



```
In [0]: import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import sqlite3
import csv
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from wordcloud import WordCloud
import re
import os
from sqlalchemy import create_engine # database connection
import datetime as dt
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem.snowball import SnowballStemmer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import SGDClassifier
from sklearn import metrics
from sklearn.metrics import f1_score, precision_score, recall_score
from sklearn import svm
from sklearn.linear_model import LogisticRegression
from skmultilearn.adapt import mlknn
from skmultilearn.problem_transform import ClassifierChain
from skmultilearn.problem_transform import BinaryRelevance
from skmultilearn.problem_transform import LabelPowerset
from sklearn.naive_bayes import GaussianNB
from datetime import datetime
```

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/>

1.2 Source / useful links

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

3.1.2 Counting the number of rows

```
In [0]: 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 : 0:01:15.750352

3.1.3 Checking for duplicates

```
In [0]: #Learn SQL: https://www.w3schools.com/sql/default.asp
if os.path.isfile('train.db'):
    start = datetime.now()
    con = sqlite3.connect('train.db')
    df_no_dup = pd.read_sql_query('SELECT Title, Body, Tags, COUNT(*) as cnt_dup FROM data GROUP BY Title, Body, Tags'
, con)
    con.close()
    print("Time taken to run this cell :", datetime.now() - start)
else:
    print("Please download the train.db file from drive or run the first to generate train.db file")
```

Time taken to run this cell : 0:04:33.560122

```
In [0]: df_no_dup.head()
# we can observe that there are duplicates
```

Out[0]:

	Title	Body	Tags	cnt_dup
0	Implementing Boundary Value Analysis of S...	<pre> <code>#include<iosstream>\n#include&...	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>...	java jdbc	2

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

number of duplicate questions : 1827881 (30.2920389063 %)

```
In [0]: # number of times each question appeared in our database
df_no_dup.cnt_dup.value_counts()
```

Out[0]:

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

```
In [0]: start = datetime.now()
df_no_dup["tag_count"] = df_no_dup["Tags"].apply(lambda text: len(text.split(" ")))
# adding a new feature number of tags per question
print("Time taken to run this cell :", datetime.now() - start)
df_no_dup.head()
```

Time taken to run this cell : 0:00:03.169523

Out[0]:

	Title	Body	Tags	cnt_dup	tag_count
0	Implementing Boundary Value Analysis of S...	<pre> <code>#include<iosstream>\n#include&...	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>...	java jdbc	2	2

```
In [0]: # distribution of number of tags per question
df_no_dup.tag_count.value_counts()
```

Out[0]:

```
3    1206157
2    1111706
4     814996
1     568298
5     505158
Name: tag_count, dtype: int64
```

```
In [0]: #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 [12]: #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 : 0:02:02.937662

3.2 Analysis of Tags

3.2.1 Total number of unique tags

```
In [0]: # 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 [14]: 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 [15]: # '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 [0]: # 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))
```



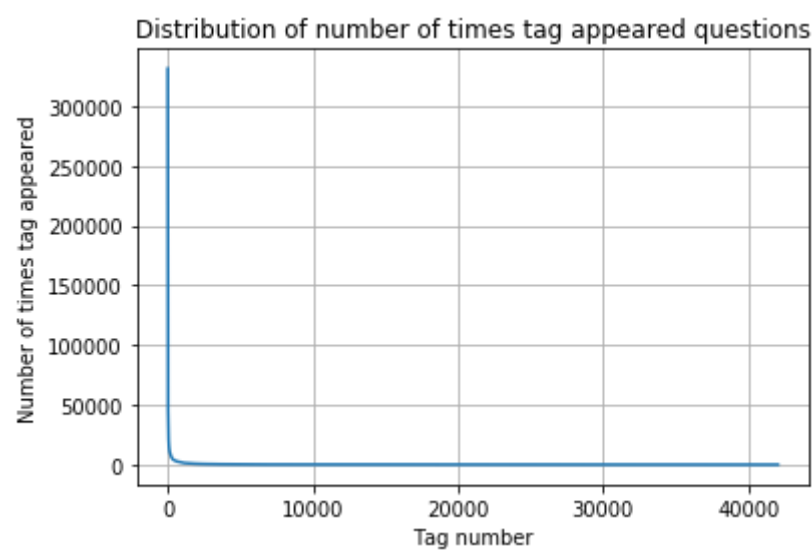
```
In [17]: #Saving this dictionary to csv files.
if not os.path.isfile('tag_counts_dict_dtm.csv'):
    with open('tag_counts_dict_dtm.csv', 'w') as csv_file:
        writer = csv.writer(csv_file)
        for key, value in result.items():
            writer.writerow([key, value])
tag_df = pd.read_csv("tag_counts_dict_dtm.csv", names=['Tags', 'Counts'])
tag_df.head()
```

Out[17]:

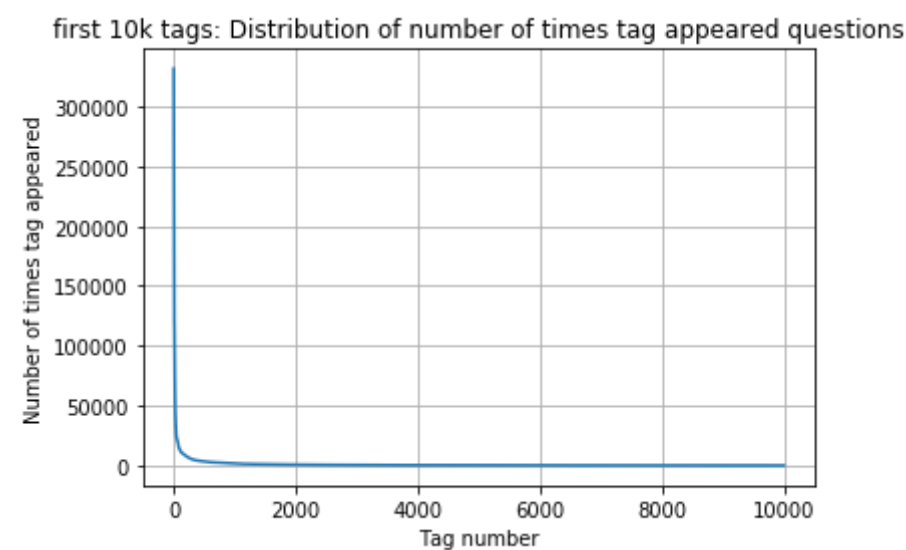
	Tags	Counts
0	.a	18
1	.app	37
2	.asp.net-mvc	1
3	.aspxauth	21
4	.bash-profile	138

```
In [0]: tag_df_sorted = tag_df.sort_values(['Counts'], ascending=False)
tag_counts = tag_df_sorted['Counts'].values
```

```
In [19]: 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()
```

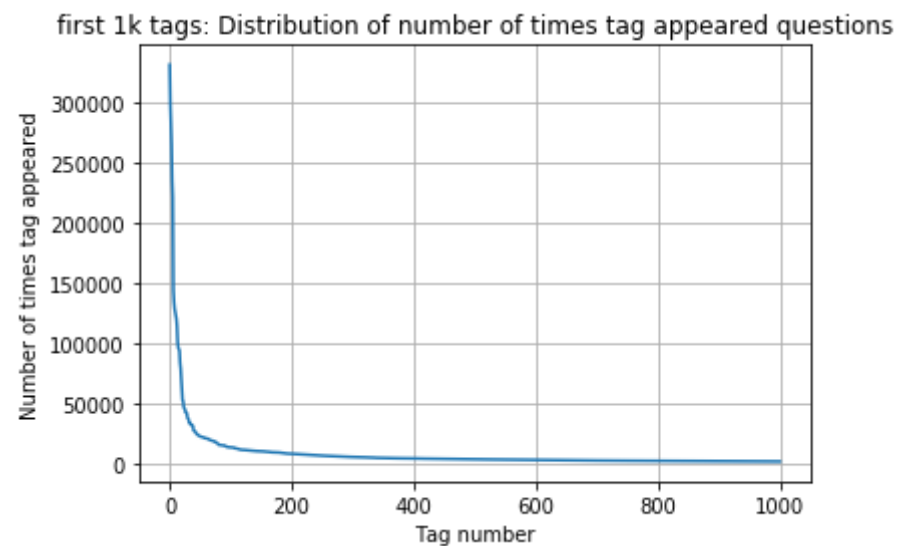


```
In [20]: plt.plot(tag_counts[0:10000])
plt.title('first 10k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:10000:25]), tag_counts[0:10000:25])
```



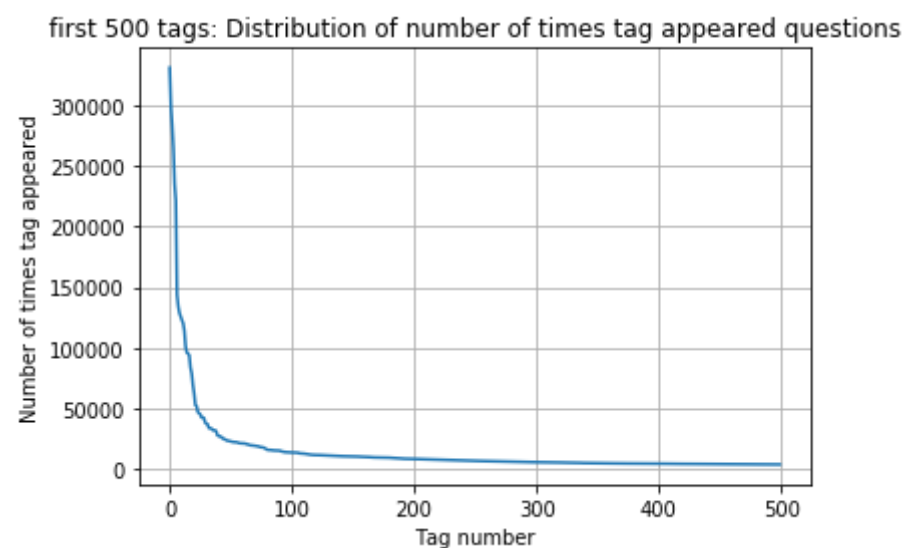
400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7151
6466	5865	5370	4983	4526	4281	4144	3929	3750	3593	
3453	3299	3123	2989	2891	2738	2647	2527	2431	2331	
2259	2186	2097	2020	1959	1900	1828	1770	1723	1673	
1631	1574	1532	1479	1448	1406	1365	1328	1300	1266	
1245	1222	1197	1181	1158	1139	1121	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 [21]: 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	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 [22]: 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])
```

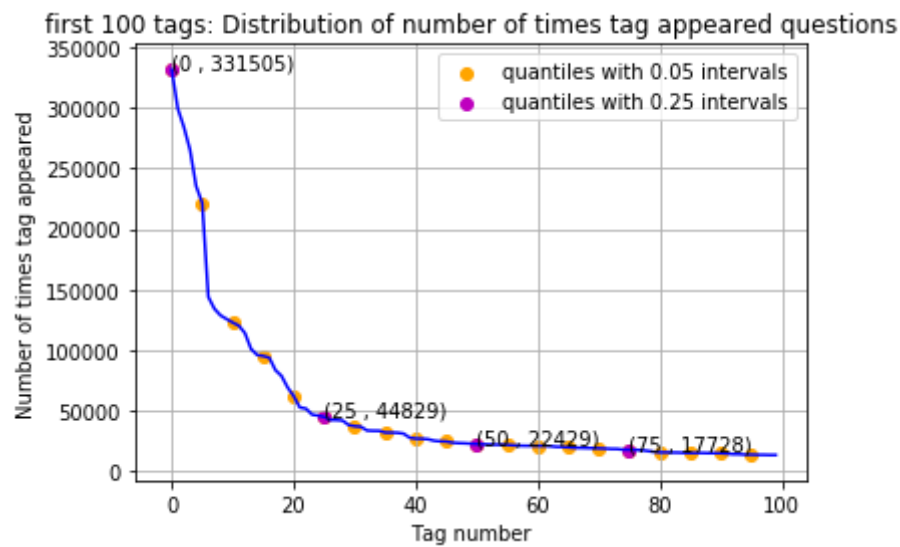


100	22429	21820	20957	19758	18905	17728	15533	15097	14884	13703
	13364	13157	12407	11658	11228	11162	10863	10600	10350	10224
	10029	9884	9719	9411	9252	9148	9040	8617	8361	8163
	8054	7867	7702	7564	7274	7151	7052	6847	6656	6553
	6466	6291	6183	6093	5971	5865	5760	5577	5490	5411
	5370	5283	5207	5107	5066	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 [23]: 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 [24]: # 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 [25]: #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 [26]: 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)))
```

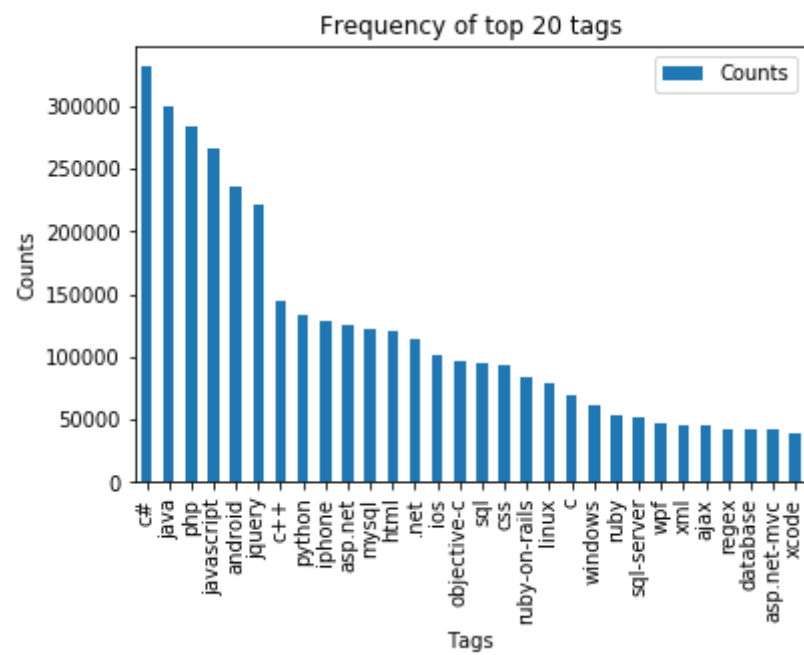
```
Maximum number of tags per question: 5
Minimum number of tags per question: 1
Avg. number of tags per question: 2.899440
```


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



Observations:

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

3.3 Cleaning and preprocessing of Questions

3.3.1 Preprocessing

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

```
In [9]: #http://www.sqlitetutorial.net/sqlite-python/create-tables/
def create_connection(db_file):
    """ create a database connection to the SQLite database
        specified by db_file
    :param db_file: database file
    :return: Connection object or None
    """
    try:
        conn = sqlite3.connect(db_file)
        return conn
    except Error as e:
        print(e)

    return None

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

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

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

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

Tables in the databse:
QuestionsProcessed

```
In [33]: # http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/select-random-row-from-a-sqlite-table
start = datetime.now()
read_db = 'train_no_dup.db'
write_db = 'Processed.db'
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        reader =conn_r.cursor()
        reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY RANDOM() LIMIT 1000000;")

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

Tables in the databse:
QuestionsProcessed
Cleared All the rows
Time taken to run this cell : 0:04:40.404914

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

In [34]: [#http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/](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=striphtml(question.encode('utf-8'))

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

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

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

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

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

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

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

```
number of questions completed= 100000
number of questions completed= 200000
number of questions completed= 300000
number of questions completed= 400000
number of questions completed= 500000
number of questions completed= 600000
number of questions completed= 700000
number of questions completed= 800000
number of questions completed= 900000
Avg. length of questions(Title+Body) before processing: 1172
Avg. length of questions(Title+Body) after processing: 327
Percent of questions containing code: 57
Time taken to run this cell : 0:23:00.589679
```

In [0]: *# 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 [36]: 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

=====

('logitech discharg batteri today replac old aaa batteri logitech previous one held month grant use mous almost everi day nso wonder discharg batteri mous mous move mous click scroll wheel',)

('bing translat languag code zh cn zh chs found zh cn zh chs use languag code microsoft tell zh chs use http msdn mic rosoft com en us librari hh aspx would make differ',)

('sql detect loop parent child relat parent child data excel get load rd parti system run ms sql server data repres d irect hope acycl graph rd parti mean complet free hand schema excel data concaten file possibl exist cross refer vari ous file someon caus loop child elsewhere write vb vba etc excel sql server db excel file almost row worri combinatori explos data set grow techniqu like creat tabl path might pretti unwieldi think simpli write program root tree travers leaf depth get greater nomin valu flag better suggest pointer previous discuss welcom',)

('unabl implement jms use apachemq iam tri implement simpl jms tradit use spring code eclips use apachemq download ap achemq apach org sampl jms sender simplequeuesend receiv simplequeueereceiv respect execut code already gone relat tut ori couldnot find answer question pleas suggest solut chang done regard classpath set activemq start info jetti ninfo activemq webconsol initi ninfo initi spring frameworkervlet dispatch ninfo activemq consol http admin ninfo activemq web demo http demo ninfo rest file access applic http fileserv ninfo start selectchannelconnector proceed next server ad eclips new server program run server program run eclips execut separ consol',)

('upgrad upgrad packag use ppm open ui show sever upgrad packag unfortun select type ctrl click first element hold sh ift select last element list realli walk element press key order select refus believ',)

('nhibern check db schema generat newbi nhibern user tri wrap brain around contempl handl deploy later inject add on web app may requir persist class think use deploy would work pretti well wonder way get nhibern tell common code base way schema export done already basic want smeth like pseudocod two function would intern use respect thank advanc pau l edit guy appreci answer far miss point bit tri set way applic allow addit remov add on may requir chang db talk ver sion code like least primari function question less deploy app add remov plug thei plugin henc pseudo code type check deploy run updat run export make sens',)

('digraph work need horizont output tri learn use php generat graph graphviz modul far got success generat graph use follow code howev show vertic graph block diagram know need use digraph show horizont get error shown second block so rri dont program background occasion purpos error get tri use digraph sure fix list sub modul recent instal imag grap hviz',)

('best way remov duplic uri param string string want creat java method remov duplic param valu alway sometim param du plic per applic function ask therefor becom program long rememb best way origin train thought went someth like grab p aram stick temp array run temp array compar array equal param name delet add back return string end loop print return string would requir length uri plus size array number param think would pretti bad consid run method around per minut incom uri better way go box java method handl overhead',)

('mssql c ef map tinyint int dozen field mark tinyint translat byte field find lot cast int code interact integ think chang entiti framework read integ instead byte implic besid chanc may pass integ rang tinyint ad addit cast may need also think use integ instead databas go high usag db edit zach comment entiti framework map sql server tinyint int lo ok like chang properti byte int ef throw error even way think',)

```
In [0]: #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 [38]: preprocessed_data.head()
```

Out[38]:

	question	tags
0	phpunit returnvaluemap yield expect result tri...	php cakephp phpunit
1	logitech discharg batteri today replac old aaa...	mouse battery logitech
2	bing translat languag code zh cn zh chs found ...	localization microsoft-translator
3	sql detect loop parent child relat parent chil...	sql loops parent-child
4	unabl implement jms use apachemq iam tri imple...	java jms activemq

```
In [39]: print("number of data points in sample :", preprocessed_data.shape[0])
print("number of dimensions :", preprocessed_data.shape[1])

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

4. Machine Learning Models

4.1 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 [0]: # binary='true' will give a binary vectorizer
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

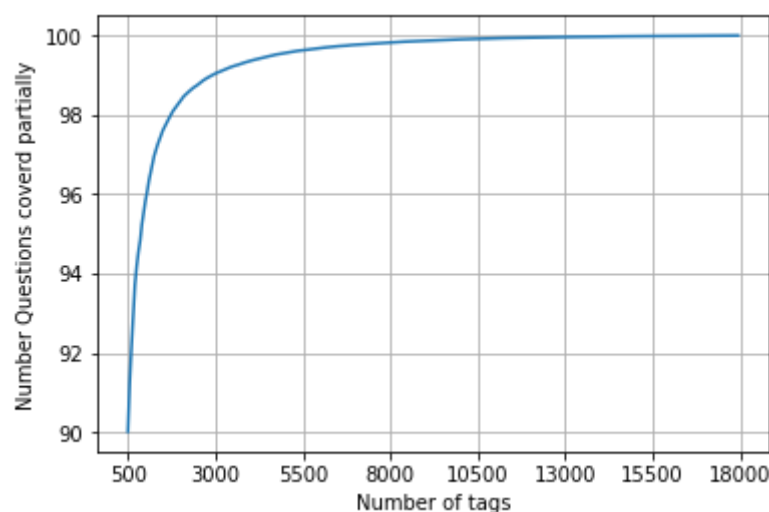
We will sample the number of tags instead considering all of them (due to limitation of computing power)

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

```
In [43]: 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 99.044 % of questions

```
In [44]: 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 : 9563 out of 999999
```

```
In [45]: 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 : 35449
number of tags taken : 5500 ( 15.515247256622189 %)
```

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

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

```
In [0]: 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 [47]: 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 : (799999, 5500)
Number of data points in test data : (200000, 5500)
```

4.3 Featurizing data

```
In [48]: 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 : 0:07:41.530373
```

```
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: (799999, 88311) Y : (799999, 5500)
Dimensions of test data X: (200000, 88311) Y: (200000, 5500)
```

```
In [0]: # 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[0]: "\nfrom skmultilearn.adapt import MLkNN\n\nclassifier = MLkNN(k=21)\n\n# train\n\nclassifier.fit(x_train_multilabel, y_train)\n\n# predict\n\npredictions = classifier.predict(x_test_multilabel)\n\nprint(accuracy_score(y_test, predictions))\n\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"
```

4.4 Applying Logistic Regression with OneVsRest Classifier

```
In [0]: # this will be taking so much time try not to run it, download the lr_with_equal_weight.pkl file and use to predict
# This takes about 6-7 hours to run.
#classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_jobs=-1)
#classifier.fit(x_train_multilabel, y_train)
from sklearn.externals import joblib
classifier=joblib.load('lr_with_equal_weight.pkl')
predictions = classifier.predict(x_test_multilabel)

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

```
In [0]: from sklearn.externals import joblib
joblib.dump(classifier, 'lr_with_equal_weight.pkl')
```

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

```
In [51]: 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 [52]: # http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/select-random-row-from-a-sqlite-table

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

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

Tables in the 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 [53]: #http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/
start = datetime.now()
preprocessed_data_list=[]
reader.fetchone()
questions_with_code=0
len_pre=0
len_post=0
questions_proccesed = 0
for row in reader:

    is_code = 0

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

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

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

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

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

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

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

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

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

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

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

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

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

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

```

```

number of questions completed= 100000
number of questions completed= 200000
number of questions completed= 300000
number of questions completed= 400000
number of questions completed= 500000
Avg. length of questions(Title+Body) before processing: 1239
Avg. length of questions(Title+Body) after processing: 424
Percent of questions containing code: 57
Time taken to run this cell : 0:17:19.462904

```

```

In [0]: # 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 [55]: 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 nth ank repli advance..',)

('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid java.lang.noclassdeffounderror javax servle t jsp tagext taglibraryvalid java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid follow guid link instal jstl got follow error tri launch jsp page java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryva lid taglib declar instal jstl 1.1 tomcat webapp tri project work also tri version 1.2 jstl still messag caus solv',)

('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index java.sql.sqlexcept microsoft odbc driver ma nag invalid descriptor index java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index use follow code display caus solv',)

('better way updat feed fb php sdk better way updat feed fb php sdk better way updat feed fb php sdk novic facebook a pi read mani tutori still confused.i find post feed api method like correct second way use curl someth like way bette r',)

('btnadd click event open two window record ad btnadd click event open two window record ad btnadd click event open t wo window record ad open window search.aspx use code hav add button search.aspx nwhen insert record btnadd click even t 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 p revent correct form submiss php check everyth think make sure input field safe type sql inject good news safe bad new s 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 se e submit noth work flawless statement though also mention script work flawless local machin use host come across prob lem state list input test mess',)

('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl subaddit lebesgu measur let lbrace rbrac e sequenc set sigma -algebra mathcal want show left bigcup right leq sum left right countabl addit measur defin set s igma algebra mathcal think use monoton properti somewher proof start appreci littl help nthank ad han answer make fol low addit construct given han answer clear bigcup bigcup cap emptyset neq left bigcup right left bigcup right sum lef t right also construct subset monoton left right leq left right final would sum leq sum result follow',)

('hql equival sql queri hql equival sql queri hql equival sql queri hql queri replac name class properti name error o ccur hql error',)

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

Saving Preprocessed data to a Database

```
In [0]: #Taking 0.5 Million entries to a dataframe.
write_db = 'Titlemoreweight.db'
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM QuestionsProcessed""", conn_r)
conn_r.commit()
conn_r.close()
```

```
In [57]: preprocessed_data.head()
```

Out[57]:

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

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

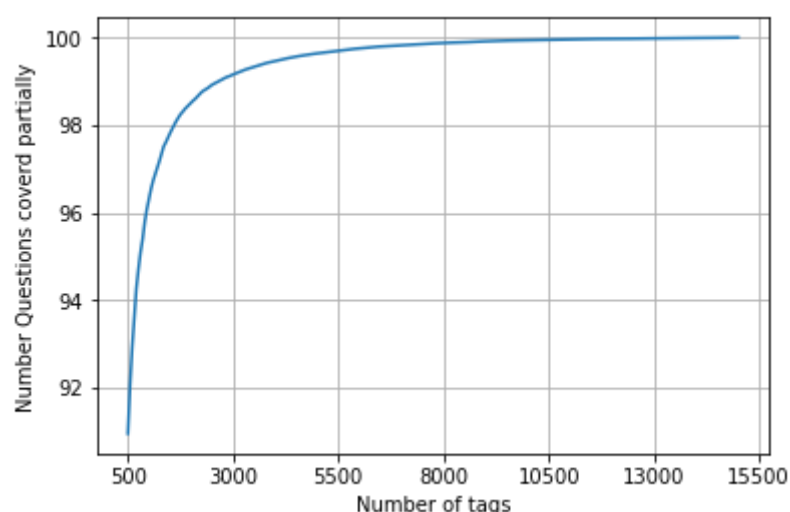
Converting string Tags to multilable output variables

```
In [0]: vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

Selecting 500 Tags

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

```
In [61]: fig, ax = plt.subplots()
ax.plot(questions_explained)
xlabel = list(500+np.array(range(-50,450,50))*50)
ax.set_xticklabels(xlabel)
plt.xlabel("Number of tags")
plt.ylabel("Number Questions covered partially")
plt.grid()
plt.show()
# you can choose any number of tags based on your computing power, minimum is 500(it covers 90% of the tags)
print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



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

```
In [62]: # we will be taking 500 tags
multilabel_yx = tags_to_choose(500)
print("number of questions that are not covered :", questions_explained_fn(500),"out of ", total_qs)
```

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

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

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

```
In [64]: print("Number of data points in train data :", y_train.shape)
print("Number of data points in test data :", y_test.shape)
```

