

Stack overflow tagging

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 knowledge, and build their careers.

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

Problem 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/>

Data Source : <https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data> (<https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data>)

Youtube : <https://youtu.be/nNDqbUhtIRg> (<https://youtu.be/nNDqbUhtIRg>)

Research paper : <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tagging-1.pdf> (<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> (<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> (<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 programming. The number of questions from each site may vary, and no filtering has been performed on the questions (such as closed questions).

Data Field 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, should not contain tabs '\t' or 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\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\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\n
    system("PAUSE");\n
    return 0;    \n
}\n

```

\n\n

The answer should come in the form of a table like
 \n\n

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

\n\n

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)

\n\n

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

\n'

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 as predicting properties of a data-point that are not mutually exclusive, such as topics that are relevant for a document. A question on Stackoverflow might be about any of C, Pointers, FileIO and/or memory-management at the same time or none of these.

__Credit__ : <http://scikit-learn.org/stable/modules/multiclass.html>

2.2.2 Performance metric

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

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

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

'Micro f1 score':

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

'Macro f1 score':

Calculate metrics for each label, and find their unweighted mean. This does not take label imbalance into account.

<https://www.kaggle.com/wiki/MeanFScore> (<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> (<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

In [2]:

```
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 sklearnlearn.adapt import mlknn
from sklearnlearn.problem_transform import ClassifierChain
from sklearnlearn.problem_transform import BinaryRelevance
from sklearnlearn.problem_transform import LabelPowerset
from sklearn.naive_bayes import GaussianNB
from datetime import datetime
```

In [3]:

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

In [4]:

```
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 file")
```

Number of rows in the database :

6034196

Time taken to count the number of rows : 0:01:57.510265

Checking for duplicates

```
In [7]:
#Learn SQL: https://www.w3schools.com/sql/default.asp
#below am considering 10k data points
if os.path.isfile('train.db'):
    start = datetime.now()
    con = sqlite3.connect('train.db')
    df_no_dup = pd.read_sql_query('SELECT Title, Body, Tags, COUNT(*) as cnt_dup FROM data GROUP BY Title, Body, Tags limit 10000 ', con)
    con.close()
    print("Time taken to run this cell :", datetime.now() - start)
else:
    print("Please download the train.db file from drive or run the first to generate train.db file")
```

Time taken to run this cell : 0:09:47.184572

```
In [8]:
df_no_dup.head()
# we can observe that there are duplicates
```

Out[8]:

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

```
In [9]:
print("number of duplicate questions :", num_rows['count(*)'].values[0]- df_no_dup.shape[0], "(", (1-((df_no_dup.s
hape[0])/(num_rows['count(*)'].values[0]))) *100, "% )")

number of duplicate questions : 6024196 ( 99.83427783916864 % )
```

```
In [10]:
# number of times each question appeared in our database
df_no_dup.cnt_dup.value_counts()
```

Out[10]:

1	6227
2	3068
3	705

Name: cnt_dup, dtype: int64

```
In [11]:
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 : 0:00:01.095063

Out[11]:

	Title	Body	Tags	cnt_dup	tag_count
0	Implementing Boundary Value Analysis of S...	<pre><code>#include<stdio.h>\n#include<...</pre>	c++ c	1	2
1	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamicall...	c# silverlight data-binding	1	3
2	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamicall...	c# silverlight data-binding columns	1	4
3	java.lang.NoClassDefFoundError: javax/serv...	<p>I followed the guide in <a href="http://sta...	jsp jstl	1	2
4	java.sql.SQLException: [Microsoft][ODBC Dri...	<p>I use the following code</p>\n\n<pre><code>...</pre>	java jdbc	2	2

In [12]:

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

Out[12]:

```
3    2918
2    2657
4    1908
1    1303
5    1214
Name: tag_count, dtype: int64
```

In [3]:

```
#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)
```

In [4]:

```
#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 file")
```

Time taken to run this cell : 0:02:10.011231

3.2 Analysis of Tags

3.2.1 Total number of unique tags

In [5]:

```
# 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'])
```

In [6]:

```
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

In [7]:

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

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

3.2.3 Number of times a tag appeared

In [8]:

```
# 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))
```

In [9]:

```
#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()
```

Out[9]:

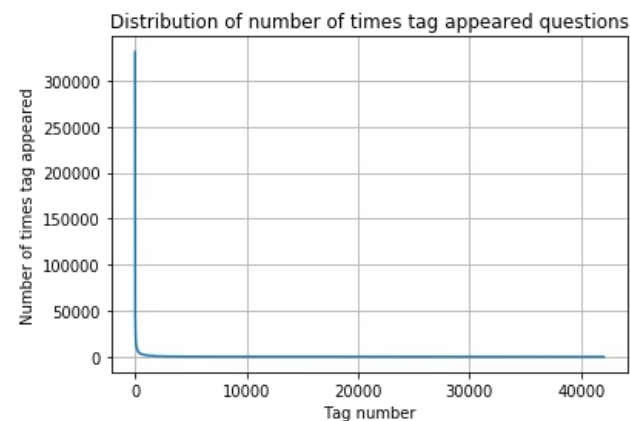
	Tags	Counts
0	.a	18
1	.app	37
2	.asp.net-mvc	1
3	.aspxauth	21
4	.bash-profile	138

In [10]:

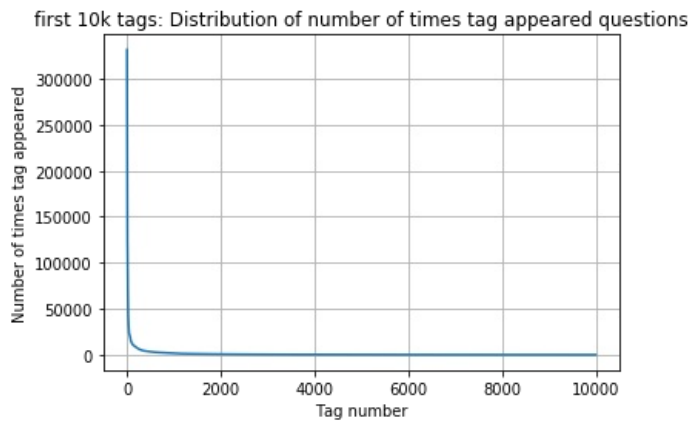
```
tag_df_sorted = tag_df.sort_values(['Counts'], ascending=False)
tag_counts = tag_df_sorted['Counts'].values
```

In [11]:

```
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()
```



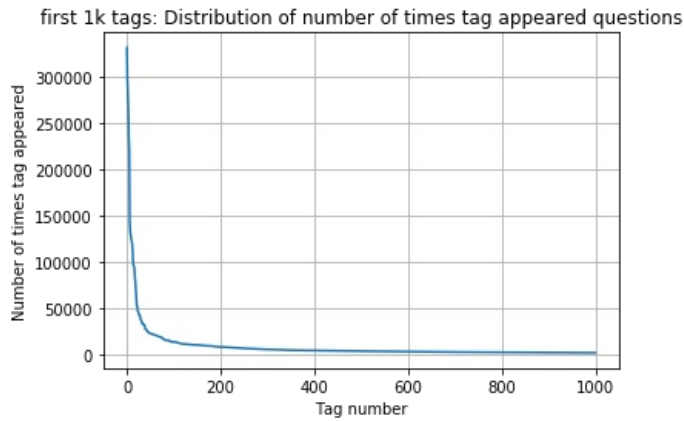

```
In [12]:
plt.plot(tag_counts[0:10000])
plt.title('first 10k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:10000:25]), tag_counts[0:10000:25])
```



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]	

In [13]:

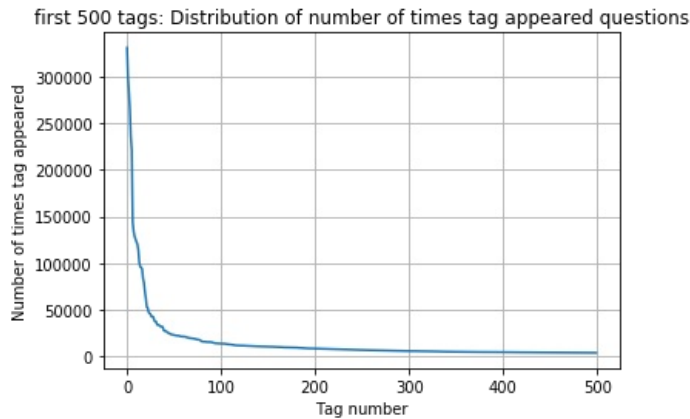
```
plt.plot(tag_counts[0:1000])
plt.title('first 1k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



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]	

In [14]:

```
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])
```



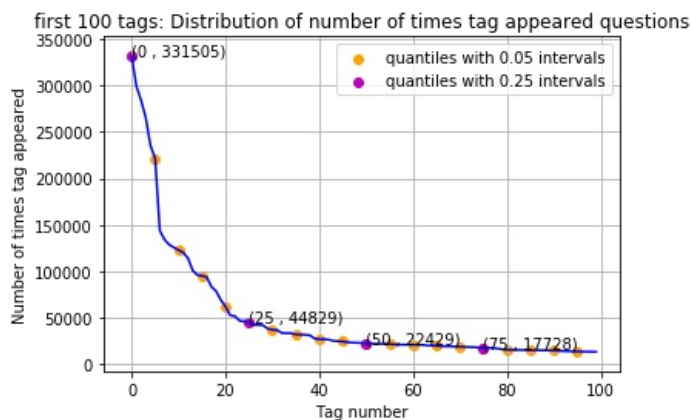
```
100 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291 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]
```

In [15]:

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

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

plt.title('first 100 tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.legend()
plt.show()
print(len(tag_counts[0:100:5]), tag_counts[0:100:5])
```



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

In [16]:

```
# 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 for this problem.

3.2.4 Tags Per Question

In [17]:

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

print(tag_quest_count[:5])
```

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

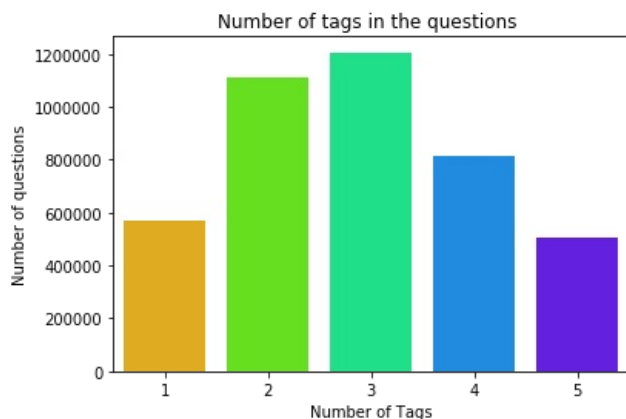
In [18]:

```
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

In [19]:

```
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

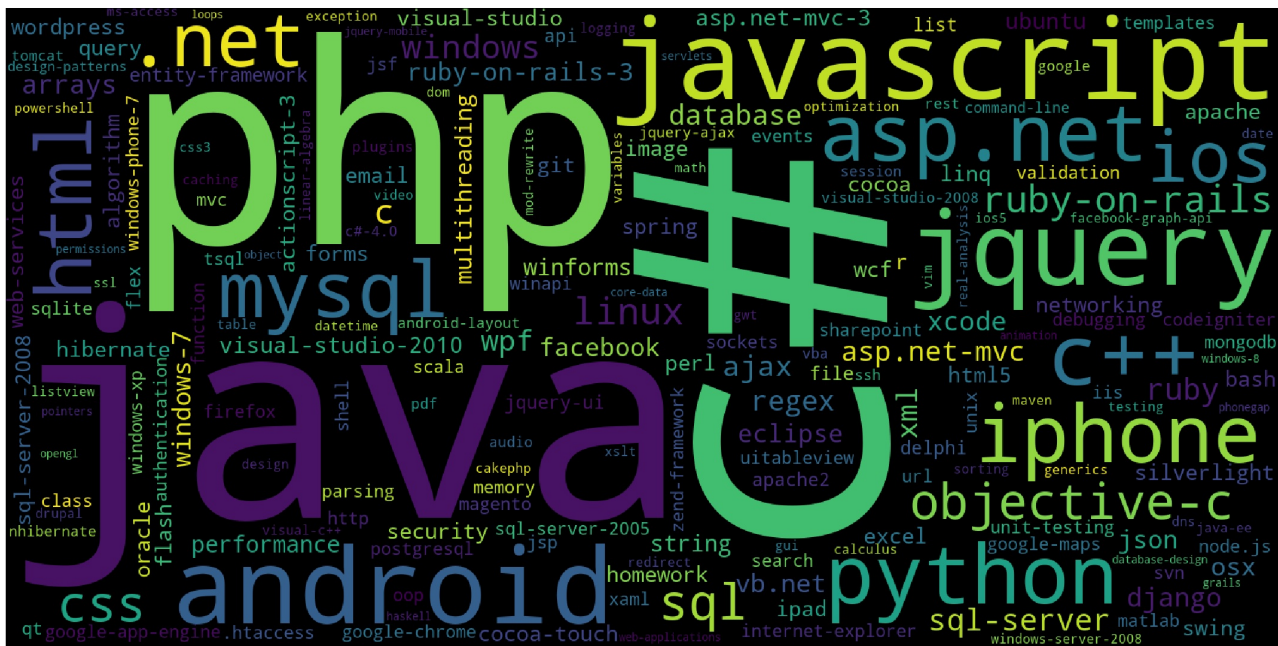
3.2.5 Most Frequent Tags

In [20]:

```
# Plotting 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 : 0:00:13.678783

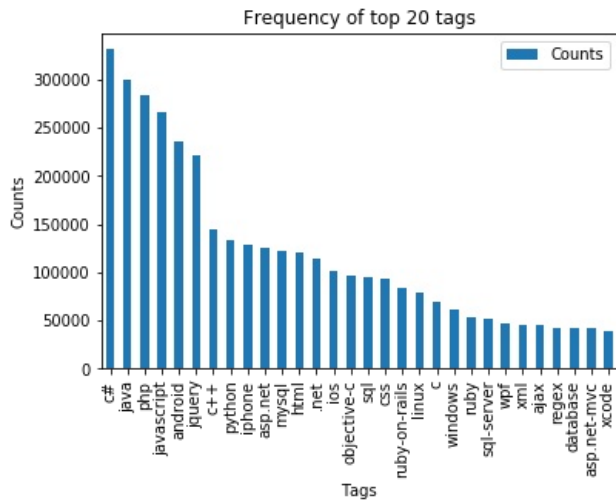
Observations:

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

In [21]:

```
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 100k 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

In [22]:

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

In [23]:

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

    return None

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

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

def create_database_table(database, query):
    conn = create_connection(database)
    if conn is not None:
        create_table(conn, query)
        checkTableExists(conn)
    else:
        print("Error! cannot create the database connection.")
    conn.close()

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

Tables in the database:
QuestionsProcessed

In [24]:

```
# 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 : 0:03:55.652026

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

In [25]:

```
#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], 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=stripthtml(question.encode('utf-8'))

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

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

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

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

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

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

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

Avg. length of questions(Title+Body) before processing: 1165
Avg. length of questions(Title+Body) after processing: 327
Percent of questions containing code: 57
Time taken to run this cell : 0:08:07.599397

In [26]:

```
# 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()
```


In [27]:

```
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

=====

('cassandra get multipl row column slice possibl execut queri get row key fake key fake key column s
lice c',)

('take java method name function arg clojur want creat function take symbol repres java method appli
object execut get result much differ expect question symbol call second arg function instead valu bo
und would actual want',)

('matlab candl chart handl use use handl chang background color chart thank',)

('cant play mpmovieplayercontrol record via avaudiorecord time tri play video small area screen reco
rd user sing time howev seem record success ni ad code point possibl blunder hope make sens flow see
m enter even audiorecord stop pleas help search sinc last day luck',)

('hibern insid glassfish tri integr hibern framework web applic run glassfish howev find practic doc
ument tutori subject pleas recommend practic document subject thank advanc',)

