.43	0.58	245
48	0.58	0.16	0.26	91
49	0.71	0.29	0.41	157
50	0.90	0.69	0.78	132
51	0.91	0.73	0.81	41
52	0.62	0.44	0.52	124
53	0.23	0.14	0.17	96
54	0.27	0.09	0.13	128
55	0.68	0.37	0.48	46
56	0.70	0.58	0.64	151
57	0.00	0.00	0.00	80
58	0.36	0.14	0.20	65
59	0.50	0.19	0.27	182
60	0.96	0.66	0.78	148
61	0.33	0.06	0.10	196
62	0.43	0.17	0.25	58
63	0.79	0.26	0.39	43
64	0.75	0.27	0.40	197
65	0.65	0.38	0.48	82
66	0.72	0.36	0.48	50
67	0.66	0.49	0.56	105
68	0.22	0.08	0.12	98
69	0.27	0.03	0.06	238
70	0.33	0.06	0.10	35
71	0.63	0.31	0.42	54
72	0.22	0.08	0.12	25
73	0.50	0.24	0.33	29
74	0.00	0.00	0.00	29
75	0.64	0.23	0.33	40
76	0.82	0.53	0.65	105
77	0.50	0.36	0.42	28
78	0.27	0.03	0.06	202
79	0.56	0.41	0.47	37
80	0.75	0.20	0.32	15
81	0.46	0.31	0.37	52
82	0.38	0.22	0.28	50
83	0.17	0.02	0.03	56
84	0.71	0.54	0.61	54
85	0.52	0.47	0.49	34
86	0.55	0.20	0.29	30
87	0.56	0.34	0.43	29
88	0.90	0.75	0.82	24
89	0.90	0.79	0.85	117
90	0.19	0.05	0.07	66
91	0.47	0.10	0.17	68
92	0.84	0.31	0.46	67
93	0.44	0.25	0.32	28
94	0.67	0.24	0.35	17
95	0.89	0.47	0.62	51
96	0.78	0.40	0.53	53
97	0.50	0.02	0.03	61
98	0.00	0.00	0.00	79
99	0.86	0.33	0.48	18
100	0.00	0.00	0.00	11
101	0.70	0.49	0.58	207
102	0.00	0.00	0.00	6
103	0.50	0.03	0.06	30
104	0.50	0.02	0.04	54

105	0.80	0.41	0.54	39
106	0.40	0.11	0.18	70
107	1.00	0.14	0.25	14
108	0.65	0.17	0.27	66
109	0.53	0.18	0.27	50
110	0.79	0.13	0.22	87
111	0.43	0.39	0.41	51
112	0.80	0.01	0.03	291
113	0.97	0.76	0.85	49
114	0.17	0.01	0.02	110
115	0.00	0.00	0.00	28
116	0.00	0.00	0.00	5
117	0.40	0.07	0.12	56
118	0.80	0.44	0.57	125
119	0.88	0.34	0.49	44
120	0.75	0.14	0.24	42
121	0.61	0.20	0.30	55
122	0.83	0.51	0.64	68
123	0.00	0.00	0.00	82
124	0.00	0.00	0.00	0
125	1.00	0.71	0.83	7
126	0.25	0.11	0.15	18
127	0.60	0.19	0.29	31
128	0.86	0.46	0.60	13
129	0.70	0.52	0.60	50
130	0.00	0.00	0.00	91
131	0.78	0.60	0.68	35
132	0.17	0.04	0.06	26
133	0.33	0.03	0.06	32
134	0.70	0.40	0.51	35
135	0.88	0.59	0.71	37
136	0.00	0.00	0.00	55
137	0.26	0.27	0.26	41
138	0.43	0.20	0.27	15
139	0.42	0.11	0.18	99
140	0.92	0.55	0.69	86
141	0.63	0.23	0.33	53
142	0.27	0.08	0.13	36
143	0.59	0.44	0.50	66
144	0.71	0.42	0.53	64
145	0.25	0.04	0.07	25
146	0.05	0.01	0.01	125
147	0.38	0.33	0.36	15
148	0.79	0.48	0.60	48
149	0.37	0.23	0.28	65
150	0.50	0.09	0.15	11
151	0.43	0.20	0.27	15
152	0.36	0.15	0.22	52
153	0.50	0.33	0.40	18
154	0.75	0.19	0.30	16
155	0.80	0.20	0.32	20
156	0.48	0.11	0.18	121
157	0.57	0.36	0.44	107
158	0.50	0.13	0.21	15
159	0.77	0.47	0.58	105
160	0.46	0.23	0.31	69
161	0.59	0.18	0.27	56
162	0.00	0.00	0.00	17

102	0.00	0.00	0.00	47
163	0.43	0.02	0.05	121
164	0.50	0.24	0.33	41
165	0.00	0.00	0.00	229
166	0.88	0.22	0.36	98
167	0.64	0.21	0.32	33
168	0.73	0.25	0.37	44
169	0.65	0.49	0.56	45
170	0.84	0.41	0.55	51
171	0.00	0.00	0.00	18
172	0.67	0.50	0.57	48
173	0.25	0.17	0.20	12
174	0.32	0.10	0.15	62
175	0.73	0.55	0.62	44
176	0.95	0.70	0.81	30
177	0.63	0.40	0.49	30
178	0.00	0.00	0.00	0
179	1.00	1.00	1.00	1
180	0.61	0.28	0.38	40
181	0.43	0.07	0.12	44
182	0.00	0.00	0.00	2
183	0.63	0.36	0.46	75
184	1.00	0.25	0.40	4
185	0.75	0.33	0.46	64
186	0.50	0.25	0.33	12
187	0.97	0.60	0.74	55
188	0.80	0.61	0.69	64
189	0.43	0.10	0.17	96
190	0.20	0.05	0.07	22
191	0.94	0.22	0.36	76
192	0.68	0.42	0.52	45
193	0.71	0.36	0.48	14
194	0.71	0.40	0.51	50
195	0.78	0.35	0.48	20
196	0.85	0.63	0.72	35
197	0.65	0.23	0.34	94
198	1.00	0.07	0.13	14
199	0.00	0.00	0.00	25
200	0.71	0.09	0.16	54
201	0.50	0.14	0.21	22
202	0.60	0.14	0.23	43
203	0.00	0.00	0.00	43
204	0.97	0.58	0.73	62
205	0.00	0.00	0.00	3
206	0.33	0.05	0.08	43
207	0.00	0.00	0.00	7
208	0.25	0.12	0.17	8
209	0.21	0.07	0.11	42
210	0.36	0.40	0.38	10
211	0.41	0.17	0.25	40
212	0.71	0.22	0.33	23
213	0.00	0.00	0.00	6
214	0.76	0.40	0.53	47
215	0.46	0.10	0.16	62
216	0.66	0.32	0.43	77
217	0.67	0.18	0.29	22
218	0.00	0.00	0.00	3
219	0.00	0.00	0.00	28