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

4.5.2 Featurizing data with Tfidf vectorizer

```
In [65]: 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 : 0:04:37.441362
```

```
In [66]: print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.shape)
        print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

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

4.5.3 Applying Logistic Regression with OneVsRest Classifier


```
In [0]: 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 : 0.23617
 Hamming loss 0.00278304
 Micro-average quality numbers
 Precision: 0.7215, Recall: 0.3248, F1-measure: 0.4480
 Macro-average quality numbers
 Precision: 0.5464, Recall: 0.2576, F1-measure: 0.3343

	precision	recall	f1-score	support
0	0.94	0.64	0.76	5519
1	0.69	0.26	0.38	8190
2	0.81	0.38	0.51	6529
3	0.81	0.43	0.56	3231
4	0.81	0.40	0.54	6430
5	0.82	0.34	0.48	2879
6	0.87	0.50	0.63	5086
7	0.88	0.54	0.67	4533
8	0.61	0.13	0.21	3000
9	0.81	0.53	0.64	2765
10	0.59	0.16	0.26	3051
11	0.70	0.33	0.45	3009
12	0.66	0.25	0.36	2630
13	0.71	0.23	0.35	1426
14	0.90	0.53	0.67	2548
15	0.67	0.18	0.28	2371
16	0.64	0.23	0.33	873
17	0.89	0.61	0.72	2151
18	0.63	0.23	0.34	2204
19	0.72	0.41	0.52	831
20	0.77	0.40	0.53	1860
21	0.27	0.07	0.11	2023
22	0.49	0.21	0.30	1513
23	0.91	0.49	0.63	1207
24	0.57	0.30	0.39	506
25	0.67	0.30	0.41	425
26	0.64	0.40	0.49	793
27	0.60	0.32	0.42	1291
28	0.73	0.36	0.48	1208
29	0.43	0.09	0.14	406
30	0.74	0.18	0.29	504
31	0.29	0.09	0.14	732
32	0.57	0.24	0.33	441
33	0.56	0.17	0.26	1645
34	0.71	0.25	0.37	1058
35	0.83	0.55	0.66	946
36	0.66	0.18	0.29	644
37	0.97	0.67	0.79	136
38	0.63	0.36	0.45	570
39	0.85	0.28	0.43	766
40	0.60	0.28	0.38	1132
41	0.45	0.18	0.26	174
42	0.79	0.51	0.62	210
43	0.80	0.40	0.54	433
44	0.66	0.51	0.58	626
45	0.74	0.32	0.44	852
46	0.75	0.42	0.54	534
47	0.33	0.13	0.18	350
48	0.74	0.50	0.60	496
49	0.79	0.61	0.69	785
50	0.18	0.05	0.07	475
51	0.32	0.10	0.15	305
52	0.53	0.03	0.06	251
53	0.68	0.40	0.50	914
54	0.45	0.16	0.23	728
55	0.18	0.01	0.01	258
56	0.46	0.19	0.26	821
57	0.45	0.09	0.15	541
58	0.78	0.28	0.41	748
59	0.94	0.62	0.75	724
60	0.33	0.06	0.11	660
61	0.84	0.18	0.29	235
62	0.91	0.71	0.80	718
63	0.83	0.63	0.72	468
64	0.54	0.32	0.41	191
65	0.35	0.12	0.18	429
66	0.29	0.05	0.09	415
67	0.75	0.49	0.59	274
68	0.82	0.52	0.64	510
69	0.67	0.45	0.54	466
70	0.30	0.07	0.11	305
71	0.46	0.15	0.23	247
72	0.78	0.48	0.59	401
73	0.98	0.73	0.84	86
74	0.75	0.37	0.49	120
75	0.89	0.68	0.77	129
76	0.29	0.00	0.01	473
77	0.38	0.27	0.31	143
78	0.79	0.44	0.57	347
79	0.73	0.23	0.35	479

80	0.54	0.33	0.41	279
81	0.78	0.17	0.28	461
82	0.22	0.01	0.03	298
83	0.77	0.45	0.57	396
84	0.55	0.34	0.42	184
85	0.66	0.20	0.31	573
86	0.46	0.05	0.09	325
87	0.51	0.27	0.35	273
88	0.42	0.20	0.27	135
89	0.31	0.07	0.12	232
90	0.55	0.30	0.39	409
91	0.63	0.24	0.35	420
92	0.76	0.53	0.62	408
93	0.69	0.49	0.58	241
94	0.31	0.04	0.07	211
95	0.33	0.07	0.12	277
96	0.26	0.03	0.06	410
97	0.89	0.30	0.45	501
98	0.74	0.60	0.66	136
99	0.52	0.28	0.36	239
100	0.56	0.13	0.21	324
101	0.93	0.61	0.74	277
102	0.92	0.70	0.79	613
103	0.49	0.17	0.25	157
104	0.22	0.06	0.09	295
105	0.83	0.34	0.48	334
106	0.77	0.12	0.21	335
107	0.76	0.48	0.59	389
108	0.57	0.23	0.32	251
109	0.53	0.40	0.46	317
110	0.65	0.08	0.14	187
111	0.48	0.07	0.12	140
112	0.58	0.25	0.35	154
113	0.65	0.18	0.29	332
114	0.45	0.28	0.34	323
115	0.47	0.22	0.30	344
116	0.76	0.49	0.60	370
117	0.56	0.22	0.32	313
118	0.78	0.68	0.72	874
119	0.46	0.20	0.28	293
120	0.00	0.00	0.00	200
121	0.77	0.48	0.59	463
122	0.36	0.08	0.14	119
123	0.75	0.01	0.02	256
124	0.91	0.70	0.79	195
125	0.45	0.14	0.21	138
126	0.80	0.49	0.61	376
127	0.14	0.03	0.05	122
128	0.14	0.03	0.05	252
129	0.42	0.10	0.16	144
130	0.40	0.08	0.13	150
131	0.27	0.01	0.03	210
132	0.66	0.26	0.37	361
133	0.94	0.54	0.68	453
134	0.89	0.73	0.81	124
135	0.21	0.03	0.06	91
136	0.67	0.27	0.38	128
137	0.57	0.33	0.42	218
138	0.75	0.15	0.25	243
139	0.38	0.18	0.24	149
140	0.76	0.43	0.55	318
141	0.28	0.12	0.17	159
142	0.66	0.36	0.47	274
143	0.87	0.72	0.78	362
144	0.61	0.17	0.26	118
145	0.67	0.37	0.47	164
146	0.60	0.28	0.38	461
147	0.66	0.42	0.51	159
148	0.34	0.14	0.20	166
149	0.98	0.46	0.63	346
150	0.63	0.08	0.14	350
151	0.90	0.65	0.76	55
152	0.79	0.46	0.58	387
153	0.50	0.11	0.18	150
154	0.59	0.12	0.20	281
155	0.23	0.04	0.07	202
156	0.76	0.62	0.68	130
157	0.29	0.07	0.12	245
158	0.89	0.58	0.70	177
159	0.49	0.25	0.33	130
160	0.50	0.13	0.20	336
161	0.93	0.59	0.72	220
162	0.16	0.03	0.05	229
163	0.89	0.40	0.55	316
164	0.76	0.35	0.48	283
165	0.63	0.32	0.42	197
166	0.51	0.27	0.35	101
167	0.47	0.18	0.26	231

168	0.59	0.23	0.33	370
169	0.42	0.19	0.26	258
170	0.26	0.05	0.08	101
171	0.38	0.22	0.28	89
172	0.51	0.35	0.41	193
173	0.42	0.22	0.29	309
174	0.52	0.15	0.23	172
175	0.93	0.71	0.80	95
176	0.94	0.59	0.72	346
177	0.95	0.44	0.60	322
178	0.63	0.46	0.53	232
179	0.29	0.06	0.09	125
180	0.55	0.27	0.36	145
181	0.39	0.09	0.15	77
182	0.15	0.02	0.04	182
183	0.61	0.31	0.41	257
184	0.08	0.01	0.02	216
185	0.31	0.07	0.11	242
186	0.38	0.15	0.22	165
187	0.75	0.56	0.64	263
188	0.30	0.09	0.14	174
189	0.71	0.31	0.43	136
190	0.88	0.50	0.63	202
191	0.42	0.15	0.22	134
192	0.73	0.40	0.52	230
193	0.43	0.18	0.25	90
194	0.57	0.48	0.52	185
195	0.18	0.04	0.06	156
196	0.42	0.09	0.15	160
197	0.64	0.07	0.12	266
198	0.38	0.05	0.09	284
199	0.41	0.06	0.11	145
200	0.94	0.69	0.80	212
201	0.67	0.22	0.33	317
202	0.78	0.54	0.64	427
203	0.28	0.08	0.12	232
204	0.51	0.21	0.30	217
205	0.48	0.44	0.46	527
206	0.14	0.02	0.03	124
207	0.47	0.09	0.15	103
208	0.89	0.49	0.63	287
209	0.34	0.09	0.14	193
210	0.70	0.31	0.43	220
211	0.78	0.20	0.32	140
212	0.16	0.02	0.03	161
213	0.50	0.22	0.31	72
214	0.61	0.46	0.52	396
215	0.86	0.33	0.48	134
216	0.45	0.04	0.08	400
217	0.51	0.24	0.33	75
218	0.96	0.75	0.85	219
219	0.77	0.36	0.49	210
220	0.91	0.59	0.72	298
221	0.97	0.59	0.73	266
222	0.77	0.41	0.54	290
223	0.08	0.01	0.01	128
224	0.80	0.40	0.53	159
225	0.58	0.29	0.39	164
226	0.63	0.35	0.45	144
227	0.59	0.31	0.41	276
228	0.16	0.02	0.03	235
229	0.42	0.02	0.04	216
230	0.35	0.18	0.23	228
231	0.72	0.48	0.58	64
232	0.44	0.07	0.12	103
233	0.71	0.31	0.43	216
234	0.71	0.09	0.15	116
235	0.56	0.39	0.46	77
236	0.96	0.64	0.77	67
237	0.54	0.06	0.11	218
238	0.27	0.06	0.09	139
239	0.17	0.01	0.02	94
240	0.57	0.30	0.39	77
241	0.50	0.08	0.14	167
242	0.83	0.29	0.43	86
243	0.42	0.14	0.21	58
244	0.60	0.17	0.26	269
245	0.18	0.06	0.09	112
246	0.95	0.74	0.83	255
247	0.46	0.21	0.29	58
248	0.25	0.02	0.04	81
249	0.00	0.00	0.00	131
250	0.40	0.20	0.27	93
251	0.68	0.29	0.40	154
252	0.33	0.04	0.07	129
253	0.62	0.29	0.39	83
254	0.39	0.09	0.14	191
255	0.15	0.02	0.04	219

256	0.22	0.03	0.05	130
257	0.45	0.28	0.34	93
258	0.68	0.41	0.51	217
259	0.29	0.10	0.15	141
260	0.95	0.13	0.22	143
261	0.53	0.11	0.18	219
262	0.55	0.29	0.38	107
263	0.38	0.21	0.27	236
264	0.27	0.17	0.21	119
265	0.37	0.15	0.22	72
266	0.00	0.00	0.00	70
267	0.30	0.13	0.18	107
268	0.67	0.44	0.53	169
269	0.32	0.11	0.16	129
270	0.73	0.52	0.61	159
271	0.80	0.35	0.49	190
272	0.59	0.22	0.32	248
273	0.91	0.70	0.79	264
274	0.89	0.65	0.75	105
275	0.62	0.08	0.14	104
276	0.14	0.02	0.03	115
277	0.83	0.60	0.70	170
278	0.67	0.25	0.36	145
279	0.92	0.62	0.74	230
280	0.56	0.44	0.49	80
281	0.68	0.56	0.61	217
282	0.74	0.47	0.58	175
283	0.33	0.06	0.10	269
284	0.62	0.24	0.35	74
285	0.86	0.50	0.63	206
286	0.90	0.59	0.71	227
287	0.85	0.31	0.45	130
288	0.38	0.06	0.11	129
289	0.40	0.03	0.05	80
290	0.15	0.07	0.10	99
291	0.77	0.31	0.44	208
292	0.29	0.03	0.05	67
293	0.83	0.40	0.54	109
294	0.39	0.24	0.30	140
295	0.25	0.08	0.12	241
296	0.23	0.10	0.14	72
297	0.22	0.04	0.06	107
298	0.80	0.39	0.53	61
299	0.91	0.38	0.53	77
300	0.19	0.06	0.10	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.56	0.33	0.41	176
304	0.96	0.72	0.82	230
305	0.96	0.60	0.74	156
306	0.52	0.37	0.43	146
307	0.29	0.08	0.13	98
308	0.00	0.00	0.00	78
309	0.71	0.05	0.10	94
310	0.76	0.36	0.49	162
311	0.81	0.53	0.64	116
312	0.48	0.26	0.34	57
313	0.75	0.05	0.09	65
314	0.48	0.35	0.40	138
315	0.55	0.21	0.30	195
316	0.43	0.23	0.30	69
317	0.35	0.10	0.16	134
318	0.50	0.34	0.41	148
319	0.85	0.45	0.59	161
320	0.21	0.14	0.17	104
321	0.85	0.54	0.66	156
322	0.58	0.31	0.40	134
323	0.57	0.38	0.45	232
324	0.42	0.16	0.23	92
325	0.45	0.31	0.37	197
326	0.13	0.02	0.04	126
327	0.45	0.04	0.08	115
328	0.98	0.64	0.78	198
329	0.61	0.31	0.41	125
330	0.80	0.20	0.32	81
331	0.40	0.06	0.11	94
332	0.50	0.02	0.03	56
333	0.15	0.03	0.05	260
334	0.20	0.03	0.06	60
335	0.29	0.07	0.12	110
336	0.63	0.41	0.50	71
337	0.19	0.05	0.07	66
338	0.45	0.31	0.37	150
339	0.00	0.00	0.00	54
340	0.85	0.54	0.66	195
341	0.88	0.19	0.31	79
342	0.40	0.16	0.23	38
343	0.71	0.40	0.51	43

344	0.53	0.24	0.33	68
345	0.68	0.37	0.48	73
346	0.30	0.03	0.05	116
347	0.88	0.32	0.47	111
348	0.30	0.10	0.14	63
349	0.82	0.57	0.67	104
350	0.62	0.45	0.53	44
351	0.78	0.17	0.29	40
352	0.95	0.40	0.57	136
353	0.44	0.20	0.28	54
354	0.42	0.04	0.07	134
355	0.57	0.28	0.37	120
356	0.52	0.21	0.30	228
357	0.66	0.26	0.38	269
358	0.70	0.35	0.47	80
359	0.87	0.46	0.60	140
360	0.37	0.13	0.19	125
361	0.90	0.62	0.73	169
362	0.11	0.04	0.05	56
363	0.94	0.66	0.77	154
364	0.50	0.07	0.12	58
365	0.26	0.13	0.17	71
366	1.00	0.65	0.79	54
367	0.36	0.04	0.08	116
368	0.33	0.02	0.04	54
369	0.00	0.00	0.00	71
370	0.20	0.03	0.06	61
371	0.50	0.08	0.14	71
372	0.65	0.46	0.54	52
373	0.79	0.35	0.49	150
374	0.38	0.13	0.19	93
375	0.15	0.03	0.05	67
376	0.00	0.00	0.00	76
377	0.74	0.16	0.26	106
378	0.11	0.01	0.02	86
379	0.33	0.07	0.12	14
380	1.00	0.39	0.56	122
381	0.18	0.03	0.05	104
382	0.28	0.08	0.12	66
383	0.50	0.27	0.35	110
384	0.00	0.00	0.00	155
385	0.45	0.10	0.16	50
386	0.27	0.11	0.16	64
387	0.31	0.05	0.09	93
388	0.61	0.27	0.38	102
389	0.07	0.01	0.02	108
390	0.96	0.64	0.77	178
391	0.58	0.16	0.25	115
392	0.77	0.40	0.53	42
393	0.00	0.00	0.00	134
394	0.50	0.02	0.03	112
395	0.43	0.12	0.19	176
396	0.42	0.09	0.15	125
397	0.69	0.24	0.36	224
398	0.88	0.57	0.69	63
399	0.00	0.00	0.00	59
400	0.50	0.33	0.40	63
401	0.45	0.17	0.25	98
402	0.57	0.16	0.25	162
403	0.41	0.14	0.21	83
404	0.73	0.84	0.78	19
405	0.26	0.07	0.10	92
406	0.86	0.15	0.25	41
407	0.64	0.33	0.43	43
408	0.80	0.32	0.46	160
409	0.20	0.12	0.15	50
410	0.00	0.00	0.00	19
411	0.36	0.10	0.15	175
412	0.27	0.06	0.09	72
413	0.56	0.05	0.10	95
414	0.19	0.03	0.05	97
415	0.32	0.17	0.22	48
416	0.45	0.30	0.36	83
417	0.50	0.07	0.13	40
418	0.33	0.07	0.11	91
419	0.49	0.28	0.35	90
420	0.29	0.22	0.25	37
421	0.00	0.00	0.00	66
422	0.62	0.33	0.43	73
423	0.48	0.25	0.33	56
424	0.93	0.82	0.87	33
425	0.00	0.00	0.00	76
426	0.25	0.05	0.08	81
427	0.99	0.67	0.80	150
428	0.95	0.66	0.78	29
429	0.99	0.70	0.82	389
430	0.65	0.36	0.46	167
431	0.53	0.08	0.14	123

432	0.45	0.36	0.40	39
433	0.29	0.16	0.20	82
434	1.00	0.65	0.79	66
435	0.65	0.46	0.54	93
436	0.52	0.26	0.35	87
437	0.27	0.07	0.11	86
438	0.77	0.48	0.59	104
439	0.62	0.13	0.21	100
440	0.20	0.01	0.01	141
441	0.42	0.25	0.31	110
442	0.40	0.14	0.20	123
443	0.50	0.13	0.20	71
444	0.44	0.07	0.13	109
445	0.36	0.17	0.23	48
446	0.43	0.25	0.32	76
447	0.28	0.13	0.18	38
448	0.68	0.53	0.60	81
449	0.58	0.17	0.26	132
450	0.47	0.28	0.35	81
451	0.88	0.28	0.42	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.91	0.44	0.60	70
455	0.48	0.07	0.12	155
456	0.47	0.16	0.24	43
457	0.48	0.19	0.28	72
458	0.26	0.08	0.12	62
459	0.60	0.13	0.21	69
460	0.07	0.01	0.02	119
461	0.80	0.15	0.26	79
462	0.69	0.23	0.35	47
463	0.25	0.05	0.08	104
464	0.64	0.32	0.43	106
465	0.54	0.11	0.18	64
466	0.58	0.28	0.38	173
467	0.82	0.38	0.52	107
468	0.83	0.12	0.21	126
469	0.00	0.00	0.00	114
470	0.94	0.79	0.86	140
471	0.91	0.25	0.40	79
472	0.40	0.29	0.34	143
473	0.69	0.30	0.42	158
474	0.38	0.07	0.11	138
475	0.00	0.00	0.00	59
476	0.57	0.31	0.40	88
477	0.86	0.57	0.69	176
478	0.94	0.71	0.81	24
479	0.09	0.01	0.02	92
480	0.82	0.45	0.58	100
481	0.49	0.17	0.26	103
482	0.49	0.23	0.31	74
483	0.83	0.57	0.68	105
484	0.25	0.02	0.04	83
485	0.14	0.01	0.02	82
486	0.41	0.13	0.19	71
487	0.45	0.21	0.29	120
488	0.33	0.02	0.04	105
489	0.69	0.29	0.41	87
490	1.00	0.81	0.90	32
491	0.00	0.00	0.00	69
492	0.00	0.00	0.00	49
493	0.00	0.00	0.00	117
494	0.50	0.16	0.25	61
495	0.99	0.44	0.61	344
496	0.39	0.23	0.29	52
497	0.63	0.20	0.30	137
498	0.31	0.04	0.07	98
499	0.68	0.16	0.27	79
micro avg	0.72	0.32	0.45	173812
macro avg	0.55	0.26	0.33	173812
weighted avg	0.66	0.32	0.42	173812
samples avg	0.41	0.31	0.33	173812

Time taken to run this cell : 0:19:45.078539

```
In [0]: joblib.dump(classifier, 'lr_with_more_title_weight.pkl')
```

```
Out[0]: ['lr_with_more_title_weight.pkl']
```