('get current file encod python file environ tell encod sourc file insid run python process even pos
sibl',)

('send packet java applet written simpl java applet send sip packet server run within eclips sun app
let viewer everyth work perfect attempt emb applet web browser use applet tag applet success run pac
ket sent verifi use wireshark kind secur set ie chrome awar guess show code necessari thank',)

('generat random number sort java goal generat random number add linkedlist object sort element code
far run problem want display sort element error get except thread main java util illegalformatconver
sionexcept java util array arraylist someon throw light problem thank enter imag descript',)

('someon use net framework client profil good reason use net framework client profil instead full ve
rision mean real life reason creat net applic sinc quit easi creat instal instal net framework client
machin bother use client profil',)

In [28]:

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

In [29]:

```
preprocessed_data.head()
```

Out[29]:

	question	tags
0	codeignit uri error live site use codeignit ho...	codeigniter routing uri http-status-code-404 host
1	cassandra get multipl row column slice possibl...	cassandra
2	take java method name function arg clojur want...	function clojure macros
3	matlab candl chart handl use use handl chang b...	matlab candlestick-chart
4	cant play mpmovieplayercontrol record via avau...	iphone mpmovieplayercontroller avaudiorecorder

In [30]:

```
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

In [38]:

```
# 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)

In [40]:

```
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

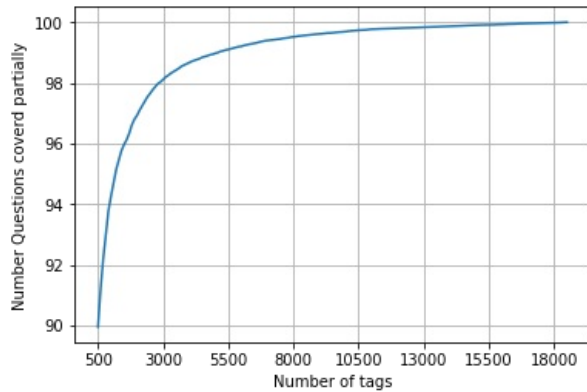
def questions_explained_fn(n):
    multilabel_yn = tags_to_choose(n)
    x= multilabel_yn.sum(axis=1)
    return (np.count_nonzero(x==0))
```

In [40]:

```
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))
```

In [41]:

```
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 50(it covers 90% of the tags)
print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
```



with 5500 tags we are covering 99.099 % of questions

In [42]:

```
multilabel_yx = tags_to_choose(5500)
print("number of questions that are not covered :", questions_explained_fn(5500),"out of ", total_qs)
```

number of questions that are not covered : 901 out of 99999

In [43]:

```
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.shape[1])*100,"%")
")
```

Number of tags in sample : 18558
number of tags taken : 5500 (29.636814311887054 %)

We consider top 15% tags which covers 99% of the questions

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

In [44]:

```
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,:]
```

In [45]:

```
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, 5500)
Number of data points in test data : (20000, 5500)

4.3 Featurizing data

In [46]:

```
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 : 0:01:37.727128

In [47]:

```
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, 90081) Y : (79999, 5500)
Dimensions of test data X: (20000, 90081) Y: (20000, 5500)

In [48]:

```
# 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)
```

Out[48]:

```
"\nfrom skmultilearn.adapt import MLkNN\nnclassifier = MLkNN(k=21)\n\n# train\nnclassifier.fit(x_train_multilabel, y_train)\n\n# predict\nnpredictions = classifier.predict(x_test_multilabel)\n\nprint(accuracy_score(y_test,predictions))\nprint(metrics.f1_score(y_test, predictions, average = 'macro'))\n\nprint(metrics.f1_score(y_test, predictions, average = 'micro'))\n\nprint(metrics.hamming_loss(y_test,predictions))\n\n"
```

4.5 Modeling with less data points (100k data points) and more weight to title and 500 tags only.

In [30]:

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

Tables in the database:
QuestionsProcessed

In [31]:

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

read_db = 'train_no_dup.db'
write_db = 'Titlemoreweight.db'
train_datasize = 90000
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 100000;")
        # 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**
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

In [32]:

```
#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=stripthtml(question.encode('utf-8'))

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

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

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

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

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

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

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

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

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

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

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 : 0:13:32.701480

In [33]:

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

Sample quesitons after preprocessing of data

In [34]:

```
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 silverlight bind
datagrid dynam code wrote code debug code block seem bind correct grid come column form come grid
column although necessari bind nthank repli advance.',)
-----
('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid java.lang.noclassdeffounde
rror javax servlet jsp tagext taglibraryvalid java.lang.noclassdeffounderror javax servlet jsp tagex
t taglibraryvalid follow guid link instal jstl got follow error tri launch jsp page java.lang.noclas
sdeffounderror javax servlet jsp tagext taglibraryvalid taglib declar instal jstl 1.1 tomcat webapp
tri project work also tri version 1.2 jstl still messag caus solv',)
-----
('java.sql.sqllexcept microsoft odbc driver manag invalid descriptor index java.sql.sqllexcept microso
ft odbc driver manag invalid descriptor index java.sql.sqllexcept microsoft odbc driver manag invalid
descriptor index use follow code display caus solv',)
-----
('better way updat feed fb php sdk better way updat feed fb php sdk better way updat feed fb php sdk
novic facebook api read mani tutori still confused.i find post feed apimethod like correct second w
ay use curl someth like way better',)
-----
('btnadd click event open two window record ad btnadd click event open two window record ad btnadd c
lick event open two window record ad open window search.aspx use code hav add button search.aspx nwh
en insert record btnadd click event open anoth window nafter insert record close window',)
-----
('sql inject issu prevent correct form submiss php sql inject issu prevent correct form submiss php
sql inject issu prevent correct form submiss php check everyth think make sure input field safe type
sql inject good news safe bad news one tag mess form submiss place even touch life figur exact html
use templat file forgiv okay entir php script get execut see data post none forum field post problem
use someth titl field none data get post current use print post see submit noth work flawless statem
ent though also mention script work flawless local machin use host come across problem state list in
put test mess',)
-----
('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl subaddit lebesgu measur
let lbrace rbrace sequenc set sigma -algebra mathcal want show left bigcup right leq sum left right
countabl addit measur defin set sigma algebra mathcal think use monoton properti somewher proof star
t appreci littl help nthank ad han answer make follow addit construct given han answer clear bigcup
bigcup cap emptyset neq left bigcup right left bigcup right sum left right also construct subset mon
oton left right leq left right final would sum leq sum result follow',)
-----
('hql equival sql queri hql equival sql queri hql equival sql queri hql queri replac name class prop
erti name error occur hql error',)
-----
('undefin symbol architectur i386 objc class skpsmtpmessag referenc error undefin symbol architectur
i386 objc class skpsmtpmessag referenc error undefin symbol architectur i386 objc class skpsmtpmessa
g referenc error import framework send email applic background import framework i.e skpsmtpmessag so
mebodi suggest get error collect2 ld return exit status import framework correct sorc taken framewor
k follow mfmcomposeviewcontrol question lock field updat answer drag drop folder project click co
pi nthat',)
-----
```

Saving Preprocessed data to a Database

In [35]:

```
#Taking 100k entries to a dataframe.
write_db = 'Titlemoreweight.db'
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM QuestionsProcessed""", conn_r)
conn_r.commit()
conn_r.close()
```

In [36]:

```
preprocessed_data.head()
```

Out[36]:

	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.sqlexcept microsoft odbc driver manag...	java jdbc
4	better way updat feed fb php sdk better way up...	facebook api facebook-php-sdk

In [37]:

```
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

Converting string Tags to multilable output variables

In [38]:

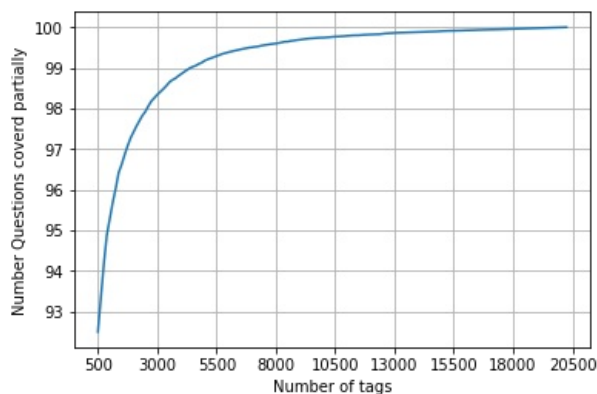
```
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

In [41]:

```
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))
```

In [42]:

```
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, minimun is 500(it covers 90% of the tags)
print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



with 5500 tags we are covering 99.481 % of questions
with 500 tags we are covering 92.5 % of questions

Selecting 500 Tags

In [43]:

```
# 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 99999

In [44]:

```
x_train=preprocessed_data.head(train_datasize)
x_test=preprocessed_data.tail(preprocessed_data.shape[0] - 90000)

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

In [45]:

```
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 : (90000, 500)
Number of data points in test data : (9999, 500)

4.5.2 Featurizing data with Tfidf vectorizer

In [46]:

```
start = datetime.now()
vectorizer = TfidfVectorizer(min_df=0.00009, max_features=20000, 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 : 1:22:00.514420

In [49]:

```
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: (90000, 20000) Y : (90000, 500)
Dimensions of test data X: (9999, 20000) Y: (9999, 500)