220	0.80	0.05	0.09	81
221	0.75	0.10	0.17	31
222	1.00	0.03	0.06	34
223	1.00	0.38	0.55	60
224	1.00	0.20	0.33	10
225	0.83	0.50	0.62	10
226	0.77	0.61	0.68	92
227	0.67	0.31	0.42	13
228	0.67	0.15	0.25	13
229	0.89	0.74	0.81	43
230	0.55	0.17	0.26	35
231	0.00	0.00	0.00	4
232	0.40	0.20	0.27	20
233	0.50	0.12	0.20	145
234	0.85	0.53	0.65	55
235	0.00	0.00	0.00	2
236	0.39	0.19	0.25	37
237	0.83	0.43	0.57	90
238	0.50	0.07	0.12	58
239	1.00	0.25	0.40	20
240	0.97	0.56	0.71	61
241	0.83	0.71	0.77	42
242	0.55	0.77	0.64	30
243	0.88	0.55	0.67	66
244	0.53	0.19	0.28	42
245	0.10	0.03	0.05	31
246	1.00	0.33	0.50	6
247	0.60	0.17	0.26	18
248	0.85	0.55	0.67	51
249	0.67	0.47	0.55	17
250	0.59	0.45	0.51	22
251	0.78	0.35	0.48	52
252	0.50	0.03	0.06	29
253	0.10	0.04	0.05	28
254	0.00	0.00	0.00	10
255	0.25	0.20	0.22	5
256	0.17	0.33	0.22	3
257	0.73	0.27	0.39	41
258	0.50	0.13	0.21	30
259	1.00	0.33	0.50	3
260	1.00	0.03	0.05	38
261	0.00	0.00	0.00	1
262	0.64	0.37	0.47	19
263	0.00	0.00	0.00	14
264	0.50	0.03	0.05	37
265	0.14	0.11	0.12	9
266	0.38	0.11	0.17	45
267	0.57	0.52	0.54	33
268	0.77	0.62	0.69	16
269	0.67	0.51	0.58	35
270	0.43	0.27	0.33	11
271	0.00	0.00	0.00	30
272	0.67	0.50	0.57	8
273	0.12	0.05	0.07	21
274	0.39	0.07	0.12	123
275	0.64	0.31	0.42	67
276	0.84	0.80	0.82	20
277	0.00	0.00	0.00	14

278	0.75	0.16	0.26	19
279	0.86	0.50	0.63	12
280	0.00	0.00	0.00	15
281	0.86	0.71	0.77	17
282	1.00	0.68	0.81	41
283	0.60	0.20	0.30	15
284	0.59	0.26	0.36	74
285	0.50	0.16	0.24	38
286	0.33	0.12	0.18	16
287	0.50	0.07	0.12	30
288	0.94	0.57	0.71	28
289	0.00	0.00	0.00	21
290	0.81	0.51	0.63	41
291	0.33	0.17	0.22	12
292	1.00	0.17	0.29	24
293	0.60	0.45	0.51	20
294	0.00	0.00	0.00	23
295	0.00	0.00	0.00	29
296	0.67	0.07	0.13	28
297	0.50	0.19	0.28	42
298	0.00	0.00	0.00	53
299	0.00	0.00	0.00	36
300	0.45	0.12	0.19	41
301	0.69	0.49	0.57	37
302	0.83	0.38	0.53	26
303	0.67	0.36	0.47	11
304	0.27	0.10	0.14	31
305	0.50	0.18	0.26	17
306	1.00	0.11	0.20	9
307	1.00	0.17	0.29	6
308	0.00	0.00	0.00	34
309	0.69	0.42	0.52	43
310	0.14	0.03	0.05	30
311	0.33	0.12	0.18	50
312	0.00	0.00	0.00	24
313	0.50	0.02	0.05	42
314	0.57	0.18	0.28	22
315	1.00	0.02	0.03	58
316	0.00	0.00	0.00	10
317	0.39	0.16	0.23	57
318	1.00	0.40	0.57	10
319	0.00	0.00	0.00	11
320	0.00	0.00	0.00	11
321	0.33	0.25	0.29	8
322	0.80	0.36	0.50	22
323	0.86	0.64	0.73	28
324	0.69	0.48	0.56	50
325	0.67	0.11	0.19	18
326	0.00	0.00	0.00	33
327	0.29	0.12	0.17	17
328	0.60	0.10	0.18	29
329	1.00	0.29	0.44	7
330	0.62	0.50	0.56	10
331	0.11	0.04	0.06	25
332	0.67	1.00	0.80	2
333	0.75	0.27	0.40	11
334	0.00	0.00	0.00	24

	१.३८	०.४०	०.२८	०
336	0.60	0.09	0.16	33
337	0.67	0.20	0.31	30
338	0.96	0.57	0.72	42
339	0.40	0.08	0.13	26
340	0.50	0.44	0.47	36
341	1.00	0.46	0.63	13
342	0.67	0.36	0.47	11
343	1.00	0.40	0.57	10
344	0.36	0.43	0.39	21
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	6
347	0.60	0.25	0.35	12
348	0.50	0.15	0.24	13
349	0.75	0.12	0.21	24
350	0.80	0.30	0.43	27
351	0.50	0.14	0.22	43
352	0.00	0.00	0.00	30
353	0.56	0.23	0.32	22
354	1.00	0.03	0.06	31
355	0.75	0.60	0.67	10
356	1.00	0.20	0.33	20
357	0.76	0.65	0.70	20
358	0.48	0.36	0.41	28
359	0.67	0.38	0.48	21
360	0.33	0.12	0.18	25
361	0.60	0.34	0.44	35
362	0.85	0.61	0.71	36
363	0.33	0.18	0.23	17
364	1.00	0.46	0.63	13
365	0.00	0.00	0.00	21
366	1.00	0.06	0.11	18
367	0.44	0.07	0.12	97
368	0.67	0.34	0.45	29
369	1.00	0.83	0.91	12
370	0.80	0.31	0.44	13
371	0.25	0.11	0.15	18
372	0.00	0.00	0.00	6
373	0.50	0.50	0.50	6
374	0.25	0.03	0.06	30
375	0.57	0.15	0.24	27
376	0.50	0.07	0.12	28
377	0.00	0.00	0.00	2
378	0.67	0.50	0.57	4
379	0.00	0.00	0.00	19
380	0.25	0.20	0.22	5
381	1.00	0.33	0.50	18
382	0.73	0.36	0.48	22
383	0.00	0.00	0.00	16
384	0.43	0.23	0.30	13
385	0.33	0.17	0.22	18
386	0.90	0.82	0.86	11
387	0.39	0.24	0.30	88
388	0.75	0.23	0.35	13
389	0.00	0.00	0.00	6
390	0.00	0.00	0.00	6
391	0.94	0.65	0.77	51
392	0.50	0.08	0.13	13

393	0.52	0.35	0.42	37
394	0.00	0.00	0.00	6
395	0.50	0.11	0.18	9
396	0.33	0.08	0.12	13
397	1.00	0.50	0.67	6
398	0.71	0.34	0.47	29
399	0.96	0.73	0.83	33
400	0.50	0.03	0.06	31
401	0.80	0.08	0.15	50
402	0.92	0.67	0.77	18
403	0.50	0.14	0.22	7
404	0.62	0.58	0.60	26
405	0.89	0.71	0.79	56
406	1.00	0.25	0.40	4
407	0.20	0.06	0.09	17
408	1.00	0.36	0.53	11
409	0.00	0.00	0.00	18
410	0.67	0.40	0.50	10
411	0.25	0.04	0.08	45
412	0.78	0.35	0.48	20
413	0.60	0.12	0.20	25
414	0.25	0.05	0.08	20
415	0.50	0.17	0.25	6
416	1.00	0.04	0.07	26
417	1.00	0.10	0.18	10
418	0.00	0.00	0.00	18
419	1.00	0.17	0.29	6
420	0.47	0.41	0.44	17
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	6
423	0.00	0.00	0.00	12
424	1.00	0.75	0.86	4
425	0.00	0.00	0.00	11
426	0.33	0.09	0.14	11
427	0.86	0.75	0.80	8
428	0.50	0.12	0.19	26
429	0.57	0.40	0.47	40
430	0.00	0.00	0.00	2
431	0.00	0.00	0.00	35
432	1.00	0.20	0.33	15
433	0.00	0.00	0.00	18
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.20	0.04	0.06	28
437	0.50	0.18	0.27	33
438	0.83	0.50	0.62	20
439	0.00	0.00	0.00	36
440	0.40	0.11	0.17	18
441	0.57	0.44	0.50	18
442	0.82	0.56	0.67	16
443	0.11	0.05	0.06	22
444	0.00	0.00	0.00	6
445	0.83	0.48	0.61	21
446	0.90	0.59	0.71	46
447	0.09	0.03	0.04	69
448	0.00	0.00	0.00	7
449	0.00	0.00	0.00	3
450	0.00	0.00	0.00	52