4.5.4 Applying Logistic Regression (directly) with OneVsRest Classifier

```
In [0]: 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 : 0.25107
 Hamming loss 0.00270298
 Micro-average quality numbers
 Precision: 0.7172, Recall: 0.3673, F1-measure: 0.4858
 Macro-average quality numbers
 Precision: 0.5570, Recall: 0.2951, F1-measure: 0.3710

	precision	recall	f1-score	support
0	0.94	0.72	0.82	5519
1	0.70	0.34	0.45	8190
2	0.80	0.42	0.55	6529
3	0.82	0.49	0.61	3231
4	0.80	0.44	0.57	6430
5	0.82	0.38	0.52	2879
6	0.86	0.53	0.66	5086
7	0.87	0.58	0.70	4533
8	0.60	0.13	0.22	3000
9	0.82	0.57	0.67	2765
10	0.60	0.20	0.30	3051
11	0.68	0.38	0.49	3009
12	0.62	0.29	0.40	2630
13	0.73	0.30	0.43	1426
14	0.89	0.57	0.70	2548
15	0.65	0.23	0.34	2371
16	0.65	0.25	0.37	873
17	0.89	0.63	0.74	2151
18	0.60	0.25	0.35	2204
19	0.71	0.41	0.52	831
20	0.76	0.47	0.58	1860
21	0.29	0.09	0.14	2023
22	0.52	0.24	0.33	1513
23	0.89	0.55	0.68	1207
24	0.56	0.28	0.38	506
25	0.69	0.34	0.45	425
26	0.65	0.43	0.52	793
27	0.62	0.38	0.47	1291
28	0.74	0.39	0.51	1208
29	0.46	0.10	0.17	406
30	0.76	0.21	0.33	504
31	0.26	0.08	0.12	732
32	0.60	0.29	0.39	441
33	0.60	0.27	0.38	1645
34	0.69	0.26	0.38	1058
35	0.83	0.58	0.68	946
36	0.65	0.24	0.35	644
37	0.98	0.65	0.78	136
38	0.62	0.38	0.47	570
39	0.84	0.31	0.45	766
40	0.59	0.35	0.44	1132
41	0.47	0.18	0.26	174
42	0.75	0.48	0.59	210
43	0.75	0.42	0.54	433
44	0.66	0.52	0.58	626
45	0.71	0.36	0.47	852
46	0.77	0.45	0.57	534
47	0.37	0.15	0.22	350
48	0.75	0.52	0.62	496
49	0.78	0.64	0.71	785
50	0.21	0.06	0.09	475
51	0.37	0.13	0.19	305
52	0.42	0.03	0.06	251
53	0.66	0.40	0.50	914
54	0.49	0.17	0.26	728
55	0.47	0.03	0.05	258
56	0.45	0.24	0.31	821
57	0.46	0.10	0.17	541
58	0.76	0.31	0.45	748
59	0.94	0.66	0.77	724
60	0.35	0.10	0.15	660
61	0.78	0.20	0.31	235
62	0.92	0.74	0.82	718
63	0.83	0.69	0.75	468
64	0.55	0.36	0.43	191
65	0.33	0.11	0.17	429
66	0.29	0.06	0.10	415
67	0.74	0.50	0.59	274
68	0.82	0.53	0.64	510
69	0.67	0.45	0.54	466
70	0.30	0.09	0.13	305
71	0.49	0.17	0.25	247
72	0.78	0.53	0.64	401
73	0.99	0.77	0.86	86
74	0.72	0.42	0.53	120
75	0.92	0.67	0.78	129
76	0.47	0.02	0.04	473
77	0.40	0.29	0.33	143
78	0.79	0.49	0.60	347
79	0.69	0.25	0.36	479

80	0.56	0.34	0.43	279
81	0.70	0.23	0.34	461
82	0.34	0.04	0.07	298
83	0.78	0.50	0.61	396
84	0.55	0.29	0.38	184
85	0.61	0.24	0.35	573
86	0.50	0.07	0.12	325
87	0.51	0.29	0.37	273
88	0.49	0.21	0.30	135
89	0.36	0.11	0.17	232
90	0.56	0.34	0.43	409
91	0.61	0.27	0.37	420
92	0.78	0.57	0.66	408
93	0.66	0.44	0.53	241
94	0.30	0.04	0.07	211
95	0.37	0.10	0.15	277
96	0.28	0.04	0.07	410
97	0.86	0.43	0.57	501
98	0.75	0.63	0.69	136
99	0.54	0.34	0.42	239
100	0.57	0.15	0.24	324
101	0.91	0.68	0.78	277
102	0.91	0.75	0.82	613
103	0.47	0.17	0.25	157
104	0.22	0.06	0.10	295
105	0.75	0.43	0.55	334
106	0.88	0.28	0.43	335
107	0.75	0.54	0.63	389
108	0.58	0.27	0.37	251
109	0.58	0.45	0.51	317
110	0.68	0.10	0.18	187
111	0.73	0.11	0.20	140
112	0.67	0.43	0.52	154
113	0.58	0.20	0.29	332
114	0.46	0.27	0.34	323
115	0.47	0.26	0.33	344
116	0.75	0.55	0.63	370
117	0.58	0.24	0.34	313
118	0.78	0.73	0.75	874
119	0.45	0.21	0.29	293
120	0.11	0.01	0.01	200
121	0.77	0.51	0.61	463
122	0.32	0.10	0.15	119
123	0.67	0.02	0.03	256
124	0.91	0.70	0.79	195
125	0.44	0.14	0.21	138
126	0.81	0.54	0.65	376
127	0.27	0.03	0.06	122
128	0.20	0.04	0.07	252
129	0.48	0.22	0.30	144
130	0.42	0.11	0.18	150
131	0.33	0.03	0.06	210
132	0.65	0.28	0.39	361
133	0.92	0.59	0.72	453
134	0.89	0.77	0.82	124
135	0.31	0.05	0.09	91
136	0.69	0.28	0.40	128
137	0.55	0.37	0.44	218
138	0.67	0.18	0.28	243
139	0.45	0.18	0.26	149
140	0.77	0.46	0.58	318
141	0.32	0.10	0.15	159
142	0.63	0.38	0.47	274
143	0.85	0.79	0.82	362
144	0.54	0.21	0.30	118
145	0.63	0.39	0.48	164
146	0.54	0.31	0.39	461
147	0.68	0.45	0.54	159
148	0.30	0.12	0.17	166
149	0.97	0.55	0.70	346
150	0.64	0.13	0.21	350
151	0.93	0.67	0.78	55
152	0.78	0.52	0.63	387
153	0.51	0.17	0.25	150
154	0.58	0.12	0.21	281
155	0.25	0.06	0.10	202
156	0.81	0.67	0.73	130
157	0.28	0.06	0.10	245
158	0.93	0.63	0.75	177
159	0.53	0.34	0.41	130
160	0.48	0.18	0.26	336
161	0.90	0.65	0.75	220
162	0.28	0.06	0.09	229
163	0.87	0.44	0.58	316
164	0.78	0.44	0.56	283
165	0.60	0.34	0.44	197
166	0.65	0.43	0.51	101
167	0.45	0.18	0.26	231

168	0.56	0.27	0.36	370
169	0.40	0.21	0.27	258
170	0.33	0.07	0.11	101
171	0.38	0.24	0.29	89
172	0.53	0.36	0.43	193
173	0.47	0.26	0.33	309
174	0.62	0.14	0.23	172
175	0.92	0.73	0.81	95
176	0.93	0.62	0.74	346
177	0.86	0.57	0.69	322
178	0.65	0.51	0.57	232
179	0.20	0.04	0.07	125
180	0.65	0.33	0.44	145
181	0.44	0.10	0.17	77
182	0.26	0.06	0.10	182
183	0.60	0.32	0.41	257
184	0.21	0.03	0.05	216
185	0.35	0.09	0.14	242
186	0.43	0.18	0.25	165
187	0.75	0.59	0.66	263
188	0.39	0.12	0.18	174
189	0.75	0.40	0.53	136
190	0.89	0.55	0.68	202
191	0.44	0.16	0.24	134
192	0.68	0.40	0.51	230
193	0.44	0.18	0.25	90
194	0.57	0.48	0.52	185
195	0.26	0.05	0.09	156
196	0.33	0.07	0.11	160
197	0.49	0.10	0.16	266
198	0.47	0.13	0.20	284
199	0.32	0.04	0.07	145
200	0.93	0.74	0.82	212
201	0.65	0.26	0.37	317
202	0.78	0.59	0.67	427
203	0.36	0.11	0.17	232
204	0.51	0.29	0.37	217
205	0.50	0.46	0.48	527
206	0.24	0.03	0.06	124
207	0.50	0.17	0.26	103
208	0.85	0.53	0.65	287
209	0.33	0.11	0.16	193
210	0.75	0.38	0.50	220
211	0.72	0.21	0.32	140
212	0.12	0.02	0.03	161
213	0.63	0.43	0.51	72
214	0.64	0.45	0.53	396
215	0.87	0.34	0.49	134
216	0.61	0.17	0.27	400
217	0.51	0.24	0.33	75
218	0.96	0.76	0.85	219
219	0.77	0.42	0.54	210
220	0.88	0.64	0.74	298
221	0.96	0.70	0.81	266
222	0.76	0.45	0.57	290
223	0.11	0.01	0.01	128
224	0.78	0.45	0.57	159
225	0.55	0.29	0.38	164
226	0.58	0.31	0.41	144
227	0.56	0.29	0.38	276
228	0.19	0.03	0.05	235
229	0.33	0.03	0.06	216
230	0.40	0.17	0.23	228
231	0.70	0.48	0.57	64
232	0.48	0.10	0.16	103
233	0.72	0.35	0.47	216
234	0.72	0.11	0.19	116
235	0.54	0.36	0.43	77
236	0.90	0.67	0.77	67
237	0.58	0.13	0.21	218
238	0.40	0.14	0.20	139
239	0.00	0.00	0.00	94
240	0.55	0.35	0.43	77
241	0.47	0.08	0.14	167
242	0.78	0.37	0.50	86
243	0.40	0.10	0.16	58
244	0.62	0.27	0.38	269
245	0.16	0.04	0.07	112
246	0.95	0.76	0.84	255
247	0.44	0.24	0.31	58
248	0.44	0.05	0.09	81
249	0.23	0.02	0.04	131
250	0.43	0.24	0.31	93
251	0.61	0.29	0.39	154
252	0.36	0.04	0.07	129
253	0.69	0.40	0.50	83
254	0.34	0.08	0.13	191
255	0.15	0.03	0.05	219

256	0.32	0.05	0.09	130
257	0.48	0.26	0.34	93
258	0.65	0.48	0.55	217
259	0.41	0.13	0.20	141
260	0.86	0.17	0.29	143
261	0.62	0.17	0.27	219
262	0.55	0.27	0.36	107
263	0.41	0.27	0.32	236
264	0.32	0.22	0.26	119
265	0.57	0.24	0.33	72
266	0.00	0.00	0.00	70
267	0.36	0.14	0.20	107
268	0.67	0.44	0.53	169
269	0.32	0.14	0.19	129
270	0.74	0.53	0.62	159
271	0.88	0.48	0.62	190
272	0.61	0.27	0.37	248
273	0.90	0.75	0.82	264
274	0.90	0.68	0.77	105
275	0.52	0.12	0.20	104
276	0.08	0.01	0.02	115
277	0.83	0.63	0.72	170
278	0.74	0.41	0.52	145
279	0.90	0.70	0.78	230
280	0.58	0.42	0.49	80
281	0.66	0.54	0.59	217
282	0.75	0.50	0.60	175
283	0.33	0.13	0.18	269
284	0.65	0.32	0.43	74
285	0.82	0.49	0.61	206
286	0.89	0.66	0.75	227
287	0.84	0.41	0.55	130
288	0.32	0.07	0.11	129
289	0.57	0.05	0.09	80
290	0.21	0.09	0.13	99
291	0.76	0.35	0.48	208
292	0.42	0.07	0.13	67
293	0.84	0.48	0.61	109
294	0.46	0.26	0.34	140
295	0.24	0.12	0.16	241
296	0.31	0.12	0.18	72
297	0.44	0.11	0.18	107
298	0.77	0.49	0.60	61
299	0.89	0.51	0.64	77
300	0.21	0.08	0.12	111
301	0.00	0.00	0.00	126
302	0.25	0.01	0.03	73
303	0.57	0.43	0.49	176
304	0.91	0.79	0.85	230
305	0.92	0.72	0.81	156
306	0.50	0.37	0.43	146
307	0.34	0.11	0.17	98
308	0.00	0.00	0.00	78
309	0.80	0.13	0.22	94
310	0.74	0.41	0.53	162
311	0.79	0.51	0.62	116
312	0.52	0.28	0.36	57
313	0.83	0.08	0.14	65
314	0.52	0.36	0.42	138
315	0.54	0.22	0.31	195
316	0.56	0.35	0.43	69
317	0.29	0.13	0.18	134
318	0.56	0.39	0.46	148
319	0.84	0.50	0.63	161
320	0.24	0.19	0.21	104
321	0.82	0.61	0.70	156
322	0.60	0.37	0.46	134
323	0.58	0.44	0.50	232
324	0.34	0.15	0.21	92
325	0.41	0.24	0.31	197
326	0.14	0.03	0.05	126
327	0.20	0.03	0.05	115
328	0.99	0.70	0.82	198
329	0.59	0.32	0.41	125
330	0.73	0.20	0.31	81
331	0.45	0.10	0.16	94
332	0.54	0.12	0.20	56
333	0.19	0.05	0.08	260
334	0.42	0.13	0.20	60
335	0.35	0.08	0.13	110
336	0.62	0.49	0.55	71
337	0.18	0.05	0.07	66
338	0.47	0.36	0.41	150
339	0.00	0.00	0.00	54
340	0.84	0.57	0.68	195
341	0.91	0.52	0.66	79
342	0.38	0.26	0.31	38
343	0.62	0.42	0.50	43

344	0.56	0.29	0.38	68
345	0.62	0.33	0.43	73
346	0.14	0.03	0.04	116
347	0.86	0.43	0.57	111
348	0.33	0.11	0.17	63
349	0.84	0.65	0.74	104
350	0.62	0.48	0.54	44
351	0.57	0.30	0.39	40
352	0.93	0.57	0.70	136
353	0.38	0.15	0.21	54
354	0.39	0.09	0.15	134
355	0.64	0.35	0.45	120
356	0.54	0.29	0.38	228
357	0.66	0.36	0.47	269
358	0.62	0.38	0.47	80
359	0.84	0.59	0.69	140
360	0.39	0.18	0.24	125
361	0.90	0.71	0.79	169
362	0.14	0.05	0.08	56
363	0.92	0.73	0.82	154
364	0.46	0.10	0.17	58
365	0.22	0.08	0.12	71
366	1.00	0.69	0.81	54
367	0.31	0.07	0.11	116
368	0.38	0.06	0.10	54
369	0.33	0.03	0.05	71
370	0.00	0.00	0.00	61
371	0.40	0.08	0.14	71
372	0.72	0.44	0.55	52
373	0.78	0.41	0.54	150
374	0.41	0.14	0.21	93
375	0.20	0.04	0.07	67
376	0.00	0.00	0.00	76
377	0.58	0.28	0.38	106
378	0.25	0.02	0.04	86
379	0.50	0.14	0.22	14
380	0.91	0.52	0.67	122
381	0.23	0.07	0.10	104
382	0.46	0.20	0.28	66
383	0.54	0.35	0.42	110
384	0.14	0.01	0.01	155
385	0.69	0.22	0.33	50
386	0.20	0.06	0.10	64
387	0.32	0.08	0.12	93
388	0.53	0.24	0.33	102
389	0.07	0.01	0.02	108
390	0.96	0.68	0.80	178
391	0.49	0.17	0.26	115
392	0.81	0.40	0.54	42
393	0.00	0.00	0.00	134
394	0.22	0.04	0.06	112
395	0.54	0.27	0.36	176
396	0.47	0.13	0.20	125
397	0.74	0.37	0.49	224
398	0.84	0.67	0.74	63
399	0.30	0.05	0.09	59
400	0.51	0.32	0.39	63
401	0.50	0.24	0.33	98
402	0.51	0.19	0.27	162
403	0.38	0.14	0.21	83
404	0.76	0.84	0.80	19
405	0.34	0.11	0.17	92
406	0.69	0.22	0.33	41
407	0.64	0.37	0.47	43
408	0.80	0.46	0.58	160
409	0.20	0.12	0.15	50
410	0.00	0.00	0.00	19
411	0.35	0.11	0.17	175
412	0.28	0.07	0.11	72
413	0.38	0.05	0.09	95
414	0.12	0.02	0.04	97
415	0.33	0.10	0.16	48
416	0.53	0.35	0.42	83
417	0.43	0.07	0.13	40
418	0.48	0.16	0.25	91
419	0.53	0.37	0.43	90
420	0.38	0.27	0.32	37
421	0.04	0.02	0.02	66
422	0.69	0.45	0.55	73
423	0.48	0.25	0.33	56
424	0.94	0.88	0.91	33
425	0.00	0.00	0.00	76
426	0.27	0.05	0.08	81
427	0.98	0.73	0.84	150
428	0.95	0.69	0.80	29
429	0.99	0.93	0.96	389
430	0.63	0.40	0.49	167
431	0.57	0.11	0.18	123

432	0.52	0.31	0.39	39
433	0.33	0.21	0.25	82
434	1.00	0.70	0.82	66
435	0.55	0.38	0.45	93
436	0.56	0.37	0.44	87
437	0.10	0.02	0.04	86
438	0.72	0.53	0.61	104
439	0.54	0.13	0.21	100
440	0.38	0.04	0.06	141
441	0.43	0.33	0.37	110
442	0.37	0.15	0.22	123
443	0.57	0.18	0.28	71
444	0.32	0.06	0.11	109
445	0.45	0.31	0.37	48
446	0.47	0.29	0.36	76
447	0.39	0.18	0.25	38
448	0.67	0.54	0.60	81
449	0.67	0.26	0.37	132
450	0.42	0.27	0.33	81
451	0.89	0.32	0.47	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.84	0.51	0.64	70
455	0.39	0.18	0.25	155
456	0.50	0.21	0.30	43
457	0.54	0.28	0.37	72
458	0.35	0.13	0.19	62
459	0.63	0.25	0.35	69
460	0.00	0.00	0.00	119
461	0.71	0.19	0.30	79
462	0.61	0.23	0.34	47
463	0.39	0.14	0.21	104
464	0.70	0.42	0.52	106
465	0.64	0.22	0.33	64
466	0.55	0.35	0.43	173
467	0.78	0.42	0.55	107
468	0.56	0.26	0.36	126
469	0.20	0.01	0.02	114
470	0.93	0.81	0.87	140
471	0.85	0.42	0.56	79
472	0.40	0.35	0.37	143
473	0.67	0.37	0.47	158
474	0.48	0.10	0.17	138
475	0.00	0.00	0.00	59
476	0.63	0.33	0.43	88
477	0.83	0.65	0.73	176
478	0.95	0.79	0.86	24
479	0.22	0.04	0.07	92
480	0.79	0.50	0.61	100
481	0.51	0.28	0.36	103
482	0.40	0.22	0.28	74
483	0.78	0.63	0.69	105
484	0.20	0.02	0.04	83
485	0.20	0.02	0.04	82
486	0.48	0.15	0.23	71
487	0.45	0.21	0.29	120
488	0.50	0.06	0.10	105
489	0.73	0.37	0.49	87
490	1.00	0.81	0.90	32
491	0.33	0.03	0.05	69
492	0.33	0.02	0.04	49
493	0.11	0.02	0.03	117
494	0.52	0.23	0.32	61
495	0.95	0.79	0.87	344
496	0.32	0.13	0.19	52
497	0.59	0.28	0.38	137
498	0.31	0.10	0.15	98
499	0.48	0.20	0.29	79
micro avg	0.72	0.37	0.49	173812
macro avg	0.56	0.30	0.37	173812
weighted avg	0.67	0.37	0.46	173812
samples avg	0.46	0.35	0.37	173812

Time taken to run this cell : 1:16:19.984150

4.5.5 Applying Linear SVM(L1 Regularization) using SGD Classifier with OneVsRest Classifier

```
In [67]: import warnings
warnings.filterwarnings("ignore")

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

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))
```

Time taken to run this cell : 0:07:25.261763

Accuracy : 0.24735

Hamming loss 0.0026982

Micro-average quality numbers

Precision: 0.8074, Recall: 0.2939, F1-measure: 0.4310

Macro-average quality numbers

Precision: 0.4206, Recall: 0.2167, F1-measure: 0.2635

	precision	recall	f1-score	support
0	0.95	0.67	0.78	5519
1	0.70	0.22	0.34	8190
2	0.84	0.35	0.50	6529
3	0.82	0.42	0.55	3231
4	0.85	0.37	0.52	6430
5	0.82	0.35	0.49	2879
6	0.87	0.52	0.65	5086
7	0.89	0.54	0.67	4533
8	0.62	0.14	0.23	3000
9	0.83	0.51	0.63	2765
10	0.62	0.01	0.02	3051
11	0.78	0.31	0.44	3009
12	0.74	0.22	0.34	2630
13	0.71	0.18	0.28	1426
14	0.90	0.56	0.69	2548
15	0.77	0.16	0.26	2371
16	0.68	0.23	0.35	873
17	0.88	0.60	0.72	2151
18	0.73	0.20	0.32	2204
19	0.68	0.49	0.57	831
20	0.76	0.46	0.57	1860
21	0.00	0.00	0.00	2023
22	0.54	0.01	0.02	1513
23	0.88	0.55	0.68	1207
24	0.58	0.01	0.03	506
25	0.73	0.36	0.48	425
26	0.65	0.39	0.48	793
27	0.68	0.26	0.37	1291
28	0.82	0.35	0.49	1208
29	0.60	0.01	0.01	406
30	0.77	0.17	0.28	504
31	0.00	0.00	0.00	732
32	0.63	0.26	0.37	441
33	0.00	0.00	0.00	1645
34	0.71	0.27	0.39	1058
35	0.82	0.60	0.69	946
36	0.76	0.15	0.24	644
37	0.96	0.80	0.87	136
38	0.66	0.34	0.44	570
39	0.86	0.27	0.41	766
40	0.69	0.16	0.26	1132
41	0.43	0.17	0.25	174
42	0.75	0.56	0.64	210
43	0.79	0.43	0.55	433
44	0.66	0.56	0.60	626
45	0.82	0.19	0.31	852
46	0.77	0.39	0.52	534
47	0.67	0.01	0.01	350
48	0.74	0.58	0.65	496
49	0.78	0.69	0.73	785
50	0.00	0.00	0.00	475
51	0.00	0.00	0.00	305
52	0.00	0.00	0.00	251
53	0.69	0.35	0.46	914
54	0.00	0.00	0.00	728
55	0.00	0.00	0.00	258
56	0.00	0.00	0.00	821
57	0.00	0.00	0.00	541
58	0.80	0.27	0.40	748
59	0.92	0.70	0.79	724
60	0.57	0.01	0.02	660
61	0.89	0.20	0.33	235
62	0.91	0.73	0.81	718
63	0.82	0.66	0.73	468
64	0.51	0.29	0.37	191
65	0.00	0.00	0.00	429
66	0.00	0.00	0.00	415
67	0.75	0.55	0.63	274
68	0.82	0.58	0.68	510
69	0.66	0.47	0.55	466
70	0.00	0.00	0.00	305
71	0.00	0.00	0.00	247
72	0.81	0.53	0.64	401
73	0.96	0.78	0.86	86
74	0.82	0.44	0.57	120
75	0.89	0.74	0.81	129
76	0.00	0.00	0.00	473
77	0.38	0.29	0.33	143
78	0.77	0.55	0.64	347

79	0.75	0.26	0.39	479
80	0.72	0.21	0.32	279
81	0.83	0.24	0.37	461
82	0.00	0.00	0.00	298
83	0.76	0.54	0.63	396
84	0.53	0.05	0.09	184
85	0.77	0.13	0.23	573
86	0.62	0.02	0.05	325
87	0.00	0.00	0.00	273
88	0.00	0.00	0.00	135
89	0.00	0.00	0.00	232
90	0.62	0.02	0.04	409
91	0.00	0.00	0.00	420
92	0.76	0.59	0.66	408
93	0.63	0.51	0.57	241
94	0.00	0.00	0.00	211
95	0.00	0.00	0.00	277
96	0.00	0.00	0.00	410
97	0.87	0.51	0.64	501
98	0.73	0.74	0.73	136
99	0.00	0.00	0.00	239
100	0.00	0.00	0.00	324
101	0.90	0.75	0.82	277
102	0.92	0.71	0.80	613
103	0.00	0.00	0.00	157
104	0.00	0.00	0.00	295
105	0.88	0.32	0.47	334
106	0.96	0.22	0.36	335
107	0.77	0.53	0.63	389
108	0.00	0.00	0.00	251
109	0.52	0.46	0.49	317
110	0.00	0.00	0.00	187
111	0.75	0.13	0.22	140
112	0.66	0.35	0.46	154
113	0.71	0.15	0.25	332
114	0.00	0.00	0.00	323
115	0.00	0.00	0.00	344
116	0.77	0.57	0.65	370
117	0.61	0.13	0.21	313
118	0.78	0.67	0.72	874
119	0.33	0.00	0.01	293
120	0.00	0.00	0.00	200
121	0.73	0.55	0.63	463
122	0.00	0.00	0.00	119
123	0.00	0.00	0.00	256
124	0.88	0.78	0.83	195
125	0.00	0.00	0.00	138
126	0.81	0.47	0.59	376
127	0.00	0.00	0.00	122
128	0.00	0.00	0.00	252
129	0.00	0.00	0.00	144
130	0.00	0.00	0.00	150
131	0.00	0.00	0.00	210
132	0.00	0.00	0.00	361
133	0.93	0.56	0.70	453
134	0.88	0.83	0.85	124
135	0.00	0.00	0.00	91
136	0.84	0.16	0.27	128
137	0.59	0.32	0.41	218
138	0.81	0.07	0.13	243
139	0.00	0.00	0.00	149
140	0.73	0.56	0.64	318
141	0.00	0.00	0.00	159
142	0.66	0.50	0.57	274
143	0.86	0.81	0.83	362
144	0.00	0.00	0.00	118
145	0.63	0.40	0.49	164
146	0.00	0.00	0.00	461
147	0.64	0.48	0.55	159
148	0.00	0.00	0.00	166
149	0.94	0.59	0.73	346
150	0.86	0.03	0.07	350
151	0.83	0.69	0.75	55
152	0.82	0.48	0.61	387
153	0.43	0.04	0.07	150
154	0.00	0.00	0.00	281
155	0.00	0.00	0.00	202
156	0.73	0.69	0.71	130
157	0.00	0.00	0.00	245
158	0.86	0.64	0.74	177
159	0.68	0.32	0.44	130
160	0.00	0.00	0.00	336
161	0.92	0.62	0.74	220
162	0.00	0.00	0.00	229
163	0.87	0.48	0.62	316
164	0.77	0.40	0.52	283
165	0.66	0.26	0.38	197
166	0.73	0.35	0.47	101

167	0.00	0.00	0.00	231
168	0.00	0.00	0.00	370
169	0.00	0.00	0.00	258
170	0.00	0.00	0.00	101
171	0.47	0.25	0.32	89
172	0.33	0.01	0.01	193
173	0.00	0.00	0.00	309
174	0.00	0.00	0.00	172
175	0.91	0.83	0.87	95
176	0.92	0.59	0.72	346
177	0.89	0.52	0.66	322
178	0.64	0.49	0.56	232
179	0.00	0.00	0.00	125
180	0.62	0.39	0.47	145
181	0.00	0.00	0.00	77
182	0.00	0.00	0.00	182
183	0.00	0.00	0.00	257
184	0.00	0.00	0.00	216
185	0.00	0.00	0.00	242
186	0.00	0.00	0.00	165
187	0.75	0.64	0.69	263
188	0.00	0.00	0.00	174
189	0.79	0.14	0.24	136
190	0.93	0.63	0.75	202
191	0.00	0.00	0.00	134
192	0.73	0.48	0.58	230
193	0.00	0.00	0.00	90
194	0.58	0.52	0.55	185
195	0.00	0.00	0.00	156
196	0.00	0.00	0.00	160
197	0.64	0.10	0.18	266
198	0.00	0.00	0.00	284
199	0.00	0.00	0.00	145
200	0.92	0.77	0.84	212
201	0.71	0.21	0.32	317
202	0.77	0.59	0.67	427
203	0.00	0.00	0.00	232
204	0.00	0.00	0.00	217
205	0.00	0.00	0.00	527
206	0.00	0.00	0.00	124
207	0.40	0.02	0.04	103
208	0.88	0.52	0.65	287
209	0.00	0.00	0.00	193
210	0.80	0.15	0.25	220
211	0.76	0.23	0.35	140
212	0.00	0.00	0.00	161
213	0.52	0.19	0.28	72
214	0.60	0.02	0.04	396
215	0.84	0.40	0.55	134
216	0.50	0.01	0.03	400
217	0.75	0.04	0.08	75
218	0.94	0.78	0.86	219
219	0.86	0.29	0.43	210
220	0.88	0.64	0.74	298
221	0.89	0.74	0.81	266
222	0.77	0.39	0.52	290
223	0.00	0.00	0.00	128
224	0.79	0.47	0.58	159
225	0.68	0.27	0.39	164
226	0.61	0.43	0.51	144
227	0.76	0.06	0.11	276
228	0.00	0.00	0.00	235
229	0.00	0.00	0.00	216
230	0.00	0.00	0.00	228
231	0.74	0.62	0.68	64
232	0.00	0.00	0.00	103
233	0.72	0.36	0.48	216
234	0.00	0.00	0.00	116
235	0.59	0.51	0.55	77
236	0.91	0.73	0.81	67
237	0.75	0.10	0.17	218
238	0.00	0.00	0.00	139
239	0.00	0.00	0.00	94
240	0.66	0.35	0.46	77
241	0.00	0.00	0.00	167
242	0.82	0.33	0.47	86
243	0.00	0.00	0.00	58
244	0.78	0.13	0.23	269
245	0.00	0.00	0.00	112
246	0.94	0.79	0.86	255
247	0.50	0.31	0.38	58
248	0.00	0.00	0.00	81
249	0.00	0.00	0.00	131
250	0.00	0.00	0.00	93
251	0.00	0.00	0.00	154
252	0.00	0.00	0.00	129
253	0.67	0.27	0.38	83
254	0.00	0.00	0.00	191

255	0.00	0.00	0.00	219
256	0.00	0.00	0.00	130
257	1.00	0.01	0.02	93
258	0.66	0.54	0.59	217
259	0.00	0.00	0.00	141
260	0.73	0.15	0.25	143
261	0.00	0.00	0.00	219
262	1.00	0.02	0.04	107
263	0.00	0.00	0.00	236
264	0.33	0.02	0.03	119
265	0.60	0.04	0.08	72
266	0.00	0.00	0.00	70
267	0.00	0.00	0.00	107
268	0.68	0.50	0.58	169
269	0.00	0.00	0.00	129
270	0.74	0.67	0.70	159
271	0.90	0.49	0.64	190
272	0.72	0.09	0.16	248
273	0.90	0.76	0.82	264
274	0.90	0.71	0.80	105
275	0.00	0.00	0.00	104
276	0.00	0.00	0.00	115
277	0.83	0.69	0.75	170
278	0.73	0.33	0.45	145
279	0.91	0.70	0.79	230
280	0.54	0.39	0.45	80
281	0.66	0.74	0.70	217
282	0.75	0.62	0.68	175
283	0.00	0.00	0.00	269
284	0.64	0.38	0.47	74
285	0.85	0.51	0.64	206
286	0.88	0.64	0.74	227
287	0.90	0.34	0.49	130
288	0.00	0.00	0.00	129
289	0.00	0.00	0.00	80
290	0.00	0.00	0.00	99
291	0.76	0.28	0.41	208
292	0.00	0.00	0.00	67
293	0.95	0.38	0.54	109
294	0.00	0.00	0.00	140
295	0.00	0.00	0.00	241
296	0.00	0.00	0.00	72
297	0.00	0.00	0.00	107
298	0.76	0.31	0.44	61
299	0.81	0.34	0.48	77
300	0.00	0.00	0.00	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.57	0.35	0.44	176
304	0.94	0.80	0.86	230
305	0.91	0.76	0.83	156
306	1.00	0.01	0.01	146
307	0.00	0.00	0.00	98
308	0.00	0.00	0.00	78
309	0.59	0.14	0.22	94
310	0.78	0.13	0.22	162
311	0.76	0.58	0.66	116
312	0.50	0.39	0.44	57
313	0.00	0.00	0.00	65
314	0.46	0.13	0.20	138
315	0.00	0.00	0.00	195
316	1.00	0.01	0.03	69
317	0.62	0.07	0.13	134
318	1.00	0.03	0.05	148
319	0.84	0.55	0.66	161
320	0.00	0.00	0.00	104
321	0.82	0.67	0.74	156
322	0.62	0.10	0.17	134
323	0.29	0.01	0.02	232
324	0.00	0.00	0.00	92
325	0.00	0.00	0.00	197
326	0.00	0.00	0.00	126
327	0.00	0.00	0.00	115
328	0.96	0.73	0.83	198
329	0.62	0.32	0.42	125
330	0.67	0.02	0.05	81
331	0.00	0.00	0.00	94
332	0.00	0.00	0.00	56
333	0.00	0.00	0.00	260
334	0.00	0.00	0.00	60
335	0.00	0.00	0.00	110
336	0.58	0.51	0.54	71
337	0.00	0.00	0.00	66
338	0.59	0.09	0.15	150
339	0.00	0.00	0.00	54
340	0.85	0.63	0.72	195
341	0.77	0.22	0.34	79
342	1.00	0.05	0.10	38

343	0.64	0.49	0.55	43
344	0.00	0.00	0.00	68
345	0.50	0.01	0.03	73
346	0.00	0.00	0.00	116
347	0.84	0.47	0.60	111
348	0.00	0.00	0.00	63
349	0.83	0.70	0.76	104
350	0.63	0.59	0.61	44
351	0.70	0.17	0.28	40
352	0.90	0.56	0.69	136
353	0.00	0.00	0.00	54
354	0.00	0.00	0.00	134
355	0.73	0.30	0.43	120
356	0.00	0.00	0.00	228
357	0.00	0.00	0.00	269
358	0.79	0.14	0.23	80
359	0.84	0.62	0.72	140
360	0.00	0.00	0.00	125
361	0.89	0.73	0.81	169
362	0.00	0.00	0.00	56
363	0.93	0.73	0.82	154
364	0.00	0.00	0.00	58
365	0.00	0.00	0.00	71
366	1.00	0.70	0.83	54
367	0.00	0.00	0.00	116
368	0.00	0.00	0.00	54
369	0.00	0.00	0.00	71
370	0.00	0.00	0.00	61
371	0.00	0.00	0.00	71
372	0.61	0.54	0.57	52
373	0.81	0.45	0.58	150
374	0.00	0.00	0.00	93
375	0.00	0.00	0.00	67
376	0.00	0.00	0.00	76
377	0.00	0.00	0.00	106
378	0.00	0.00	0.00	86
379	0.75	0.21	0.33	14
380	0.85	0.61	0.71	122
381	0.00	0.00	0.00	104
382	0.00	0.00	0.00	66
383	1.00	0.05	0.09	110
384	0.00	0.00	0.00	155
385	0.36	0.08	0.13	50
386	0.00	0.00	0.00	64
387	0.00	0.00	0.00	93
388	1.00	0.01	0.02	102
389	0.00	0.00	0.00	108
390	0.94	0.70	0.80	178
391	0.00	0.00	0.00	115
392	0.86	0.57	0.69	42
393	0.00	0.00	0.00	134
394	0.00	0.00	0.00	112
395	0.00	0.00	0.00	176
396	0.00	0.00	0.00	125
397	0.67	0.26	0.37	224
398	0.77	0.70	0.73	63
399	0.00	0.00	0.00	59
400	0.40	0.03	0.06	63
401	0.00	0.00	0.00	98
402	0.00	0.00	0.00	162
403	0.00	0.00	0.00	83
404	0.67	0.84	0.74	19
405	0.00	0.00	0.00	92
406	0.70	0.17	0.27	41
407	0.73	0.26	0.38	43
408	0.76	0.33	0.46	160
409	0.00	0.00	0.00	50
410	0.00	0.00	0.00	19
411	0.00	0.00	0.00	175
412	0.00	0.00	0.00	72
413	0.00	0.00	0.00	95
414	0.00	0.00	0.00	97
415	0.00	0.00	0.00	48
416	0.38	0.06	0.10	83
417	0.00	0.00	0.00	40
418	0.00	0.00	0.00	91
419	0.48	0.16	0.24	90
420	0.00	0.00	0.00	37
421	0.00	0.00	0.00	66
422	0.00	0.00	0.00	73
423	0.47	0.32	0.38	56
424	0.91	0.88	0.89	33
425	0.00	0.00	0.00	76
426	0.00	0.00	0.00	81
427	0.97	0.75	0.85	150
428	0.91	0.72	0.81	29
429	0.99	0.92	0.96	389
430	0.67	0.37	0.47	167

431	0.00	0.00	0.00	123
432	0.00	0.00	0.00	39
433	0.00	0.00	0.00	82
434	0.96	0.73	0.83	66
435	0.59	0.39	0.47	93
436	0.62	0.23	0.34	87
437	0.00	0.00	0.00	86
438	0.74	0.52	0.61	104
439	0.00	0.00	0.00	100
440	0.00	0.00	0.00	141
441	0.00	0.00	0.00	110
442	0.00	0.00	0.00	123
443	0.25	0.03	0.05	71
444	0.00	0.00	0.00	109
445	0.69	0.23	0.34	48
446	0.21	0.04	0.07	76
447	0.00	0.00	0.00	38
448	0.66	0.70	0.68	81
449	0.00	0.00	0.00	132
450	0.00	0.00	0.00	81
451	0.86	0.32	0.46	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.81	0.49	0.61	70
455	0.00	0.00	0.00	155
456	0.67	0.05	0.09	43
457	0.00	0.00	0.00	72
458	0.00	0.00	0.00	62
459	1.00	0.10	0.18	69
460	0.00	0.00	0.00	119
461	0.00	0.00	0.00	79
462	1.00	0.06	0.12	47
463	1.00	0.12	0.22	104
464	0.00	0.00	0.00	106
465	0.00	0.00	0.00	64
466	0.00	0.00	0.00	173
467	0.72	0.36	0.47	107
468	0.85	0.09	0.16	126
469	0.00	0.00	0.00	114
470	0.91	0.81	0.86	140
471	0.94	0.39	0.55	79
472	0.36	0.13	0.19	143
473	0.72	0.37	0.49	158
474	0.00	0.00	0.00	138
475	0.00	0.00	0.00	59
476	0.65	0.27	0.38	88
477	0.84	0.69	0.76	176
478	0.95	0.75	0.84	24
479	0.00	0.00	0.00	92
480	0.79	0.57	0.66	100
481	0.53	0.08	0.14	103
482	0.00	0.00	0.00	74
483	0.81	0.63	0.71	105
484	0.00	0.00	0.00	83
485	0.00	0.00	0.00	82
486	0.00	0.00	0.00	71
487	0.00	0.00	0.00	120
488	0.67	0.08	0.14	105
489	0.77	0.38	0.51	87
490	1.00	0.81	0.90	32
491	0.00	0.00	0.00	69
492	0.00	0.00	0.00	49
493	0.00	0.00	0.00	117
494	0.47	0.30	0.36	61
495	0.97	0.75	0.85	344
496	1.00	0.04	0.07	52
497	0.67	0.12	0.20	137
498	0.00	0.00	0.00	98
499	0.58	0.18	0.27	79
micro avg	0.81	0.29	0.43	173812
macro avg	0.42	0.22	0.26	173812
weighted avg	0.60	0.29	0.37	173812
samples avg	0.41	0.28	0.32	173812

4.5.6 Applying Linear SVM(L2 Regularization) using SGD Classifier with OneVsRest Classifier