4.5.3 Applying Logistic Regression with OneVsRest Classifier

In [50]:

```
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 : 0.19091909190919093
Hamming loss 0.0030715071507150713
Micro-average quality numbers
Precision: 0.7155, Recall: 0.3213, F1-measure: 0.4435

Macro-average quality numbers
Precision: 0.4468, Recall: 0.2155, F1-measure: 0.2710

	precision	recall	f1-score	support
--	-----------	--------	----------	---------

0	0.84	0.38	0.53	374
1	0.69	0.21	0.33	902
2	0.57	0.10	0.17	229
3	0.63	0.15	0.24	129
4	0.78	0.50	0.61	683
5	0.90	0.48	0.62	517
6	0.78	0.38	0.51	704
7	0.87	0.64	0.74	385
8	0.91	0.58	0.71	902
9	0.72	0.59	0.65	945
10	0.83	0.51	0.63	335
11	0.74	0.42	0.54	66
12	1.00	0.42	0.59	12
13	0.78	0.42	0.55	308
14	0.32	0.12	0.17	143
15	0.92	0.14	0.24	81
16	0.65	0.23	0.34	199
17	0.63	0.25	0.36	252
18	0.63	0.49	0.55	69
19	0.89	0.57	0.70	292
20	0.55	0.21	0.31	472
21	0.78	0.56	0.65	116
22	0.82	0.42	0.56	354
23	0.61	0.31	0.41	107
24	0.81	0.36	0.50	138
25	0.19	0.09	0.12	67
26	0.90	0.20	0.32	96
27	0.28	0.05	0.09	176
28	0.62	0.14	0.23	35
29	1.00	0.15	0.26	20
30	0.31	0.67	0.42	27
31	0.00	0.00	0.00	6
32	0.71	0.57	0.64	96
33	0.65	0.42	0.51	71
34	0.81	0.26	0.40	184
35	0.84	0.61	0.71	124
36	0.00	0.00	0.00	22
37	0.34	0.18	0.23	73
38	0.69	0.41	0.51	134
39	0.66	0.27	0.38	109
40	0.00	0.00	0.00	7
41	0.53	0.12	0.20	80
42	0.00	0.00	0.00	6
43	0.40	0.10	0.16	416
44	0.44	0.30	0.36	27
45	0.52	0.25	0.33	142
46	0.70	0.22	0.34	72
47	0.79	0.51	0.62	121
48	0.67	0.29	0.41	34
49	0.73	0.27	0.40	88
50	0.84	0.72	0.77	60
51	0.83	0.62	0.71	16
52	0.66	0.37	0.48	67
53	0.37	0.17	0.23	59
54	0.20	0.09	0.12	35
55	0.67	0.42	0.52	19
56	0.68	0.55	0.61	74
57	0.00	0.00	0.00	46
58	0.58	0.21	0.31	33
59	0.41	0.15	0.22	81
60	0.93	0.73	0.82	73
61	0.30	0.04	0.08	135
62	0.36	0.19	0.25	21
63	0.88	0.29	0.44	24
64	0.68	0.25	0.36	153
65	0.71	0.38	0.49	40
66	0.50	0.14	0.22	7
67	0.66	0.53	0.58	59
68	0.21	0.10	0.14	29
69	0.47	0.15	0.22	131
70	0.50	0.05	0.09	21
71	0.62	0.40	0.48	20
72	0.00	0.00	0.00	9
73	0.00	0.00	0.00	1
74	1.00	0.07	0.12	15
75	0.71	0.29	0.42	17
76	0.87	0.62	0.72	55
77	0.00	0.00	0.00	2
78	0.44	0.04	0.07	103

79	0.58	0.27	0.37	26
80	1.00	0.21	0.35	14
81	0.43	0.26	0.32	23
82	0.43	0.18	0.25	17
83	0.00	0.00	0.00	35
84	0.80	0.43	0.56	28
85	0.44	0.36	0.40	11
86	0.00	0.00	0.00	11
87	1.00	0.40	0.57	5
88	0.79	0.85	0.81	13
89	0.93	0.90	0.91	87
90	0.00	0.00	0.00	39
91	0.57	0.12	0.21	32
92	0.67	0.14	0.24	28
93	0.55	0.33	0.41	18
94	0.67	0.25	0.36	8
95	0.82	0.47	0.60	38
96	0.69	0.36	0.47	25
97	1.00	0.03	0.07	29
98	0.00	0.00	0.00	57
99	0.00	0.00	0.00	0
100	0.00	0.00	0.00	4
101	0.83	0.43	0.57	23
102	0.00	0.00	0.00	3
103	0.50	0.07	0.12	15
104	0.00	0.00	0.00	19
105	0.91	0.43	0.59	23
106	0.54	0.14	0.22	50
107	0.00	0.00	0.00	10
108	0.71	0.14	0.23	37
109	0.50	0.25	0.33	28
110	0.90	0.18	0.30	51
111	0.38	0.38	0.38	21
112	0.00	0.00	0.00	286
113	0.96	0.87	0.91	30
114	0.00	0.00	0.00	9
115	0.17	0.06	0.08	18
116	0.00	0.00	0.00	2
117	0.75	0.13	0.22	23
118	0.81	0.44	0.57	66
119	0.78	0.29	0.42	24
120	0.00	0.00	0.00	1
121	0.56	0.21	0.30	24
122	0.92	0.46	0.62	26
123	0.00	0.00	0.00	34
124	0.00	0.00	0.00	0
125	1.00	0.67	0.80	3
126	0.00	0.00	0.00	8
127	1.00	0.19	0.32	21
128	1.00	0.62	0.77	8
129	0.71	0.52	0.60	23
130	0.00	0.00	0.00	53
131	0.00	0.00	0.00	1
132	0.00	0.00	0.00	9
133	0.00	0.00	0.00	6
134	1.00	0.20	0.33	5
135	1.00	0.60	0.75	10
136	0.00	0.00	0.00	27
137	0.10	0.14	0.12	21
138	1.00	0.33	0.50	3
139	0.33	0.07	0.11	29
140	0.93	0.58	0.72	24
141	0.73	0.31	0.43	36
142	0.43	0.20	0.27	15
143	0.44	0.30	0.36	27
144	0.68	0.50	0.58	34
145	0.00	0.00	0.00	8
146	0.17	0.01	0.02	86
147	0.40	0.25	0.31	8
148	0.88	0.37	0.52	19
149	0.47	0.20	0.28	35
150	0.00	0.00	0.00	8
151	0.25	0.33	0.29	3
152	0.20	0.07	0.11	28
153	0.50	0.50	0.50	2
154	0.00	0.00	0.00	1
155	0.33	0.09	0.14	11
156	0.67	0.15	0.25	52
157	0.53	0.40	0.45	43
158	1.00	0.14	0.25	7
159	0.68	0.42	0.52	55
160	0.56	0.27	0.36	37
161	1.00	0.25	0.40	16

162	0.00	0.00	0.00	21
163	0.40	0.03	0.05	75
164	0.25	0.05	0.08	20
165	0.00	0.00	0.00	227
166	0.88	0.09	0.17	75
167	0.57	0.29	0.38	14
168	0.62	0.16	0.26	31
169	0.73	0.44	0.55	18
170	0.85	0.48	0.61	23
171	0.00	0.00	0.00	5
172	0.75	0.56	0.64	27
173	0.33	0.50	0.40	2
174	0.33	0.08	0.13	36
175	0.82	0.47	0.60	19
176	1.00	0.80	0.89	10
177	0.17	0.08	0.11	12
178	0.00	0.00	0.00	0
179	0.00	0.00	0.00	0
180	0.50	0.21	0.30	14
181	0.00	0.00	0.00	12
182	0.00	0.00	0.00	1
183	0.65	0.34	0.45	38
184	0.00	0.00	0.00	2
185	0.88	0.25	0.39	28
186	0.33	0.50	0.40	4
187	1.00	0.70	0.82	30
188	0.75	0.56	0.64	27
189	0.43	0.12	0.19	48
190	0.00	0.00	0.00	12
191	0.73	0.22	0.34	50
192	0.54	0.32	0.40	22
193	0.00	0.00	0.00	4
194	0.71	0.48	0.57	25
195	1.00	0.14	0.25	7
196	0.86	0.63	0.73	19
197	0.62	0.25	0.36	52
198	1.00	0.11	0.20	9
199	0.00	0.00	0.00	13
200	1.00	0.04	0.07	28
201	0.50	0.17	0.25	6
202	0.80	0.24	0.36	17
203	0.00	0.00	0.00	21
204	1.00	0.47	0.64	34
205	0.00	0.00	0.00	1
206	0.25	0.03	0.05	35
207	0.00	0.00	0.00	3
208	0.00	0.00	0.00	4
209	0.25	0.11	0.15	28
210	0.00	0.00	0.00	1
211	0.75	0.15	0.25	20
212	1.00	0.50	0.67	6
213	0.00	0.00	0.00	2
214	0.50	0.27	0.35	15
215	0.57	0.13	0.22	30
216	0.74	0.37	0.49	38
217	0.33	0.08	0.13	12
218	0.00	0.00	0.00	1
219	0.00	0.00	0.00	16
220	0.80	0.05	0.10	79
221	1.00	0.13	0.24	15
222	0.00	0.00	0.00	15
223	1.00	0.32	0.49	34
224	0.00	0.00	0.00	5
225	1.00	0.67	0.80	3
226	0.77	0.71	0.74	48
227	0.00	0.00	0.00	0
228	0.50	0.20	0.29	5
229	0.87	0.77	0.82	26
230	0.50	0.12	0.19	26
231	0.00	0.00	0.00	0
232	0.71	0.36	0.48	14
233	0.00	0.00	0.00	2
234	0.93	0.54	0.68	24
235	0.00	0.00	0.00	1
236	0.43	0.20	0.27	15
237	0.76	0.39	0.52	41
238	0.67	0.18	0.29	22
239	1.00	0.20	0.33	10
240	1.00	0.50	0.67	26
241	0.92	0.73	0.81	15
242	0.58	0.73	0.65	15
243	0.83	0.69	0.75	29
244	0.67	0.21	0.32	29

245	0.00	0.00	0.00	6
246	0.00	0.00	0.00	2
247	1.00	0.08	0.14	13
248	0.83	0.50	0.62	30
249	0.71	0.45	0.56	11
250	0.20	0.10	0.13	10
251	0.82	0.38	0.51	24
252	0.20	0.08	0.12	12
253	0.25	0.08	0.12	12
254	0.00	0.00	0.00	1
255	1.00	0.50	0.67	2
256	0.00	0.00	0.00	1
257	0.67	0.25	0.36	16
258	0.50	0.06	0.11	16
259	0.00	0.00	0.00	2
260	0.00	0.00	0.00	17
261	0.00	0.00	0.00	0
262	0.75	0.27	0.40	11
263	0.00	0.00	0.00	1
264	0.33	0.05	0.09	20
265	0.00	0.00	0.00	3
266	0.25	0.04	0.06	28
267	0.50	0.47	0.48	17
268	0.71	0.50	0.59	10
269	0.79	0.48	0.59	23
270	0.38	0.38	0.38	8
271	0.00	0.00	0.00	20
272	0.00	0.00	0.00	0
273	0.00	0.00	0.00	6
274	0.25	0.17	0.20	6
275	0.71	0.31	0.43	39
276	1.00	0.78	0.88	9
277	0.00	0.00	0.00	8
278	0.00	0.00	0.00	6
279	0.80	0.80	0.80	5
280	0.00	0.00	0.00	4
281	1.00	0.67	0.80	3
282	1.00	0.67	0.80	15
283	0.00	0.00	0.00	0
284	0.67	0.32	0.44	37
285	0.50	0.14	0.22	21
286	0.50	0.09	0.15	11
287	0.00	0.00	0.00	18
288	0.88	0.44	0.58	16
289	0.00	0.00	0.00	11
290	0.74	0.58	0.65	24
291	0.50	0.25	0.33	4
292	0.50	0.22	0.31	9
293	0.50	0.55	0.52	11
294	0.00	0.00	0.00	14
295	0.00	0.00	0.00	13
296	0.50	0.25	0.33	8
297	0.43	0.19	0.26	16
298	0.00	0.00	0.00	34
299	0.00	0.00	0.00	16
300	0.50	0.10	0.16	21
301	0.81	0.57	0.67	23
302	0.80	0.36	0.50	11
303	0.00	0.00	0.00	3
304	0.29	0.12	0.17	16
305	0.00	0.00	0.00	6
306	0.00	0.00	0.00	3
307	0.00	0.00	0.00	2
308	0.00	0.00	0.00	14
309	0.67	0.40	0.50	25
310	1.00	0.06	0.11	17
311	0.25	0.07	0.11	30
312	0.00	0.00	0.00	11
313	1.00	0.07	0.13	14
314	0.33	0.14	0.20	14
315	0.00	0.00	0.00	38
316	1.00	0.14	0.25	7
317	0.38	0.23	0.29	26
318	0.50	1.00	0.67	1
319	0.00	0.00	0.00	5
320	0.00	0.00	0.00	0
321	0.00	0.00	0.00	0
322	0.80	0.50	0.62	8
323	0.89	0.67	0.76	12
324	0.64	0.55	0.59	29
325	0.00	0.00	0.00	4
326	0.20	0.10	0.13	10
327	0.00	0.00	0.00	8

328	0.00	0.00	0.00	10
329	0.00	0.00	0.00	0
330	1.00	0.33	0.50	3
331	1.00	0.08	0.15	12
332	0.00	0.00	0.00	0
333	0.67	0.29	0.40	7
334	0.00	0.00	0.00	14
335	0.00	0.00	0.00	0
336	0.00	0.00	0.00	17
337	0.89	0.40	0.55	20
338	1.00	0.06	0.12	16
339	0.67	0.31	0.42	13
340	0.59	0.50	0.54	26
341	1.00	0.83	0.91	6
342	0.00	0.00	0.00	0
343	0.00	0.00	0.00	0
344	0.25	0.07	0.11	15
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	1
347	1.00	0.25	0.40	4
348	0.00	0.00	0.00	5
349	0.67	0.14	0.24	14
350	0.50	0.25	0.33	12
351	0.33	0.15	0.21	13
352	0.00	0.00	0.00	15
353	0.60	0.33	0.43	9
354	0.00	0.00	0.00	19
355	0.50	0.50	0.50	2
356	1.00	0.17	0.29	6
357	0.71	0.62	0.67	8
358	0.45	0.26	0.33	19
359	0.50	0.17	0.25	12
360	0.50	0.05	0.09	20
361	0.62	0.26	0.37	19
362	1.00	0.56	0.72	16
363	1.00	0.10	0.18	10
364	1.00	0.38	0.55	8
365	0.00	0.00	0.00	11
366	0.00	0.00	0.00	4
367	0.00	0.00	0.00	1
368	0.44	0.44	0.44	9
369	1.00	0.83	0.91	6
370	1.00	0.14	0.25	7
371	0.50	0.20	0.29	10
372	0.00	0.00	0.00	1
373	0.00	0.00	0.00	2
374	0.50	0.06	0.10	18
375	0.00	0.00	0.00	12
376	0.00	0.00	0.00	16
377	0.00	0.00	0.00	1
378	0.00	0.00	0.00	0
379	0.00	0.00	0.00	14
380	1.00	0.33	0.50	3
381	1.00	0.57	0.73	7
382	0.80	0.40	0.53	10
383	0.00	0.00	0.00	9
384	0.50	0.22	0.31	9
385	0.00	0.00	0.00	7
386	0.83	0.83	0.83	6
387	0.00	0.00	0.00	4
388	0.00	0.00	0.00	8
389	0.00	0.00	0.00	4
390	0.00	0.00	0.00	3
391	0.92	0.46	0.62	26
392	0.00	0.00	0.00	3
393	0.45	0.25	0.32	20
394	0.00	0.00	0.00	3
395	0.00	0.00	0.00	1
396	0.00	0.00	0.00	3
397	1.00	0.67	0.80	3
398	0.73	0.42	0.53	19
399	0.94	0.65	0.77	23
400	0.00	0.00	0.00	13
401	0.88	0.33	0.48	21
402	0.86	0.55	0.67	11
403	1.00	0.17	0.29	6
404	0.54	0.54	0.54	13
405	0.86	0.72	0.78	25
406	0.00	0.00	0.00	2
407	0.25	0.14	0.18	7
408	0.00	0.00	0.00	1
409	0.00	0.00	0.00	6
410	0.00	0.00	0.00	0

411	0.00	0.00	0.00	18
412	0.88	0.41	0.56	17
413	1.00	0.29	0.44	7
414	0.00	0.00	0.00	14
415	0.00	0.00	0.00	3
416	1.00	0.06	0.11	17
417	0.00	0.00	0.00	0
418	0.00	0.00	0.00	8
419	0.00	0.00	0.00	0
420	0.71	0.50	0.59	10
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	2
423	0.00	0.00	0.00	6
424	0.00	0.00	0.00	1
425	0.00	0.00	0.00	4
426	0.00	0.00	0.00	8
427	0.67	1.00	0.80	2
428	0.40	0.11	0.17	18
429	0.61	0.54	0.57	26
430	0.00	0.00	0.00	0
431	0.00	0.00	0.00	14
432	1.00	0.30	0.46	10
433	0.00	0.00	0.00	7
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.00	0.00	0.00	18
437	0.57	0.14	0.23	28
438	0.75	0.50	0.60	12
439	0.00	0.00	0.00	35
440	1.00	0.30	0.46	10
441	0.50	0.25	0.33	12
442	0.75	0.50	0.60	6
443	0.00	0.00	0.00	9
444	0.00	0.00	0.00	2
445	1.00	0.57	0.73	7
446	0.85	0.52	0.65	21
447	0.00	0.00	0.00	6
448	0.00	0.00	0.00	6
449	1.00	0.50	0.67	2
450	0.00	0.00	0.00	30
451	0.00	0.00	0.00	9
452	1.00	0.80	0.89	10
453	0.00	0.00	0.00	3
454	0.67	0.67	0.67	3
455	0.00	0.00	0.00	9
456	0.00	0.00	0.00	1
457	0.00	0.00	0.00	3
458	0.00	0.00	0.00	4
459	1.00	0.27	0.42	15
460	0.00	0.00	0.00	13
461	1.00	0.29	0.44	7
462	1.00	0.23	0.38	13
463	0.00	0.00	0.00	8
464	1.00	0.68	0.81	22
465	0.33	0.17	0.22	6
466	0.75	0.23	0.35	13
467	0.20	0.11	0.14	19
468	0.40	0.06	0.10	35
469	0.00	0.00	0.00	1
470	0.00	0.00	0.00	2
471	0.78	0.41	0.54	17
472	0.00	0.00	0.00	44
473	1.00	0.20	0.33	10
474	1.00	0.64	0.78	11
475	0.60	0.52	0.55	66
476	0.00	0.00	0.00	3
477	0.33	0.10	0.15	10
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	9
480	0.00	0.00	0.00	5
481	0.00	0.00	0.00	13
482	0.00	0.00	0.00	1
483	0.50	0.20	0.29	5
484	1.00	0.22	0.36	9
485	1.00	0.17	0.29	12
486	0.50	0.29	0.36	7
487	1.00	0.08	0.14	13
488	0.00	0.00	0.00	0
489	1.00	0.25	0.40	4
490	1.00	0.33	0.50	3
491	0.00	0.00	0.00	3
492	0.00	0.00	0.00	1
493	0.50	0.07	0.12	14

494	0.33	0.33	0.33	3
495	1.00	0.40	0.57	5
496	1.00	0.17	0.29	12
497	0.00	0.00	0.00	0
498	0.00	0.00	0.00	6
499	0.50	0.27	0.35	11
micro avg	0.72	0.32	0.44	19042
macro avg	0.45	0.22	0.27	19042
weighted avg	0.63	0.32	0.41	19042
samples avg	0.45	0.32	0.35	19042

Time taken to run this cell : 0:04:41.032406

In [52]:

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

Out[52]:

```
['lr_with_more_title_weight.pkl']
```

In [53]:

```
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)
```

```
Accuracy : 0.18771877187718772
Hamming loss  0.0030895089508950896
Micro-average quality numbers
Precision: 0.7030, Recall: 0.3270, F1-measure: 0.4464
Macro-average quality numbers
Precision: 0.4477, Recall: 0.2272, F1-measure: 0.2819
```

	precision	recall	f1-score	support
0	0.84	0.39	0.53	374
1	0.69	0.23	0.34	902
2	0.53	0.08	0.14	229
3	0.65	0.16	0.25	129
4	0.78	0.51	0.62	683
5	0.91	0.47	0.62	517
6	0.79	0.38	0.51	704
7	0.88	0.63	0.73	385
8	0.91	0.58	0.71	902
9	0.71	0.60	0.65	945
10	0.83	0.53	0.65	335
11	0.72	0.39	0.51	66
12	1.00	0.50	0.67	12
13	0.78	0.42	0.54	308
14	0.35	0.15	0.21	143
15	0.92	0.15	0.26	81
16	0.62	0.23	0.33	199
17	0.64	0.26	0.37	252
18	0.67	0.52	0.59	69
19	0.89	0.57	0.69	292
20	0.55	0.21	0.31	472
21	0.78	0.56	0.65	116
22	0.81	0.42	0.55	354
23	0.58	0.30	0.40	107
24	0.75	0.36	0.48	138
25	0.24	0.12	0.16	67

26	0.91	0.22	0.35	96
27	0.29	0.05	0.09	176
28	0.50	0.11	0.19	35
29	1.00	0.15	0.26	20
30	0.29	0.63	0.40	27
31	1.00	0.17	0.29	6
32	0.72	0.55	0.62	96
33	0.67	0.39	0.50	71
34	0.80	0.28	0.41	184
35	0.84	0.58	0.69	124
36	0.00	0.00	0.00	22
37	0.37	0.22	0.28	73
38	0.67	0.43	0.52	134
39	0.64	0.26	0.37	109
40	0.00	0.00	0.00	7
41	0.50	0.12	0.20	80
42	0.00	0.00	0.00	6
43	0.39	0.12	0.19	416
44	0.47	0.30	0.36	27
45	0.47	0.26	0.34	142
46	0.71	0.21	0.32	72
47	0.80	0.53	0.64	121
48	0.71	0.29	0.42	34
49	0.74	0.26	0.39	88
50	0.84	0.70	0.76	60
51	0.83	0.62	0.71	16
52	0.64	0.37	0.47	67
53	0.36	0.17	0.23	59
54	0.17	0.09	0.11	35
55	0.64	0.47	0.55	19
56	0.69	0.55	0.62	74
57	0.00	0.00	0.00	46
58	0.58	0.21	0.31	33
59	0.39	0.15	0.21	81
60	0.91	0.73	0.81	73
61	0.30	0.07	0.12	135
62	0.40	0.19	0.26	21
63	0.89	0.33	0.48	24
64	0.68	0.25	0.36	153
65	0.67	0.40	0.50	40
66	0.67	0.29	0.40	7
67	0.67	0.49	0.57	59
68	0.19	0.10	0.13	29
69	0.42	0.11	0.17	131
70	0.50	0.05	0.09	21
71	0.57	0.40	0.47	20
72	0.00	0.00	0.00	9
73	0.00	0.00	0.00	1
74	0.50	0.07	0.12	15
75	0.67	0.35	0.46	17
76	0.85	0.60	0.70	55
77	0.00	0.00	0.00	2
78	0.36	0.05	0.09	103
79	0.64	0.35	0.45	26
80	0.75	0.21	0.33	14
81	0.43	0.26	0.32	23
82	0.33	0.18	0.23	17
83	0.00	0.00	0.00	35
84	0.73	0.39	0.51	28
85	0.36	0.36	0.36	11
86	0.00	0.00	0.00	11
87	1.00	0.40	0.57	5
88	0.83	0.77	0.80	13
89	0.94	0.90	0.92	87
90	0.10	0.03	0.04	39
91	0.56	0.16	0.24	32
92	0.56	0.18	0.27	28
93	0.55	0.33	0.41	18
94	0.50	0.12	0.20	8
95	0.86	0.47	0.61	38
96	0.64	0.36	0.46	25
97	1.00	0.03	0.07	29
98	0.00	0.00	0.00	57
99	0.00	0.00	0.00	0
100	0.00	0.00	0.00	4
101	0.75	0.39	0.51	23
102	0.00	0.00	0.00	3
103	0.50	0.07	0.12	15
104	0.00	0.00	0.00	19
105	0.91	0.43	0.59	23
106	0.50	0.12	0.19	50
107	0.00	0.00	0.00	10
108	0.75	0.16	0.27	37

109	0.47	0.29	0.36	28
110	0.90	0.18	0.30	51
111	0.33	0.38	0.36	21
112	0.00	0.00	0.00	286
113	0.96	0.87	0.91	30
114	0.00	0.00	0.00	9
115	0.00	0.00	0.00	18
116	0.00	0.00	0.00	2
117	0.33	0.09	0.14	23
118	0.79	0.47	0.59	66
119	0.80	0.33	0.47	24
120	0.00	0.00	0.00	1
121	0.50	0.21	0.29	24
122	0.86	0.46	0.60	26
123	0.25	0.03	0.05	34
124	0.00	0.00	0.00	0
125	1.00	0.67	0.80	3
126	0.00	0.00	0.00	8
127	1.00	0.19	0.32	21
128	1.00	0.62	0.77	8
129	0.71	0.52	0.60	23
130	0.00	0.00	0.00	53
131	0.00	0.00	0.00	1
132	0.00	0.00	0.00	9
133	0.00	0.00	0.00	6
134	1.00	0.20	0.33	5
135	0.86	0.60	0.71	10
136	0.00	0.00	0.00	27
137	0.11	0.14	0.12	21
138	1.00	0.33	0.50	3
139	0.33	0.07	0.11	29
140	0.89	0.67	0.76	24
141	0.75	0.33	0.46	36
142	0.33	0.13	0.19	15
143	0.47	0.30	0.36	27
144	0.64	0.47	0.54	34
145	0.00	0.00	0.00	8
146	0.14	0.01	0.02	86
147	0.40	0.25	0.31	8
148	0.88	0.37	0.52	19
149	0.44	0.20	0.27	35
150	0.00	0.00	0.00	8
151	0.20	0.33	0.25	3
152	0.17	0.07	0.10	28
153	1.00	0.50	0.67	2
154	0.00	0.00	0.00	1
155	0.33	0.09	0.14	11
156	0.53	0.15	0.24	52
157	0.56	0.42	0.48	43
158	1.00	0.14	0.25	7
159	0.63	0.47	0.54	55
160	0.53	0.27	0.36	37
161	1.00	0.25	0.40	16
162	0.00	0.00	0.00	21
163	0.25	0.03	0.05	75
164	0.20	0.05	0.08	20
165	0.00	0.00	0.00	227
166	0.89	0.11	0.19	75
167	0.40	0.14	0.21	14
168	0.56	0.16	0.25	31
169	0.71	0.56	0.63	18
170	0.85	0.48	0.61	23
171	0.00	0.00	0.00	5
172	0.74	0.52	0.61	27
173	0.50	0.50	0.50	2
174	0.30	0.08	0.13	36
175	0.82	0.47	0.60	19
176	1.00	0.80	0.89	10
177	0.17	0.08	0.11	12
178	0.00	0.00	0.00	0
179	0.00	0.00	0.00	0
180	0.33	0.14	0.20	14
181	0.00	0.00	0.00	12
182	0.00	0.00	0.00	1
183	0.64	0.37	0.47	38
184	0.00	0.00	0.00	2
185	0.89	0.29	0.43	28
186	0.33	0.50	0.40	4
187	1.00	0.70	0.82	30
188	0.68	0.56	0.61	27
189	0.39	0.15	0.21	48
190	0.00	0.00	0.00	12
191	0.71	0.20	0.31	50

192	0.54	0.32	0.40	22
193	0.00	0.00	0.00	4
194	0.69	0.44	0.54	25
195	1.00	0.14	0.25	7
196	0.86	0.63	0.73	19
197	0.64	0.27	0.38	52
198	1.00	0.11	0.20	9
199	0.00	0.00	0.00	13
200	0.00	0.00	0.00	28
201	0.50	0.17	0.25	6
202	0.80	0.24	0.36	17
203	0.00	0.00	0.00	21
204	1.00	0.56	0.72	34
205	0.00	0.00	0.00	1
206	0.40	0.06	0.10	35
207	0.00	0.00	0.00	3
208	0.00	0.00	0.00	4
209	0.25	0.11	0.15	28
210	0.00	0.00	0.00	1
211	0.60	0.15	0.24	20
212	1.00	0.50	0.67	6
213	0.00	0.00	0.00	2
214	0.56	0.33	0.42	15
215	0.67	0.13	0.22	30
216	0.70	0.37	0.48	38
217	0.33	0.08	0.13	12
218	0.00	0.00	0.00	1
219	0.00	0.00	0.00	16
220	0.80	0.05	0.10	79
221	1.00	0.13	0.24	15
222	0.00	0.00	0.00	15
223	1.00	0.32	0.49	34
224	0.00	0.00	0.00	5
225	1.00	0.67	0.80	3
226	0.77	0.69	0.73	48
227	0.00	0.00	0.00	0
228	0.50	0.20	0.29	5
229	0.87	0.77	0.82	26
230	0.43	0.12	0.18	26
231	0.00	0.00	0.00	0
232	0.50	0.21	0.30	14
233	0.00	0.00	0.00	2
234	0.87	0.54	0.67	24
235	0.00	0.00	0.00	1
236	0.29	0.13	0.18	15
237	0.71	0.41	0.52	41
238	0.80	0.18	0.30	22
239	1.00	0.20	0.33	10
240	1.00	0.54	0.70	26
241	0.92	0.73	0.81	15
242	0.57	0.80	0.67	15
243	0.83	0.69	0.75	29
244	0.50	0.24	0.33	29
245	0.00	0.00	0.00	6
246	0.00	0.00	0.00	2
247	1.00	0.08	0.14	13
248	0.84	0.53	0.65	30
249	0.67	0.55	0.60	11
250	0.33	0.20	0.25	10
251	0.83	0.42	0.56	24
252	0.25	0.08	0.12	12
253	0.25	0.08	0.12	12
254	0.00	0.00	0.00	1
255	1.00	0.50	0.67	2
256	0.00	0.00	0.00	1
257	0.67	0.38	0.48	16
258	0.67	0.12	0.21	16
259	0.00	0.00	0.00	2
260	0.00	0.00	0.00	17
261	0.00	0.00	0.00	0
262	0.75	0.27	0.40	11
263	0.00	0.00	0.00	1
264	0.50	0.10	0.17	20
265	0.00	0.00	0.00	3
266	0.33	0.04	0.06	28
267	0.53	0.47	0.50	17
268	0.71	0.50	0.59	10
269	0.75	0.39	0.51	23
270	0.38	0.38	0.38	8
271	0.00	0.00	0.00	20
272	0.00	0.00	0.00	0
273	0.00	0.00	0.00	6
274	0.20	0.17	0.18	6

275	0.67	0.31	0.42	39
276	0.88	0.78	0.82	9
277	0.00	0.00	0.00	8
278	1.00	0.17	0.29	6
279	0.80	0.80	0.80	5
280	0.00	0.00	0.00	4
281	0.67	0.67	0.67	3
282	1.00	0.67	0.80	15
283	0.00	0.00	0.00	0
284	0.57	0.35	0.43	37
285	0.50	0.14	0.22	21
286	0.33	0.09	0.14	11
287	0.00	0.00	0.00	18
288	0.88	0.44	0.58	16
289	0.00	0.00	0.00	11
290	0.75	0.62	0.68	24
291	0.50	0.25	0.33	4
292	0.50	0.22	0.31	9
293	0.55	0.55	0.55	11
294	0.00	0.00	0.00	14
295	0.00	0.00	0.00	13
296	0.50	0.25	0.33	8
297	0.43	0.19	0.26	16
298	0.00	0.00	0.00	34
299	0.00	0.00	0.00	16
300	0.60	0.14	0.23	21
301	0.80	0.52	0.63	23
302	0.80	0.36	0.50	11
303	0.00	0.00	0.00	3
304	0.38	0.19	0.25	16
305	0.00	0.00	0.00	6
306	0.00	0.00	0.00	3
307	0.00	0.00	0.00	2
308	0.00	0.00	0.00	14
309	0.64	0.36	0.46	25
310	1.00	0.06	0.11	17
311	0.33	0.13	0.19	30
312	0.00	0.00	0.00	11
313	0.67	0.14	0.24	14
314	0.40	0.14	0.21	14
315	0.50	0.03	0.05	38
316	1.00	0.14	0.25	7
317	0.33	0.19	0.24	26
318	0.50	1.00	0.67	1
319	0.00	0.00	0.00	5
320	0.00	0.00	0.00	0
321	0.00	0.00	0.00	0
322	0.80	0.50	0.62	8
323	0.89	0.67	0.76	12
324	0.64	0.55	0.59	29
325	0.00	0.00	0.00	4
326	0.17	0.10	0.12	10
327	0.00	0.00	0.00	8
328	1.00	0.10	0.18	10
329	0.00	0.00	0.00	0
330	0.50	0.33	0.40	3
331	0.67	0.17	0.27	12
332	0.00	0.00	0.00	0
333	0.75	0.43	0.55	7
334	0.00	0.00	0.00	14
335	0.00	0.00	0.00	0
336	0.00	0.00	0.00	17
337	0.79	0.55	0.65	20
338	1.00	0.06	0.12	16
339	0.67	0.31	0.42	13
340	0.57	0.46	0.51	26
341	1.00	0.83	0.91	6
342	0.00	0.00	0.00	0
343	0.00	0.00	0.00	0
344	0.67	0.13	0.22	15
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	1
347	1.00	0.25	0.40	4
348	0.00	0.00	0.00	5
349	0.75	0.21	0.33	14
350	0.50	0.25	0.33	12
351	0.38	0.23	0.29	13
352	0.00	0.00	0.00	15
353	0.60	0.33	0.43	9
354	0.00	0.00	0.00	19
355	0.50	0.50	0.50	2
356	1.00	0.17	0.29	6
357	0.83	0.62	0.71	8

358	0.55	0.32	0.40	19
359	0.56	0.42	0.48	12
360	0.25	0.05	0.08	20
361	0.56	0.26	0.36	19
362	1.00	0.62	0.77	16
363	1.00	0.20	0.33	10
364	1.00	0.50	0.67	8
365	0.00	0.00	0.00	11
366	0.00	0.00	0.00	4
367	0.00	0.00	0.00	1
368	0.44	0.44	0.44	9
369	1.00	1.00	1.00	6
370	1.00	0.14	0.25	7
371	0.33	0.10	0.15	10
372	0.00	0.00	0.00	1
373	1.00	0.50	0.67	2
374	0.75	0.17	0.27	18
375	1.00	0.08	0.15	12
376	0.00	0.00	0.00	16
377	0.00	0.00	0.00	1
378	0.00	0.00	0.00	0
379	0.00	0.00	0.00	14
380	0.00	0.00	0.00	3
381	0.80	0.57	0.67	7
382	0.71	0.50	0.59	10
383	0.00	0.00	0.00	9
384	0.50	0.22	0.31	9
385	0.00	0.00	0.00	7
386	0.83	0.83	0.83	6
387	0.00	0.00	0.00	4
388	0.00	0.00	0.00	8
389	0.00	0.00	0.00	4
390	0.00	0.00	0.00	3
391	0.94	0.58	0.71	26
392	0.00	0.00	0.00	3
393	0.50	0.30	0.37	20
394	0.00	0.00	0.00	3
395	0.00	0.00	0.00	1
396	0.00	0.00	0.00	3
397	1.00	0.67	0.80	3
398	0.73	0.42	0.53	19
399	0.94	0.65	0.77	23
400	0.50	0.08	0.13	13
401	0.70	0.33	0.45	21
402	0.88	0.64	0.74	11
403	1.00	0.17	0.29	6
404	0.58	0.54	0.56	13
405	0.86	0.72	0.78	25
406	0.00	0.00	0.00	2
407	0.33	0.14	0.20	7
408	0.00	0.00	0.00	1
409	0.00	0.00	0.00	6
410	0.00	0.00	0.00	0
411	0.00	0.00	0.00	18
412	0.78	0.41	0.54	17
413	1.00	0.29	0.44	7
414	0.00	0.00	0.00	14
415	0.33	0.33	0.33	3
416	1.00	0.06	0.11	17
417	0.00	0.00	0.00	0
418	0.00	0.00	0.00	8
419	0.00	0.00	0.00	0
420	0.71	0.50	0.59	10
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	2
423	0.00	0.00	0.00	6
424	1.00	1.00	1.00	1
425	0.00	0.00	0.00	4
426	0.00	0.00	0.00	8
427	0.67	1.00	0.80	2
428	0.60	0.17	0.26	18
429	0.62	0.58	0.60	26
430	0.00	0.00	0.00	0
431	0.00	0.00	0.00	14
432	1.00	0.30	0.46	10
433	0.00	0.00	0.00	7
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.00	0.00	0.00	18
437	0.56	0.18	0.27	28
438	0.78	0.58	0.67	12
439	0.00	0.00	0.00	35
440	1.00	0.30	0.46	10

441	0.50	0.25	0.33	12
442	0.75	0.50	0.60	6
443	0.00	0.00	0.00	9
444	0.00	0.00	0.00	2
445	1.00	0.43	0.60	7
446	0.80	0.57	0.67	21
447	0.00	0.00	0.00	6
448	0.00	0.00	0.00	6
449	1.00	0.50	0.67	2
450	1.00	0.03	0.06	30
451	0.00	0.00	0.00	9
452	1.00	0.90	0.95	10
453	0.00	0.00	0.00	3
454	0.67	0.67	0.67	3
455	0.00	0.00	0.00	9
456	0.00	0.00	0.00	1
457	0.00	0.00	0.00	3
458	0.00	0.00	0.00	4
459	1.00	0.53	0.70	15
460	0.00	0.00	0.00	13
461	1.00	0.29	0.44	7
462	1.00	0.23	0.38	13
463	0.00	0.00	0.00	8
464	1.00	0.68	0.81	22
465	0.25	0.17	0.20	6
466	0.50	0.23	0.32	13
467	0.20	0.16	0.18	19
468	0.33	0.03	0.05	35
469	0.00	0.00	0.00	1
470	0.00	0.00	0.00	2
471	0.70	0.41	0.52	17
472	0.12	0.02	0.04	44
473	1.00	0.20	0.33	10
474	1.00	0.64	0.78	11
475	0.58	0.55	0.56	66
476	0.00	0.00	0.00	3
477	0.40	0.20	0.27	10
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	9
480	0.00	0.00	0.00	5
481	0.33	0.08	0.12	13
482	0.00	0.00	0.00	1
483	0.50	0.20	0.29	5
484	1.00	0.33	0.50	9
485	0.80	0.33	0.47	12
486	0.50	0.29	0.36	7
487	1.00	0.08	0.14	13
488	0.00	0.00	0.00	0
489	1.00	0.25	0.40	4
490	0.50	0.33	0.40	3
491	0.00	0.00	0.00	3
492	0.00	0.00	0.00	1
493	0.29	0.14	0.19	14
494	0.50	0.33	0.40	3
495	0.67	0.40	0.50	5
496	0.00	0.00	0.00	12
497	0.00	0.00	0.00	0
498	0.00	0.00	0.00	6
499	0.67	0.36	0.47	11
micro avg	0.70	0.33	0.45	19042
macro avg	0.45	0.23	0.28	19042
weighted avg	0.63	0.33	0.41	19042
samples avg	0.45	0.32	0.35	19042

Time taken to run this cell : 0:18:32.439164

4.5.4 Featurizing data with BOW vectorizer and up to 4 grams

In [47]:

```
start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009, max_features=20000, \
                             tokenizer = lambda x: x.split(), ngram_range=(1,4))
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 : 0:36:09.898222

In [48]:

```
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: (90000, 20000) Y : (90000, 500)
Dimensions of test data X: (9999, 20000) Y: (9999, 500)

Hyper parameter tuning on alpha

In [49]:

```
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import GridSearchCV
```

In [50]:

```
start=datetime.now()
param_grid = {"estimator__alpha": [10**-5, 10**-3, 10**-1, 10**1, 10**2]}

clf = OneVsRestClassifier(SGDClassifier(loss='log',penalty='l1'))

model = GridSearchCV(clf,param_grid, scoring = 'f1_micro', cv=2,n_jobs=-1)

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

Time taken to run this cell : 7:56:27.496709

In [57]:

```
best_param=print(model.best_params_)

{'estimator__alpha': 0.001}
```

Applying Logistic Regression with OneVsRest Classifier

In [64]:

```
start = datetime.now()
classifier_2 = OneVsRestClassifier(LogisticRegression(C=0.001,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)
```

Accuracy : 0.11421142114211422
Hamming loss 0.003558755875587559
Micro-average quality numbers
Precision: 0.8559, Recall: 0.0789, F1-measure: 0.1445
Macro-average quality numbers
Precision: 0.0546, Recall: 0.0102, F1-measure: 0.0156

	precision	recall	f1-score	support
0	0.78	0.20	0.32	374
1	0.38	0.00	0.01	902
2	0.47	0.07	0.12	229
3	0.53	0.12	0.20	129
4	0.88	0.27	0.41	683
5	0.95	0.38	0.55	517
6	0.86	0.26	0.40	704
7	0.90	0.46	0.61	385

8	0.97	0.26	0.42	902
9	0.90	0.02	0.04	945
10	0.90	0.18	0.30	335
11	0.72	0.39	0.51	66
12	1.00	0.33	0.50	12
13	0.93	0.18	0.30	308
14	1.00	0.01	0.01	143
15	0.00	0.00	0.00	81
16	1.00	0.01	0.01	199
17	0.00	0.00	0.00	252
18	0.71	0.17	0.28	69
19	0.94	0.35	0.51	292
20	0.00	0.00	0.00	472
21	0.78	0.45	0.57	116
22	0.81	0.06	0.12	354
23	0.67	0.02	0.04	107
24	0.00	0.00	0.00	138
25	0.00	0.00	0.00	67
26	0.00	0.00	0.00	96
27	1.00	0.01	0.02	176
28	0.00	0.00	0.00	35
29	0.00	0.00	0.00	20
30	0.00	0.00	0.00	27
31	0.00	0.00	0.00	6
32	0.88	0.07	0.13	96
33	0.25	0.01	0.03	71
34	0.00	0.00	0.00	184
35	0.83	0.08	0.15	124
36	0.00	0.00	0.00	22
37	0.00	0.00	0.00	73
38	0.50	0.04	0.08	134
39	0.00	0.00	0.00	109
40	0.00	0.00	0.00	7
41	0.00	0.00	0.00	80
42	0.00	0.00	0.00	6
43	0.00	0.00	0.00	416
44	0.00	0.00	0.00	27
45	0.36	0.04	0.06	142
46	0.00	0.00	0.00	72
47	0.00	0.00	0.00	121
48	0.00	0.00	0.00	34
49	0.00	0.00	0.00	88
50	0.88	0.12	0.21	60
51	0.00	0.00	0.00	16
52	0.67	0.06	0.11	67
53	0.00	0.00	0.00	59
54	0.00	0.00	0.00	35
55	0.25	0.05	0.09	19
56	0.62	0.07	0.12	74
57	0.00	0.00	0.00	46
58	0.00	0.00	0.00	33
59	0.00	0.00	0.00	81
60	0.00	0.00	0.00	73
61	0.00	0.00	0.00	135
62	0.00	0.00	0.00	21
63	0.00	0.00	0.00	24
64	0.00	0.00	0.00	153
65	1.00	0.05	0.10	40
66	0.00	0.00	0.00	7
67	0.00	0.00	0.00	59
68	0.00	0.00	0.00	29
69	0.00	0.00	0.00	131
70	0.00	0.00	0.00	21
71	0.00	0.00	0.00	20
72	0.00	0.00	0.00	9
73	0.00	0.00	0.00	1
74	0.00	0.00	0.00	15
75	0.00	0.00	0.00	17
76	0.00	0.00	0.00	55
77	0.00	0.00	0.00	2
78	0.00	0.00	0.00	103
79	0.50	0.04	0.07	26
80	0.00	0.00	0.00	14
81	1.00	0.04	0.08	23
82	0.00	0.00	0.00	17
83	0.00	0.00	0.00	35
84	0.00	0.00	0.00	28
85	0.00	0.00	0.00	11
86	0.00	0.00	0.00	11
87	0.00	0.00	0.00	5
88	0.00	0.00	0.00	13
89	0.00	0.00	0.00	87
90	0.00	0.00	0.00	39

91	0.00	0.00	0.00	32
92	0.00	0.00	0.00	28
93	0.50	0.06	0.10	18
94	1.00	0.12	0.22	8
95	0.00	0.00	0.00	38
96	0.00	0.00	0.00	25
97	0.00	0.00	0.00	29
98	0.00	0.00	0.00	57
99	0.00	0.00	0.00	0
100	0.00	0.00	0.00	4
101	0.00	0.00	0.00	23
102	0.00	0.00	0.00	3
103	0.00	0.00	0.00	15
104	0.00	0.00	0.00	19
105	0.00	0.00	0.00	23
106	0.00	0.00	0.00	50
107	0.00	0.00	0.00	10
108	0.00	0.00	0.00	37
109	0.00	0.00	0.00	28
110	0.00	0.00	0.00	51
111	0.00	0.00	0.00	21
112	0.00	0.00	0.00	286
113	0.00	0.00	0.00	30
114	0.00	0.00	0.00	9
115	0.00	0.00	0.00	18
116	0.00	0.00	0.00	2
117	0.00	0.00	0.00	23
118	0.00	0.00	0.00	66
119	0.00	0.00	0.00	24
120	0.00	0.00	0.00	1
121	0.00	0.00	0.00	24
122	0.00	0.00	0.00	26
123	0.00	0.00	0.00	34
124	0.00	0.00	0.00	0
125	0.00	0.00	0.00	3
126	0.00	0.00	0.00	8
127	0.00	0.00	0.00	21
128	0.00	0.00	0.00	8
129	0.00	0.00	0.00	23
130	0.00	0.00	0.00	53
131	0.00	0.00	0.00	1
132	0.00	0.00	0.00	9
133	0.00	0.00	0.00	6
134	0.00	0.00	0.00	5
135	0.00	0.00	0.00	10
136	0.00	0.00	0.00	27
137	0.00	0.00	0.00	21
138	0.00	0.00	0.00	3
139	0.00	0.00	0.00	29
140	0.00	0.00	0.00	24
141	0.00	0.00	0.00	36
142	0.00	0.00	0.00	15
143	0.00	0.00	0.00	27
144	0.00	0.00	0.00	34
145	0.00	0.00	0.00	8
146	0.00	0.00	0.00	86
147	0.00	0.00	0.00	8
148	0.00	0.00	0.00	19
149	0.00	0.00	0.00	35
150	0.00	0.00	0.00	8
151	0.00	0.00	0.00	3
152	0.00	0.00	0.00	28
153	0.00	0.00	0.00	2
154	0.00	0.00	0.00	1
155	0.00	0.00	0.00	11
156	0.00	0.00	0.00	52
157	0.00	0.00	0.00	43
158	0.00	0.00	0.00	7
159	0.00	0.00	0.00	55
160	0.00	0.00	0.00	37
161	0.00	0.00	0.00	16
162	0.00	0.00	0.00	21
163	0.00	0.00	0.00	75
164	0.00	0.00	0.00	20
165	0.00	0.00	0.00	227
166	0.00	0.00	0.00	75
167	0.00	0.00	0.00	14
168	0.00	0.00	0.00	31
169	0.00	0.00	0.00	18
170	0.00	0.00	0.00	23
171	0.00	0.00	0.00	5
172	0.00	0.00	0.00	27
173	0.00	0.00	0.00	2

174	0.00	0.00	0.00	36
175	0.00	0.00	0.00	19
176	0.00	0.00	0.00	10
177	0.00	0.00	0.00	12
178	0.00	0.00	0.00	0
179	0.00	0.00	0.00	0
180	0.00	0.00	0.00	14
181	0.00	0.00	0.00	12
182	0.00	0.00	0.00	1
183	0.00	0.00	0.00	38
184	0.00	0.00	0.00	2
185	0.00	0.00	0.00	28
186	0.00	0.00	0.00	4
187	0.00	0.00	0.00	30
188	0.00	0.00	0.00	27
189	0.00	0.00	0.00	48
190	0.00	0.00	0.00	12
191	0.00	0.00	0.00	50
192	0.00	0.00	0.00	22
193	0.00	0.00	0.00	4
194	0.00	0.00	0.00	25
195	0.00	0.00	0.00	7
196	0.00	0.00	0.00	19
197	0.00	0.00	0.00	52
198	0.00	0.00	0.00	9
199	0.00	0.00	0.00	13
200	0.00	0.00	0.00	28
201	0.00	0.00	0.00	6
202	0.00	0.00	0.00	17
203	0.00	0.00	0.00	21
204	0.00	0.00	0.00	34
205	0.00	0.00	0.00	1
206	0.00	0.00	0.00	35
207	0.00	0.00	0.00	3
208	0.00	0.00	0.00	4
209	0.00	0.00	0.00	28
210	0.00	0.00	0.00	1
211	0.00	0.00	0.00	20
212	0.00	0.00	0.00	6
213	0.00	0.00	0.00	2
214	0.00	0.00	0.00	15
215	0.00	0.00	0.00	30
216	0.00	0.00	0.00	38
217	0.00	0.00	0.00	12
218	0.00	0.00	0.00	1
219	0.00	0.00	0.00	16
220	0.00	0.00	0.00	79
221	0.00	0.00	0.00	15
222	0.00	0.00	0.00	15
223	0.00	0.00	0.00	34
224	0.00	0.00	0.00	5
225	0.00	0.00	0.00	3
226	0.00	0.00	0.00	48
227	0.00	0.00	0.00	0
228	0.00	0.00	0.00	5
229	0.00	0.00	0.00	26
230	0.00	0.00	0.00	26
231	0.00	0.00	0.00	0
232	0.00	0.00	0.00	14
233	0.00	0.00	0.00	2
234	0.00	0.00	0.00	24
235	0.00	0.00	0.00	1
236	0.00	0.00	0.00	15
237	0.00	0.00	0.00	41
238	0.00	0.00	0.00	22
239	0.00	0.00	0.00	10
240	0.00	0.00	0.00	26
241	0.00	0.00	0.00	15
242	0.00	0.00	0.00	15
243	0.00	0.00	0.00	29
244	0.00	0.00	0.00	29
245	0.00	0.00	0.00	6
246	0.00	0.00	0.00	2
247	0.00	0.00	0.00	13
248	0.00	0.00	0.00	30
249	0.00	0.00	0.00	11
250	0.00	0.00	0.00	10
251	0.00	0.00	0.00	24
252	0.00	0.00	0.00	12
253	0.00	0.00	0.00	12
254	0.00	0.00	0.00	1
255	0.00	0.00	0.00	2
256	0.00	0.00	0.00	1

257	0.00	0.00	0.00	16
258	0.00	0.00	0.00	16
259	0.00	0.00	0.00	2
260	0.00	0.00	0.00	17
261	0.00	0.00	0.00	0
262	0.00	0.00	0.00	11
263	0.00	0.00	0.00	1
264	0.00	0.00	0.00	20
265	0.00	0.00	0.00	3
266	0.00	0.00	0.00	28
267	0.00	0.00	0.00	17
268	0.00	0.00	0.00	10
269	0.00	0.00	0.00	23
270	0.00	0.00	0.00	8
271	0.00	0.00	0.00	20
272	0.00	0.00	0.00	0
273	0.00	0.00	0.00	6
274	0.00	0.00	0.00	6
275	0.00	0.00	0.00	39
276	0.00	0.00	0.00	9
277	0.00	0.00	0.00	8
278	0.00	0.00	0.00	6
279	0.00	0.00	0.00	5
280	0.00	0.00	0.00	4
281	0.00	0.00	0.00	3
282	0.00	0.00	0.00	15
283	0.00	0.00	0.00	0
284	0.00	0.00	0.00	37
285	0.00	0.00	0.00	21
286	0.00	0.00	0.00	11
287	0.00	0.00	0.00	18
288	0.00	0.00	0.00	16
289	0.00	0.00	0.00	11
290	0.00	0.00	0.00	24
291	0.00	0.00	0.00	4
292	0.00	0.00	0.00	9
293	0.00	0.00	0.00	11
294	0.00	0.00	0.00	14
295	0.00	0.00	0.00	13
296	0.00	0.00	0.00	8
297	0.00	0.00	0.00	16
298	0.00	0.00	0.00	34
299	0.00	0.00	0.00	16
300	0.00	0.00	0.00	21
301	0.00	0.00	0.00	23
302	0.00	0.00	0.00	11
303	0.00	0.00	0.00	3
304	0.00	0.00	0.00	16
305	0.00	0.00	0.00	6
306	0.00	0.00	0.00	3
307	0.00	0.00	0.00	2
308	0.00	0.00	0.00	14
309	0.00	0.00	0.00	25
310	0.00	0.00	0.00	17
311	0.00	0.00	0.00	30
312	0.00	0.00	0.00	11
313	0.00	0.00	0.00	14
314	0.00	0.00	0.00	14
315	0.00	0.00	0.00	38
316	0.00	0.00	0.00	7
317	0.00	0.00	0.00	26
318	0.00	0.00	0.00	1
319	0.00	0.00	0.00	5
320	0.00	0.00	0.00	0
321	0.00	0.00	0.00	0
322	0.00	0.00	0.00	8
323	0.00	0.00	0.00	12
324	0.00	0.00	0.00	29
325	0.00	0.00	0.00	4
326	0.00	0.00	0.00	10
327	0.00	0.00	0.00	8
328	0.00	0.00	0.00	10
329	0.00	0.00	0.00	0
330	0.00	0.00	0.00	3
331	0.00	0.00	0.00	12
332	0.00	0.00	0.00	0
333	0.00	0.00	0.00	7
334	0.00	0.00	0.00	14
335	0.00	0.00	0.00	0
336	0.00	0.00	0.00	17
337	0.00	0.00	0.00	20
338	0.00	0.00	0.00	16
339	0.00	0.00	0.00	13