451	0.00	0.00	0.00	16
452	1.00	0.71	0.83	17
453	0.00	0.00	0.00	13
454	0.80	0.36	0.50	11
455	0.00	0.00	0.00	12
456	0.00	0.00	0.00	6
457	0.14	0.06	0.08	18
458	0.00	0.00	0.00	15
459	0.94	0.57	0.71	28
460	0.00	0.00	0.00	18
461	1.00	0.30	0.46	10
462	0.50	0.12	0.20	24
463	0.75	0.17	0.27	18
464	0.95	0.49	0.64	39
465	0.33	0.36	0.35	11
466	0.10	0.03	0.04	35
467	0.08	0.05	0.06	21
468	0.33	0.03	0.05	37
469	0.67	0.40	0.50	5
470	1.00	0.12	0.22	8
471	0.79	0.30	0.43	37
472	0.11	0.02	0.04	47
473	0.43	0.21	0.29	14
474	1.00	0.61	0.76	23
475	0.60	0.64	0.62	66
476	0.00	0.00	0.00	3
477	0.43	0.32	0.36	19
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	23
480	1.00	0.03	0.06	60
481	0.43	0.12	0.18	26
482	0.67	0.50	0.57	4
483	0.50	0.12	0.20	8
484	0.89	0.35	0.50	23
485	0.62	0.44	0.52	18
486	0.40	0.17	0.24	12
487	0.85	0.38	0.52	29
488	1.00	1.00	1.00	1
489	0.50	0.17	0.25	6
490	0.33	0.29	0.31	7
491	0.00	0.00	0.00	3
492	0.40	0.20	0.27	10
493	0.50	0.42	0.46	19
494	0.00	0.00	0.00	7
495	0.80	0.50	0.62	8
496	0.50	0.50	0.50	18
497	0.00	0.00	0.00	72
498	0.00	0.00	0.00	8
499	0.50	0.25	0.33	32
<hr/>				
micro avg	0.71	0.31	0.43	37472
macro avg	0.52	0.25	0.32	37472
weighted avg	0.64	0.31	0.40	37472
samples avg	0.42	0.30	0.33	37472

Time taken to run this cell : 0:12:43.233552

C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me

```
precision , predicted , average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
```

## 5. Assignments

1. Use bag of words upto 4 grams and compute the micro f1 score with Logistic regression(OvR)
  2. Perform hyperparam tuning on alpha (or lambda) for Logistic regression to improve the performance using G
  3. Try OneVsRestClassifier with Linear-SVM (SGDClassifier with loss-hinge)

# Featurizing data with BOW vectorizer

```
start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009, max_features=20000, tokenizer = lambda x: x.split(), ngram_range=(1,2))
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 : 0:02:19.083209

```
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: (80000, 20000) Y : (80000, 500)  
Dimensions of test data X: (20000, 20000) Y: (20000, 500)

Applying Logistic Regression with OneVsRest Classifier(Hyperparameter Tuning)

```
from sklearn.model_selection import GridSearchCV
start=datetime.now()
parameters=[10**-2,10**-1,1,10]
params = {'estimator__C':parameters}
model=OneVsRestClassifier(LogisticRegression())
clf=GridSearchCV(model, params, cv=3, n_jobs=-1)

clf.fit(x_train_multilabel,y_train)
print('Time to train'.datetime.now()-start)
```

C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\sv  
"the number of iterations.", ConvergenceWarning)  
Time to train 3:01:57.431302

```
best_c=clf.best_params_['estimator__C']
print(best_c)
```

0.1

```
log=OneVsRestClassifier(LogisticRegression(C=best_c), n_jobs=-1)
log.fit(x_train_multilabel,y_train)
```

```
pred=log.predict(x_test_multilabel)
print("accuracy :",metrics.accuracy_score(y_test,pred))
print("macro f1 score :",metrics.f1_score(y_test, pred, average = 'macro'))
print("micro f1 score :",metrics.f1_score(y_test, pred, average = 'micro'))
print("hamming loss :",metrics.hamming_loss(y_test,pred))
print("Precision recall report :\n",metrics.classification_report(y_test, pred))
```

accuracy : 0.17885  
macro f1 score : 0.30207194081161576  
micro f1 scoore : 0.4222632760993986  
hamming loss : 0.0031898

Precision recall report :

	precision	recall	f1-score	support
0	0.78	0.36	0.49	820
1	0.63	0.26	0.37	1931
2	0.48	0.14	0.21	544
3	0.65	0.18	0.29	222
4	0.76	0.48	0.59	1311
5	0.85	0.48	0.61	1014
6	0.74	0.39	0.51	1374
7	0.82	0.58	0.68	702
8	0.91	0.57	0.70	1424
9	0.76	0.44	0.55	1037
10	0.78	0.53	0.63	797
11	0.68	0.38	0.49	156
12	0.68	0.42	0.52	36
13	0.74	0.38	0.50	610
14	0.47	0.22	0.30	405
15	0.64	0.16	0.26	144
16	0.58	0.20	0.30	425
17	0.65	0.22	0.33	485
18	0.83	0.66	0.73	269
19	0.90	0.56	0.69	518
20	0.52	0.23	0.32	529
21	0.84	0.57	0.68	294
22	0.80	0.36	0.50	520
23	0.60	0.27	0.37	246
24	0.63	0.35	0.45	312
25	0.53	0.30	0.38	314
26	0.68	0.25	0.37	190
27	0.23	0.07	0.10	342
28	0.45	0.20	0.28	96
29	0.14	0.03	0.05	32
30	0.65	0.36	0.46	747
31	0.57	0.29	0.38	14
32	0.62	0.60	0.61	166
33	0.63	0.27	0.38	171
34	0.70	0.24	0.36	256
35	0.85	0.52	0.65	199
36	0.07	0.02	0.03	60
37	0.32	0.19	0.24	203
38	0.70	0.44	0.54	201
39	0.53	0.32	0.40	208
40	0.14	0.08	0.10	13
41	0.56	0.13	0.21	154
42	0.41	0.29	0.34	69
43	0.31	0.11	0.16	426
44	0.44	0.22	0.29	77
45	0.52	0.26	0.34	223
46	0.56	0.28	0.37	144
47	0.86	0.40	0.55	245
48	0.63	0.19	0.29	91
49	0.63	0.33	0.44	157

50	0.90	0.63	0.74	132
51	0.89	0.61	0.72	41
52	0.59	0.39	0.47	124
53	0.26	0.23	0.24	96
54	0.26	0.12	0.16	128
55	0.67	0.30	0.42	46
56	0.69	0.59	0.64	151
57	0.00	0.00	0.00	80
58	0.43	0.15	0.23	65
59	0.49	0.12	0.19	182
60	0.92	0.65	0.76	148
61	0.36	0.10	0.15	196
62	0.31	0.16	0.21	58
63	0.79	0.26	0.39	43
64	0.62	0.23	0.33	197
65	0.60	0.32	0.42	82
66	0.69	0.36	0.47	50
67	0.68	0.50	0.57	105
68	0.36	0.09	0.15	98
69	0.23	0.06	0.10	238
70	0.67	0.06	0.11	35
71	0.54	0.39	0.45	54
72	0.40	0.08	0.13	25
73	0.36	0.14	0.20	29
74	0.22	0.07	0.11	29
75	0.45	0.25	0.32	40
76	0.77	0.50	0.61	105
77	0.70	0.50	0.58	28
78	0.21	0.07	0.10	202
79	0.60	0.41	0.48	37
80	1.00	0.33	0.50	15
81	0.46	0.37	0.41	52
82	0.42	0.20	0.27	50
83	0.20	0.04	0.06	56
84	0.70	0.48	0.57	54
85	0.59	0.59	0.59	34
86	0.29	0.13	0.18	30
87	0.61	0.48	0.54	29
88	0.86	0.75	0.80	24
89	0.86	0.78	0.82	117
90	0.21	0.09	0.13	66
91	0.60	0.22	0.32	68
92	0.82	0.27	0.40	67
93	0.58	0.39	0.47	28
94	0.40	0.24	0.30	17
95	0.86	0.49	0.62	51
96	0.68	0.36	0.47	53
97	0.14	0.02	0.03	61
98	0.20	0.01	0.02	79
99	0.80	0.44	0.57	18
100	0.00	0.00	0.00	11
101	0.68	0.41	0.51	207
102	0.00	0.00	0.00	6
103	0.33	0.03	0.06	30
104	0.36	0.07	0.12	54
105	0.79	0.38	0.52	39
106	0.31	0.11	0.17	70
107	0.67	0.14	0.24	14