```
In [68]: import warnings
warnings.filterwarnings("ignore")

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.00001, penalty='l2'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict (x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:05:19.925545

Accuracy : 0.24899

Hamming loss 0.0026806

Micro-average quality numbers

Precision: 0.8247, Recall: 0.2907, F1-measure: 0.4298

Macro-average quality numbers

Precision: 0.4562, Recall: 0.2073, F1-measure: 0.2587

	precision	recall	f1-score	support
0	0.95	0.69	0.80	5519
1	0.74	0.24	0.36	8190
2	0.86	0.35	0.50	6529
3	0.84	0.43	0.57	3231
4	0.87	0.37	0.52	6430
5	0.84	0.36	0.51	2879
6	0.90	0.49	0.64	5086
7	0.90	0.54	0.67	4533
8	0.62	0.14	0.23	3000
9	0.83	0.56	0.67	2765
10	0.79	0.01	0.03	3051
11	0.80	0.31	0.44	3009
12	0.75	0.22	0.34	2630
13	0.81	0.22	0.35	1426
14	0.92	0.53	0.67	2548
15	0.85	0.12	0.22	2371
16	0.69	0.26	0.38	873
17	0.90	0.59	0.72	2151
18	0.73	0.19	0.31	2204
19	0.69	0.47	0.56	831
20	0.77	0.47	0.59	1860
21	0.00	0.00	0.00	2023
22	0.68	0.11	0.19	1513
23	0.91	0.52	0.66	1207
24	0.72	0.09	0.17	506
25	0.73	0.33	0.45	425
26	0.66	0.42	0.52	793
27	0.70	0.23	0.35	1291
28	0.84	0.35	0.49	1208
29	0.50	0.01	0.01	406
30	0.79	0.16	0.26	504
31	0.00	0.00	0.00	732
32	0.68	0.27	0.39	441
33	1.00	0.00	0.01	1645
34	0.73	0.26	0.38	1058
35	0.83	0.61	0.70	946
36	0.73	0.16	0.26	644
37	0.96	0.75	0.84	136
38	0.65	0.38	0.48	570
39	0.90	0.25	0.39	766
40	0.72	0.18	0.29	1132
41	0.52	0.22	0.31	174
42	0.79	0.55	0.65	210
43	0.83	0.41	0.55	433
44	0.69	0.54	0.60	626
45	0.83	0.24	0.37	852
46	0.81	0.40	0.54	534
47	1.00	0.00	0.01	350
48	0.74	0.57	0.64	496
49	0.78	0.70	0.74	785
50	0.00	0.00	0.00	475
51	0.00	0.00	0.00	305
52	0.00	0.00	0.00	251
53	0.70	0.34	0.46	914
54	0.00	0.00	0.00	728
55	0.00	0.00	0.00	258
56	0.00	0.00	0.00	821
57	1.00	0.00	0.00	541
58	0.81	0.26	0.39	748
59	0.93	0.64	0.76	724
60	0.50	0.00	0.00	660
61	0.90	0.20	0.32	235
62	0.92	0.73	0.81	718
63	0.83	0.66	0.73	468
64	0.53	0.41	0.46	191
65	0.00	0.00	0.00	429
66	0.00	0.00	0.00	415
67	0.77	0.54	0.63	274
68	0.84	0.57	0.68	510
69	0.67	0.49	0.57	466
70	0.00	0.00	0.00	305
71	0.00	0.00	0.00	247
72	0.80	0.51	0.63	401
73	0.99	0.77	0.86	86
74	0.84	0.41	0.55	120
75	0.90	0.74	0.81	129
76	0.00	0.00	0.00	473
77	0.44	0.24	0.31	143
78	0.78	0.52	0.62	347

79	0.77	0.25	0.37	479
80	0.70	0.23	0.35	279
81	0.85	0.19	0.31	461
82	0.00	0.00	0.00	298
83	0.80	0.48	0.60	396
84	0.60	0.18	0.28	184
85	0.78	0.06	0.12	573
86	0.78	0.02	0.04	325
87	0.71	0.04	0.08	273
88	0.00	0.00	0.00	135
89	0.00	0.00	0.00	232
90	0.69	0.06	0.11	409
91	0.00	0.00	0.00	420
92	0.77	0.56	0.65	408
93	0.68	0.50	0.58	241
94	0.00	0.00	0.00	211
95	0.00	0.00	0.00	277
96	0.00	0.00	0.00	410
97	0.89	0.43	0.58	501
98	0.75	0.71	0.73	136
99	0.69	0.08	0.15	239
100	0.00	0.00	0.00	324
101	0.93	0.65	0.76	277
102	0.93	0.72	0.81	613
103	0.14	0.01	0.01	157
104	0.00	0.00	0.00	295
105	0.89	0.34	0.49	334
106	0.97	0.17	0.28	335
107	0.79	0.50	0.61	389
108	0.00	0.00	0.00	251
109	0.57	0.43	0.49	317
110	0.00	0.00	0.00	187
111	0.70	0.10	0.17	140
112	0.68	0.32	0.44	154
113	0.67	0.15	0.24	332
114	0.00	0.00	0.00	323
115	1.00	0.00	0.01	344
116	0.75	0.57	0.65	370
117	0.61	0.12	0.20	313
118	0.79	0.61	0.69	874
119	0.75	0.01	0.02	293
120	0.00	0.00	0.00	200
121	0.75	0.54	0.63	463
122	0.00	0.00	0.00	119
123	0.00	0.00	0.00	256
124	0.89	0.78	0.83	195
125	0.00	0.00	0.00	138
126	0.83	0.49	0.61	376
127	0.00	0.00	0.00	122
128	0.00	0.00	0.00	252
129	0.00	0.00	0.00	144
130	1.00	0.01	0.01	150
131	0.00	0.00	0.00	210
132	1.00	0.00	0.01	361
133	0.95	0.55	0.70	453
134	0.87	0.81	0.84	124
135	0.00	0.00	0.00	91
136	0.95	0.15	0.26	128
137	0.62	0.30	0.41	218
138	1.00	0.01	0.02	243
139	0.00	0.00	0.00	149
140	0.75	0.52	0.61	318
141	0.00	0.00	0.00	159
142	0.65	0.48	0.55	274
143	0.86	0.79	0.83	362
144	0.50	0.01	0.02	118
145	0.65	0.35	0.45	164
146	0.71	0.09	0.16	461
147	0.69	0.58	0.63	159
148	0.00	0.00	0.00	166
149	0.98	0.47	0.63	346
150	1.00	0.04	0.07	350
151	0.88	0.69	0.78	55
152	0.83	0.47	0.60	387
153	0.57	0.03	0.05	150
154	0.00	0.00	0.00	281
155	0.00	0.00	0.00	202
156	0.80	0.69	0.74	130
157	0.00	0.00	0.00	245
158	0.89	0.66	0.75	177
159	0.68	0.32	0.44	130
160	0.00	0.00	0.00	336
161	0.91	0.63	0.74	220
162	0.00	0.00	0.00	229
163	0.89	0.43	0.58	316
164	0.79	0.34	0.47	283
165	0.64	0.30	0.41	197
166	0.72	0.49	0.58	101

167	0.00	0.00	0.00	231
168	0.83	0.01	0.03	370
169	0.00	0.00	0.00	258
170	0.00	0.00	0.00	101
171	0.48	0.26	0.34	89
172	0.00	0.00	0.00	193
173	0.25	0.00	0.01	309
174	0.00	0.00	0.00	172
175	0.92	0.84	0.88	95
176	0.94	0.58	0.71	346
177	0.93	0.46	0.62	322
178	0.66	0.46	0.54	232
179	0.00	0.00	0.00	125
180	0.54	0.36	0.43	145
181	0.00	0.00	0.00	77
182	0.00	0.00	0.00	182
183	0.00	0.00	0.00	257
184	0.00	0.00	0.00	216
185	0.00	0.00	0.00	242
186	0.00	0.00	0.00	165
187	0.75	0.59	0.66	263
188	0.00	0.00	0.00	174
189	0.90	0.19	0.32	136
190	0.95	0.50	0.66	202
191	0.00	0.00	0.00	134
192	0.74	0.45	0.56	230
193	0.00	0.00	0.00	90
194	0.60	0.44	0.50	185
195	0.00	0.00	0.00	156
196	0.00	0.00	0.00	160
197	0.61	0.06	0.12	266
198	0.00	0.00	0.00	284
199	0.00	0.00	0.00	145
200	0.93	0.70	0.80	212
201	0.67	0.10	0.17	317
202	0.80	0.57	0.67	427
203	0.00	0.00	0.00	232
204	0.00	0.00	0.00	217
205	0.50	0.00	0.00	527
206	0.00	0.00	0.00	124
207	1.00	0.02	0.04	103
208	0.90	0.50	0.64	287
209	0.00	0.00	0.00	193
210	0.76	0.25	0.38	220
211	0.83	0.17	0.28	140
212	0.00	0.00	0.00	161
213	0.61	0.31	0.41	72
214	0.67	0.01	0.01	396
215	0.84	0.38	0.52	134
216	0.00	0.00	0.00	400
217	0.40	0.05	0.09	75
218	0.97	0.76	0.85	219
219	0.85	0.27	0.41	210
220	0.91	0.60	0.73	298
221	0.97	0.65	0.78	266
222	0.77	0.38	0.51	290
223	0.00	0.00	0.00	128
224	0.78	0.46	0.58	159
225	0.73	0.26	0.39	164
226	0.64	0.38	0.48	144
227	0.73	0.03	0.06	276
228	0.00	0.00	0.00	235
229	0.00	0.00	0.00	216
230	0.00	0.00	0.00	228
231	0.75	0.61	0.67	64
232	0.00	0.00	0.00	103
233	0.73	0.31	0.44	216
234	0.00	0.00	0.00	116
235	0.56	0.42	0.48	77
236	0.96	0.66	0.78	67
237	0.79	0.07	0.13	218
238	0.00	0.00	0.00	139
239	0.00	0.00	0.00	94
240	0.67	0.21	0.32	77
241	0.00	0.00	0.00	167
242	0.87	0.31	0.46	86
243	0.00	0.00	0.00	58
244	0.79	0.10	0.18	269
245	0.00	0.00	0.00	112
246	0.94	0.78	0.85	255
247	0.48	0.24	0.32	58
248	0.00	0.00	0.00	81
249	0.00	0.00	0.00	131
250	0.50	0.02	0.04	93
251	0.83	0.06	0.12	154
252	0.00	0.00	0.00	129
253	0.74	0.28	0.40	83
254	0.00	0.00	0.00	191

255	0.00	0.00	0.00	219
256	0.00	0.00	0.00	130
257	1.00	0.03	0.06	93
258	0.66	0.52	0.58	217
259	0.00	0.00	0.00	141
260	0.90	0.13	0.23	143
261	0.00	0.00	0.00	219
262	0.60	0.03	0.05	107
263	0.00	0.00	0.00	236
264	0.00	0.00	0.00	119
265	0.50	0.03	0.05	72
266	0.00	0.00	0.00	70
267	0.00	0.00	0.00	107
268	0.69	0.46	0.55	169
269	0.00	0.00	0.00	129
270	0.73	0.67	0.70	159
271	0.89	0.40	0.55	190
272	0.78	0.12	0.22	248
273	0.91	0.73	0.81	264
274	0.89	0.71	0.79	105
275	0.00	0.00	0.00	104
276	0.00	0.00	0.00	115
277	0.84	0.68	0.75	170
278	0.75	0.25	0.37	145
279	0.92	0.66	0.77	230
280	0.57	0.34	0.43	80
281	0.67	0.70	0.68	217
282	0.76	0.59	0.66	175
283	0.00	0.00	0.00	269
284	0.58	0.26	0.36	74
285	0.84	0.52	0.65	206
286	0.89	0.63	0.73	227
287	0.93	0.32	0.47	130
288	0.00	0.00	0.00	129
289	0.00	0.00	0.00	80
290	0.00	0.00	0.00	99
291	0.79	0.26	0.39	208
292	0.00	0.00	0.00	67
293	0.96	0.40	0.57	109
294	0.00	0.00	0.00	140
295	0.00	0.00	0.00	241
296	0.00	0.00	0.00	72
297	0.00	0.00	0.00	107
298	0.85	0.38	0.52	61
299	0.92	0.44	0.60	77
300	0.00	0.00	0.00	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.68	0.31	0.42	176
304	0.96	0.76	0.84	230
305	0.96	0.62	0.75	156
306	1.00	0.01	0.01	146
307	0.00	0.00	0.00	98
308	0.00	0.00	0.00	78
309	0.78	0.07	0.14	94
310	0.79	0.14	0.23	162
311	0.77	0.57	0.65	116
312	0.51	0.33	0.40	57
313	0.00	0.00	0.00	65
314	0.51	0.22	0.30	138
315	0.00	0.00	0.00	195
316	1.00	0.06	0.11	69
317	0.58	0.05	0.10	134
318	0.00	0.00	0.00	148
319	0.84	0.47	0.60	161
320	0.00	0.00	0.00	104
321	0.81	0.64	0.72	156
322	0.72	0.25	0.37	134
323	0.61	0.11	0.18	232
324	0.00	0.00	0.00	92
325	0.00	0.00	0.00	197
326	0.00	0.00	0.00	126
327	0.00	0.00	0.00	115
328	0.99	0.67	0.80	198
329	0.65	0.24	0.35	125
330	0.67	0.07	0.13	81
331	0.00	0.00	0.00	94
332	0.00	0.00	0.00	56
333	0.00	0.00	0.00	260
334	0.00	0.00	0.00	60
335	0.00	0.00	0.00	110
336	0.63	0.51	0.56	71
337	0.00	0.00	0.00	66
338	0.80	0.05	0.10	150
339	0.00	0.00	0.00	54
340	0.84	0.58	0.69	195
341	0.93	0.18	0.30	79
342	1.00	0.03	0.05	38

343	0.63	0.44	0.52	43
344	0.75	0.04	0.08	68
345	0.61	0.15	0.24	73
346	0.00	0.00	0.00	116
347	0.87	0.36	0.51	111
348	0.00	0.00	0.00	63
349	0.83	0.69	0.75	104
350	0.62	0.59	0.60	44
351	0.75	0.15	0.25	40
352	0.96	0.48	0.64	136
353	0.00	0.00	0.00	54
354	0.00	0.00	0.00	134
355	0.80	0.28	0.41	120
356	0.00	0.00	0.00	228
357	0.75	0.02	0.04	269
358	0.80	0.15	0.25	80
359	0.85	0.54	0.66	140
360	0.00	0.00	0.00	125
361	0.90	0.71	0.79	169
362	0.00	0.00	0.00	56
363	0.93	0.68	0.78	154
364	0.00	0.00	0.00	58
365	0.00	0.00	0.00	71
366	1.00	0.70	0.83	54
367	0.00	0.00	0.00	116
368	0.00	0.00	0.00	54
369	0.00	0.00	0.00	71
370	0.00	0.00	0.00	61
371	0.00	0.00	0.00	71
372	0.67	0.50	0.57	52
373	0.83	0.39	0.53	150
374	0.00	0.00	0.00	93
375	0.00	0.00	0.00	67
376	0.00	0.00	0.00	76
377	0.00	0.00	0.00	106
378	0.00	0.00	0.00	86
379	0.50	0.07	0.12	14
380	1.00	0.45	0.62	122
381	0.00	0.00	0.00	104
382	0.00	0.00	0.00	66
383	0.92	0.10	0.18	110
384	0.00	0.00	0.00	155
385	0.44	0.08	0.14	50
386	0.00	0.00	0.00	64
387	0.00	0.00	0.00	93
388	1.00	0.02	0.04	102
389	0.00	0.00	0.00	108
390	0.95	0.68	0.79	178
391	0.00	0.00	0.00	115
392	0.89	0.40	0.56	42
393	0.00	0.00	0.00	134
394	0.00	0.00	0.00	112
395	0.00	0.00	0.00	176
396	0.00	0.00	0.00	125
397	0.80	0.33	0.46	224
398	0.85	0.70	0.77	63
399	0.00	0.00	0.00	59
400	0.50	0.03	0.06	63
401	0.00	0.00	0.00	98
402	0.00	0.00	0.00	162
403	0.00	0.00	0.00	83
404	0.67	0.84	0.74	19
405	0.00	0.00	0.00	92
406	0.78	0.17	0.28	41
407	0.76	0.30	0.43	43
408	0.79	0.28	0.41	160
409	0.00	0.00	0.00	50
410	0.00	0.00	0.00	19
411	0.00	0.00	0.00	175
412	0.00	0.00	0.00	72
413	0.00	0.00	0.00	95
414	0.00	0.00	0.00	97
415	0.00	0.00	0.00	48
416	0.50	0.04	0.07	83
417	0.00	0.00	0.00	40
418	0.00	0.00	0.00	91
419	0.56	0.10	0.17	90
420	0.00	0.00	0.00	37
421	0.00	0.00	0.00	66
422	0.00	0.00	0.00	73
423	0.47	0.25	0.33	56
424	0.90	0.85	0.88	33
425	0.00	0.00	0.00	76
426	0.00	0.00	0.00	81
427	0.98	0.72	0.83	150
428	0.95	0.69	0.80	29
429	0.99	0.88	0.93	389
430	0.83	0.26	0.39	167

431	0.00	0.00	0.00	123
432	0.67	0.05	0.10	39
433	0.00	0.00	0.00	82
434	1.00	0.68	0.81	66
435	0.65	0.30	0.41	93
436	0.61	0.25	0.36	87
437	0.00	0.00	0.00	86
438	0.70	0.41	0.52	104
439	0.00	0.00	0.00	100
440	0.00	0.00	0.00	141
441	0.00	0.00	0.00	110
442	0.00	0.00	0.00	123
443	0.25	0.01	0.03	71
444	0.00	0.00	0.00	109
445	0.75	0.06	0.12	48
446	0.17	0.03	0.05	76
447	0.00	0.00	0.00	38
448	0.66	0.74	0.70	81
449	1.00	0.01	0.02	132
450	0.00	0.00	0.00	81
451	0.88	0.28	0.42	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.92	0.49	0.64	70
455	0.00	0.00	0.00	155
456	1.00	0.02	0.05	43
457	1.00	0.01	0.03	72
458	0.00	0.00	0.00	62
459	1.00	0.12	0.21	69
460	0.00	0.00	0.00	119
461	0.00	0.00	0.00	79
462	0.50	0.09	0.15	47
463	0.00	0.00	0.00	104
464	0.00	0.00	0.00	106
465	0.00	0.00	0.00	64
466	0.00	0.00	0.00	173
467	0.82	0.35	0.49	107
468	0.83	0.04	0.08	126
469	0.00	0.00	0.00	114
470	0.92	0.81	0.86	140
471	1.00	0.37	0.54	79
472	0.20	0.01	0.01	143
473	0.77	0.28	0.41	158
474	0.00	0.00	0.00	138
475	0.00	0.00	0.00	59
476	0.65	0.19	0.30	88
477	0.85	0.67	0.75	176
478	0.90	0.79	0.84	24
479	0.00	0.00	0.00	92
480	0.80	0.53	0.64	100
481	0.50	0.05	0.09	103
482	0.00	0.00	0.00	74
483	0.86	0.56	0.68	105
484	0.00	0.00	0.00	83
485	0.00	0.00	0.00	82
486	0.00	0.00	0.00	71
487	0.00	0.00	0.00	120
488	0.40	0.02	0.04	105
489	0.79	0.26	0.40	87
490	1.00	0.81	0.90	32
491	0.00	0.00	0.00	69
492	0.00	0.00	0.00	49
493	0.00	0.00	0.00	117
494	0.59	0.21	0.31	61
495	0.97	0.75	0.85	344
496	0.00	0.00	0.00	52
497	0.60	0.02	0.04	137
498	0.00	0.00	0.00	98
499	0.68	0.16	0.27	79
micro avg	0.82	0.29	0.43	173812
macro avg	0.46	0.21	0.26	173812
weighted avg	0.65	0.29	0.38	173812
samples avg	0.41	0.28	0.32	173812