340	0.00	0.00	0.00	26
341	0.00	0.00	0.00	6
342	0.00	0.00	0.00	0
343	0.00	0.00	0.00	0
344	0.00	0.00	0.00	15
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	1
347	0.00	0.00	0.00	4
348	0.00	0.00	0.00	5
349	0.00	0.00	0.00	14
350	0.00	0.00	0.00	12
351	0.00	0.00	0.00	13
352	0.00	0.00	0.00	15
353	0.00	0.00	0.00	9
354	0.00	0.00	0.00	19
355	0.00	0.00	0.00	2
356	0.00	0.00	0.00	6
357	0.00	0.00	0.00	8
358	0.00	0.00	0.00	19
359	0.00	0.00	0.00	12
360	0.00	0.00	0.00	20
361	0.00	0.00	0.00	19
362	0.00	0.00	0.00	16
363	0.00	0.00	0.00	10
364	0.00	0.00	0.00	8
365	0.00	0.00	0.00	11
366	0.00	0.00	0.00	4
367	0.00	0.00	0.00	1
368	0.00	0.00	0.00	9
369	0.00	0.00	0.00	6
370	0.00	0.00	0.00	7
371	0.00	0.00	0.00	10
372	0.00	0.00	0.00	1
373	0.00	0.00	0.00	2
374	0.00	0.00	0.00	18
375	0.00	0.00	0.00	12
376	0.00	0.00	0.00	16
377	0.00	0.00	0.00	1
378	0.00	0.00	0.00	0
379	0.00	0.00	0.00	14
380	0.00	0.00	0.00	3
381	0.00	0.00	0.00	7
382	0.00	0.00	0.00	10
383	0.00	0.00	0.00	9
384	0.00	0.00	0.00	9
385	0.00	0.00	0.00	7
386	0.00	0.00	0.00	6
387	0.00	0.00	0.00	4
388	0.00	0.00	0.00	8
389	0.00	0.00	0.00	4
390	0.00	0.00	0.00	3
391	0.00	0.00	0.00	26
392	0.00	0.00	0.00	3
393	0.00	0.00	0.00	20
394	0.00	0.00	0.00	3
395	0.00	0.00	0.00	1
396	0.00	0.00	0.00	3
397	0.00	0.00	0.00	3
398	0.00	0.00	0.00	19
399	0.00	0.00	0.00	23
400	0.00	0.00	0.00	13
401	0.00	0.00	0.00	21
402	0.00	0.00	0.00	11
403	0.00	0.00	0.00	6
404	0.00	0.00	0.00	13
405	0.00	0.00	0.00	25
406	0.00	0.00	0.00	2
407	0.00	0.00	0.00	7
408	0.00	0.00	0.00	1
409	0.00	0.00	0.00	6
410	0.00	0.00	0.00	0
411	0.00	0.00	0.00	18
412	0.00	0.00	0.00	17
413	0.00	0.00	0.00	7
414	0.00	0.00	0.00	14
415	0.00	0.00	0.00	3
416	0.00	0.00	0.00	17
417	0.00	0.00	0.00	0
418	0.00	0.00	0.00	8
419	0.00	0.00	0.00	0
420	0.00	0.00	0.00	10
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	2

423	0.00	0.00	0.00	6
424	0.00	0.00	0.00	1
425	0.00	0.00	0.00	4
426	0.00	0.00	0.00	8
427	0.00	0.00	0.00	2
428	0.00	0.00	0.00	18
429	0.00	0.00	0.00	26
430	0.00	0.00	0.00	0
431	0.00	0.00	0.00	14
432	0.00	0.00	0.00	10
433	0.00	0.00	0.00	7
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.00	0.00	0.00	18
437	0.00	0.00	0.00	28
438	0.00	0.00	0.00	12
439	0.00	0.00	0.00	35
440	0.00	0.00	0.00	10
441	0.00	0.00	0.00	12
442	0.00	0.00	0.00	6
443	0.00	0.00	0.00	9
444	0.00	0.00	0.00	2
445	0.00	0.00	0.00	7
446	0.00	0.00	0.00	21
447	0.00	0.00	0.00	6
448	0.00	0.00	0.00	6
449	0.00	0.00	0.00	2
450	0.00	0.00	0.00	30
451	0.00	0.00	0.00	9
452	0.00	0.00	0.00	10
453	0.00	0.00	0.00	3
454	0.00	0.00	0.00	3
455	0.00	0.00	0.00	9
456	0.00	0.00	0.00	1
457	0.00	0.00	0.00	3
458	0.00	0.00	0.00	4
459	0.00	0.00	0.00	15
460	0.00	0.00	0.00	13
461	0.00	0.00	0.00	7
462	0.00	0.00	0.00	13
463	0.00	0.00	0.00	8
464	0.00	0.00	0.00	22
465	0.00	0.00	0.00	6
466	0.00	0.00	0.00	13
467	0.00	0.00	0.00	19
468	0.00	0.00	0.00	35
469	0.00	0.00	0.00	1
470	0.00	0.00	0.00	2
471	0.00	0.00	0.00	17
472	0.00	0.00	0.00	44
473	0.00	0.00	0.00	10
474	0.00	0.00	0.00	11
475	0.00	0.00	0.00	66
476	0.00	0.00	0.00	3
477	0.00	0.00	0.00	10
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	9
480	0.00	0.00	0.00	5
481	0.00	0.00	0.00	13
482	0.00	0.00	0.00	1
483	0.00	0.00	0.00	5
484	0.00	0.00	0.00	9
485	0.00	0.00	0.00	12
486	0.00	0.00	0.00	7
487	0.00	0.00	0.00	13
488	0.00	0.00	0.00	0
489	0.00	0.00	0.00	4
490	0.00	0.00	0.00	3
491	0.00	0.00	0.00	3
492	0.00	0.00	0.00	1
493	0.00	0.00	0.00	14
494	0.00	0.00	0.00	3
495	0.00	0.00	0.00	5
496	0.00	0.00	0.00	12
497	0.00	0.00	0.00	0
498	0.00	0.00	0.00	6
499	0.00	0.00	0.00	11
micro avg	0.86	0.08	0.14	19042
macro avg	0.05	0.01	0.02	19042
weighted avg	0.37	0.08	0.12	19042
samples avg	0.14	0.08	0.10	19042

Applying Logistic Regression with OneVsRest Classifier and loss='log'

In [59]:

```

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.001, 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 : 0.1415141514151415
Hamming loss  0.003614161416141614
Micro-average quality numbers
Precision: 0.5458, Recall: 0.3044, F1-measure: 0.3909
Macro-average quality numbers
Precision: 0.3243, Recall: 0.1876, F1-measure: 0.2169

```

	precision	recall	f1-score	support
0	0.38	0.47	0.42	374
1	0.67	0.16	0.25	902
2	0.26	0.12	0.16	229
3	0.34	0.22	0.27	129
4	0.69	0.48	0.57	683
5	0.69	0.52	0.59	517
6	0.79	0.29	0.42	704
7	0.63	0.70	0.66	385
8	0.69	0.51	0.58	902
9	0.71	0.62	0.66	945
10	0.72	0.49	0.59	335
11	0.59	0.44	0.50	66
12	0.47	0.67	0.55	12
13	0.72	0.44	0.54	308
14	0.27	0.17	0.21	143
15	0.48	0.30	0.37	81
16	0.49	0.31	0.38	199
17	0.56	0.23	0.33	252
18	0.41	0.36	0.38	69
19	0.76	0.63	0.69	292
20	0.47	0.27	0.35	472
21	0.69	0.68	0.69	116
22	0.80	0.41	0.54	354
23	0.50	0.24	0.33	107
24	0.86	0.09	0.16	138
25	0.21	0.15	0.17	67
26	0.81	0.27	0.41	96
27	0.21	0.18	0.19	176
28	0.64	0.20	0.30	35
29	1.00	0.25	0.40	20
30	0.06	0.30	0.10	27
31	1.00	0.33	0.50	6
32	0.56	0.65	0.60	96
33	0.57	0.39	0.47	71
34	0.68	0.21	0.32	184
35	0.78	0.56	0.65	124
36	0.00	0.00	0.00	22
37	0.17	0.04	0.07	73
38	0.53	0.51	0.52	134
39	0.17	0.06	0.09	109
40	0.00	0.00	0.00	7

41	0.69	0.14	0.23	80
42	0.00	0.00	0.00	6
43	0.00	0.00	0.00	416
44	0.42	0.30	0.35	27
45	0.38	0.32	0.35	142
46	0.67	0.19	0.30	72
47	0.88	0.42	0.57	121
48	0.41	0.26	0.32	34
49	0.64	0.16	0.25	88
50	0.78	0.82	0.80	60
51	0.59	0.81	0.68	16
52	0.60	0.40	0.48	67
53	0.22	0.10	0.14	59
54	0.09	0.06	0.07	35
55	0.64	0.37	0.47	19
56	0.58	0.62	0.60	74
57	1.00	0.04	0.08	46
58	0.50	0.27	0.35	33
59	0.35	0.20	0.25	81
60	0.89	0.74	0.81	73
61	0.28	0.12	0.17	135
62	0.13	0.19	0.16	21
63	0.67	0.25	0.36	24
64	0.74	0.32	0.45	153
65	0.45	0.53	0.48	40
66	0.20	0.14	0.17	7
67	0.65	0.59	0.62	59
68	0.11	0.07	0.09	29
69	0.29	0.04	0.07	131
70	0.03	0.05	0.04	21
71	0.16	0.20	0.18	20
72	0.17	0.11	0.13	9
73	0.00	0.00	0.00	1
74	0.00	0.00	0.00	15
75	0.27	0.18	0.21	17
76	0.77	0.73	0.75	55
77	0.00	0.00	0.00	2
78	0.08	0.01	0.02	103
79	0.58	0.27	0.37	26
80	1.00	0.29	0.44	14
81	0.26	0.30	0.28	23
82	0.19	0.18	0.18	17
83	0.10	0.06	0.07	35
84	0.33	0.54	0.41	28
85	0.36	0.36	0.36	11
86	0.09	0.09	0.09	11
87	0.25	0.20	0.22	5
88	0.80	0.62	0.70	13
89	0.86	0.80	0.83	87
90	0.30	0.18	0.23	39
91	0.25	0.06	0.10	32
92	0.80	0.14	0.24	28
93	0.42	0.28	0.33	18
94	0.50	0.38	0.43	8
95	0.81	0.55	0.66	38
96	0.47	0.36	0.41	25
97	0.50	0.03	0.06	29
98	0.00	0.00	0.00	57
99	0.00	0.00	0.00	0
100	0.00	0.00	0.00	4
101	0.82	0.39	0.53	23
102	0.00	0.00	0.00	3
103	0.33	0.07	0.11	15
104	0.00	0.00	0.00	19
105	0.48	0.48	0.48	23
106	0.53	0.34	0.41	50
107	0.00	0.00	0.00	10
108	0.80	0.11	0.19	37
109	0.50	0.29	0.36	28
110	1.00	0.04	0.08	51
111	0.37	0.33	0.35	21
112	0.00	0.00	0.00	286
113	0.91	0.70	0.79	30
114	0.00	0.00	0.00	9
115	0.14	0.11	0.12	18
116	0.00	0.00	0.00	2
117	0.12	0.09	0.10	23
118	0.79	0.50	0.61	66
119	0.67	0.17	0.27	24
120	0.00	0.00	0.00	1
121	0.14	0.08	0.11	24
122	0.48	0.58	0.53	26
123	0.00	0.00	0.00	34

124	0.00	0.00	0.00	0
125	1.00	0.67	0.80	3
126	0.33	0.12	0.18	8
127	0.12	0.05	0.07	21
128	0.71	0.62	0.67	8
129	0.68	0.57	0.62	23
130	0.08	0.02	0.03	53
131	0.00	0.00	0.00	1
132	0.25	0.11	0.15	9
133	0.00	0.00	0.00	6
134	0.50	0.20	0.29	5
135	0.88	0.70	0.78	10
136	0.00	0.00	0.00	27
137	0.30	0.38	0.33	21
138	0.00	0.00	0.00	3
139	0.19	0.21	0.20	29
140	0.93	0.58	0.72	24
141	0.50	0.11	0.18	36
142	0.23	0.20	0.21	15
143	0.41	0.26	0.32	27
144	0.54	0.38	0.45	34
145	0.00	0.00	0.00	8
146	0.00	0.00	0.00	86
147	0.24	0.50	0.32	8
148	0.62	0.26	0.37	19
149	0.24	0.14	0.18	35
150	0.00	0.00	0.00	8
151	0.33	1.00	0.50	3
152	0.25	0.21	0.23	28
153	0.07	0.50	0.12	2
154	0.00	0.00	0.00	1
155	0.00	0.00	0.00	11
156	0.39	0.31	0.34	52
157	0.45	0.35	0.39	43
158	0.50	0.14	0.22	7
159	0.70	0.38	0.49	55
160	0.44	0.19	0.26	37
161	0.23	0.31	0.26	16
162	0.27	0.14	0.19	21
163	0.18	0.04	0.07	75
164	0.27	0.20	0.23	20
165	0.00	0.00	0.00	227
166	0.00	0.00	0.00	75
167	0.00	0.00	0.00	14
168	0.43	0.10	0.16	31
169	0.25	0.11	0.15	18
170	0.80	0.52	0.63	23
171	0.00	0.00	0.00	5
172	0.78	0.67	0.72	27
173	0.17	0.50	0.25	2
174	0.43	0.08	0.14	36
175	0.63	0.63	0.63	19
176	0.90	0.90	0.90	10
177	0.09	0.17	0.12	12
178	0.00	0.00	0.00	0
179	0.00	0.00	0.00	0
180	0.33	0.21	0.26	14
181	0.00	0.00	0.00	12
182	1.00	1.00	1.00	1
183	0.67	0.21	0.32	38
184	0.00	0.00	0.00	2
185	0.50	0.14	0.22	28
186	0.33	0.50	0.40	4
187	0.92	0.73	0.81	30
188	0.71	0.56	0.63	27
189	0.16	0.15	0.15	48
190	0.00	0.00	0.00	12
191	0.75	0.18	0.29	50
192	0.29	0.18	0.22	22
193	0.00	0.00	0.00	4
194	0.50	0.56	0.53	25
195	1.00	0.14	0.25	7
196	1.00	0.58	0.73	19
197	0.61	0.21	0.31	52
198	0.00	0.00	0.00	9
199	0.25	0.08	0.12	13
200	0.00	0.00	0.00	28
201	0.09	0.17	0.12	6
202	0.00	0.00	0.00	17
203	0.00	0.00	0.00	21
204	1.00	0.21	0.34	34
205	0.00	0.00	0.00	1
206	0.18	0.09	0.12	35