108	0.50	0.09	0.15	66
109	0.54	0.28	0.37	50
110	0.84	0.18	0.30	87
111	0.44	0.39	0.42	51
112	0.00	0.00	0.00	291
113	0.97	0.76	0.85	49
114	0.34	0.09	0.14	110
115	0.17	0.04	0.06	28
116	0.00	0.00	0.00	5
117	0.36	0.09	0.14	56
118	0.76	0.38	0.51	125
119	0.88	0.34	0.49	44
120	0.83	0.12	0.21	42
121	0.42	0.18	0.25	55
122	0.73	0.44	0.55	68
123	0.12	0.05	0.07	82
124	0.00	0.00	0.00	0
125	0.83	0.71	0.77	7
126	0.11	0.06	0.07	18
127	0.50	0.10	0.16	31
128	0.86	0.46	0.60	13
129	0.74	0.46	0.57	50
130	0.21	0.05	0.09	91
131	0.73	0.54	0.62	35
132	0.38	0.19	0.26	26
133	0.20	0.03	0.05	32
134	0.65	0.37	0.47	35
135	0.95	0.51	0.67	37
136	0.00	0.00	0.00	55
137	0.30	0.34	0.32	41
138	0.31	0.27	0.29	15
139	0.33	0.12	0.18	99
140	0.94	0.52	0.67	86
141	0.62	0.28	0.39	53
142	0.56	0.14	0.22	36
143	0.54	0.45	0.49	66
144	0.63	0.34	0.44	64
145	0.43	0.12	0.19	25
146	0.14	0.06	0.08	125
147	0.25	0.27	0.26	15
148	0.75	0.44	0.55	48
149	0.42	0.26	0.32	65
150	0.33	0.09	0.14	11
151	0.36	0.27	0.31	15
152	0.30	0.15	0.20	52
153	0.44	0.39	0.41	18
154	0.60	0.19	0.29	16
155	0.80	0.20	0.32	20
156	0.48	0.17	0.25	121
157	0.50	0.34	0.40	107
158	0.00	0.00	0.00	15
159	0.71	0.43	0.54	105
160	0.58	0.28	0.37	69
161	0.65	0.36	0.46	56
162	0.00	0.00	0.00	47
163	0.23	0.05	0.08	121
164	0.48	0.24	0.32	41
165	0.00	0.00	0.00	220

	v.vv	v.vv	v.vv	vvv
165	0.88	0.07	0.13	98
166	0.67	0.24	0.36	33
167	0.57	0.18	0.28	44
168	0.64	0.47	0.54	45
169	0.84	0.31	0.46	51
170	0.00	0.00	0.00	18
171	0.58	0.40	0.47	48
172	0.44	0.33	0.38	12
173	0.35	0.15	0.20	62
174	0.75	0.41	0.53	44
175	0.95	0.70	0.81	30
176	0.52	0.37	0.43	30
177	0.00	0.00	0.00	0
178	1.00	1.00	1.00	1
179	0.65	0.33	0.43	40
180	0.25	0.09	0.13	44
181	0.00	0.00	0.00	2
182	0.59	0.35	0.44	75
183	0.67	0.50	0.57	4
184	0.57	0.20	0.30	64
185	0.43	0.25	0.32	12
186	0.97	0.55	0.70	55
187	0.83	0.59	0.69	64
188	0.41	0.19	0.26	96
189	0.00	0.00	0.00	22
190	0.82	0.24	0.37	76
191	0.65	0.49	0.56	45
192	0.86	0.43	0.57	14
193	0.67	0.40	0.50	50
194	1.00	0.30	0.46	20
195	0.91	0.57	0.70	35
196	0.67	0.26	0.37	94
197	0.00	0.00	0.00	14
198	0.00	0.00	0.00	25
199	0.40	0.04	0.07	54
200	0.25	0.05	0.08	22
201	0.33	0.12	0.17	43
202	1.00	0.02	0.05	43
203	0.97	0.52	0.67	62
204	0.00	0.00	0.00	3
205	0.35	0.14	0.20	43
206	0.50	0.14	0.22	7
207	0.25	0.12	0.17	8
208	0.43	0.07	0.12	42
209	0.40	0.40	0.40	10
210	0.35	0.15	0.21	40
211	0.70	0.30	0.42	23
212	0.00	0.00	0.00	6
213	0.78	0.45	0.57	47
214	0.33	0.06	0.11	62
215	0.71	0.32	0.45	77
216	0.22	0.09	0.13	22
217	0.00	0.00	0.00	3
218	0.10	0.04	0.05	28
219	0.80	0.05	0.09	81
220	0.50	0.19	0.28	31
221	0.50	0.03	0.06	34

223	1.00	0.30	0.46	60
224	0.50	0.20	0.29	10
225	0.86	0.60	0.71	10
226	0.75	0.70	0.72	92
227	0.67	0.31	0.42	13
228	0.40	0.15	0.22	13
229	0.86	0.72	0.78	43
230	0.38	0.14	0.21	35
231	0.00	0.00	0.00	4
232	0.33	0.10	0.15	20
233	0.39	0.14	0.20	145
234	0.83	0.45	0.59	55
235	0.00	0.00	0.00	2
236	0.57	0.11	0.18	37
237	0.68	0.33	0.45	90
238	0.17	0.02	0.03	58
239	0.62	0.25	0.36	20
240	0.96	0.44	0.61	61
241	0.86	0.57	0.69	42
242	0.60	0.70	0.65	30
243	0.79	0.47	0.59	66
244	0.58	0.17	0.26	42
245	0.00	0.00	0.00	31
246	1.00	0.33	0.50	6
247	0.50	0.17	0.25	18
248	0.79	0.45	0.58	51
249	0.78	0.41	0.54	17
250	0.54	0.32	0.40	22
251	0.71	0.29	0.41	52
252	0.29	0.07	0.11	29
253	0.07	0.04	0.05	28
254	0.00	0.00	0.00	10
255	0.20	0.20	0.20	5
256	0.25	0.33	0.29	3
257	0.67	0.29	0.41	41
258	0.33	0.13	0.19	30
259	1.00	0.33	0.50	3
260	0.00	0.00	0.00	38
261	0.00	0.00	0.00	1
262	0.88	0.37	0.52	19
263	0.00	0.00	0.00	14
264	0.25	0.11	0.15	37
265	0.12	0.11	0.12	9
266	0.24	0.20	0.22	45
267	0.59	0.52	0.55	33
268	0.73	0.50	0.59	16
269	0.52	0.34	0.41	35
270	0.43	0.27	0.33	11
271	0.00	0.00	0.00	30
272	0.29	0.25	0.27	8
273	0.09	0.05	0.06	21
274	0.48	0.08	0.14	123
275	0.47	0.25	0.33	67
276	0.84	0.80	0.82	20
277	0.00	0.00	0.00	14
278	0.25	0.05	0.09	19
279	0.75	0.50	0.60	12
280	0.00	0.00	0.00	15