4.5.7 Logistic Regression Hyper Parameter Tuning (1 to 3grams)

```
In [71]: from sklearn.model_selection import GridSearchCV

start = datetime.now()

classifier = OneVsRestClassifier(SGDClassifier(loss='log', penalty='l1'), n_jobs=-1)

parameters = {'estimator__alpha':[2*10**-7, 4*10**-7, 6*10**-7, 8*10**-7]}

clf = GridSearchCV(classifier, parameters, cv= 3, return_train_score=True, scoring='f1_micro', n_jobs=-1)

clf.fit(x_train_multilabel, y_train)

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

Time taken to run this cell : 1:53:21.364263

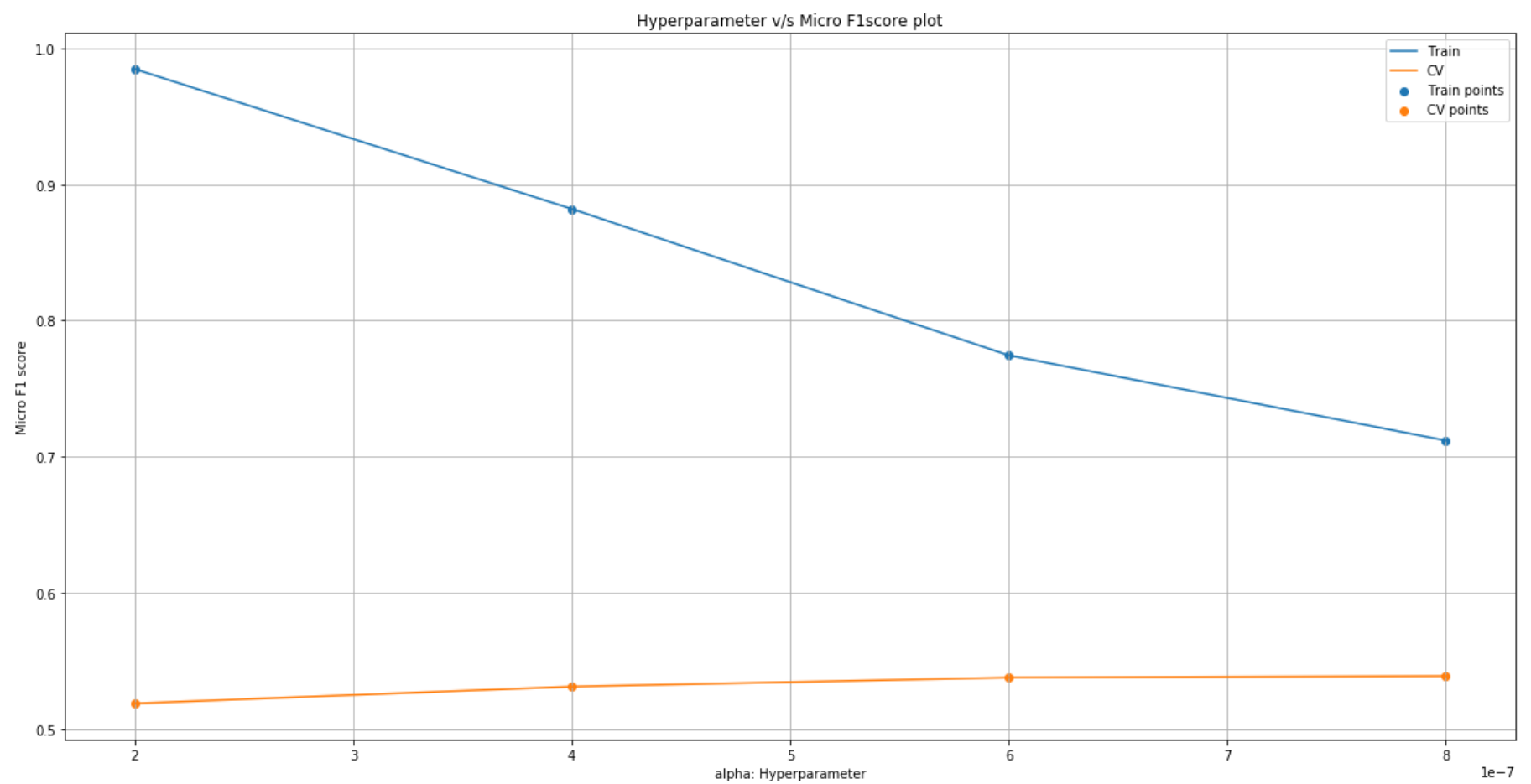
```
In [72]: plt.figure(figsize=(20,10))

train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']

plt.plot(parameters['estimator__alpha'], train_auc, label='Train')
plt.plot(parameters['estimator__alpha'], cv_auc, label='CV')

plt.scatter(parameters['estimator__alpha'], train_auc, label='Train points')
plt.scatter(parameters['estimator__alpha'], cv_auc, label='CV points')

plt.legend()
plt.xlabel("alpha: Hyperparameter")
plt.ylabel("Micro F1 score")
plt.title("Hyperparameter v/s Micro F1score plot")
plt.grid()
plt.show()
```



```
In [73]: import warnings
warnings.filterwarnings("ignore")

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=8*10**-7, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:11:15.409339

Hamming loss 0.00272156

Micro-average quality numbers

Precision: 0.6923, Recall: 0.3908, F1-measure: 0.4996

Macro-average quality numbers

Precision: 0.5528, Recall: 0.3126, F1-measure: 0.3855

	precision	recall	f1-score	support
0	0.94	0.76	0.84	5519
1	0.64	0.39	0.48	8190
2	0.77	0.44	0.56	6529
3	0.79	0.53	0.63	3231
4	0.75	0.48	0.58	6430
5	0.80	0.41	0.54	2879
6	0.84	0.56	0.67	5086
7	0.85	0.60	0.70	4533
8	0.56	0.14	0.22	3000
9	0.79	0.60	0.68	2765
10	0.56	0.25	0.35	3051
11	0.66	0.40	0.50	3009
12	0.58	0.33	0.42	2630
13	0.70	0.35	0.47	1426
14	0.89	0.59	0.71	2548
15	0.60	0.27	0.37	2371
16	0.66	0.29	0.40	873
17	0.87	0.63	0.73	2151
18	0.56	0.26	0.35	2204
19	0.70	0.42	0.53	831
20	0.75	0.50	0.60	1860
21	0.31	0.14	0.19	2023
22	0.51	0.26	0.34	1513
23	0.88	0.59	0.70	1207
24	0.54	0.29	0.38	506
25	0.63	0.36	0.46	425
26	0.64	0.41	0.50	793
27	0.61	0.41	0.49	1291
28	0.72	0.42	0.53	1208
29	0.41	0.12	0.18	406
30	0.67	0.21	0.32	504
31	0.26	0.08	0.13	732
32	0.60	0.30	0.40	441
33	0.60	0.31	0.40	1645
34	0.66	0.26	0.37	1058
35	0.83	0.58	0.68	946
36	0.65	0.27	0.38	644
37	0.98	0.66	0.79	136
38	0.61	0.40	0.48	570
39	0.81	0.32	0.46	766
40	0.59	0.40	0.48	1132
41	0.47	0.21	0.29	174
42	0.74	0.51	0.60	210
43	0.73	0.44	0.55	433
44	0.66	0.50	0.57	626
45	0.69	0.38	0.49	852
46	0.75	0.48	0.59	534
47	0.39	0.21	0.27	350
48	0.75	0.53	0.62	496
49	0.78	0.67	0.72	785
50	0.20	0.11	0.14	475
51	0.40	0.16	0.23	305
52	0.38	0.04	0.07	251
53	0.63	0.40	0.49	914
54	0.49	0.19	0.28	728
55	0.38	0.05	0.08	258
56	0.43	0.24	0.31	821
57	0.48	0.12	0.19	541
58	0.72	0.34	0.46	748
59	0.93	0.68	0.78	724
60	0.35	0.13	0.19	660
61	0.70	0.21	0.32	235
62	0.92	0.74	0.82	718
63	0.82	0.70	0.75	468
64	0.53	0.32	0.40	191
65	0.29	0.11	0.16	429
66	0.28	0.07	0.11	415
67	0.72	0.51	0.60	274
68	0.82	0.53	0.64	510
69	0.67	0.45	0.54	466
70	0.32	0.12	0.18	305
71	0.50	0.17	0.25	247
72	0.78	0.51	0.62	401
73	0.94	0.79	0.86	86
74	0.82	0.44	0.57	120
75	0.92	0.67	0.78	129
76	0.34	0.03	0.05	473
77	0.45	0.30	0.36	143
78	0.78	0.50	0.61	347
79	0.67	0.26	0.37	479

80	0.54	0.39	0.45	279
81	0.70	0.25	0.37	461
82	0.36	0.04	0.07	298
83	0.77	0.48	0.59	396
84	0.53	0.38	0.44	184
85	0.59	0.25	0.36	573
86	0.46	0.07	0.12	325
87	0.51	0.38	0.44	273
88	0.43	0.20	0.27	135
89	0.36	0.12	0.19	232
90	0.52	0.39	0.45	409
91	0.61	0.29	0.39	420
92	0.78	0.59	0.67	408
93	0.67	0.46	0.55	241
94	0.38	0.08	0.13	211
95	0.43	0.13	0.20	277
96	0.23	0.08	0.11	410
97	0.85	0.48	0.61	501
98	0.77	0.60	0.67	136
99	0.50	0.34	0.41	239
100	0.53	0.15	0.24	324
101	0.91	0.71	0.80	277
102	0.90	0.76	0.82	613
103	0.47	0.18	0.26	157
104	0.30	0.11	0.16	295
105	0.76	0.45	0.56	334
106	0.87	0.33	0.48	335
107	0.75	0.56	0.64	389
108	0.62	0.29	0.40	251
109	0.56	0.47	0.51	317
110	0.70	0.11	0.19	187
111	0.62	0.13	0.21	140
112	0.71	0.48	0.57	154
113	0.55	0.21	0.31	332
114	0.50	0.31	0.38	323
115	0.47	0.29	0.36	344
116	0.74	0.55	0.63	370
117	0.56	0.24	0.34	313
118	0.77	0.74	0.76	874
119	0.41	0.24	0.30	293
120	0.30	0.04	0.08	200
121	0.75	0.50	0.60	463
122	0.34	0.11	0.17	119
123	0.36	0.02	0.03	256
124	0.91	0.72	0.80	195
125	0.37	0.13	0.19	138
126	0.80	0.58	0.67	376
127	0.26	0.04	0.07	122
128	0.22	0.03	0.06	252
129	0.50	0.29	0.37	144
130	0.41	0.13	0.19	150
131	0.30	0.06	0.10	210
132	0.65	0.31	0.42	361
133	0.91	0.63	0.75	453
134	0.89	0.77	0.83	124
135	0.27	0.07	0.11	91
136	0.71	0.36	0.48	128
137	0.56	0.40	0.47	218
138	0.57	0.21	0.31	243
139	0.36	0.21	0.27	149
140	0.75	0.49	0.59	318
141	0.37	0.13	0.19	159
142	0.61	0.40	0.48	274
143	0.86	0.81	0.83	362
144	0.53	0.25	0.34	118
145	0.60	0.35	0.45	164
146	0.56	0.33	0.42	461
147	0.70	0.50	0.58	159
148	0.35	0.14	0.20	166
149	0.96	0.57	0.71	346
150	0.67	0.15	0.25	350
151	0.93	0.71	0.80	55
152	0.77	0.55	0.64	387
153	0.47	0.20	0.28	150
154	0.56	0.13	0.21	281
155	0.39	0.13	0.20	202
156	0.80	0.66	0.72	130
157	0.41	0.07	0.12	245
158	0.93	0.65	0.76	177
159	0.54	0.33	0.41	130
160	0.47	0.21	0.29	336
161	0.88	0.65	0.75	220
162	0.23	0.06	0.10	229
163	0.88	0.45	0.59	316
164	0.75	0.46	0.57	283
165	0.56	0.31	0.40	197
166	0.69	0.50	0.58	101
167	0.45	0.16	0.24	231

168	0.58	0.30	0.39	370
169	0.40	0.21	0.28	258
170	0.58	0.14	0.22	101
171	0.38	0.27	0.31	89
172	0.53	0.35	0.42	193
173	0.48	0.29	0.36	309
174	0.54	0.12	0.20	172
175	0.93	0.71	0.80	95
176	0.91	0.62	0.74	346
177	0.86	0.60	0.71	322
178	0.63	0.52	0.57	232
179	0.38	0.09	0.14	125
180	0.63	0.36	0.46	145
181	0.55	0.14	0.23	77
182	0.23	0.08	0.12	182
183	0.60	0.35	0.44	257
184	0.25	0.08	0.12	216
185	0.37	0.17	0.23	242
186	0.44	0.17	0.25	165
187	0.76	0.58	0.66	263
188	0.45	0.13	0.20	174
189	0.75	0.43	0.55	136
190	0.89	0.57	0.69	202
191	0.44	0.18	0.25	134
192	0.62	0.42	0.50	230
193	0.38	0.14	0.21	90
194	0.56	0.57	0.56	185
195	0.27	0.08	0.12	156
196	0.31	0.10	0.15	160
197	0.41	0.11	0.17	266
198	0.38	0.12	0.18	284
199	0.33	0.07	0.11	145
200	0.93	0.77	0.84	212
201	0.65	0.25	0.36	317
202	0.77	0.60	0.68	427
203	0.33	0.14	0.20	232
204	0.45	0.30	0.36	217
205	0.50	0.50	0.50	527
206	0.29	0.06	0.11	124
207	0.43	0.20	0.28	103
208	0.85	0.54	0.66	287
209	0.33	0.10	0.16	193
210	0.70	0.40	0.51	220
211	0.68	0.19	0.30	140
212	0.18	0.04	0.07	161
213	0.60	0.44	0.51	72
214	0.63	0.46	0.53	396
215	0.79	0.36	0.49	134
216	0.54	0.22	0.31	400
217	0.50	0.27	0.35	75
218	0.97	0.76	0.85	219
219	0.74	0.41	0.53	210
220	0.88	0.66	0.76	298
221	0.95	0.71	0.81	266
222	0.75	0.42	0.54	290
223	0.10	0.01	0.01	128
224	0.76	0.46	0.57	159
225	0.56	0.31	0.40	164
226	0.55	0.32	0.40	144
227	0.56	0.38	0.46	276
228	0.14	0.03	0.04	235
229	0.23	0.04	0.06	216
230	0.38	0.19	0.26	228
231	0.76	0.55	0.64	64
232	0.30	0.15	0.20	103
233	0.71	0.37	0.48	216
234	0.60	0.21	0.31	116
235	0.52	0.34	0.41	77
236	0.90	0.69	0.78	67
237	0.53	0.15	0.24	218
238	0.41	0.15	0.22	139
239	0.33	0.01	0.02	94
240	0.62	0.30	0.40	77
241	0.47	0.10	0.16	167
242	0.79	0.36	0.50	86
243	0.41	0.12	0.19	58
244	0.60	0.31	0.41	269
245	0.33	0.12	0.18	112
246	0.95	0.78	0.86	255
247	0.36	0.16	0.22	58
248	0.36	0.05	0.09	81
249	0.17	0.02	0.04	131
250	0.44	0.26	0.32	93
251	0.62	0.31	0.41	154
252	0.33	0.04	0.07	129
253	0.62	0.37	0.47	83
254	0.24	0.08	0.12	191
255	0.13	0.05	0.07	219

256	0.28	0.05	0.09	130
257	0.48	0.25	0.33	93
258	0.65	0.48	0.55	217
259	0.33	0.11	0.16	141
260	0.76	0.20	0.32	143
261	0.57	0.17	0.27	219
262	0.59	0.35	0.44	107
263	0.38	0.19	0.25	236
264	0.29	0.21	0.24	119
265	0.55	0.22	0.32	72
266	0.10	0.01	0.02	70
267	0.35	0.17	0.23	107
268	0.66	0.47	0.55	169
269	0.42	0.12	0.19	129
270	0.72	0.56	0.63	159
271	0.85	0.50	0.63	190
272	0.61	0.28	0.38	248
273	0.90	0.75	0.82	264
274	0.89	0.70	0.78	105
275	0.46	0.11	0.17	104
276	0.06	0.01	0.02	115
277	0.84	0.62	0.71	170
278	0.74	0.41	0.53	145
279	0.90	0.72	0.80	230
280	0.56	0.41	0.47	80
281	0.65	0.53	0.59	217
282	0.75	0.52	0.61	175
283	0.32	0.13	0.18	269
284	0.62	0.34	0.44	74
285	0.83	0.49	0.62	206
286	0.89	0.67	0.76	227
287	0.81	0.42	0.55	130
288	0.31	0.07	0.11	129
289	0.38	0.06	0.11	80
290	0.29	0.13	0.18	99
291	0.76	0.38	0.50	208
292	0.27	0.04	0.08	67
293	0.82	0.50	0.62	109
294	0.45	0.25	0.32	140
295	0.21	0.15	0.17	241
296	0.39	0.15	0.22	72
297	0.43	0.14	0.21	107
298	0.64	0.49	0.56	61
299	0.82	0.55	0.66	77
300	0.17	0.09	0.12	111
301	1.00	0.01	0.02	126
302	0.14	0.01	0.03	73
303	0.57	0.40	0.47	176
304	0.91	0.79	0.85	230
305	0.88	0.74	0.81	156
306	0.48	0.42	0.45	146
307	0.41	0.13	0.20	98
308	0.17	0.01	0.02	78
309	0.71	0.13	0.22	94
310	0.71	0.38	0.50	162
311	0.76	0.52	0.62	116
312	0.58	0.25	0.35	57
313	0.86	0.09	0.17	65
314	0.49	0.32	0.39	138
315	0.54	0.30	0.39	195
316	0.45	0.36	0.40	69
317	0.32	0.15	0.20	134
318	0.52	0.39	0.44	148
319	0.84	0.53	0.65	161
320	0.25	0.17	0.20	104
321	0.80	0.62	0.70	156
322	0.61	0.40	0.48	134
323	0.53	0.41	0.46	232
324	0.28	0.16	0.21	92
325	0.46	0.27	0.34	197
326	0.17	0.04	0.06	126
327	0.18	0.03	0.05	115
328	0.98	0.71	0.82	198
329	0.58	0.30	0.39	125
330	0.73	0.23	0.36	81
331	0.41	0.12	0.18	94
332	0.53	0.18	0.27	56
333	0.23	0.06	0.09	260
334	0.40	0.13	0.20	60
335	0.50	0.08	0.14	110
336	0.66	0.49	0.56	71
337	0.20	0.08	0.11	66
338	0.44	0.37	0.40	150
339	0.00	0.00	0.00	54
340	0.82	0.61	0.70	195
341	0.90	0.58	0.71	79
342	0.52	0.37	0.43	38
343	0.61	0.44	0.51	43

344	0.57	0.31	0.40	68
345	0.68	0.36	0.47	73
346	0.17	0.03	0.06	116
347	0.82	0.46	0.59	111
348	0.30	0.13	0.18	63
349	0.84	0.64	0.73	104
350	0.62	0.55	0.58	44
351	0.58	0.35	0.44	40
352	0.93	0.59	0.72	136
353	0.50	0.22	0.31	54
354	0.42	0.11	0.18	134
355	0.63	0.36	0.46	120
356	0.51	0.34	0.41	228
357	0.64	0.42	0.50	269
358	0.63	0.36	0.46	80
359	0.83	0.61	0.71	140
360	0.39	0.19	0.26	125
361	0.91	0.73	0.81	169
362	0.11	0.04	0.05	56
363	0.90	0.73	0.81	154
364	0.43	0.16	0.23	58
365	0.21	0.08	0.12	71
366	0.97	0.67	0.79	54
367	0.30	0.09	0.14	116
368	0.36	0.07	0.12	54
369	0.09	0.01	0.02	71
370	0.33	0.02	0.03	61
371	0.35	0.08	0.14	71
372	0.70	0.44	0.54	52
373	0.74	0.49	0.59	150
374	0.38	0.14	0.20	93
375	0.21	0.06	0.09	67
376	0.17	0.01	0.02	76
377	0.49	0.35	0.41	106
378	0.25	0.02	0.04	86
379	0.40	0.14	0.21	14
380	0.91	0.55	0.68	122
381	0.23	0.07	0.10	104
382	0.41	0.14	0.20	66
383	0.58	0.37	0.45	110
384	0.16	0.02	0.03	155
385	0.75	0.24	0.36	50
386	0.23	0.09	0.13	64
387	0.35	0.10	0.15	93
388	0.48	0.24	0.32	102
389	0.06	0.01	0.02	108
390	0.93	0.69	0.79	178
391	0.47	0.20	0.28	115
392	0.75	0.43	0.55	42
393	0.00	0.00	0.00	134
394	0.33	0.07	0.12	112
395	0.49	0.28	0.36	176
396	0.35	0.12	0.18	125
397	0.72	0.43	0.54	224
398	0.82	0.67	0.74	63
399	0.23	0.05	0.08	59
400	0.52	0.38	0.44	63
401	0.51	0.32	0.39	98
402	0.53	0.22	0.31	162
403	0.36	0.19	0.25	83
404	0.68	0.79	0.73	19
405	0.34	0.11	0.17	92
406	0.82	0.34	0.48	41
407	0.62	0.35	0.45	43
408	0.80	0.49	0.61	160
409	0.15	0.10	0.12	50
410	0.00	0.00	0.00	19
411	0.32	0.19	0.24	175
412	0.40	0.08	0.14	72
413	0.46	0.06	0.11	95
414	0.27	0.07	0.11	97
415	0.41	0.15	0.22	48
416	0.52	0.41	0.46	83
417	0.36	0.12	0.19	40
418	0.48	0.16	0.25	91
419	0.57	0.33	0.42	90
420	0.35	0.22	0.27	37
421	0.08	0.02	0.03	66
422	0.63	0.52	0.57	73
423	0.48	0.25	0.33	56
424	0.90	0.85	0.88	33
425	0.00	0.00	0.00	76
426	0.17	0.02	0.04	81
427	0.96	0.73	0.83	150
428	0.95	0.69	0.80	29
429	0.99	0.94	0.97	389
430	0.63	0.40	0.49	167
431	0.55	0.13	0.21	123

432	0.46	0.31	0.37	39
433	0.31	0.28	0.29	82
434	1.00	0.71	0.83	66
435	0.60	0.41	0.49	93
436	0.57	0.37	0.45	87
437	0.15	0.03	0.06	86
438	0.71	0.51	0.59	104
439	0.56	0.15	0.24	100
440	0.44	0.06	0.10	141
441	0.48	0.37	0.42	110
442	0.35	0.15	0.21	123
443	0.50	0.14	0.22	71
444	0.30	0.07	0.12	109
445	0.53	0.42	0.47	48
446	0.44	0.37	0.40	76
447	0.33	0.24	0.28	38
448	0.63	0.52	0.57	81
449	0.60	0.33	0.42	132
450	0.44	0.32	0.37	81
451	0.87	0.34	0.49	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.83	0.57	0.68	70
455	0.37	0.18	0.24	155
456	0.48	0.23	0.31	43
457	0.56	0.31	0.40	72
458	0.40	0.13	0.20	62
459	0.58	0.26	0.36	69
460	0.10	0.03	0.04	119
461	0.72	0.23	0.35	79
462	0.55	0.23	0.33	47
463	0.44	0.16	0.24	104
464	0.63	0.41	0.49	106
465	0.58	0.28	0.38	64
466	0.54	0.44	0.48	173
467	0.72	0.45	0.55	107
468	0.56	0.25	0.35	126
469	0.57	0.04	0.07	114
470	0.93	0.82	0.87	140
471	0.79	0.43	0.56	79
472	0.40	0.37	0.38	143
473	0.65	0.42	0.51	158
474	0.50	0.11	0.18	138
475	0.00	0.00	0.00	59
476	0.66	0.38	0.48	88
477	0.81	0.69	0.74	176
478	0.95	0.79	0.86	24
479	0.18	0.08	0.11	92
480	0.79	0.54	0.64	100
481	0.49	0.30	0.37	103
482	0.37	0.22	0.27	74
483	0.78	0.62	0.69	105
484	0.28	0.06	0.10	83
485	0.24	0.05	0.08	82
486	0.45	0.18	0.26	71
487	0.48	0.20	0.28	120
488	0.50	0.06	0.10	105
489	0.70	0.37	0.48	87
490	1.00	0.84	0.92	32
491	0.14	0.01	0.03	69
492	0.50	0.02	0.04	49
493	0.17	0.04	0.07	117
494	0.52	0.23	0.32	61
495	0.93	0.82	0.87	344
496	0.37	0.19	0.25	52
497	0.59	0.28	0.38	137
498	0.33	0.15	0.21	98
499	0.50	0.22	0.30	79
micro avg	0.69	0.39	0.50	173812
macro avg	0.55	0.31	0.39	173812
weighted avg	0.65	0.39	0.48	173812
samples avg	0.47	0.37	0.39	173812

4.5.8 Linear SVM Hyper Parameter Tuning (1 to 3grams)

```

In [74]: start = datetime.now()

classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', penalty='l1'), n_jobs=-1)

parameters = {'estimator__alpha':[2*10**-6, 4*10**-6, 6*10**-6, 8*10**-6]}

clf = GridSearchCV(classifier, parameters, cv= 3, return_train_score=True, scoring='f1_micro', n_jobs=-1)

clf.fit(x_train_multilabel, y_train)

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

```

Time taken to run this cell : 1:17:52.986729

```

In [75]: plt.figure(figsize=(20,10))

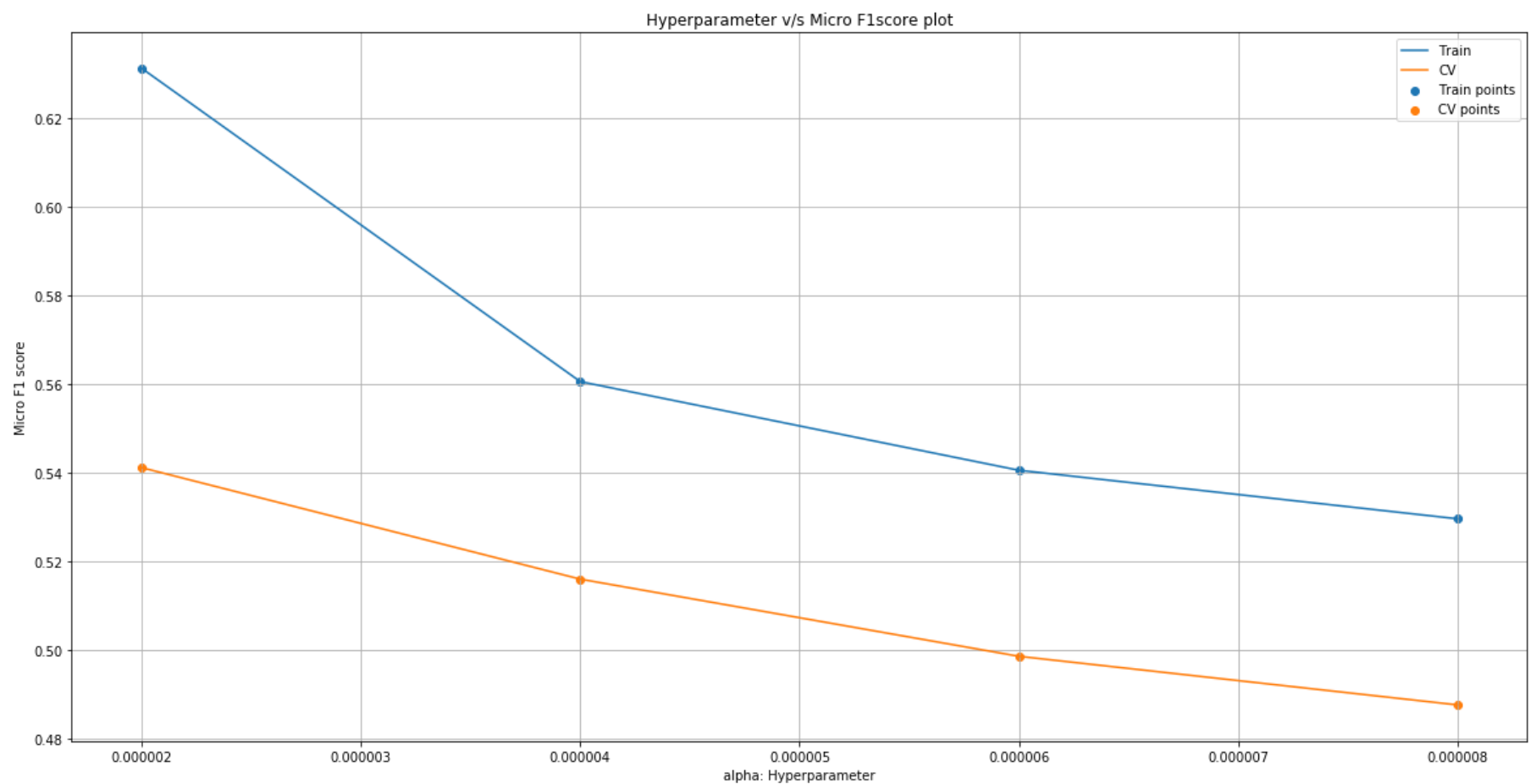
train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']

plt.plot(parameters['estimator__alpha'], train_auc, label='Train')
plt.plot(parameters['estimator__alpha'], cv_auc, label='CV')

plt.scatter(parameters['estimator__alpha'], train_auc, label='Train points')
plt.scatter(parameters['estimator__alpha'], cv_auc, label='CV points')

plt.legend()
plt.xlabel("alpha: Hyperparameter")
plt.ylabel("Micro F1 score")
plt.title("Hyperparameter v/s Micro F1score plot")
plt.grid()
plt.show()