207	0.00	0.00	0.00	3
208	0.00	0.00	0.00	4
209	0.09	0.07	0.08	28
210	0.00	0.00	0.00	1
211	0.12	0.05	0.07	20
212	0.20	0.17	0.18	6
213	0.00	0.00	0.00	2
214	0.56	0.33	0.42	15
215	0.00	0.00	0.00	30
216	1.00	0.13	0.23	38
217	0.29	0.17	0.21	12
218	0.00	0.00	0.00	1
219	0.20	0.06	0.10	16
220	0.75	0.04	0.07	79
221	0.33	0.07	0.11	15
222	0.00	0.00	0.00	15
223	1.00	0.24	0.38	34
224	0.00	0.00	0.00	5
225	0.50	0.67	0.57	3
226	0.81	0.71	0.76	48
227	0.00	0.00	0.00	0
228	0.00	0.00	0.00	5
229	0.85	0.88	0.87	26
230	0.00	0.00	0.00	26
231	0.00	0.00	0.00	0
232	0.44	0.29	0.35	14
233	0.00	0.00	0.00	2
234	0.52	0.71	0.60	24
235	0.00	0.00	0.00	1
236	0.14	0.13	0.14	15
237	0.65	0.41	0.51	41
238	0.00	0.00	0.00	22
239	0.67	0.20	0.31	10
240	1.00	0.62	0.76	26
241	0.81	0.87	0.84	15
242	0.44	0.47	0.45	15
243	0.75	0.52	0.61	29
244	0.43	0.21	0.28	29
245	0.00	0.00	0.00	6
246	0.33	0.50	0.40	2
247	0.33	0.08	0.12	13
248	0.87	0.43	0.58	30
249	0.57	0.36	0.44	11
250	0.00	0.00	0.00	10
251	0.71	0.21	0.32	24
252	0.44	0.33	0.38	12
253	0.06	0.08	0.07	12
254	0.00	0.00	0.00	1
255	0.20	0.50	0.29	2
256	0.00	0.00	0.00	1
257	0.50	0.12	0.20	16
258	0.00	0.00	0.00	16
259	0.00	0.00	0.00	2
260	0.00	0.00	0.00	17
261	0.00	0.00	0.00	0
262	0.67	0.18	0.29	11
263	0.00	0.00	0.00	1
264	0.00	0.00	0.00	20
265	0.00	0.00	0.00	3
266	0.20	0.04	0.06	28
267	0.33	0.24	0.28	17
268	0.75	0.90	0.82	10
269	0.88	0.30	0.45	23
270	0.40	0.25	0.31	8
271	0.00	0.00	0.00	20
272	0.00	0.00	0.00	0
273	0.00	0.00	0.00	6
274	0.00	0.00	0.00	6
275	0.46	0.33	0.39	39
276	0.75	0.67	0.71	9
277	0.00	0.00	0.00	8
278	0.00	0.00	0.00	6
279	0.00	0.00	0.00	5
280	0.00	0.00	0.00	4
281	0.50	0.67	0.57	3
282	1.00	0.60	0.75	15
283	0.00	0.00	0.00	0
284	0.41	0.24	0.31	37
285	0.00	0.00	0.00	21
286	0.50	0.09	0.15	11
287	0.00	0.00	0.00	18
288	0.50	0.19	0.27	16
289	0.00	0.00	0.00	11

290	0.76	0.54	0.63	24
291	0.00	0.00	0.00	4
292	0.00	0.00	0.00	9
293	0.60	0.27	0.37	11
294	0.17	0.07	0.10	14
295	0.00	0.00	0.00	13
296	0.00	0.00	0.00	8
297	0.00	0.00	0.00	16
298	0.00	0.00	0.00	34
299	0.00	0.00	0.00	16
300	0.00	0.00	0.00	21
301	0.56	0.22	0.31	23
302	0.67	0.36	0.47	11
303	0.00	0.00	0.00	3
304	0.00	0.00	0.00	16
305	0.50	0.17	0.25	6
306	0.14	0.33	0.20	3
307	0.00	0.00	0.00	2
308	0.00	0.00	0.00	14
309	0.65	0.44	0.52	25
310	0.00	0.00	0.00	17
311	0.00	0.00	0.00	30
312	0.00	0.00	0.00	11
313	0.00	0.00	0.00	14
314	0.20	0.07	0.11	14
315	0.00	0.00	0.00	38
316	0.40	0.29	0.33	7
317	0.27	0.12	0.16	26
318	0.00	0.00	0.00	1
319	0.00	0.00	0.00	5
320	0.00	0.00	0.00	0
321	0.00	0.00	0.00	0
322	0.75	0.38	0.50	8
323	1.00	0.42	0.59	12
324	0.62	0.34	0.44	29
325	0.33	0.25	0.29	4
326	0.00	0.00	0.00	10
327	0.00	0.00	0.00	8
328	1.00	0.10	0.18	10
329	0.00	0.00	0.00	0
330	1.00	0.33	0.50	3
331	0.14	0.08	0.11	12
332	0.00	0.00	0.00	0
333	0.36	0.57	0.44	7
334	0.00	0.00	0.00	14
335	0.00	0.00	0.00	0
336	0.00	0.00	0.00	17
337	0.90	0.45	0.60	20
338	0.00	0.00	0.00	16
339	0.00	0.00	0.00	13
340	0.33	0.15	0.21	26
341	1.00	0.67	0.80	6
342	0.00	0.00	0.00	0
343	0.00	0.00	0.00	0
344	0.00	0.00	0.00	15
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	1
347	0.20	0.25	0.22	4
348	0.00	0.00	0.00	5
349	1.00	0.07	0.13	14
350	0.50	0.08	0.14	12
351	0.33	0.08	0.12	13
352	0.00	0.00	0.00	15
353	1.00	0.44	0.62	9
354	0.00	0.00	0.00	19
355	1.00	0.50	0.67	2
356	1.00	0.17	0.29	6
357	0.42	0.62	0.50	8
358	0.56	0.26	0.36	19
359	0.75	0.25	0.38	12
360	0.00	0.00	0.00	20
361	0.56	0.47	0.51	19
362	0.57	0.50	0.53	16
363	0.00	0.00	0.00	10
364	0.00	0.00	0.00	8
365	0.00	0.00	0.00	11
366	0.00	0.00	0.00	4
367	0.00	0.00	0.00	1
368	0.29	0.22	0.25	9
369	1.00	0.17	0.29	6
370	0.00	0.00	0.00	7
371	0.00	0.00	0.00	10
372	0.00	0.00	0.00	1

373	0.00	0.00	0.00	2
374	0.00	0.00	0.00	18
375	0.00	0.00	0.00	12
376	0.00	0.00	0.00	16
377	0.00	0.00	0.00	1
378	0.00	0.00	0.00	0
379	0.00	0.00	0.00	14
380	0.33	0.33	0.33	3
381	1.00	0.29	0.44	7
382	0.67	0.40	0.50	10
383	0.00	0.00	0.00	9
384	0.83	0.56	0.67	9
385	0.00	0.00	0.00	7
386	0.62	0.83	0.71	6
387	0.00	0.00	0.00	4
388	0.00	0.00	0.00	8
389	1.00	0.25	0.40	4
390	0.00	0.00	0.00	3
391	0.94	0.58	0.71	26
392	0.00	0.00	0.00	3
393	1.00	0.15	0.26	20
394	0.00	0.00	0.00	3
395	0.00	0.00	0.00	1
396	0.00	0.00	0.00	3
397	1.00	0.33	0.50	3
398	1.00	0.11	0.19	19
399	0.95	0.78	0.86	23
400	0.00	0.00	0.00	13
401	0.17	0.10	0.12	21
402	1.00	0.27	0.43	11
403	0.50	0.17	0.25	6
404	0.57	0.31	0.40	13
405	0.91	0.40	0.56	25
406	1.00	0.50	0.67	2
407	0.00	0.00	0.00	7
408	0.00	0.00	0.00	1
409	0.00	0.00	0.00	6
410	0.00	0.00	0.00	0
411	0.00	0.00	0.00	18
412	0.83	0.29	0.43	17
413	0.00	0.00	0.00	7
414	0.00	0.00	0.00	14
415	0.00	0.00	0.00	3
416	0.40	0.12	0.18	17
417	0.00	0.00	0.00	0
418	1.00	0.12	0.22	8
419	0.00	0.00	0.00	0
420	0.67	0.20	0.31	10
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	2
423	0.00	0.00	0.00	6
424	0.00	0.00	0.00	1
425	1.00	0.25	0.40	4
426	0.00	0.00	0.00	8
427	1.00	1.00	1.00	2
428	0.62	0.28	0.38	18
429	0.39	0.46	0.42	26
430	0.00	0.00	0.00	0
431	0.00	0.00	0.00	14
432	0.00	0.00	0.00	10
433	0.00	0.00	0.00	7
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.00	0.00	0.00	18
437	0.38	0.11	0.17	28
438	0.50	0.08	0.14	12
439	0.00	0.00	0.00	35
440	0.00	0.00	0.00	10
441	1.00	0.25	0.40	12
442	0.67	0.67	0.67	6
443	0.00	0.00	0.00	9
444	0.00	0.00	0.00	2
445	1.00	0.29	0.44	7
446	0.83	0.24	0.37	21
447	0.00	0.00	0.00	6
448	0.00	0.00	0.00	6
449	0.00	0.00	0.00	2
450	0.00	0.00	0.00	30
451	0.00	0.00	0.00	9
452	1.00	0.40	0.57	10
453	0.00	0.00	0.00	3
454	0.00	0.00	0.00	3
455	0.00	0.00	0.00	9

456	0.00	0.00	0.00	1
457	0.00	0.00	0.00	3
458	0.00	0.00	0.00	4
459	1.00	0.27	0.42	15
460	0.00	0.00	0.00	13
461	0.00	0.00	0.00	7
462	0.00	0.00	0.00	13
463	0.00	0.00	0.00	8
464	0.00	0.00	0.00	22
465	0.00	0.00	0.00	6
466	0.50	0.08	0.13	13
467	0.12	0.11	0.11	19
468	0.25	0.06	0.09	35
469	0.00	0.00	0.00	1
470	0.00	0.00	0.00	2
471	0.86	0.35	0.50	17
472	0.00	0.00	0.00	44
473	1.00	0.10	0.18	10
474	1.00	0.64	0.78	11
475	0.00	0.00	0.00	66
476	0.00	0.00	0.00	3
477	0.40	0.20	0.27	10
478	0.00	0.00	0.00	1
479	0.00	0.00	0.00	9
480	0.00	0.00	0.00	5
481	0.00	0.00	0.00	13
482	0.00	0.00	0.00	1
483	1.00	0.20	0.33	5
484	0.33	0.22	0.27	9
485	0.40	0.17	0.24	12
486	1.00	0.14	0.25	7
487	0.00	0.00	0.00	13
488	0.00	0.00	0.00	0
489	1.00	0.50	0.67	4
490	0.50	0.33	0.40	3
491	0.00	0.00	0.00	3
492	0.00	0.00	0.00	1
493	0.14	0.07	0.10	14
494	0.00	0.00	0.00	3
495	0.50	0.20	0.29	5
496	0.00	0.00	0.00	12
497	0.00	0.00	0.00	0
498	0.00	0.00	0.00	6
499	0.00	0.00	0.00	11
micro avg	0.55	0.30	0.39	19042
macro avg	0.32	0.19	0.22	19042
weighted avg	0.50	0.30	0.36	19042
samples avg	0.39	0.30	0.31	19042

Time taken to run this cell : 0:03:01.978203

Applying Logistic Regression with OneVsRest Classifier loss='Hinge'

In [61]:

```
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 : 0.0968096809680968
Hamming loss  0.006205020502050205
Micro-average quality numbers
Precision: 0.2915, Recall: 0.4398, F1-measure: 0.3506
Macro-average quality numbers
Precision: 0.1702, Recall: 0.2992, F1-measure: 0.2044
```

	precision	recall	f1-score	support
0	0.40	0.52	0.45	374
1	0.43	0.42	0.42	902
2	0.14	0.21	0.17	229
3	0.20	0.19	0.19	129
4	0.51	0.64	0.57	683
5	0.59	0.56	0.58	517
6	0.50	0.48	0.49	704
7	0.51	0.70	0.59	385
8	0.66	0.72	0.69	902
9	0.69	0.69	0.69	945
10	0.50	0.56	0.53	335
11	0.24	0.53	0.33	66
12	0.17	0.67	0.27	12
13	0.50	0.54	0.52	308
14	0.13	0.22	0.16	143
15	0.21	0.35	0.26	81
16	0.27	0.36	0.31	199
17	0.39	0.45	0.42	252
18	0.27	0.57	0.36	69
19	0.66	0.65	0.66	292
20	0.36	0.62	0.46	472
21	0.48	0.60	0.54	116
22	0.56	0.57	0.57	354
23	0.22	0.36	0.27	107
24	0.34	0.41	0.37	138
25	0.18	0.27	0.22	67
26	0.26	0.32	0.29	96
27	0.12	0.16	0.14	176
28	0.14	0.29	0.19	35
29	0.12	0.30	0.17	20
30	0.11	0.67	0.18	27
31	0.11	0.50	0.18	6
32	0.44	0.56	0.49	96
33	0.33	0.44	0.38	71
34	0.34	0.36	0.35	184
35	0.56	0.54	0.55	124
36	0.07	0.14	0.09	22
37	0.13	0.21	0.16	73
38	0.41	0.49	0.45	134
39	0.28	0.41	0.33	109
40	0.04	0.14	0.07	7
41	0.17	0.31	0.22	80
42	0.00	0.00	0.00	6
43	0.33	0.58	0.42	416
44	0.17	0.41	0.24	27
45	0.31	0.33	0.32	142

46	0.25	0.25	0.25	72
47	0.46	0.64	0.54	121
48	0.16	0.41	0.23	34
49	0.29	0.36	0.32	88
50	0.45	0.70	0.55	60
51	0.42	0.81	0.55	16
52	0.37	0.51	0.43	67
53	0.19	0.24	0.21	59
54	0.04	0.17	0.07	35
55	0.21	0.42	0.28	19
56	0.42	0.57	0.48	74
57	0.11	0.22	0.15	46
58	0.12	0.30	0.17	33
59	0.24	0.44	0.31	81
60	0.57	0.75	0.65	73
61	0.24	0.21	0.23	135
62	0.07	0.19	0.10	21
63	0.16	0.38	0.22	24
64	0.27	0.44	0.34	153
65	0.17	0.40	0.24	40
66	0.21	0.57	0.31	7
67	0.34	0.54	0.42	59
68	0.06	0.17	0.09	29
69	0.20	0.17	0.18	131
70	0.06	0.10	0.07	21
71	0.14	0.25	0.18	20
72	0.05	0.11	0.07	9
73	0.00	0.00	0.00	1
74	0.08	0.27	0.12	15
75	0.07	0.18	0.10	17
76	0.44	0.71	0.55	55
77	0.00	0.00	0.00	2
78	0.18	0.27	0.22	103
79	0.19	0.38	0.25	26
80	0.24	0.50	0.33	14
81	0.20	0.48	0.28	23
82	0.08	0.18	0.11	17
83	0.10	0.11	0.11	35
84	0.35	0.57	0.43	28
85	0.26	0.64	0.37	11
86	0.00	0.00	0.00	11
87	0.36	0.80	0.50	5
88	0.29	0.77	0.43	13
89	0.66	0.89	0.76	87
90	0.14	0.23	0.17	39
91	0.16	0.28	0.20	32
92	0.06	0.29	0.10	28
93	0.14	0.33	0.20	18
94	0.13	0.25	0.17	8
95	0.35	0.66	0.46	38
96	0.20	0.44	0.28	25
97	0.03	0.03	0.03	29
98	0.05	0.05	0.05	57
99	0.00	0.00	0.00	0
100	0.00	0.00	0.00	4
101	0.27	0.61	0.38	23
102	0.00	0.00	0.00	3
103	0.00	0.00	0.00	15
104	0.00	0.00	0.00	19
105	0.27	0.43	0.33	23
106	0.17	0.18	0.18	50
107	0.08	0.10	0.09	10
108	0.20	0.30	0.24	37
109	0.25	0.46	0.32	28
110	0.19	0.29	0.23	51
111	0.19	0.43	0.26	21
112	0.32	0.04	0.07	286
113	0.61	0.83	0.70	30
114	0.00	0.00	0.00	9
115	0.03	0.06	0.04	18
116	0.00	0.00	0.00	2
117	0.04	0.13	0.06	23
118	0.52	0.53	0.53	66
119	0.35	0.38	0.36	24
120	0.00	0.00	0.00	1
121	0.13	0.29	0.18	24
122	0.32	0.50	0.39	26
123	0.04	0.09	0.06	34
124	0.00	0.00	0.00	0
125	0.15	0.67	0.25	3
126	0.03	0.12	0.05	8
127	0.20	0.24	0.22	21
128	0.23	0.62	0.33	8