281	0.91	0.59	0.71	17
282	1.00	0.63	0.78	41
283	0.71	0.33	0.45	15
284	0.62	0.18	0.27	74
285	0.40	0.05	0.09	38
286	0.25	0.12	0.17	16
287	0.40	0.07	0.11	30
288	0.93	0.46	0.62	28
289	0.00	0.00	0.00	21
290	0.84	0.51	0.64	41
291	0.25	0.17	0.20	12
292	0.50	0.08	0.14	24
293	0.44	0.35	0.39	20
294	0.17	0.09	0.11	23
295	0.33	0.03	0.06	29
296	0.40	0.14	0.21	28
297	0.36	0.10	0.15	42
298	0.08	0.02	0.03	53
299	0.20	0.03	0.05	36
300	0.46	0.15	0.22	41
301	0.65	0.41	0.50	37
302	0.83	0.38	0.53	26
303	0.12	0.09	0.11	11
304	0.31	0.13	0.18	31
305	0.45	0.29	0.36	17
306	0.50	0.22	0.31	9
307	0.50	0.17	0.25	6
308	0.00	0.00	0.00	34
309	0.71	0.35	0.47	43
310	0.08	0.03	0.05	30
311	0.26	0.12	0.16	50
312	0.33	0.04	0.07	24
313	0.00	0.00	0.00	42
314	0.38	0.14	0.20	22
315	0.33	0.02	0.03	58
316	0.00	0.00	0.00	10
317	0.33	0.26	0.29	57
318	0.50	0.30	0.37	10
319	0.00	0.00	0.00	11
320	0.67	0.18	0.29	11
321	0.60	0.38	0.46	8
322	0.89	0.36	0.52	22
323	0.93	0.50	0.65	28
324	0.71	0.48	0.57	50
325	0.60	0.17	0.26	18
326	0.18	0.06	0.09	33
327	0.18	0.12	0.14	17
328	0.67	0.14	0.23	29
329	0.67	0.29	0.40	7
330	0.56	0.50	0.53	10
331	0.29	0.16	0.21	25
332	1.00	0.50	0.67	2
333	0.62	0.45	0.53	11
334	0.00	0.00	0.00	24
335	1.00	0.20	0.33	5
336	0.00	0.00	0.00	33
337	0.62	0.17	0.26	30
338	0.05	0.50	0.66	12

330	0.25	0.04	0.07	26
339	0.50	0.31	0.38	36
340	1.00	0.46	0.63	13
341	0.57	0.36	0.44	11
342	0.80	0.40	0.53	10
343	0.29	0.10	0.14	21
344	0.00	0.00	0.00	0
345	0.00	0.00	0.00	6
346	0.50	0.08	0.14	12
347	0.25	0.08	0.12	13
348	0.60	0.12	0.21	24
349	0.82	0.33	0.47	27
350	0.50	0.23	0.32	43
351	0.00	0.00	0.00	30
352	0.50	0.27	0.35	22
353	0.33	0.06	0.11	31
354	0.58	0.70	0.64	10
355	0.50	0.15	0.23	20
356	0.72	0.65	0.68	20
357	0.42	0.29	0.34	28
358	0.56	0.43	0.49	21
359	0.00	0.00	0.00	25
360	0.57	0.46	0.51	35
361	0.86	0.53	0.66	36
362	0.43	0.18	0.25	17
363	1.00	0.23	0.38	13
364	1.00	0.10	0.17	21
365	0.00	0.00	0.00	18
366	0.36	0.04	0.07	97
367	0.65	0.52	0.58	29
368	1.00	0.50	0.67	12
369	0.33	0.08	0.12	13
370	0.20	0.11	0.14	18
371	0.00	0.00	0.00	6
372	0.50	0.33	0.40	6
373	0.60	0.20	0.30	30
374	0.19	0.19	0.19	27
375	0.25	0.04	0.06	28
376	0.00	0.00	0.00	2
377	0.33	0.50	0.40	4
378	0.00	0.00	0.00	19
379	0.29	0.40	0.33	5
380	1.00	0.39	0.56	18
381	0.46	0.27	0.34	22
382	0.00	0.00	0.00	16
383	0.83	0.38	0.53	13
384	0.33	0.11	0.17	18
385	0.89	0.73	0.80	11
386	0.43	0.33	0.37	88
387	0.00	0.00	0.00	13
388	1.00	0.17	0.29	6
389	0.00	0.00	0.00	6
390	1.00	0.35	0.52	51
391	0.00	0.00	0.00	13
392	0.59	0.27	0.37	37
393	0.00	0.00	0.00	6
394	0.00	0.00	0.00	9
395	0.00	0.00	0.00	

396	0.00	0.00	0.00	13
397	1.00	0.50	0.67	6
398	0.42	0.34	0.38	29
399	0.94	0.52	0.67	33
400	0.50	0.03	0.06	31
401	0.50	0.06	0.11	50
402	0.86	0.33	0.48	18
403	0.33	0.14	0.20	7
404	0.76	0.50	0.60	26
405	0.83	0.62	0.71	56
406	1.00	0.75	0.86	4
407	0.33	0.24	0.28	17
408	0.50	0.09	0.15	11
409	0.33	0.06	0.10	18
410	0.40	0.20	0.27	10
411	0.64	0.20	0.31	45
412	0.80	0.40	0.53	20
413	0.50	0.04	0.07	25
414	0.43	0.15	0.22	20
415	0.00	0.00	0.00	6
416	0.20	0.08	0.11	26
417	1.00	0.40	0.57	10
418	0.00	0.00	0.00	18
419	1.00	0.17	0.29	6
420	0.62	0.47	0.53	17
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	6
423	0.00	0.00	0.00	12
424	1.00	0.25	0.40	4
425	1.00	0.36	0.53	11
426	0.50	0.09	0.15	11
427	1.00	0.75	0.86	8
428	0.71	0.19	0.30	26
429	0.56	0.38	0.45	40
430	0.00	0.00	0.00	2
431	0.00	0.00	0.00	35
432	0.50	0.07	0.12	15
433	0.00	0.00	0.00	18
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.30	0.11	0.16	28
437	0.36	0.12	0.18	33
438	0.80	0.40	0.53	20
439	0.00	0.00	0.00	36
440	1.00	0.11	0.20	18
441	0.47	0.44	0.46	18
442	0.75	0.56	0.64	16
443	0.33	0.05	0.08	22
444	0.00	0.00	0.00	6
445	0.86	0.29	0.43	21
446	0.78	0.30	0.44	46
447	0.26	0.09	0.13	69
448	0.00	0.00	0.00	7
449	0.00	0.00	0.00	3
450	0.22	0.04	0.07	52
451	0.00	0.00	0.00	16
452	1.00	0.71	0.83	17
453	0.00	0.00	0.00	13