```



```
In [76]: import warnings
warnings.filterwarnings("ignore")

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=2*10**-6, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:08:38.912048
Hamming loss 0.00259098
Micro-average quality numbers
Precision: 0.7847, Recall: 0.3510, F1-measure: 0.4850
Macro-average quality numbers
Precision: 0.5278, Recall: 0.2689, F1-measure: 0.3264

	precision	recall	f1-score	support
0	0.94	0.77	0.84	5519
1	0.72	0.33	0.45	8190
2	0.82	0.42	0.55	6529
3	0.83	0.50	0.62	3231
4	0.83	0.43	0.57	6430
5	0.81	0.40	0.54	2879
6	0.87	0.55	0.68	5086
7	0.88	0.61	0.72	4533
8	0.62	0.15	0.24	3000
9	0.82	0.59	0.68	2765
10	0.64	0.09	0.15	3051
11	0.73	0.38	0.50	3009
12	0.71	0.26	0.38	2630
13	0.72	0.33	0.45	1426
14	0.89	0.63	0.74	2548
15	0.76	0.20	0.31	2371
16	0.67	0.30	0.41	873
17	0.89	0.64	0.74	2151
18	0.72	0.21	0.32	2204
19	0.70	0.46	0.55	831
20	0.75	0.56	0.64	1860
21	0.40	0.00	0.01	2023
22	0.59	0.19	0.28	1513
23	0.90	0.60	0.72	1207
24	0.62	0.18	0.28	506
25	0.73	0.38	0.50	425
26	0.64	0.41	0.50	793
27	0.69	0.32	0.44	1291
28	0.82	0.39	0.53	1208
29	0.34	0.03	0.06	406
30	0.72	0.20	0.31	504
31	0.00	0.00	0.00	732
32	0.65	0.29	0.40	441
33	0.70	0.13	0.22	1645
34	0.74	0.27	0.39	1058
35	0.81	0.64	0.71	946
36	0.68	0.21	0.32	644
37	0.95	0.77	0.85	136
38	0.66	0.41	0.50	570
39	0.83	0.34	0.48	766
40	0.65	0.28	0.39	1132
41	0.54	0.26	0.35	174
42	0.76	0.55	0.64	210
43	0.82	0.42	0.56	433
44	0.66	0.57	0.61	626
45	0.80	0.29	0.42	852
46	0.77	0.50	0.61	534
47	0.38	0.02	0.03	350
48	0.75	0.58	0.65	496
49	0.77	0.74	0.75	785
50	0.00	0.00	0.00	475
51	0.00	0.00	0.00	305
52	0.00	0.00	0.00	251
53	0.68	0.38	0.48	914
54	0.00	0.00	0.00	728
55	0.00	0.00	0.00	258
56	0.50	0.02	0.03	821
57	0.50	0.01	0.02	541
58	0.81	0.26	0.39	748
59	0.92	0.72	0.81	724
60	0.53	0.02	0.03	660
61	0.83	0.23	0.35	235
62	0.92	0.76	0.83	718
63	0.82	0.66	0.73	468
64	0.54	0.31	0.40	191
65	0.00	0.00	0.00	429
66	0.00	0.00	0.00	415
67	0.76	0.55	0.64	274
68	0.82	0.57	0.67	510
69	0.67	0.51	0.58	466
70	0.33	0.00	0.01	305
71	0.33	0.00	0.01	247
72	0.81	0.52	0.64	401
73	0.96	0.81	0.88	86
74	0.74	0.42	0.53	120
75	0.91	0.73	0.81	129
76	0.91	0.02	0.04	473
77	0.40	0.34	0.37	143
78	0.77	0.56	0.65	347
79	0.73	0.27	0.39	479

80	0.73	0.28	0.40	279
81	0.81	0.24	0.37	461
82	0.00	0.00	0.00	298
83	0.77	0.56	0.65	396
84	0.55	0.29	0.38	184
85	0.75	0.20	0.31	573
86	0.57	0.05	0.09	325
87	0.54	0.30	0.38	273
88	0.42	0.06	0.10	135
89	1.00	0.00	0.01	232
90	0.57	0.25	0.35	409
91	0.62	0.04	0.08	420
92	0.78	0.59	0.67	408
93	0.69	0.50	0.58	241
94	0.00	0.00	0.00	211
95	0.64	0.03	0.05	277
96	0.00	0.00	0.00	410
97	0.86	0.53	0.66	501
98	0.74	0.71	0.72	136
99	0.58	0.28	0.38	239
100	0.00	0.00	0.00	324
101	0.89	0.78	0.83	277
102	0.91	0.76	0.83	613
103	0.50	0.04	0.07	157
104	0.00	0.00	0.00	295
105	0.80	0.43	0.56	334
106	0.94	0.30	0.46	335
107	0.77	0.62	0.69	389
108	0.90	0.08	0.14	251
109	0.55	0.44	0.49	317
110	1.00	0.01	0.02	187
111	0.74	0.16	0.27	140
112	0.72	0.52	0.60	154
113	0.65	0.16	0.26	332
114	0.73	0.02	0.05	323
115	0.42	0.01	0.03	344
116	0.73	0.61	0.67	370
117	0.61	0.15	0.24	313
118	0.78	0.74	0.76	874
119	0.56	0.05	0.09	293
120	0.00	0.00	0.00	200
121	0.73	0.56	0.64	463
122	0.00	0.00	0.00	119
123	0.57	0.02	0.03	256
124	0.88	0.75	0.81	195
125	0.62	0.07	0.13	138
126	0.81	0.57	0.67	376
127	0.00	0.00	0.00	122
128	0.00	0.00	0.00	252
129	0.53	0.22	0.31	144
130	0.57	0.11	0.18	150
131	0.00	0.00	0.00	210
132	0.67	0.04	0.08	361
133	0.92	0.61	0.73	453
134	0.88	0.81	0.85	124
135	0.00	0.00	0.00	91
136	0.88	0.17	0.29	128
137	0.61	0.32	0.42	218
138	0.85	0.18	0.30	243
139	0.44	0.03	0.05	149
140	0.75	0.55	0.63	318
141	0.36	0.03	0.05	159
142	0.64	0.51	0.57	274
143	0.84	0.83	0.83	362
144	0.59	0.08	0.15	118
145	0.60	0.42	0.49	164
146	0.62	0.23	0.33	461
147	0.69	0.58	0.63	159
148	0.00	0.00	0.00	166
149	0.93	0.64	0.76	346
150	0.80	0.05	0.09	350
151	0.86	0.69	0.77	55
152	0.79	0.55	0.65	387
153	0.61	0.07	0.13	150
154	0.62	0.06	0.10	281
155	0.00	0.00	0.00	202
156	0.81	0.68	0.74	130
157	0.00	0.00	0.00	245
158	0.88	0.71	0.78	177
159	0.64	0.32	0.43	130
160	0.43	0.01	0.02	336
161	0.88	0.70	0.78	220
162	0.00	0.00	0.00	229
163	0.88	0.48	0.62	316
164	0.78	0.42	0.54	283
165	0.63	0.38	0.47	197
166	0.69	0.55	0.62	101
167	1.00	0.00	0.01	231

168	0.72	0.16	0.26	370
169	0.45	0.03	0.06	258
170	0.00	0.00	0.00	101
171	0.46	0.26	0.33	89
172	0.48	0.10	0.17	193
173	0.49	0.16	0.24	309
174	1.00	0.01	0.01	172
175	0.91	0.84	0.87	95
176	0.91	0.64	0.75	346
177	0.82	0.65	0.72	322
178	0.63	0.54	0.58	232
179	0.00	0.00	0.00	125
180	0.72	0.36	0.48	145
181	0.00	0.00	0.00	77
182	0.00	0.00	0.00	182
183	0.61	0.19	0.29	257
184	0.00	0.00	0.00	216
185	0.00	0.00	0.00	242
186	1.00	0.01	0.01	165
187	0.75	0.66	0.70	263
188	0.00	0.00	0.00	174
189	0.79	0.44	0.57	136
190	0.87	0.64	0.74	202
191	0.00	0.00	0.00	134
192	0.72	0.47	0.57	230
193	0.36	0.09	0.14	90
194	0.56	0.55	0.55	185
195	0.00	0.00	0.00	156
196	0.00	0.00	0.00	160
197	0.65	0.13	0.21	266
198	0.00	0.00	0.00	284
199	0.00	0.00	0.00	145
200	0.93	0.77	0.85	212
201	0.71	0.23	0.35	317
202	0.77	0.62	0.69	427
203	0.00	0.00	0.00	232
204	0.61	0.18	0.28	217
205	0.49	0.30	0.37	527
206	0.00	0.00	0.00	124
207	0.42	0.11	0.17	103
208	0.86	0.55	0.67	287
209	0.00	0.00	0.00	193
210	0.79	0.32	0.45	220
211	0.78	0.23	0.35	140
212	0.00	0.00	0.00	161
213	0.55	0.50	0.52	72
214	0.68	0.31	0.42	396
215	0.82	0.40	0.54	134
216	0.67	0.07	0.13	400
217	0.49	0.27	0.34	75
218	0.95	0.78	0.86	219
219	0.80	0.34	0.47	210
220	0.89	0.70	0.78	298
221	0.90	0.79	0.84	266
222	0.78	0.43	0.55	290
223	0.00	0.00	0.00	128
224	0.76	0.48	0.59	159
225	0.71	0.29	0.41	164
226	0.63	0.42	0.50	144
227	0.63	0.22	0.32	276
228	0.00	0.00	0.00	235
229	0.33	0.00	0.01	216
230	0.00	0.00	0.00	228
231	0.73	0.59	0.66	64
232	1.00	0.01	0.02	103
233	0.75	0.39	0.51	216
234	0.57	0.07	0.12	116
235	0.57	0.60	0.59	77
236	0.91	0.73	0.81	67
237	0.74	0.11	0.20	218
238	0.00	0.00	0.00	139
239	0.00	0.00	0.00	94
240	0.62	0.42	0.50	77
241	0.00	0.00	0.00	167
242	0.84	0.44	0.58	86
243	0.67	0.03	0.07	58
244	0.75	0.16	0.26	269
245	0.00	0.00	0.00	112
246	0.94	0.82	0.88	255
247	0.44	0.29	0.35	58
248	0.00	0.00	0.00	81
249	0.00	0.00	0.00	131
250	0.63	0.18	0.28	93
251	0.66	0.23	0.34	154
252	0.00	0.00	0.00	129
253	0.69	0.30	0.42	83
254	0.00	0.00	0.00	191
255	0.00	0.00	0.00	219

256	0.00	0.00	0.00	130
257	0.50	0.23	0.31	93
258	0.69	0.63	0.66	217
259	1.00	0.01	0.01	141
260	0.85	0.20	0.33	143
261	1.00	0.01	0.02	219
262	0.64	0.25	0.36	107
263	0.33	0.06	0.10	236
264	0.45	0.04	0.08	119
265	0.50	0.12	0.20	72
266	0.00	0.00	0.00	70
267	0.67	0.02	0.04	107
268	0.70	0.51	0.59	169
269	0.50	0.01	0.02	129
270	0.72	0.64	0.67	159
271	0.87	0.50	0.64	190
272	0.64	0.24	0.35	248
273	0.89	0.78	0.83	264
274	0.88	0.71	0.79	105
275	0.78	0.07	0.12	104
276	0.00	0.00	0.00	115
277	0.83	0.68	0.75	170
278	0.78	0.48	0.59	145
279	0.91	0.71	0.80	230
280	0.53	0.34	0.41	80
281	0.67	0.66	0.66	217
282	0.78	0.62	0.69	175
283	0.00	0.00	0.00	269
284	0.67	0.54	0.60	74
285	0.85	0.52	0.64	206
286	0.89	0.68	0.77	227
287	0.84	0.42	0.56	130
288	0.00	0.00	0.00	129
289	0.67	0.03	0.05	80
290	0.00	0.00	0.00	99
291	0.77	0.33	0.46	208
292	0.00	0.00	0.00	67
293	0.86	0.47	0.61	109
294	0.80	0.03	0.06	140
295	0.00	0.00	0.00	241
296	0.00	0.00	0.00	72
297	0.40	0.02	0.04	107
298	0.73	0.52	0.61	61
299	0.81	0.60	0.69	77
300	0.00	0.00	0.00	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.62	0.45	0.52	176
304	0.92	0.82	0.87	230
305	0.90	0.78	0.83	156
306	0.57	0.27	0.36	146
307	0.33	0.01	0.02	98
308	0.00	0.00	0.00	78
309	0.67	0.11	0.18	94
310	0.70	0.40	0.51	162
311	0.76	0.59	0.66	116
312	0.48	0.37	0.42	57
313	0.00	0.00	0.00	65
314	0.51	0.36	0.42	138
315	0.68	0.15	0.25	195
316	0.51	0.26	0.35	69
317	0.37	0.08	0.13	134
318	0.51	0.20	0.28	148
319	0.85	0.60	0.70	161
320	0.15	0.03	0.05	104
321	0.80	0.63	0.71	156
322	0.61	0.40	0.48	134
323	0.55	0.56	0.55	232
324	0.00	0.00	0.00	92
325	0.00	0.00	0.00	197
326	0.00	0.00	0.00	126
327	0.00	0.00	0.00	115
328	0.96	0.74	0.84	198
329	0.60	0.33	0.42	125
330	0.74	0.17	0.28	81
331	1.00	0.01	0.02	94
332	0.40	0.07	0.12	56
333	0.00	0.00	0.00	260
334	0.50	0.02	0.03	60
335	0.00	0.00	0.00	110
336	0.68	0.48	0.56	71
337	0.50	0.02	0.03	66
338	0.49	0.44	0.46	150
339	0.00	0.00	0.00	54
340	0.82	0.62	0.70	195
341	0.82	0.51	0.62	79
342	0.52	0.29	0.37	38
343	0.67	0.51	0.58	43

344	0.59	0.19	0.29	68
345	0.65	0.27	0.38	73
346	0.00	0.00	0.00	116
347	0.82	0.50	0.63	111
348	0.25	0.02	0.03	63
349	0.84	0.70	0.76	104
350	0.59	0.59	0.59	44
351	0.61	0.42	0.50	40
352	0.91	0.64	0.75	136
353	0.20	0.02	0.03	54
354	0.00	0.00	0.00	134
355	0.75	0.33	0.46	120
356	0.67	0.17	0.27	228
357	0.64	0.41	0.50	269
358	0.73	0.34	0.46	80
359	0.84	0.67	0.75	140
360	0.67	0.02	0.03	125
361	0.90	0.76	0.83	169
362	1.00	0.02	0.04	56
363	0.89	0.76	0.82	154
364	0.57	0.07	0.12	58
365	0.80	0.06	0.11	71
366	0.97	0.70	0.82	54
367	0.00	0.00	0.00	116
368	0.50	0.02	0.04	54
369	0.00	0.00	0.00	71
370	0.00	0.00	0.00	61
371	0.00	0.00	0.00	71
372	0.70	0.50	0.58	52
373	0.77	0.46	0.57	150
374	0.00	0.00	0.00	93
375	0.00	0.00	0.00	67
376	0.00	0.00	0.00	76
377	0.77	0.09	0.17	106
378	0.00	0.00	0.00	86
379	0.75	0.21	0.33	14
380	0.84	0.64	0.73	122
381	0.00	0.00	0.00	104
382	0.00	0.00	0.00	66
383	0.57	0.23	0.32	110
384	0.00	0.00	0.00	155
385	0.71	0.24	0.36	50
386	0.43	0.05	0.08	64
387	0.00	0.00	0.00	93
388	0.40	0.12	0.18	102
389	0.00	0.00	0.00	108
390	0.94	0.71	0.81	178
391	0.52	0.20	0.29	115
392	0.85	0.52	0.65	42
393	0.00	0.00	0.00	134
394	0.00	0.00	0.00	112
395	0.00	0.00	0.00	176
396	0.57	0.03	0.06	125
397	0.72	0.40	0.52	224
398	0.81	0.70	0.75	63
399	0.00	0.00	0.00	59
400	0.56	0.32	0.40	63
401	0.64	0.18	0.29	98
402	0.00	0.00	0.00	162
403	0.00	0.00	0.00	83
404	0.67	0.84	0.74	19
405	0.67	0.02	0.04	92
406	0.78	0.34	0.47	41
407	0.63	0.44	0.52	43
408	0.78	0.44	0.56	160
409	0.00	0.00	0.00	50
410	0.00	0.00	0.00	19
411	0.00	0.00	0.00	175
412	0.00	0.00	0.00	72
413	0.50	0.05	0.10	95
414	0.00	0.00	0.00	97
415	0.20	0.02	0.04	48
416	0.64	0.17	0.27	83
417	1.00	0.03	0.05	40
418	0.00	0.00	0.00	91
419	0.60	0.28	0.38	90
420	0.25	0.03	0.05	37
421	0.00	0.00	0.00	66
422	0.64	0.44	0.52	73
423	0.45	0.27	0.34	56
424	0.88	0.88	0.88	33
425	0.00	0.00	0.00	76
426	0.00	0.00	0.00	81
427	0.97	0.79	0.87	150
428	0.91	0.69	0.78	29
429	0.99	0.93	0.96	389
430	0.67	0.37	0.48	167
431	0.00	0.00	0.00	123

432	0.47	0.23	0.31	39
433	0.50	0.04	0.07	82
434	0.96	0.71	0.82	66
435	0.60	0.44	0.51	93
436	0.67	0.53	0.59	87
437	0.50	0.01	0.02	86
438	0.72	0.69	0.71	104
439	0.00	0.00	0.00	100
440	0.00	0.00	0.00	141
441	0.62	0.21	0.31	110
442	0.00	0.00	0.00	123
443	0.47	0.24	0.32	71
444	0.00	0.00	0.00	109
445	0.64	0.33	0.44	48
446	0.48	0.29	0.36	76
447	0.00	0.00	0.00	38
448	0.66	0.65	0.66	81
449	0.84	0.20	0.33	132
450	0.64	0.09	0.15	81
451	0.84	0.34	0.49	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.80	0.51	0.63	70
455	0.50	0.04	0.07	155
456	0.60	0.21	0.31	43
457	0.52	0.17	0.25	72
458	0.00	0.00	0.00	62
459	0.76	0.23	0.36	69
460	0.00	0.00	0.00	119
461	0.86	0.08	0.14	79
462	0.62	0.17	0.27	47
463	1.00	0.13	0.24	104
464	0.64	0.36	0.46	106
465	0.62	0.08	0.14	64
466	0.61	0.30	0.40	173
467	0.72	0.38	0.50	107
468	0.60	0.17	0.26	126
469	0.00	0.00	0.00	114
470	0.91	0.82	0.86	140
471	0.87	0.43	0.58	79
472	0.42	0.32	0.36	143
473	0.71	0.37	0.48	158
474	0.89	0.06	0.11	138
475	0.00	0.00	0.00	59
476	0.69	0.40	0.50	88
477	0.81	0.71	0.76	176
478	0.91	0.88	0.89	24
479	0.00	0.00	0.00	92
480	0.79	0.55	0.65	100
481	0.53	0.32	0.40	103
482	0.00	0.00	0.00	74
483	0.83	0.62	0.71	105
484	0.00	0.00	0.00	83
485	0.00	0.00	0.00	82
486	0.00	0.00	0.00	71
487	0.00	0.00	0.00	120
488	0.64	0.09	0.15	105
489	0.71	0.46	0.56	87
490	1.00	0.81	0.90	32
491	0.00	0.00	0.00	69
492	0.00	0.00	0.00	49
493	0.00	0.00	0.00	117
494	0.54	0.31	0.40	61
495	0.96	0.81	0.88	344
496	0.50	0.10	0.16	52
497	0.64	0.25	0.36	137
498	0.41	0.09	0.15	98
499	0.56	0.23	0.32	79
micro avg	0.78	0.35	0.48	173812
macro avg	0.53	0.27	0.33	173812
weighted avg	0.67	0.35	0.44	173812
samples avg	0.47	0.34	0.37	173812

4.5.9 Featurizing data with BOW vectorizer (1 to 4 grams)

Reduce Datasize to 0.1M due to Time complexity calculating 1 to 4 grams(was taking more than 10 hours)

```
In [10]: 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("Titlemoreweight100k.db", sql_create_table)
```

Tables in the database:
QuestionsProcessed

```
In [11]: # 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 = 'Titlemoreweight100k.db'

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 0.1M rows
        reader.execute("SELECT Title, Body, Tags From no_dup_train LIMIT 100001;")
        # for selecting random points
        #reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY RANDOM() LIMIT 100001;")

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

In [12]: [#http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/](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=striphtml(question.encode('utf-8'))

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

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

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

    len_post+=len(question)
    tup = (question,code,tags,x,len(question),is_code)
    questions_proccesed += 1
    writer.execute("insert into QuestionsProcessed(question,code,tags,words_pre,words_post,is_code) values
(?,?,?,?,?,?)",tup)
    if (questions_proccesed%10000==0):
        print("number of questions completed=",questions_proccesed)

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

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

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

```
number of questions completed= 10000
number of questions completed= 20000
number of questions completed= 30000
number of questions completed= 40000
number of questions completed= 50000
number of questions completed= 60000
number of questions completed= 70000
number of questions completed= 80000
number of questions completed= 90000
number of questions completed= 100000
Avg. length of questions(Title+Body) before processing: 1232
Avg. length of questions(Title+Body) after processing: 355
Percent of questions containing code: 57
Time taken to run this cell : 0:02:32.358797
```

In [0]: *# 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()
```

In [0]: *#Taking 0.1 Million entries to a dataframe.*

```
write_db = 'Titlmoreweight100k.db'
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        preprocessed_data_100k = pd.read_sql_query("""SELECT question, Tags FROM QuestionsProcessed""", conn_r)
conn_r.commit()
conn_r.close()
```

```
In [15]: preprocessed_data_100k.head()
```

Out[15]:

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

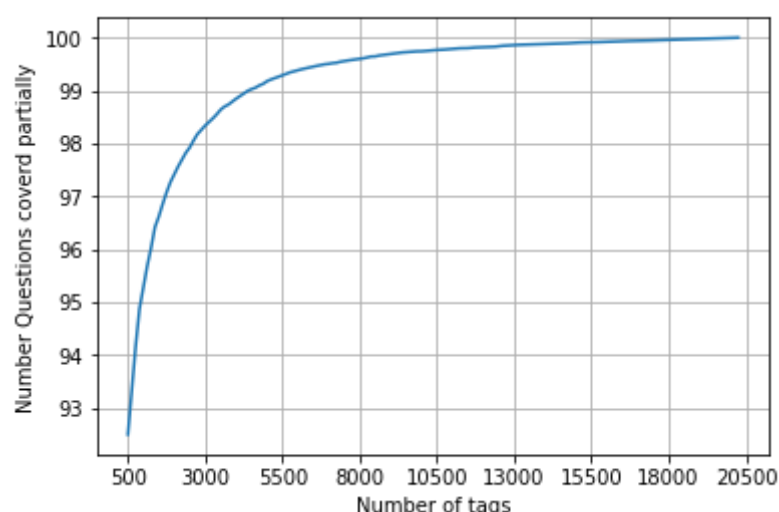
```
In [16]: print("number of data points in sample :", preprocessed_data_100k.shape[0])
print("number of dimensions :", preprocessed_data_100k.shape[1])
```

```
number of data points in sample : 100000
number of dimensions : 2
```

```
In [0]: vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data_100k['tags'])
```

```
In [0]: questions_explained = []
total_tags=multilabel_y.shape[1]
total_qs=preprocessed_data_100k.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 [22]: fig, ax = plt.subplots()
ax.plot(questions_explained)
xlabel = list(500+np.array(range(-50,450,50))*50)
ax.set_xticklabels(xlabel)
plt.xlabel("Number of tags")
plt.ylabel("Number Questions coverd partially")
plt.grid()
plt.show()
# you can choose any number of tags based on your computing power, minimun is 500(it covers 90% of the tags)
print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



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

```
In [23]: # we will be taking 5500 tags
multilabel_yx = tags_to_choose(500)
print("number of questions that are not covered :", questions_explained_fn(500),"out of ", total_qs)
```

```
number of questions that are not covered : 7500 out of 100000
```

```
In [0]: x_train=preprocessed_data_100k.head(70000)
x_test=preprocessed_data_100k.tail(preprocessed_data_100k.shape[0] - 70000)

y_train = multilabel_yx[0:70000,:]
y_test = multilabel_yx[70000:preprocessed_data_100k.shape[0],:]
```

```
In [25]: 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 : (70000, 500)
Number of data points in test data : (30000, 500)
```

```
In [26]: start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009, max_features=200000, tokenizer = lambda x: str(x).split(), ngram_range=(1,4), lowercase = False)
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 : 0:01:12.103851
```

```
In [27]: 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: (70000, 105469) Y : (70000, 500)
Dimensions of test data X: (30000, 105469) Y: (30000, 500)
```

4.5.10 Applying Logistic Regression (SGD with Log Loss) with OneVsRest Classifier(1 to 4grams)

Hyper parameter Tuning

```
In [29]: from sklearn.model_selection import GridSearchCV
start = datetime.now()

classifier = OneVsRestClassifier(SGDClassifier(loss='log', penalty='l1'), n_jobs= 1)

parameters = {'estimator__alpha':[2*10**-5, 4*10**-5, 6*10**-5, 8*10**-5]}

clf = GridSearchCV(classifier, parameters, cv= 3, return_train_score=True, scoring='f1_micro')

clf.fit(x_train_multilabel, y_train)

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

Time taken to run this cell : 3:30:02.857675
```