129	0.46	0.57	0.51	23
130	0.09	0.15	0.11	53
131	0.00	0.00	0.00	1
132	0.03	0.11	0.05	9
133	0.12	0.33	0.17	6
134	1.00	0.20	0.33	5
135	0.39	0.70	0.50	10
136	0.00	0.00	0.00	27
137	0.12	0.33	0.18	21
138	0.00	0.00	0.00	3
139	0.09	0.21	0.13	29
140	0.38	0.62	0.47	24
141	0.35	0.61	0.45	36
142	0.09	0.20	0.13	15
143	0.26	0.37	0.30	27
144	0.32	0.62	0.42	34
145	0.00	0.00	0.00	8
146	0.11	0.10	0.11	86
147	0.14	0.38	0.21	8
148	0.38	0.53	0.44	19
149	0.15	0.31	0.21	35
150	0.02	0.12	0.04	8
151	0.06	0.67	0.11	3
152	0.10	0.14	0.12	28
153	0.08	0.50	0.13	2
154	0.00	0.00	0.00	1
155	0.08	0.18	0.11	11
156	0.19	0.29	0.23	52
157	0.31	0.60	0.41	43
158	0.03	0.14	0.06	7
159	0.41	0.62	0.50	55
160	0.27	0.41	0.33	37
161	0.09	0.25	0.14	16
162	0.02	0.05	0.03	21
163	0.10	0.15	0.12	75
164	0.12	0.15	0.13	20
165	0.40	0.02	0.03	227
166	0.24	0.13	0.17	75
167	0.12	0.36	0.18	14
168	0.19	0.19	0.19	31
169	0.18	0.33	0.23	18
170	0.41	0.52	0.46	23
171	0.00	0.00	0.00	5
172	0.45	0.70	0.55	27
173	0.07	0.50	0.12	2
174	0.11	0.25	0.16	36
175	0.21	0.42	0.28	19
176	0.42	0.80	0.55	10
177	0.10	0.17	0.12	12
178	0.00	0.00	0.00	0
179	0.00	0.00	0.00	0
180	0.21	0.50	0.29	14
181	0.00	0.00	0.00	12
182	0.00	0.00	0.00	1
183	0.30	0.50	0.38	38
184	0.00	0.00	0.00	2
185	0.28	0.32	0.30	28
186	0.12	0.75	0.20	4
187	0.59	0.73	0.66	30
188	0.43	0.56	0.48	27
189	0.16	0.19	0.17	48
190	0.06	0.17	0.08	12
191	0.37	0.32	0.34	50
192	0.15	0.27	0.19	22
193	0.00	0.00	0.00	4
194	0.37	0.56	0.44	25
195	0.21	0.57	0.31	7
196	0.45	0.68	0.54	19
197	0.30	0.42	0.35	52
198	0.11	0.33	0.16	9
199	0.06	0.15	0.08	13
200	0.10	0.14	0.12	28
201	0.05	0.17	0.08	6
202	0.10	0.12	0.11	17
203	0.03	0.05	0.04	21
204	0.66	0.74	0.69	34
205	0.00	0.00	0.00	1
206	0.07	0.11	0.09	35
207	0.00	0.00	0.00	3
208	0.00	0.00	0.00	4
209	0.16	0.21	0.18	28
210	0.00	0.00	0.00	1
211	0.08	0.20	0.11	20

212	0.11	0.50	0.18	6
213	0.00	0.00	0.00	2
214	0.21	0.40	0.27	15
215	0.08	0.17	0.11	30
216	0.26	0.37	0.31	38
217	0.05	0.17	0.07	12
218	0.00	0.00	0.00	1
219	0.07	0.12	0.09	16
220	0.14	0.24	0.18	79
221	0.20	0.33	0.25	15
222	0.07	0.20	0.11	15
223	0.36	0.38	0.37	34
224	0.16	0.60	0.25	5
225	0.40	0.67	0.50	3
226	0.63	0.71	0.67	48
227	0.00	0.00	0.00	0
228	0.08	0.40	0.13	5
229	0.62	0.77	0.69	26
230	0.20	0.27	0.23	26
231	0.00	0.00	0.00	0
232	0.22	0.29	0.25	14
233	0.00	0.00	0.00	2
234	0.54	0.54	0.54	24
235	0.00	0.00	0.00	1
236	0.04	0.13	0.06	15
237	0.27	0.44	0.33	41
238	0.08	0.18	0.11	22
239	0.06	0.20	0.10	10
240	0.50	0.58	0.54	26
241	0.52	0.73	0.61	15
242	0.20	0.33	0.25	15
243	0.31	0.55	0.40	29
244	0.23	0.24	0.23	29
245	0.04	0.17	0.07	6
246	0.00	0.00	0.00	2
247	0.04	0.08	0.05	13
248	0.26	0.67	0.37	30
249	0.33	0.55	0.41	11
250	0.17	0.20	0.18	10
251	0.42	0.58	0.49	24
252	0.08	0.25	0.12	12
253	0.07	0.25	0.10	12
254	0.25	1.00	0.40	1
255	0.07	0.50	0.12	2
256	0.00	0.00	0.00	1
257	0.21	0.56	0.31	16
258	0.07	0.19	0.10	16
259	0.00	0.00	0.00	2
260	0.00	0.00	0.00	17
261	0.00	0.00	0.00	0
262	0.29	0.55	0.37	11
263	0.00	0.00	0.00	1
264	0.10	0.20	0.13	20
265	0.00	0.00	0.00	3
266	0.12	0.21	0.15	28
267	0.45	0.59	0.51	17
268	0.41	0.70	0.52	10
269	0.46	0.26	0.33	23
270	0.09	0.25	0.13	8
271	0.00	0.00	0.00	20
272	0.00	0.00	0.00	0
273	0.14	0.33	0.20	6
274	0.04	0.17	0.07	6
275	0.23	0.46	0.31	39
276	0.47	0.89	0.62	9
277	0.00	0.00	0.00	8
278	0.00	0.00	0.00	6
279	0.19	0.80	0.31	5
280	0.00	0.00	0.00	4
281	0.22	0.67	0.33	3
282	0.50	0.80	0.62	15
283	0.00	0.00	0.00	0
284	0.27	0.43	0.33	37
285	0.10	0.24	0.14	21
286	0.12	0.27	0.16	11
287	0.07	0.11	0.09	18
288	0.40	0.38	0.39	16
289	0.04	0.09	0.05	11
290	0.35	0.67	0.46	24
291	0.05	0.25	0.08	4
292	0.06	0.22	0.10	9
293	0.14	0.27	0.18	11
294	0.03	0.07	0.05	14

295	0.00	0.00	0.00	13
296	0.06	0.25	0.09	8
297	0.10	0.31	0.16	16
298	0.07	0.12	0.09	34
299	0.16	0.25	0.20	16
300	0.15	0.19	0.17	21
301	0.32	0.43	0.37	23
302	0.29	0.55	0.37	11
303	0.06	0.67	0.11	3
304	0.07	0.19	0.10	16
305	0.06	0.17	0.09	6
306	0.04	0.33	0.07	3
307	0.00	0.00	0.00	2
308	0.03	0.14	0.05	14
309	0.31	0.44	0.37	25
310	0.02	0.06	0.03	17
311	0.14	0.20	0.16	30
312	0.10	0.18	0.12	11
313	0.04	0.14	0.06	14
314	0.07	0.36	0.11	14
315	0.02	0.03	0.02	38
316	0.19	0.43	0.26	7
317	0.14	0.27	0.19	26
318	0.06	1.00	0.11	1
319	0.00	0.00	0.00	5
320	0.00	0.00	0.00	0
321	0.00	0.00	0.00	0
322	0.19	0.50	0.28	8
323	0.57	0.67	0.62	12
324	0.42	0.55	0.48	29
325	0.10	0.50	0.16	4
326	0.11	0.50	0.18	10
327	0.00	0.00	0.00	8
328	0.12	0.30	0.17	10
329	0.00	0.00	0.00	0
330	0.00	0.00	0.00	3
331	0.08	0.25	0.12	12
332	0.00	0.00	0.00	0
333	0.12	0.43	0.19	7
334	0.03	0.07	0.04	14
335	0.00	0.00	0.00	0
336	0.03	0.06	0.04	17
337	0.31	0.45	0.37	20
338	0.17	0.25	0.20	16
339	0.06	0.15	0.08	13
340	0.31	0.38	0.34	26
341	0.26	0.83	0.40	6
342	0.00	0.00	0.00	0
343	0.00	0.00	0.00	0
344	0.11	0.20	0.14	15
345	0.00	0.00	0.00	0
346	0.00	0.00	0.00	1
347	0.04	0.25	0.06	4
348	0.06	0.20	0.09	5
349	0.10	0.21	0.14	14
350	0.12	0.33	0.18	12
351	0.09	0.31	0.14	13
352	0.00	0.00	0.00	15
353	0.38	0.56	0.45	9
354	0.08	0.11	0.09	19
355	0.10	0.50	0.17	2
356	0.05	0.17	0.07	6
357	0.43	0.75	0.55	8
358	0.32	0.47	0.38	19
359	0.14	0.17	0.15	12
360	0.03	0.05	0.04	20
361	0.40	0.32	0.35	19
362	0.67	0.62	0.65	16
363	0.15	0.30	0.20	10
364	0.21	0.38	0.27	8
365	0.12	0.27	0.16	11
366	0.05	0.25	0.09	4
367	0.00	0.00	0.00	1
368	0.12	0.44	0.19	9
369	0.32	1.00	0.48	6
370	0.00	0.00	0.00	7
371	0.07	0.20	0.11	10
372	0.00	0.00	0.00	1
373	0.00	0.00	0.00	2
374	0.12	0.17	0.14	18
375	0.12	0.17	0.14	12
376	0.05	0.12	0.07	16
377	0.00	0.00	0.00	1

378	0.00	0.00	0.00	0
379	0.04	0.07	0.05	14
380	0.20	0.33	0.25	3
381	0.25	0.57	0.35	7
382	0.22	0.40	0.29	10
383	0.00	0.00	0.00	9
384	0.22	0.56	0.31	9
385	0.00	0.00	0.00	7
386	0.40	0.67	0.50	6
387	0.04	0.25	0.07	4
388	0.00	0.00	0.00	8
389	0.00	0.00	0.00	4
390	0.00	0.00	0.00	3
391	0.48	0.62	0.54	26
392	0.00	0.00	0.00	3
393	0.30	0.30	0.30	20
394	0.00	0.00	0.00	3
395	0.00	0.00	0.00	1
396	0.00	0.00	0.00	3
397	0.20	0.67	0.31	3
398	0.30	0.58	0.39	19
399	0.78	0.78	0.78	23
400	0.09	0.15	0.11	13
401	0.31	0.48	0.38	21
402	0.43	0.82	0.56	11
403	0.14	0.17	0.15	6
404	0.33	0.62	0.43	13
405	0.71	0.88	0.79	25
406	0.00	0.00	0.00	2
407	0.02	0.14	0.03	7
408	0.14	1.00	0.25	1
409	0.00	0.00	0.00	6
410	0.00	0.00	0.00	0
411	0.15	0.22	0.18	18
412	0.28	0.41	0.33	17
413	0.10	0.29	0.14	7
414	0.11	0.14	0.12	14
415	0.07	0.33	0.12	3
416	0.17	0.18	0.17	17
417	0.00	0.00	0.00	0
418	0.17	0.38	0.23	8
419	0.00	0.00	0.00	0
420	0.27	0.40	0.32	10
421	0.00	0.00	0.00	1
422	0.00	0.00	0.00	2
423	0.00	0.00	0.00	6
424	0.00	0.00	0.00	1
425	0.15	0.50	0.24	4
426	0.00	0.00	0.00	8
427	0.12	1.00	0.21	2
428	0.22	0.44	0.30	18
429	0.40	0.65	0.50	26
430	0.00	0.00	0.00	0
431	0.06	0.14	0.08	14
432	0.17	0.20	0.18	10
433	0.00	0.00	0.00	7
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	0
436	0.15	0.22	0.18	18
437	0.22	0.29	0.25	28
438	0.26	0.50	0.34	12
439	0.03	0.14	0.04	35
440	0.14	0.20	0.17	10
441	0.19	0.33	0.24	12
442	0.20	0.50	0.29	6
443	0.00	0.00	0.00	9
444	0.00	0.00	0.00	2
445	0.33	0.57	0.42	7
446	0.45	0.62	0.52	21
447	0.04	0.17	0.07	6
448	0.00	0.00	0.00	6
449	0.09	0.50	0.15	2
450	0.12	0.13	0.12	30
451	0.00	0.00	0.00	9
452	0.47	0.90	0.62	10
453	0.00	0.00	0.00	3
454	0.00	0.00	0.00	3
455	0.09	0.22	0.12	9
456	0.00	0.00	0.00	1
457	0.00	0.00	0.00	3
458	0.05	0.25	0.08	4
459	0.29	0.40	0.33	15
460	0.00	0.00	0.00	13

461	0.17	0.43	0.24	7
462	0.33	0.31	0.32	13
463	0.17	0.25	0.20	8
464	0.42	0.73	0.53	22
465	0.10	0.17	0.12	6
466	0.12	0.31	0.17	13
467	0.11	0.21	0.15	19
468	0.24	0.11	0.15	35
469	0.00	0.00	0.00	1
470	0.00	0.00	0.00	2
471	0.21	0.24	0.22	17
472	0.12	0.23	0.16	44
473	0.12	0.30	0.17	10
474	0.32	0.73	0.44	11
475	0.57	0.67	0.62	66
476	0.05	0.33	0.08	3
477	0.24	0.50	0.32	10
478	0.00	0.00	0.00	1
479	0.14	0.44	0.22	9
480	0.01	0.20	0.02	5
481	0.05	0.15	0.07	13
482	0.08	1.00	0.15	1
483	0.15	0.40	0.22	5
484	0.18	0.44	0.26	9
485	0.17	0.17	0.17	12
486	0.25	0.29	0.27	7
487	0.06	0.08	0.07	13
488	0.00	0.00	0.00	0
489	0.17	0.50	0.25	4
490	0.12	0.33	0.18	3
491	0.00	0.00	0.00	3
492	0.00	0.00	0.00	1
493	0.13	0.21	0.16	14
494	0.06	0.33	0.11	3
495	0.16	0.60	0.25	5
496	0.12	0.17	0.14	12
497	0.00	0.00	0.00	0
498	0.00	0.00	0.00	6
499	0.12	0.27	0.17	11
micro avg	0.29	0.44	0.35	19042
macro avg	0.17	0.30	0.20	19042
weighted avg	0.36	0.44	0.38	19042
samples avg	0.38	0.43	0.36	19042

Time taken to run this cell : 0:02:32.796529

Summary

In [4]:

```

from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Sr.No", "MODEL","FEATURIZATION","ALPHA",'micro-f1-score','macro-f1']

x.add_row(['1', 'SGD classifier+ one vs rest(log)','Tfidf', '0.00001', '0.4435', '0.2710' ])
x.add_row(['2', 'Logistic Regression+ one vs rest', 'BOW', '0.001', '0.1445', '0.0156'])
x.add_row(['3', 'SGD classifier+one vs rest(log)', 'BOW', '0.001', '0.3909', '0.2169'])
x.add_row(['4', 'SGD classifier+one vs rest(hinge)', 'BOW', '0.001', '0.3506', '0.2044'])
print(x)

```

Sr.No	MODEL	FEATURIZATION	ALPHA	micro-f1-score	macro-f1
1	SGD classifier+ one vs rest(log)	Tfidf	0.00001	0.4435	0.2710
2	Logistic Regression+ one vs rest	BOW	0.001	0.1445	0.0156
3	SGD classifier+one vs rest(log)	BOW	0.001	0.3909	0.2169
4	SGD classifier+one vs rest(hinge)	BOW	0.001	0.3506	0.2044

In []:

```


```