454	0.67	0.18	0.29	11
455	1.00	0.08	0.15	12
456	0.40	0.33	0.36	6
457	0.33	0.11	0.17	18
458	0.20	0.07	0.10	15
459	0.90	0.32	0.47	28
460	0.00	0.00	0.00	18
461	0.56	0.50	0.53	10
462	0.50	0.12	0.20	24
463	0.00	0.00	0.00	18
464	0.87	0.33	0.48	39
465	0.14	0.09	0.11	11
466	0.20	0.09	0.12	35
467	0.11	0.05	0.07	21
468	0.00	0.00	0.00	37
469	1.00	0.20	0.33	5
470	0.25	0.12	0.17	8
471	0.82	0.24	0.38	37
472	0.00	0.00	0.00	47
473	0.50	0.29	0.36	14
474	1.00	0.52	0.69	23
475	0.54	0.30	0.39	66
476	0.00	0.00	0.00	3
477	0.62	0.26	0.37	19
478	0.00	0.00	0.00	1
479	0.22	0.09	0.12	23
480	0.25	0.02	0.03	60
481	0.20	0.04	0.06	26
482	0.67	0.50	0.57	4
483	0.00	0.00	0.00	8
484	1.00	0.30	0.47	23
485	0.83	0.28	0.42	18
486	0.67	0.17	0.27	12
487	0.75	0.10	0.18	29
488	1.00	1.00	1.00	1
489	0.67	0.33	0.44	6
490	0.40	0.29	0.33	7
491	0.00	0.00	0.00	3
492	0.22	0.20	0.21	10
493	0.50	0.32	0.39	19
494	0.00	0.00	0.00	7
495	0.67	0.25	0.36	8
496	0.14	0.06	0.08	18
497	0.00	0.00	0.00	72
498	0.00	0.00	0.00	8
499	0.38	0.09	0.15	32
micro avg	0.66	0.31	0.42	37472
macro avg	0.48	0.24	0.30	37472
weighted avg	0.60	0.31	0.40	37472
samples avg	0.41	0.31	0.33	37472

```
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
  'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
```

```

precision , predicted , average , warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'precision', 'predicted', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
    ...