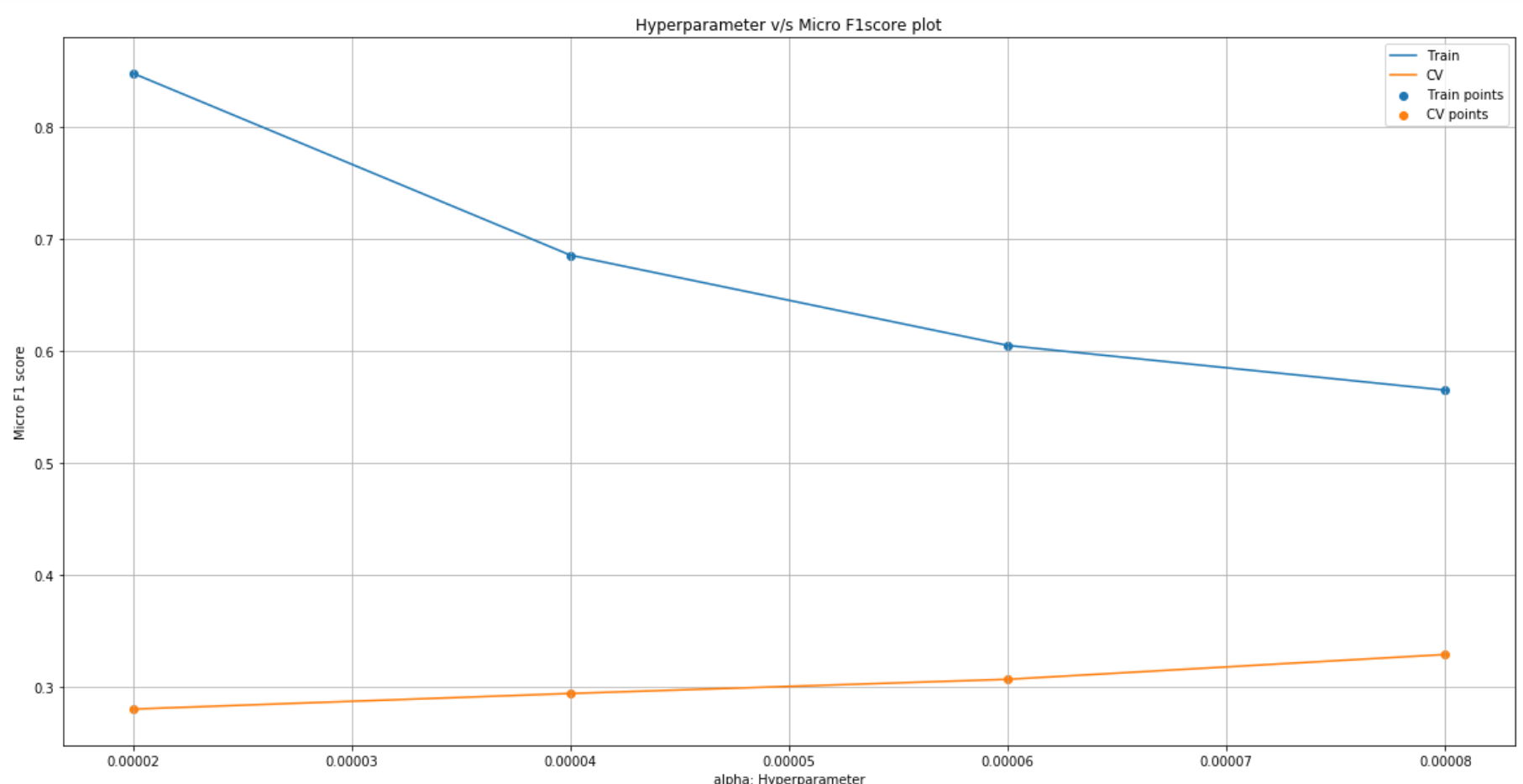
```
In [30]: plt.figure(figsize=(20,10))

train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']

plt.plot(parameters['estimator__alpha'], train_auc, label='Train')
plt.plot(parameters['estimator__alpha'], cv_auc, label='CV')

plt.scatter(parameters['estimator__alpha'], train_auc, label='Train points')
plt.scatter(parameters['estimator__alpha'], cv_auc, label='CV points')

plt.legend()
plt.xlabel("alpha: Hyperparameter")
plt.ylabel("Micro F1 score")
plt.title("Hyperparameter v/s Micro F1score plot")
plt.grid()
plt.show()
```



Model

```
In [31]: start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=8*10**-5, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict (x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:08:28.749626
 Accuracy : 0.0852
 Hamming loss 0.005211866666666667
 Micro-average quality numbers
 Precision: 0.3667, Recall: 0.4082, F1-measure: 0.3863
 Macro-average quality numbers
 Precision: 0.2508, Recall: 0.3127, F1-measure: 0.2588

	precision	recall	f1-score	support
0	0.83	0.73	0.78	6668
1	0.41	0.27	0.32	3659
2	0.15	0.15	0.15	971
3	0.49	0.56	0.52	1506
4	0.52	0.49	0.50	1649
5	0.55	0.49	0.52	1113
6	0.59	0.45	0.51	1482
7	0.53	0.63	0.58	980
8	0.81	0.68	0.74	1520
9	0.50	0.86	0.63	1041
10	0.49	0.41	0.45	861
11	0.36	0.38	0.37	386
12	0.13	0.65	0.22	37
13	0.50	0.46	0.48	917
14	0.27	0.25	0.26	519
15	0.31	0.26	0.28	656
16	0.39	0.26	0.31	794
17	0.39	0.37	0.38	700
18	0.44	0.66	0.53	363
19	0.55	0.62	0.58	541
20	0.35	0.49	0.40	540
21	0.47	0.67	0.55	362
22	0.64	0.55	0.59	551
23	0.28	0.31	0.30	309
24	0.33	0.40	0.36	331
25	0.26	0.38	0.31	424
26	0.32	0.30	0.31	465
27	0.18	0.11	0.14	386
28	0.19	0.26	0.22	107
29	0.06	0.12	0.08	195
30	0.29	0.34	0.31	758
31	0.09	0.47	0.14	15
32	0.45	0.61	0.52	323
33	0.29	0.33	0.31	279
34	0.37	0.33	0.35	275
35	0.67	0.53	0.59	268
36	0.11	0.12	0.11	76
37	0.27	0.13	0.17	269
38	0.39	0.48	0.43	255
39	0.25	0.36	0.30	249
40	0.05	0.09	0.06	66
41	0.18	0.23	0.20	209
42	0.40	0.43	0.42	72
43	0.33	0.21	0.26	430
44	0.20	0.25	0.22	279
45	0.32	0.34	0.33	240
46	0.22	0.43	0.29	157
47	0.58	0.60	0.59	249
48	0.09	0.08	0.09	198
49	0.40	0.32	0.35	171
50	0.57	0.74	0.64	200
51	0.52	0.72	0.60	85
52	0.37	0.42	0.39	175
53	0.11	0.30	0.16	114
54	0.15	0.19	0.17	223
55	0.30	0.42	0.35	122
56	0.38	0.57	0.46	168
57	0.04	0.09	0.06	176
58	0.12	0.26	0.16	140
59	0.21	0.21	0.21	191
60	0.58	0.68	0.63	152
61	0.25	0.32	0.28	208
62	0.12	0.17	0.14	136
63	0.42	0.50	0.46	158
64	0.36	0.46	0.40	203
65	0.29	0.49	0.36	105
66	0.23	0.48	0.31	58
67	0.36	0.57	0.44	128
68	0.15	0.07	0.09	158
69	0.17	0.21	0.18	248
70	0.10	0.11	0.11	201
71	0.21	0.36	0.26	89
72	0.33	0.38	0.35	157
73	0.20	0.38	0.26	29
74	0.09	0.03	0.05	58
75	0.13	0.23	0.17	158
76	0.65	0.56	0.60	110
77	0.27	0.64	0.38	33
78	0.13	0.23	0.16	210

79	0.40	0.56	0.46	169
80	0.04	0.20	0.07	15
81	0.30	0.43	0.36	214
82	0.14	0.25	0.18	65
83	0.15	0.22	0.18	156
84	0.26	0.41	0.32	59
85	0.45	0.62	0.52	55
86	0.06	0.22	0.10	36
87	0.26	0.52	0.35	29
88	0.31	0.67	0.42	54
89	0.46	0.77	0.58	137
90	0.11	0.16	0.13	103
91	0.23	0.18	0.20	79
92	0.17	0.19	0.18	84
93	0.35	0.48	0.41	133
94	0.73	0.70	0.72	318
95	0.33	0.67	0.44	51
96	0.35	0.50	0.41	82
97	0.08	0.11	0.09	75
98	0.06	0.02	0.03	120
99	0.20	0.44	0.28	18
100	0.38	0.42	0.40	196
101	0.53	0.49	0.51	208
102	0.17	0.13	0.15	122
103	0.01	0.02	0.01	62
104	0.10	0.09	0.10	88
105	0.49	0.49	0.49	65
106	0.17	0.22	0.19	115
107	0.03	0.14	0.05	29
108	0.25	0.27	0.26	109
109	0.19	0.34	0.24	73
110	0.46	0.26	0.34	102
111	0.37	0.51	0.43	180
112	0.00	0.00	0.00	292
113	0.59	0.87	0.71	54
114	0.07	0.09	0.08	120
115	0.18	0.31	0.22	107
116	0.39	0.25	0.31	52
117	0.16	0.15	0.16	72
118	0.59	0.58	0.58	139
119	0.47	0.49	0.48	57
120	0.22	0.32	0.26	44
121	0.08	0.13	0.10	85
122	0.37	0.57	0.45	82
123	0.15	0.08	0.10	100
124	0.13	0.75	0.22	4
125	0.13	0.67	0.22	9
126	0.11	0.22	0.15	46
127	0.18	0.15	0.16	54
128	0.84	0.78	0.81	195
129	0.30	0.50	0.38	54
130	0.14	0.05	0.08	96
131	0.31	0.71	0.43	35
132	0.10	0.14	0.12	58
133	0.13	0.14	0.14	36
134	0.22	0.42	0.29	36
135	0.36	0.62	0.46	39
136	0.02	0.01	0.01	97
137	0.14	0.24	0.18	70
138	0.07	0.12	0.09	17
139	0.12	0.18	0.15	119
140	0.74	0.71	0.73	101
141	0.22	0.39	0.28	115
142	0.38	0.36	0.37	94
143	0.34	0.57	0.43	84
144	0.36	0.52	0.42	64
145	0.03	0.07	0.04	61
146	0.14	0.05	0.08	132
147	0.28	0.29	0.29	119
148	0.36	0.60	0.45	62
149	0.31	0.22	0.25	83
150	0.12	0.26	0.16	72
151	0.08	0.39	0.13	23
152	0.10	0.22	0.14	76
153	0.11	0.28	0.16	18
154	0.05	0.12	0.07	17
155	0.07	0.12	0.09	24
156	0.27	0.18	0.22	136
157	0.24	0.33	0.28	129
158	0.36	0.25	0.30	143
159	0.46	0.53	0.49	107
160	0.16	0.41	0.23	78
161	0.14	0.30	0.19	73
162	0.09	0.15	0.11	106
163	0.14	0.12	0.13	126
164	0.23	0.37	0.28	63
165	0.00	0.00	0.00	229
166	0.38	0.31	0.34	115

167	0.55	0.13	0.21	46
168	0.22	0.23	0.23	69
169	0.37	0.54	0.44	70
170	0.56	0.56	0.56	54
171	0.00	0.00	0.00	43
172	0.27	0.41	0.32	76
173	0.21	0.50	0.30	12
174	0.12	0.21	0.15	76
175	0.47	0.57	0.52	91
176	0.67	0.64	0.65	157
177	0.23	0.27	0.25	41
178	0.00	0.00	0.00	0
179	0.05	1.00	0.10	1
180	0.18	0.29	0.23	55
181	0.05	0.08	0.06	62
182	0.03	0.50	0.06	2
183	0.27	0.46	0.34	80
184	0.00	0.00	0.00	206
185	0.33	0.23	0.27	86
186	0.20	0.44	0.28	66
187	0.52	0.63	0.57	59
188	0.31	0.54	0.40	68
189	0.18	0.16	0.17	108
190	0.08	0.13	0.10	85
191	0.24	0.40	0.30	86
192	0.17	0.57	0.26	46
193	0.14	0.33	0.19	18
194	0.23	0.32	0.27	74
195	0.25	0.44	0.32	55
196	0.68	0.71	0.69	38
197	0.32	0.37	0.34	95
198	0.17	0.38	0.24	16
199	0.05	0.08	0.06	39
200	0.28	0.16	0.20	58
201	0.06	0.22	0.10	55
202	0.26	0.29	0.27	58
203	0.18	0.05	0.07	66
204	0.74	0.75	0.74	64
205	0.00	0.00	0.00	10
206	0.06	0.20	0.09	66
207	0.20	0.22	0.21	73
208	0.12	0.13	0.13	54
209	0.10	0.13	0.12	61
210	0.12	0.25	0.17	12
211	0.07	0.10	0.08	59
212	0.29	0.54	0.38	26
213	0.19	0.32	0.24	105
214	0.20	0.46	0.28	50
215	0.29	0.17	0.21	65
216	0.35	0.33	0.34	79
217	0.18	0.25	0.21	55
218	0.06	0.67	0.11	3
219	0.11	0.18	0.13	62
220	0.11	0.17	0.13	81
221	0.18	0.21	0.19	34
222	0.01	0.02	0.01	64
223	0.54	0.48	0.50	61
224	0.13	0.39	0.19	18
225	0.18	0.70	0.29	10
226	0.62	0.66	0.64	99
227	0.25	0.69	0.37	13
228	0.09	0.11	0.10	74
229	0.63	0.74	0.68	50
230	0.23	0.24	0.24	74
231	0.00	0.00	0.00	4
232	0.26	0.31	0.28	26
233	0.03	0.08	0.04	146
234	0.52	0.72	0.61	61
235	0.02	0.08	0.03	13
236	0.23	0.16	0.19	49
237	0.55	0.40	0.46	90
238	0.16	0.10	0.12	58
239	0.07	0.25	0.10	24
240	0.61	0.58	0.59	64
241	0.58	0.83	0.68	75
242	0.44	0.59	0.50	63
243	0.45	0.58	0.51	76
244	0.35	0.33	0.34	63
245	0.08	0.20	0.11	41
246	0.90	0.11	0.20	162
247	0.17	0.32	0.22	22
248	0.43	0.63	0.52	52
249	0.16	0.37	0.23	19
250	0.39	0.65	0.49	23
251	0.47	0.51	0.49	57
252	0.21	0.25	0.23	36
253	0.02	0.02	0.02	41
254	0.05	0.20	0.08	10

255	0.12	0.14	0.13	22
256	0.07	0.38	0.12	8
257	0.09	0.13	0.10	62
258	0.14	0.21	0.17	43
259	0.52	0.55	0.54	87
260	0.00	0.00	0.00	56
261	0.00	0.00	0.00	3
262	0.19	0.40	0.25	20
263	0.17	0.07	0.10	15
264	0.03	0.06	0.04	50
265	0.14	0.20	0.17	25
266	0.10	0.19	0.13	47
267	0.46	0.48	0.47	97
268	0.68	0.72	0.70	36
269	0.44	0.41	0.43	56
270	0.30	0.58	0.40	38
271	0.03	0.05	0.04	58
272	0.15	0.50	0.24	8
273	0.08	0.15	0.10	27
274	0.08	0.19	0.11	123
275	0.22	0.29	0.25	69
276	0.39	0.25	0.30	112
277	0.11	0.10	0.10	31
278	0.15	0.21	0.18	29
279	0.16	0.29	0.21	38
280	0.30	0.44	0.36	50
281	0.41	0.45	0.43	20
282	0.74	0.87	0.80	45
283	0.06	0.13	0.08	15
284	0.23	0.22	0.22	74
285	0.21	0.24	0.22	46
286	0.15	0.07	0.10	29
287	0.04	0.04	0.04	54
288	0.35	0.52	0.42	33
289	0.00	0.00	0.00	26
290	0.49	0.71	0.58	41
291	0.07	0.25	0.11	24
292	0.13	0.17	0.15	40
293	0.23	0.45	0.31	33
294	0.11	0.06	0.08	31
295	0.06	0.02	0.03	47
296	0.09	0.21	0.13	33
297	0.13	0.20	0.16	45
298	0.05	0.07	0.05	59
299	0.07	0.14	0.09	51
300	0.15	0.24	0.19	49
301	0.62	0.47	0.54	38
302	0.34	0.39	0.37	28
303	0.14	0.50	0.22	16
304	0.11	0.22	0.15	32
305	0.09	0.21	0.12	24
306	0.18	0.05	0.07	44
307	0.07	0.50	0.12	6
308	0.00	0.00	0.00	48
309	0.56	0.47	0.51	49
310	0.14	0.11	0.12	38
311	0.27	0.13	0.17	62
312	0.09	0.04	0.05	27
313	0.09	0.06	0.07	49
314	0.20	0.25	0.22	24
315	0.33	0.03	0.06	59
316	0.07	0.30	0.12	10
317	0.20	0.37	0.26	67
318	0.14	0.42	0.21	12
319	0.00	0.00	0.00	14
320	0.06	0.25	0.10	12
321	0.10	0.33	0.16	9
322	0.52	0.52	0.52	23
323	0.64	0.55	0.59	33
324	0.41	0.49	0.44	57
325	0.13	0.16	0.15	25
326	0.06	0.02	0.03	44
327	0.10	0.30	0.15	27
328	0.07	0.09	0.08	34
329	0.31	0.57	0.40	7
330	0.20	0.05	0.07	22
331	0.18	0.32	0.23	25
332	0.93	0.25	0.39	106
333	0.42	0.31	0.36	84
334	0.00	0.00	0.00	36
335	0.12	0.54	0.20	13
336	0.00	0.00	0.00	37
337	0.07	0.11	0.09	38
338	0.70	0.84	0.76	44
339	0.04	0.15	0.06	34
340	0.23	0.50	0.31	40
341	0.24	0.43	0.31	23
342	0.12	0.36	0.18	11

343	0.28	0.58	0.38	12
344	0.10	0.24	0.14	25
345	0.03	1.00	0.05	1
346	0.21	0.15	0.17	41
347	0.05	0.17	0.08	46
348	0.02	0.11	0.04	19
349	0.28	0.42	0.34	38
350	0.31	0.39	0.35	33
351	0.29	0.19	0.23	53
352	0.00	0.00	0.00	49
353	0.24	0.19	0.21	27
354	0.12	0.16	0.14	31
355	0.42	0.42	0.42	12
356	0.11	0.24	0.15	33
357	0.47	0.71	0.57	24
358	0.39	0.35	0.37	34
359	0.33	0.42	0.37	33
360	0.13	0.11	0.12	47
361	0.13	0.54	0.22	39
362	0.69	0.66	0.68	38
363	0.08	0.18	0.11	17
364	0.01	0.03	0.01	33
365	0.29	0.08	0.12	26
366	0.10	0.21	0.13	19
367	0.30	0.08	0.13	98
368	0.24	0.47	0.32	38
369	0.59	0.57	0.58	28
370	0.33	0.27	0.30	15
371	0.14	0.18	0.16	22
372	0.57	0.33	0.42	12
373	0.33	0.50	0.40	6
374	0.24	0.16	0.19	31
375	0.19	0.08	0.11	38
376	0.11	0.02	0.04	42
377	0.06	0.17	0.09	23
378	0.08	0.50	0.14	4
379	0.14	0.19	0.16	37
380	0.05	0.33	0.09	6
381	0.41	0.50	0.45	18
382	0.30	0.28	0.29	40
383	0.01	0.02	0.01	53
384	0.29	0.36	0.32	25
385	0.19	0.28	0.22	53
386	0.70	1.00	0.82	14
387	0.34	0.39	0.36	88
388	0.00	0.00	0.00	16
389	0.00	0.00	0.00	8
390	0.00	0.00	0.00	37
391	0.89	0.65	0.76	52
392	0.10	0.12	0.11	17
393	0.56	0.41	0.47	37
394	0.00	0.00	0.00	19
395	0.00	0.00	0.00	9
396	0.04	0.14	0.06	14
397	0.73	0.55	0.63	29
398	0.30	0.42	0.35	38
399	0.56	0.66	0.60	38
400	0.07	0.11	0.09	36
401	0.26	0.18	0.21	56
402	0.81	0.65	0.72	20
403	0.02	0.09	0.03	11
404	0.37	0.56	0.44	27
405	0.68	0.75	0.72	57
406	0.00	0.00	0.00	95
407	0.10	0.08	0.09	25
408	0.17	0.27	0.21	11
409	0.03	0.07	0.04	27
410	0.14	0.45	0.22	11
411	0.20	0.17	0.18	53
412	0.58	0.35	0.44	31
413	0.50	0.21	0.29	29
414	0.10	0.11	0.11	27
415	0.03	0.10	0.05	30
416	0.23	0.10	0.14	31
417	0.23	0.30	0.26	10
418	0.00	0.00	0.00	23
419	0.30	0.50	0.37	6
420	0.16	0.41	0.23	22
421	0.00	0.00	0.00	1
422	0.02	0.03	0.03	59
423	0.00	0.00	0.00	38
424	0.67	0.03	0.05	76
425	0.07	0.26	0.11	19
426	0.06	0.13	0.08	15
427	0.38	0.67	0.48	48
428	0.31	0.32	0.32	28
429	0.47	0.42	0.45	40
430	0.59	0.45	0.51	29

431	0.00	0.00	0.00	43
432	0.24	0.21	0.22	19
433	0.02	0.03	0.02	34
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	2
436	0.21	0.07	0.11	40
437	0.30	0.24	0.26	38
438	0.41	0.62	0.49	26
439	0.12	0.11	0.12	36
440	0.57	0.30	0.39	27
441	0.48	0.53	0.50	19
442	0.50	0.67	0.57	21
443	0.16	0.11	0.13	35
444	0.10	0.17	0.12	18
445	0.26	0.28	0.27	25
446	0.72	0.69	0.71	49
447	0.13	0.13	0.13	71
448	0.05	0.05	0.05	19
449	0.23	0.24	0.23	55
450	0.12	0.06	0.08	52
451	0.00	0.00	0.00	25
452	0.81	0.62	0.70	40
453	0.00	0.00	0.00	14
454	0.12	0.20	0.15	15
455	0.12	0.06	0.08	18
456	0.12	0.67	0.20	6
457	0.22	0.23	0.22	22
458	0.12	0.17	0.14	18
459	0.24	0.55	0.34	29
460	0.00	0.00	0.00	24
461	0.13	0.36	0.19	14
462	0.55	0.23	0.32	26
463	0.50	0.09	0.15	22
464	0.91	0.75	0.82	40
465	0.28	0.20	0.23	41
466	0.19	0.14	0.16	42
467	0.21	0.20	0.20	51
468	0.20	0.27	0.23	37
469	0.03	0.20	0.04	5
470	0.20	0.21	0.21	19
471	0.34	0.26	0.29	43
472	0.00	0.00	0.00	55
473	0.21	0.24	0.22	29
474	0.82	0.75	0.78	24
475	0.59	0.66	0.62	68
476	0.15	0.11	0.12	38
477	0.25	0.27	0.26	22
478	0.16	0.17	0.16	53
479	0.12	0.08	0.10	26
480	0.00	0.00	0.00	64
481	0.15	0.12	0.13	26
482	0.13	0.57	0.21	7
483	0.14	0.15	0.15	13
484	0.15	0.26	0.19	23
485	0.45	0.31	0.37	29
486	0.47	0.39	0.43	23
487	0.47	0.29	0.36	31
488	0.31	0.27	0.29	30
489	0.13	0.14	0.14	36
490	0.07	0.12	0.09	16
491	0.01	0.03	0.01	39
492	0.09	0.27	0.13	11
493	0.21	0.44	0.28	25
494	0.03	0.07	0.05	15
495	0.23	0.67	0.34	9
496	0.33	0.53	0.41	19
497	0.00	0.00	0.00	72
498	0.38	0.42	0.40	19
499	0.38	0.28	0.32	32
micro avg	0.37	0.41	0.39	60294
macro avg	0.25	0.31	0.26	60294
weighted avg	0.41	0.41	0.40	60294
samples avg	0.39	0.41	0.36	60294

4.5.11 Applying Linear SVM (SGD with Hinge Loss) with OneVsRest Classifier(1 to 4grams)

Hyper parameter Tuning


```
In [33]: start = datetime.now()

classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', penalty='l1'), n_jobs=1)

parameters = {'estimator__alpha':[2*10**-5, 4*10**-5, 6*10**-5, 8*10**-5]}

clf = GridSearchCV(classifier, parameters, cv= 3, return_train_score=True, scoring='f1_micro', n_jobs=-1)

clf.fit(x_train_multilabel, y_train)

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

Time taken to run this cell : 1:31:54.491088

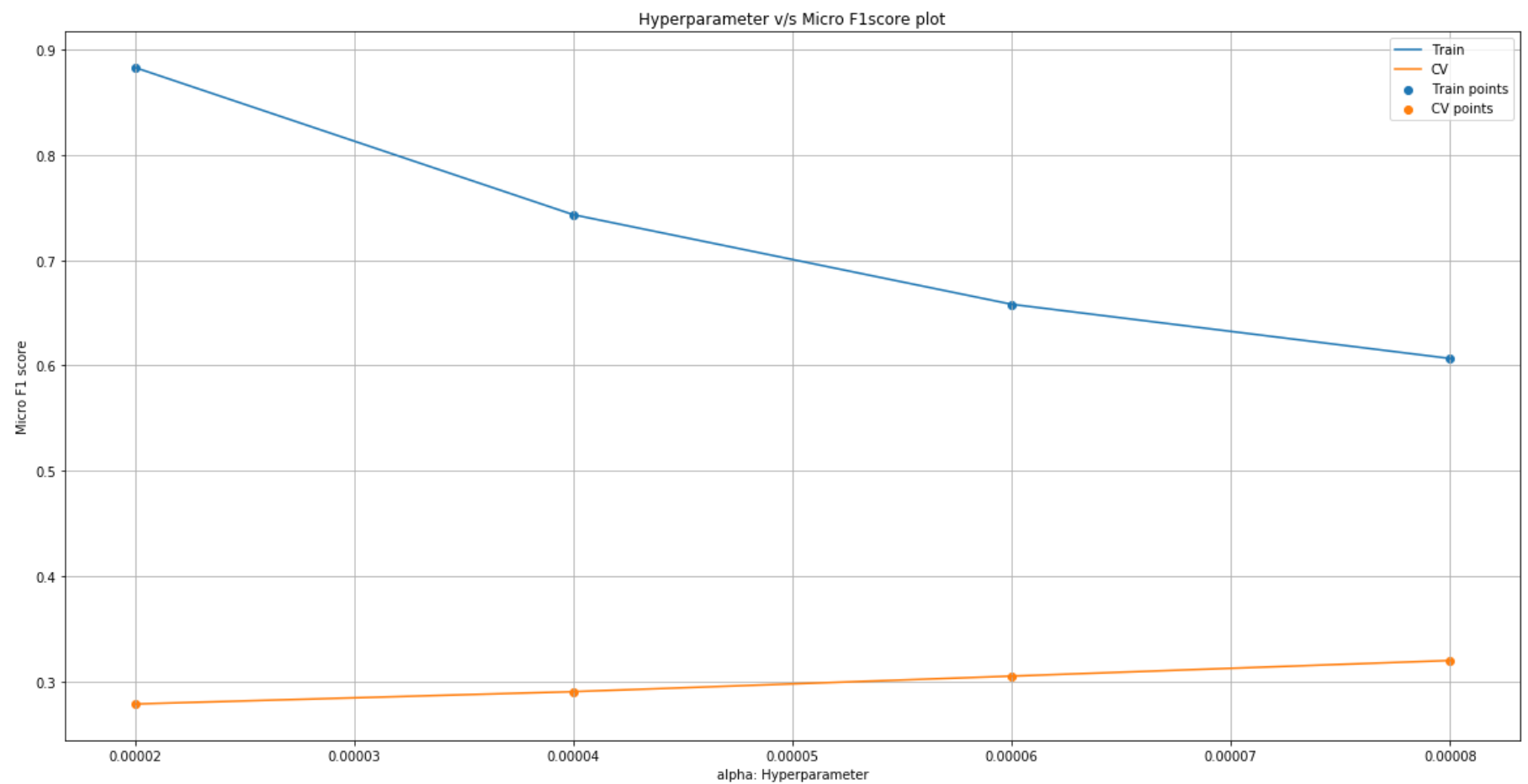
```
In [34]: plt.figure(figsize=(20,10))

train_auc= clf.cv_results_['mean_train_score']
cv_auc = clf.cv_results_['mean_test_score']

plt.plot(parameters['estimator__alpha'], train_auc, label='Train')
plt.plot(parameters['estimator__alpha'], cv_auc, label='CV')

plt.scatter(parameters['estimator__alpha'], train_auc, label='Train points')
plt.scatter(parameters['estimator__alpha'], cv_auc, label='CV points')

plt.legend()
plt.xlabel("alpha: Hyperparameter")
plt.ylabel("Micro F1 score")
plt.title("Hyperparameter v/s Micro F1score plot")
plt.grid()
plt.show()
```