```

## OneVsRestClassifier with Linear-SVM (SGDClassifier with loss-hinge)

```

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.00001, penalty='l1'), n_jobs=-1
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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 : 0.0621

Hamming loss 0.0084522

Micro-average quality numbers

Precision: 0.2038, Recall: 0.4318, F1-measure: 0.2769

Macro-average quality numbers

Precision: 0.1427, Recall: 0.3458, F1-measure: 0.1880

	precision	recall	f1-score	support
0	0.38	0.45	0.41	820
1	0.37	0.40	0.39	1931
2	0.12	0.27	0.17	544
3	0.17	0.30	0.21	222
4	0.47	0.58	0.52	1311
5	0.47	0.55	0.51	1014
6	0.45	0.51	0.48	1374
7	0.42	0.64	0.51	702
8	0.61	0.70	0.65	1424
9	0.66	0.73	0.69	1037
10	0.45	0.64	0.52	797
11	0.12	0.43	0.19	156
12	0.08	0.53	0.15	36
13	0.44	0.49	0.47	610
14	0.28	0.35	0.31	405
15	0.12	0.32	0.18	144
16	0.22	0.32	0.26	425
17	0.34	0.41	0.37	485
18	0.43	0.72	0.54	269
19	0.51	0.64	0.57	518
20	0.27	0.57	0.37	529
21	0.52	0.62	0.57	294
22	0.52	0.48	0.50	520
23	0.18	0.33	0.23	246
24	0.34	0.44	0.38	312
25	0.25	0.36	0.30	314
26	0.13	0.38	0.19	190
27	0.14	0.20	0.16	342
28	0.11	0.30	0.16	96
29	0.04	0.22	0.07	32
30	0.29	0.43	0.35	747
31	0.06	0.43	0.10	14
32	0.27	0.66	0.39	166
33	0.23	0.33	0.27	171
34	0.28	0.39	0.33	256
35	0.39	0.62	0.48	199
36	0.05	0.13	0.07	60
37	0.12	0.34	0.18	203
38	0.31	0.50	0.38	201
39	0.24	0.43	0.31	208
40	0.02	0.15	0.03	13
41	0.12	0.23	0.15	154
42	0.24	0.51	0.33	69
43	0.24	0.45	0.31	426
44	0.15	0.47	0.23	77
45	0.24	0.49	0.33	223
46	0.20	0.43	0.27	144
47	0.41	0.57	0.48	245
48	0.11	0.38	0.17	91

49	0.21	0.34	0.26	157
50	0.38	0.71	0.49	132
51	0.37	0.76	0.50	41
52	0.19	0.52	0.28	124
53	0.06	0.33	0.11	96
54	0.04	0.23	0.07	128
55	0.16	0.41	0.23	46
56	0.38	0.60	0.47	151
57	0.04	0.11	0.06	80
58	0.07	0.23	0.11	65
59	0.20	0.24	0.21	182
60	0.38	0.72	0.50	148
61	0.19	0.18	0.19	196
62	0.07	0.33	0.11	58
63	0.15	0.37	0.21	43
64	0.20	0.40	0.27	197
65	0.22	0.41	0.28	82
66	0.28	0.68	0.40	50
67	0.19	0.51	0.28	105
68	0.06	0.07	0.07	98
69	0.16	0.27	0.20	238
70	0.04	0.14	0.06	35
71	0.19	0.50	0.27	54
72	0.10	0.36	0.16	25
73	0.14	0.34	0.20	29
74	0.02	0.24	0.04	29
75	0.05	0.28	0.08	40
76	0.43	0.54	0.48	105
77	0.23	0.57	0.33	28
78	0.10	0.22	0.14	202
79	0.15	0.54	0.24	37
80	0.16	0.47	0.24	15
81	0.10	0.38	0.16	52
82	0.10	0.24	0.15	50
83	0.02	0.09	0.04	56
84	0.26	0.50	0.34	54
85	0.33	0.71	0.45	34
86	0.03	0.17	0.05	30
87	0.28	0.52	0.36	29
88	0.14	0.79	0.24	24
89	0.30	0.87	0.45	117
90	0.05	0.14	0.07	66
91	0.11	0.28	0.15	68
92	0.12	0.42	0.19	67
93	0.05	0.43	0.09	28
94	0.11	0.35	0.16	17
95	0.17	0.55	0.26	51
96	0.13	0.55	0.22	53
97	0.03	0.08	0.05	61
98	0.05	0.13	0.07	79
99	0.16	0.39	0.23	18
100	0.04	0.36	0.07	11
101	0.44	0.57	0.49	207
102	0.01	0.17	0.02	6
103	0.01	0.03	0.01	30
104	0.11	0.20	0.14	54
105	0.30	0.49	0.37	39
106	0.11	0.21	0.14	70

107	0.01	0.14	0.02	14
108	0.12	0.27	0.17	66
109	0.11	0.34	0.17	50
110	0.09	0.31	0.14	87
111	0.23	0.49	0.31	51
112	0.19	0.01	0.03	291
113	0.51	0.86	0.64	49
114	0.12	0.20	0.15	110
115	0.02	0.11	0.03	28
116	0.00	0.00	0.00	5
117	0.05	0.14	0.08	56
118	0.30	0.53	0.38	125
119	0.30	0.59	0.40	44
120	0.19	0.31	0.24	42
121	0.15	0.35	0.21	55
122	0.38	0.49	0.42	68
123	0.03	0.13	0.05	82
124	0.00	0.00	0.00	0
125	0.12	0.86	0.21	7
126	0.03	0.22	0.06	18
127	0.06	0.23	0.09	31
128	0.15	0.31	0.20	13
129	0.23	0.54	0.32	50
130	0.05	0.09	0.06	91
131	0.42	0.66	0.51	35
132	0.03	0.19	0.05	26
133	0.01	0.03	0.02	32
134	0.16	0.31	0.21	35
135	0.32	0.78	0.45	37
136	0.00	0.00	0.00	55
137	0.10	0.51	0.17	41
138	0.07	0.27	0.11	15
139	0.07	0.29	0.12	99
140	0.32	0.63	0.42	86
141	0.14	0.43	0.21	53
142	0.10	0.22	0.14	36
143	0.27	0.61	0.37	66
144	0.26	0.44	0.33	64
145	0.04	0.20	0.06	25
146	0.05	0.18	0.07	125
147	0.05	0.47	0.08	15
148	0.25	0.62	0.36	48
149	0.10	0.38	0.15	65
150	0.03	0.36	0.06	11
151	0.10	0.53	0.16	15
152	0.07	0.25	0.11	52
153	0.17	0.56	0.26	18
154	0.11	0.25	0.16	16
155	0.09	0.35	0.15	20
156	0.24	0.38	0.30	121
157	0.28	0.45	0.34	107
158	0.02	0.20	0.04	15
159	0.35	0.48	0.41	105
160	0.18	0.39	0.25	69
161	0.10	0.34	0.15	56
162	0.04	0.19	0.07	47
163	0.05	0.17	0.08	121
164	0.12	0.27	0.12	11

104	0.12	0.31	0.10	41
165	0.00	0.00	0.00	229
166	0.21	0.23	0.22	98
167	0.10	0.24	0.14	33
168	0.16	0.27	0.20	44
169	0.27	0.49	0.35	45
170	0.42	0.43	0.42	51
171	0.00	0.00	0.00	18
172	0.18	0.52	0.26	48
173	0.09	0.50	0.15	12
174	0.09	0.34	0.14	62
175	0.26	0.59	0.36	44
176	0.37	0.87	0.51	30
177	0.22	0.40	0.28	30
178	0.00	0.00	0.00	0
179	0.03	1.00	0.05	1
180	0.16	0.40	0.23	40
181	0.05	0.16	0.08	44
182	0.00	0.00	0.00	2
183	0.31	0.44	0.37	75
184	0.05	0.50	0.09	4
185	0.19	0.36	0.25	64
186	0.06	0.58	0.11	12
187	0.54	0.67	0.60	55
188	0.28	0.67	0.40	64
189	0.12	0.27	0.17	96
190	0.03	0.23	0.06	22
191	0.24	0.38	0.29	76
192	0.17	0.40	0.24	45
193	0.10	0.29	0.15	14
194	0.22	0.56	0.31	50
195	0.08	0.45	0.14	20
196	0.24	0.63	0.35	35
197	0.24	0.46	0.32	94
198	0.04	0.21	0.07	14
199	0.00	0.00	0.00	25
200	0.06	0.09	0.07	54
201	0.03	0.18	0.06	22
202	0.06	0.26	0.10	43
203	0.03	0.09	0.05	43
204	0.46	0.77	0.58	62
205	0.00	0.00	0.00	3
206	0.03	0.21	0.05	43
207	0.02	0.14	0.04	7
208	0.02	0.12	0.03	8
209	0.14	0.17	0.15	42
210	0.11	0.50	0.18	10
211	0.07	0.25	0.11	40
212	0.16	0.43	0.24	23
213	0.00	0.00	0.00	6
214	0.12	0.51	0.20	47
215	0.09	0.10	0.09	62
216	0.24	0.44	0.31	77
217	0.05	0.27	0.08	22
218	0.00	0.00	0.00	3
219	0.00	0.00	0.00	28
220	0.10	0.27	0.14	81
221	0.10	0.32	0.15	31

222	0.04	0.15	0.06	34
223	0.42	0.42	0.42	60
224	0.07	0.50	0.12	10
225	0.18	0.70	0.29	10
226	0.37	0.78	0.50	92
227	0.15	0.38	0.21	13
228	0.03	0.23	0.05	13
229	0.51	0.84	0.63	43
230	0.09	0.20	0.13	35
231	0.02	0.25	0.03	4
232	0.12	0.30	0.17	20
233	0.09	0.27	0.14	145
234	0.37	0.45	0.41	55
235	0.00	0.00	0.00	2
236	0.06	0.11	0.08	37
237	0.46	0.52	0.49	90
238	0.10	0.16	0.12	58
239	0.05	0.35	0.09	20
240	0.47	0.59	0.53	61
241	0.38	0.74	0.50	42
242	0.15	0.63	0.25	30
243	0.34	0.50	0.40	66
244	0.19	0.26	0.22	42
245	0.05	0.13	0.07	31
246	0.12	0.50	0.19	6
247	0.05	0.17	0.08	18
248	0.32	0.59	0.42	51
249	0.09	0.47	0.14	17
250	0.19	0.55	0.28	22
251	0.25	0.50	0.33	52
252	0.15	0.24	0.18	29
253	0.04	0.11	0.06	28
254	0.04	0.20	0.06	10
255	0.01	0.20	0.02	5
256	0.06	0.67	0.12	3
257	0.18	0.39	0.25	41
258	0.06	0.27	0.10	30
259	0.04	0.33	0.07	3
260	0.01	0.03	0.01	38
261	0.00	0.00	0.00	1
262	0.15	0.42	0.23	19
263	0.03	0.14	0.04	14
264	0.03	0.19	0.05	37
265	0.03	0.22	0.06	9
266	0.05	0.20	0.08	45
267	0.19	0.61	0.29	33
268	0.15	0.88	0.26	16
269	0.17	0.60	0.27	35
270	0.05	0.36	0.09	11
271	0.02	0.17	0.03	30
272	0.06	0.25	0.10	8
273	0.05	0.24	0.08	21
274	0.14	0.16	0.15	123
275	0.15	0.37	0.22	67
276	0.37	0.80	0.51	20
277	0.00	0.00	0.00	14
278	0.02	0.05	0.03	19
279	0.08	0.58	0.14	12

280	0.00	0.00	0.00	15
281	0.55	0.71	0.62	17
282	0.53	0.76	0.62	41
283	0.12	0.33	0.17	15
284	0.28	0.32	0.30	74
285	0.06	0.13	0.09	38
286	0.02	0.12	0.03	16
287	0.03	0.07	0.04	30
288	0.21	0.71	0.33	28
289	0.02	0.05	0.03	21
290	0.48	0.63	0.55	41
291	0.04	0.42	0.07	12
292	0.07	0.29	0.12	24
293	0.11	0.60	0.19	20
294	0.03	0.13	0.05	23
295	0.05	0.21	0.08	29
296	0.04	0.21	0.07	28
297	0.08	0.24	0.12	42
298	0.04	0.11	0.06	53
299	0.04	0.22	0.07	36
300	0.14	0.17	0.15	41
301	0.13	0.54	0.20	37
302	0.27	0.42	0.33	26
303	0.09	0.45	0.15	11
304	0.08	0.23	0.12	31
305	0.11	0.41	0.18	17
306	0.03	0.22	0.06	9
307	0.03	0.17	0.05	6
308	0.01	0.03	0.01	34
309	0.31	0.44	0.36	43
310	0.01	0.07	0.02	30
311	0.12	0.30	0.18	50
312	0.01	0.12	0.02	24
313	0.02	0.05	0.03	42
314	0.08	0.27	0.13	22
315	0.03	0.03	0.03	58
316	0.03	0.10	0.05	10
317	0.22	0.32	0.26	57
318	0.14	0.50	0.22	10
319	0.00	0.00	0.00	11
320	0.03	0.27	0.06	11
321	0.08	0.38	0.13	8
322	0.21	0.41	0.28	22
323	0.35	0.68	0.46	28
324	0.32	0.48	0.38	50
325	0.04	0.33	0.08	18
326	0.04	0.15	0.07	33
327	0.01	0.24	0.02	17
328	0.05	0.24	0.08	29
329	0.13	0.43	0.20	7
330	0.10	0.40	0.16	10
331	0.07	0.20	0.11	25
332	0.08	1.00	0.15	2
333	0.07	0.36	0.11	11
334	0.00	0.00	0.00	24
335	0.01	0.20	0.02	5
336	0.05	0.09	0.07	33
337	0.11	0.25	0.17	29

337	0.11	0.25	0.17	28
338	0.74	0.76	0.75	42
339	0.02	0.04	0.02	26
340	0.21	0.47	0.29	36
341	0.15	0.46	0.23	13
342	0.14	0.45	0.21	11
343	0.16	0.70	0.26	10
344	0.14	0.43	0.21	21
345	0.00	0.00	0.00	0
346	0.01	0.33	0.02	6
347	0.01	0.17	0.01	12
348	0.02	0.08	0.03	13
349	0.06	0.29	0.10	24
350	0.23	0.44	0.30	27
351	0.06	0.21	0.09	43
352	0.00	0.03	0.01	30
353	0.12	0.36	0.18	22
354	0.04	0.06	0.05	31
355	0.05	0.70	0.09	10
356	0.05	0.20	0.08	20
357	0.37	0.70	0.48	20
358	0.13	0.36	0.19	28
359	0.17	0.38	0.24	21
360	0.05	0.24	0.08	25
361	0.18	0.46	0.26	35
362	0.53	0.67	0.59	36
363	0.03	0.24	0.05	17
364	0.09	0.38	0.15	13
365	0.07	0.14	0.09	21
366	0.16	0.39	0.23	18
367	0.12	0.04	0.06	97
368	0.13	0.55	0.21	29
369	0.14	0.92	0.24	12
370	0.05	0.15	0.08	13
371	0.04	0.33	0.07	18
372	0.03	0.17	0.05	6
373	0.05	0.33	0.08	6
374	0.07	0.23	0.10	30
375	0.05	0.19	0.08	27
376	0.02	0.07	0.04	28
377	0.00	0.00	0.00	2
378	0.03	0.25	0.06	4
379	0.09	0.21	0.12	19
380	0.08	0.60	0.14	5
381	0.14	0.33	0.19	18
382	0.16	0.36	0.23	22
383	0.02	0.12	0.03	16
384	0.16	0.62	0.25	13
385	0.06	0.28	0.11	18
386	0.34	0.91	0.50	11
387	0.35	0.35	0.35	88
388	0.01	0.15	0.02	13
389	0.03	0.17	0.05	6
390	0.00	0.00	0.00	6
391	0.43	0.69	0.53	51
392	0.02	0.08	0.03	13
393	0.33	0.49	0.39	37
394	0.00	0.00	0.00	6

395	0.02	0.11	0.03	9
396	0.03	0.08	0.04	13
397	0.11	0.67	0.19	6
398	0.16	0.48	0.24	29
399	0.45	0.85	0.59	33
400	0.02	0.06	0.03	31
401	0.17	0.32	0.22	50
402	0.37	0.61	0.46	18
403	0.01	0.14	0.03	7
404	0.30	0.54	0.38	26
405	0.71	0.64	0.67	56
406	0.05	0.50	0.08	4
407	0.02	0.06	0.03	17
408	0.08	0.27	0.12	11
409	0.00	0.00	0.00	18
410	0.07	0.40	0.12	10
411	0.11	0.33	0.17	45
412	0.09	0.20	0.12	20
413	0.14	0.40	0.21	25
414	0.07	0.25	0.11	20
415	0.02	0.17	0.04	6
416	0.05	0.19	0.08	26
417	0.17	0.40	0.24	10
418	0.01	0.06	0.02	18
419	0.13	0.67	0.22	6
420	0.18	0.41	0.25	17
421	0.00	0.00	0.00	1
422	0.02	0.17	0.03	6
423	0.00	0.00	0.00	12
424	0.03	0.25	0.05	4
425	0.05	0.45	0.09	11
426	0.02	0.09	0.03	11
427	0.07	0.62	0.13	8
428	0.08	0.19	0.11	26
429	0.21	0.57	0.31	40
430	0.00	0.00	0.00	2
431	0.01	0.03	0.01	35
432	0.07	0.33	0.11	15
433	0.00	0.00	0.00	18
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.07	0.07	0.07	28
437	0.12	0.45	0.18	33
438	0.31	0.55	0.40	20
439	0.03	0.25	0.06	36
440	0.03	0.11	0.05	18
441	0.17	0.56	0.26	18
442	0.38	0.62	0.48	16
443	0.07	0.14	0.10	22
444	0.00	0.00	0.00	6
445	0.16	0.62	0.26	21
446	0.56	0.70	0.62	46
447	0.07	0.20	0.11	69
448	0.00	0.00	0.00	7
449	0.02	0.33	0.04	3
450	0.06	0.19	0.10	52
451	0.01	0.12	0.02	16
452	0.23	0.88	0.37	17

453	0.01	0.15	0.02	13
454	0.10	0.36	0.16	11
455	0.01	0.08	0.02	12
456	0.11	0.33	0.17	6
457	0.06	0.22	0.09	18
458	0.01	0.07	0.01	15
459	0.27	0.50	0.35	28
460	0.00	0.00	0.00	18
461	0.04	0.20	0.07	10
462	0.09	0.17	0.12	24
463	0.11	0.28	0.16	18
464	0.61	0.51	0.56	39
465	0.09	0.27	0.13	11
466	0.06	0.14	0.09	35
467	0.07	0.19	0.10	21
468	0.23	0.24	0.23	37
469	0.03	0.20	0.05	5
470	0.06	0.25	0.09	8
471	0.31	0.38	0.34	37
472	0.07	0.23	0.11	47
473	0.14	0.43	0.21	14
474	0.34	0.74	0.47	23
475	0.50	0.65	0.57	66
476	0.00	0.00	0.00	3
477	0.15	0.47	0.22	19
478	0.03	1.00	0.07	1
479	0.06	0.17	0.09	23
480	0.02	0.22	0.04	60
481	0.10	0.31	0.15	26
482	0.05	0.75	0.09	4
483	0.14	0.50	0.22	8
484	0.13	0.35	0.19	23
485	0.12	0.17	0.14	18
486	0.04	0.25	0.07	12
487	0.18	0.34	0.24	29
488	0.01	1.00	0.03	1
489	0.04	0.50	0.07	6
490	0.03	0.29	0.05	7
491	0.00	0.00	0.00	3
492	0.07	0.50	0.12	10
493	0.20	0.47	0.29	19
494	0.02	0.14	0.03	7
495	0.23	0.75	0.35	8
496	0.12	0.33	0.18	18
497	0.06	0.06	0.06	72
498	0.02	0.12	0.03	8
499	0.12	0.47	0.19	32
micro avg	0.20	0.43	0.28	37472
macro avg	0.14	0.35	0.19	37472
weighted avg	0.29	0.43	0.33	37472
samples avg	0.31	0.42	0.30	37472

Time taken to run this cell : 0:09:05.564627

```
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'recall', 'true', average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
'precision', 'recall', average, warn_for)
```

```
recall , true , average, warn_for)
C:\Users\Sudharshan.Reddy\AppData\Local\Continuum\anaconda3\lib\site-packages\sklearn\me
```

```
from prettytable import PrettyTable
x = PrettyTable(['Model','Accuracy','Macro f1-score','Micro f1-score','Hamming Loss'])
x.add_row(['LogReg_OvR_MoreWeights_HyperparameterTune_BOW','17.88','30.20','42.22','0.0031'])
x.add_row(['LinearSVM_OvR_MoreWeights_BOW','6.21','18.80','27.69','0.008'])
print(x)
```



Model	Accuracy	Macro f1-score	Micro f1-s
LogReg_OvR_MoreWeights_HyperparameterTune_BOW	17.88	30.20	42.22
LinearSVM_OvR_MoreWeights_BOW	6.21	18.80	27.69

1) For all the models I trained I used Only 100K data points due to limited resources.

2) I used Bag of Words to convert title text into vectors.

3) I trained two models

- Logistic Regression with Hyperparameter tuning
- Linear\_SVM(SGDClf with hing loss)

4) Considering the Micro f1-score values of two models It is clear that Logistic Regression with Hyperparameter tuning is better than SVM.