Model A with L1 Loss

```
In [35]: start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=8*10**-5, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:07:16.262753

Hamming loss 0.0053706

Micro-average quality numbers

Precision: 0.3535, Recall: 0.4053, F1-measure: 0.3776

Macro-average quality numbers

Precision: 0.2280, Recall: 0.3106, F1-measure: 0.2475

	precision	recall	f1-score	support
0	0.82	0.72	0.77	6668
1	0.40	0.28	0.33	3659
2	0.14	0.16	0.15	971
3	0.51	0.55	0.53	1506
4	0.55	0.47	0.50	1649
5	0.57	0.48	0.52	1113
6	0.53	0.45	0.49	1482
7	0.56	0.60	0.58	980
8	0.76	0.77	0.77	1520
9	0.52	0.84	0.64	1041
10	0.45	0.43	0.44	861
11	0.23	0.40	0.30	386
12	0.07	0.54	0.13	37
13	0.55	0.43	0.48	917
14	0.28	0.22	0.24	519
15	0.25	0.27	0.26	656
16	0.42	0.24	0.31	794
17	0.40	0.32	0.36	700
18	0.40	0.55	0.46	363
19	0.58	0.64	0.61	541
20	0.35	0.48	0.41	540
21	0.54	0.64	0.58	362
22	0.63	0.55	0.59	551
23	0.25	0.34	0.29	309
24	0.29	0.40	0.34	331
25	0.26	0.36	0.30	424
26	0.35	0.29	0.32	465
27	0.17	0.07	0.10	386
28	0.15	0.18	0.16	107
29	0.06	0.08	0.07	195
30	0.29	0.36	0.32	758
31	0.11	0.53	0.18	15
32	0.41	0.58	0.48	323
33	0.30	0.34	0.32	279
34	0.37	0.37	0.37	275
35	0.59	0.49	0.54	268
36	0.09	0.11	0.10	76
37	0.10	0.23	0.14	269
38	0.41	0.51	0.45	255
39	0.24	0.32	0.28	249
40	0.16	0.26	0.20	66
41	0.20	0.20	0.20	209
42	0.30	0.38	0.33	72
43	0.28	0.42	0.34	430
44	0.17	0.21	0.18	279
45	0.33	0.35	0.34	240
46	0.26	0.38	0.31	157
47	0.48	0.58	0.53	249
48	0.13	0.11	0.12	198
49	0.30	0.34	0.32	171
50	0.55	0.68	0.61	200
51	0.60	0.75	0.67	85
52	0.32	0.33	0.32	175
53	0.08	0.30	0.13	114
54	0.07	0.18	0.10	223
55	0.21	0.32	0.25	122
56	0.48	0.58	0.52	168
57	0.10	0.05	0.06	176
58	0.14	0.23	0.17	140
59	0.21	0.21	0.21	191
60	0.64	0.75	0.69	152
61	0.27	0.23	0.24	208
62	0.16	0.15	0.15	136
63	0.44	0.47	0.45	158
64	0.61	0.50	0.55	203
65	0.28	0.42	0.33	105
66	0.22	0.50	0.30	58
67	0.43	0.54	0.48	128
68	0.14	0.12	0.13	158
69	0.21	0.15	0.17	248
70	0.19	0.14	0.16	201
71	0.24	0.28	0.26	89
72	0.40	0.38	0.39	157
73	0.19	0.48	0.27	29
74	0.09	0.17	0.12	58
75	0.14	0.16	0.15	158
76	0.57	0.65	0.61	110
77	0.25	0.45	0.32	33
78	0.12	0.23	0.16	210
79	0.37	0.54	0.44	169

80	0.05	0.40	0.08	15
81	0.27	0.39	0.32	214
82	0.15	0.29	0.20	65
83	0.15	0.27	0.20	156
84	0.33	0.49	0.40	59
85	0.37	0.65	0.47	55
86	0.13	0.28	0.18	36
87	0.22	0.41	0.29	29
88	0.38	0.63	0.47	54
89	0.45	0.76	0.56	137
90	0.13	0.19	0.15	103
91	0.14	0.23	0.17	79
92	0.08	0.14	0.10	84
93	0.37	0.50	0.42	133
94	0.67	0.41	0.51	318
95	0.34	0.75	0.47	51
96	0.39	0.49	0.43	82
97	0.11	0.16	0.13	75
98	0.05	0.03	0.03	120
99	0.18	0.39	0.25	18
100	0.54	0.39	0.46	196
101	0.44	0.44	0.44	208
102	0.08	0.16	0.11	122
103	0.04	0.03	0.03	62
104	0.07	0.06	0.06	88
105	0.31	0.43	0.36	65
106	0.19	0.20	0.19	115
107	0.02	0.07	0.03	29
108	0.31	0.18	0.23	109
109	0.28	0.23	0.25	73
110	0.33	0.25	0.29	102
111	0.40	0.41	0.40	180
112	0.76	0.09	0.15	292
113	0.55	0.85	0.67	54
114	0.12	0.11	0.12	120
115	0.25	0.36	0.30	107
116	0.30	0.37	0.33	52
117	0.06	0.12	0.08	72
118	0.56	0.56	0.56	139
119	0.47	0.44	0.45	57
120	0.38	0.52	0.44	44
121	0.07	0.18	0.10	85
122	0.33	0.57	0.42	82
123	0.04	0.07	0.05	100
124	0.16	0.75	0.26	4
125	0.14	0.67	0.23	9
126	0.07	0.22	0.11	46
127	0.13	0.17	0.15	54
128	0.87	0.76	0.81	195
129	0.29	0.46	0.36	54
130	0.11	0.16	0.13	96
131	0.39	0.66	0.49	35
132	0.08	0.14	0.10	58
133	0.16	0.14	0.15	36
134	0.19	0.39	0.26	36
135	0.42	0.69	0.52	39
136	0.04	0.01	0.02	97
137	0.20	0.41	0.27	70
138	0.16	0.29	0.20	17
139	0.14	0.22	0.17	119
140	0.48	0.65	0.55	101
141	0.24	0.46	0.32	115
142	0.47	0.32	0.38	94
143	0.35	0.57	0.43	84
144	0.45	0.62	0.52	64
145	0.04	0.03	0.03	61
146	0.14	0.16	0.15	132
147	0.25	0.32	0.28	119
148	0.50	0.60	0.54	62
149	0.21	0.16	0.18	83
150	0.07	0.18	0.11	72
151	0.09	0.35	0.14	23
152	0.06	0.14	0.08	76
153	0.10	0.44	0.16	18
154	0.05	0.24	0.08	17
155	0.00	0.00	0.00	24
156	0.31	0.15	0.20	136
157	0.30	0.40	0.34	129
158	0.14	0.15	0.15	143
159	0.53	0.60	0.56	107
160	0.24	0.37	0.29	78
161	0.15	0.29	0.20	73
162	0.07	0.08	0.07	106
163	0.12	0.13	0.12	126
164	0.26	0.49	0.34	63
165	0.00	0.00	0.00	229
166	0.33	0.27	0.30	115
167	0.17	0.17	0.17	46

168	0.25	0.22	0.23	69
169	0.38	0.47	0.42	70
170	0.62	0.54	0.57	54
171	0.00	0.00	0.00	43
172	0.28	0.34	0.31	76
173	0.06	0.25	0.10	12
174	0.18	0.16	0.17	76
175	0.34	0.57	0.43	91
176	0.59	0.64	0.61	157
177	0.20	0.27	0.23	41
178	0.00	0.00	0.00	0
179	0.05	1.00	0.10	1
180	0.23	0.36	0.28	55
181	0.07	0.11	0.09	62
182	0.09	0.50	0.15	2
183	0.37	0.45	0.41	80
184	0.00	0.00	0.00	206
185	0.32	0.23	0.27	86
186	0.34	0.33	0.34	66
187	0.65	0.68	0.66	59
188	0.39	0.59	0.47	68
189	0.20	0.16	0.18	108
190	0.17	0.25	0.20	85
191	0.31	0.36	0.33	86
192	0.15	0.57	0.24	46
193	0.21	0.28	0.24	18
194	0.21	0.35	0.27	74
195	0.13	0.38	0.19	55
196	0.45	0.71	0.55	38
197	0.33	0.32	0.32	95
198	0.09	0.31	0.14	16
199	0.06	0.13	0.08	39
200	0.29	0.09	0.13	58
201	0.06	0.25	0.10	55
202	0.15	0.22	0.18	58
203	0.08	0.05	0.06	66
204	0.49	0.66	0.56	64
205	0.00	0.00	0.00	10
206	0.05	0.15	0.08	66
207	0.13	0.21	0.16	73
208	0.12	0.15	0.13	54
209	0.09	0.21	0.12	61
210	0.21	0.58	0.31	12
211	0.13	0.25	0.17	59
212	0.33	0.50	0.40	26
213	0.21	0.28	0.24	105
214	0.32	0.52	0.39	50
215	0.11	0.11	0.11	65
216	0.41	0.48	0.44	79
217	0.21	0.35	0.26	55
218	0.06	0.33	0.11	3
219	0.05	0.13	0.07	62
220	0.07	0.15	0.10	81
221	0.14	0.24	0.17	34
222	0.04	0.06	0.05	64
223	0.53	0.48	0.50	61
224	0.11	0.33	0.16	18
225	0.18	0.70	0.29	10
226	0.61	0.70	0.65	99
227	0.28	0.54	0.37	13
228	0.06	0.12	0.08	74
229	0.51	0.72	0.60	50
230	0.15	0.18	0.16	74
231	0.00	0.00	0.00	4
232	0.23	0.35	0.28	26
233	0.04	0.01	0.02	146
234	0.44	0.70	0.54	61
235	0.15	0.15	0.15	13
236	0.16	0.18	0.17	49
237	0.45	0.42	0.43	90
238	0.05	0.05	0.05	58
239	0.11	0.25	0.15	24
240	0.75	0.66	0.70	64
241	0.50	0.72	0.59	75
242	0.40	0.46	0.43	63
243	0.51	0.41	0.45	76
244	0.25	0.35	0.29	63
245	0.06	0.05	0.05	41
246	0.00	0.00	0.00	162
247	0.06	0.05	0.05	22
248	0.62	0.60	0.61	52
249	0.18	0.53	0.26	19
250	0.36	0.65	0.46	23
251	0.30	0.47	0.37	57
252	0.09	0.22	0.13	36
253	0.10	0.10	0.10	41
254	0.03	0.10	0.05	10
255	0.05	0.18	0.08	22

256	0.09	0.50	0.15	8
257	0.06	0.11	0.07	62
258	0.30	0.26	0.28	43
259	0.39	0.54	0.46	87
260	0.00	0.00	0.00	56
261	0.00	0.00	0.00	3
262	0.15	0.35	0.21	20
263	0.05	0.20	0.07	15
264	0.04	0.08	0.05	50
265	0.08	0.20	0.11	25
266	0.11	0.21	0.15	47
267	0.36	0.44	0.40	97
268	0.43	0.69	0.53	36
269	0.32	0.45	0.37	56
270	0.42	0.61	0.49	38
271	0.02	0.03	0.03	58
272	0.08	0.25	0.12	8
273	0.19	0.22	0.20	27
274	0.14	0.14	0.14	123
275	0.29	0.32	0.31	69
276	0.36	0.17	0.23	112
277	0.10	0.10	0.10	31
278	0.13	0.14	0.13	29
279	0.13	0.21	0.16	38
280	0.23	0.32	0.27	50
281	0.42	0.55	0.48	20
282	0.74	0.76	0.75	45
283	0.12	0.27	0.17	15
284	0.20	0.34	0.26	74
285	0.11	0.22	0.15	46
286	0.04	0.03	0.04	29
287	0.08	0.11	0.09	54
288	0.22	0.64	0.33	33
289	0.03	0.08	0.04	26
290	0.49	0.63	0.55	41
291	0.04	0.17	0.07	24
292	0.12	0.17	0.14	40
293	0.18	0.45	0.26	33
294	0.04	0.06	0.05	31
295	0.00	0.00	0.00	47
296	0.05	0.21	0.08	33
297	0.09	0.16	0.12	45
298	0.15	0.03	0.06	59
299	0.17	0.18	0.17	51
300	0.19	0.22	0.21	49
301	0.29	0.29	0.29	38
302	0.38	0.54	0.45	28
303	0.18	0.50	0.26	16
304	0.16	0.22	0.18	32
305	0.12	0.25	0.16	24
306	0.04	0.11	0.06	44
307	0.04	0.17	0.06	6
308	0.00	0.00	0.00	48
309	0.47	0.41	0.43	49
310	0.03	0.08	0.04	38
311	0.13	0.13	0.13	62
312	0.02	0.04	0.03	27
313	0.08	0.08	0.08	49
314	0.14	0.29	0.19	24
315	0.10	0.07	0.08	59
316	0.11	0.20	0.14	10
317	0.21	0.24	0.22	67
318	0.21	0.58	0.30	12
319	0.00	0.00	0.00	14
320	0.05	0.17	0.07	12
321	0.19	0.56	0.28	9
322	0.37	0.43	0.40	23
323	0.43	0.61	0.51	33
324	0.41	0.53	0.46	57
325	0.20	0.20	0.20	25
326	0.14	0.07	0.09	44
327	0.12	0.33	0.17	27
328	0.10	0.21	0.14	34
329	0.11	0.29	0.16	7
330	0.05	0.09	0.06	22
331	0.06	0.24	0.09	25
332	0.92	0.42	0.58	106
333	0.48	0.38	0.42	84
334	0.04	0.03	0.03	36
335	0.05	0.15	0.08	13
336	0.03	0.03	0.03	37
337	0.00	0.00	0.00	38
338	0.80	0.75	0.78	44
339	0.06	0.06	0.06	34
340	0.25	0.45	0.32	40
341	0.40	0.52	0.45	23
342	0.14	0.45	0.21	11
343	0.21	0.50	0.30	12

344	0.12	0.32	0.18	25
345	0.03	1.00	0.06	1
346	0.08	0.12	0.10	41
347	0.04	0.20	0.06	46
348	0.05	0.11	0.07	19
349	0.17	0.45	0.24	38
350	0.16	0.55	0.24	33
351	0.23	0.34	0.27	53
352	0.02	0.02	0.02	49
353	0.25	0.22	0.24	27
354	0.03	0.03	0.03	31
355	0.11	0.58	0.18	12
356	0.09	0.18	0.12	33
357	0.33	0.67	0.44	24
358	0.26	0.38	0.31	34
359	0.26	0.48	0.34	33
360	0.14	0.11	0.12	47
361	0.28	0.49	0.36	39
362	0.68	0.74	0.71	38
363	0.11	0.24	0.15	17
364	0.06	0.03	0.04	33
365	0.02	0.04	0.03	26
366	0.14	0.21	0.17	19
367	0.32	0.12	0.18	98
368	0.29	0.53	0.37	38
369	0.34	0.50	0.41	28
370	0.06	0.13	0.09	15
371	0.11	0.23	0.15	22
372	0.17	0.08	0.11	12
373	0.14	0.50	0.21	6
374	0.13	0.16	0.14	31
375	0.11	0.11	0.11	38
376	0.10	0.07	0.08	42
377	0.04	0.09	0.06	23
378	0.14	0.50	0.22	4
379	0.05	0.05	0.05	37
380	0.07	0.33	0.12	6
381	0.33	0.56	0.42	18
382	0.27	0.35	0.31	40
383	0.03	0.06	0.04	53
384	0.16	0.44	0.23	25
385	0.16	0.26	0.20	53
386	0.53	0.57	0.55	14
387	0.34	0.42	0.38	88
388	0.06	0.12	0.08	16
389	0.07	0.25	0.11	8
390	0.01	0.03	0.01	37
391	0.38	0.48	0.42	52
392	0.03	0.06	0.04	17
393	0.55	0.59	0.57	37
394	0.00	0.00	0.00	19
395	0.03	0.11	0.05	9
396	0.12	0.21	0.15	14
397	0.59	0.66	0.62	29
398	0.24	0.39	0.30	38
399	0.54	0.71	0.61	38
400	0.33	0.17	0.22	36
401	0.23	0.05	0.09	56
402	0.79	0.75	0.77	20
403	0.03	0.18	0.05	11
404	0.47	0.59	0.52	27
405	0.70	0.75	0.73	57
406	0.00	0.00	0.00	95
407	0.07	0.12	0.08	25
408	0.04	0.09	0.06	11
409	0.04	0.07	0.05	27
410	0.19	0.55	0.29	11
411	0.17	0.09	0.12	53
412	0.38	0.45	0.41	31
413	0.18	0.24	0.21	29
414	0.04	0.11	0.05	27
415	0.11	0.23	0.15	30
416	0.08	0.06	0.07	31
417	0.43	0.30	0.35	10
418	0.03	0.09	0.05	23
419	0.20	0.50	0.29	6
420	0.15	0.45	0.22	22
421	0.00	0.00	0.00	1
422	0.02	0.02	0.02	59
423	0.03	0.05	0.04	38
424	0.46	0.08	0.13	76
425	0.12	0.21	0.15	19
426	0.02	0.07	0.03	15
427	0.39	0.56	0.46	48
428	0.27	0.32	0.30	28
429	0.53	0.40	0.46	40
430	0.32	0.34	0.33	29
431	0.00	0.00	0.00	43

432	0.27	0.16	0.20	19
433	0.03	0.03	0.03	34
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	2
436	0.12	0.17	0.14	40
437	0.23	0.29	0.26	38
438	0.31	0.62	0.41	26
439	0.08	0.11	0.09	36
440	0.25	0.22	0.24	27
441	0.23	0.42	0.30	19
442	0.48	0.71	0.58	21
443	0.23	0.20	0.21	35
444	0.09	0.22	0.13	18
445	0.14	0.32	0.19	25
446	0.31	0.76	0.44	49
447	0.18	0.17	0.17	71
448	0.02	0.11	0.04	19
449	0.24	0.24	0.24	55
450	0.03	0.08	0.04	52
451	0.00	0.00	0.00	25
452	0.63	0.55	0.59	40
453	0.00	0.00	0.00	14
454	0.00	0.00	0.00	15
455	0.02	0.06	0.03	18
456	0.11	0.50	0.18	6
457	0.19	0.18	0.19	22
458	0.02	0.06	0.03	18
459	0.56	0.52	0.54	29
460	0.05	0.04	0.04	24
461	0.15	0.36	0.21	14
462	0.35	0.42	0.39	26
463	0.25	0.05	0.08	22
464	0.85	0.70	0.77	40
465	0.20	0.22	0.21	41
466	0.13	0.24	0.17	42
467	0.27	0.29	0.28	51
468	0.22	0.16	0.19	37
469	0.06	0.40	0.10	5
470	0.09	0.21	0.13	19
471	0.47	0.40	0.43	43
472	0.04	0.04	0.04	55
473	0.23	0.31	0.26	29
474	0.59	0.71	0.64	24
475	0.63	0.71	0.67	68
476	0.25	0.34	0.29	38
477	0.25	0.45	0.32	22
478	0.18	0.25	0.21	53
479	0.07	0.12	0.09	26
480	0.04	0.09	0.05	64
481	0.05	0.15	0.07	26
482	0.07	0.43	0.12	7
483	0.33	0.38	0.36	13
484	0.22	0.39	0.28	23
485	0.23	0.21	0.22	29
486	0.26	0.30	0.28	23
487	0.16	0.23	0.19	31
488	0.12	0.20	0.15	30
489	0.29	0.31	0.30	36
490	0.03	0.06	0.04	16
491	0.00	0.00	0.00	39
492	0.06	0.18	0.09	11
493	0.22	0.48	0.30	25
494	0.00	0.00	0.00	15
495	0.10	0.33	0.15	9
496	0.07	0.11	0.09	19
497	0.00	0.00	0.00	72
498	0.26	0.42	0.32	19
499	0.50	0.31	0.38	32
micro avg	0.35	0.41	0.38	60294
macro avg	0.23	0.31	0.25	60294
weighted avg	0.40	0.41	0.39	60294
samples avg	0.38	0.41	0.35	60294

Model B with L2 Loss


```
In [36]: start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=8*10**-5, penalty='l2'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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

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))
```

Time taken to run this cell : 0:01:48.287269

Hamming loss 0.0038696

Micro-average quality numbers

Precision: 0.5286, Recall: 0.3451, F1-measure: 0.4175

Macro-average quality numbers

Precision: 0.3901, Recall: 0.2320, F1-measure: 0.2749

	precision	recall	f1-score	support
0	0.83	0.71	0.76	6668
1	0.38	0.27	0.31	3659
2	0.13	0.16	0.14	971
3	0.60	0.52	0.56	1506
4	0.57	0.48	0.52	1649
5	0.54	0.48	0.51	1113
6	0.53	0.35	0.42	1482
7	0.57	0.56	0.57	980
8	0.82	0.59	0.69	1520
9	0.57	0.73	0.64	1041
10	0.49	0.41	0.45	861
11	0.38	0.37	0.37	386
12	0.34	0.57	0.43	37
13	0.60	0.34	0.44	917
14	0.30	0.16	0.21	519
15	0.34	0.28	0.31	656
16	0.46	0.19	0.27	794
17	0.42	0.25	0.31	700
18	0.67	0.56	0.61	363
19	0.77	0.52	0.62	541
20	0.43	0.37	0.40	540
21	0.68	0.55	0.61	362
22	0.75	0.40	0.53	551
23	0.40	0.27	0.32	309
24	0.39	0.33	0.36	331
25	0.35	0.31	0.33	424
26	0.46	0.29	0.36	465
27	0.20	0.12	0.15	386
28	0.21	0.20	0.20	107
29	0.07	0.07	0.07	195
30	0.68	0.29	0.40	758
31	0.12	0.13	0.13	15
32	0.64	0.56	0.60	323
33	0.39	0.35	0.37	279
34	0.47	0.31	0.37	275
35	0.70	0.43	0.53	268
36	0.06	0.04	0.05	76
37	0.21	0.09	0.13	269
38	0.54	0.40	0.46	255
39	0.38	0.27	0.32	249
40	0.31	0.23	0.26	66
41	0.26	0.08	0.12	209
42	0.42	0.39	0.40	72
43	0.34	0.12	0.18	430
44	0.31	0.17	0.22	279
45	0.47	0.19	0.27	240
46	0.47	0.30	0.37	157
47	0.75	0.46	0.57	249
48	0.19	0.08	0.11	198
49	0.47	0.32	0.38	171
50	0.83	0.65	0.73	200
51	0.74	0.66	0.70	85
52	0.52	0.35	0.42	175
53	0.18	0.32	0.23	114
54	0.19	0.12	0.14	223
55	0.34	0.21	0.26	122
56	0.60	0.45	0.52	168
57	0.15	0.02	0.04	176
58	0.31	0.26	0.29	140
59	0.44	0.13	0.20	191
60	0.90	0.62	0.74	152
61	0.32	0.15	0.20	208
62	0.17	0.07	0.10	136
63	0.62	0.41	0.49	158
64	0.59	0.30	0.40	203
65	0.41	0.34	0.37	105
66	0.30	0.28	0.29	58
67	0.51	0.45	0.48	128
68	0.13	0.05	0.07	158
69	0.20	0.06	0.10	248
70	0.24	0.08	0.13	201
71	0.37	0.25	0.30	89
72	0.47	0.39	0.43	157
73	0.29	0.31	0.30	29
74	0.19	0.07	0.10	58
75	0.29	0.13	0.18	158
76	0.79	0.52	0.63	110
77	0.47	0.45	0.46	33
78	0.16	0.10	0.13	210
79	0.61	0.48	0.54	169

80	0.06	0.07	0.06	15
81	0.48	0.35	0.40	214
82	0.34	0.28	0.31	65
83	0.27	0.19	0.22	156
84	0.56	0.46	0.50	59
85	0.66	0.60	0.63	55
86	0.24	0.19	0.22	36
87	0.44	0.48	0.46	29
88	0.55	0.56	0.55	54
89	0.69	0.70	0.69	137
90	0.16	0.12	0.13	103
91	0.33	0.18	0.23	79
92	0.32	0.08	0.13	84
93	0.60	0.53	0.56	133
94	0.77	0.34	0.47	318
95	0.60	0.55	0.57	51
96	0.52	0.29	0.38	82
97	0.18	0.04	0.07	75
98	0.00	0.00	0.00	120
99	0.41	0.39	0.40	18
100	0.50	0.43	0.46	196
101	0.66	0.30	0.42	208
102	0.33	0.14	0.20	122
103	0.00	0.00	0.00	62
104	0.08	0.02	0.04	88
105	0.64	0.35	0.46	65
106	0.33	0.15	0.20	115
107	0.07	0.07	0.07	29
108	0.37	0.12	0.18	109
109	0.58	0.29	0.39	73
110	0.56	0.19	0.28	102
111	0.60	0.29	0.39	180
112	0.67	0.01	0.01	292
113	0.86	0.69	0.76	54
114	0.24	0.03	0.06	120
115	0.36	0.22	0.28	107
116	0.58	0.21	0.31	52
117	0.26	0.07	0.11	72
118	0.65	0.37	0.47	139
119	0.67	0.25	0.36	57
120	0.67	0.36	0.47	44
121	0.06	0.01	0.02	85
122	0.65	0.41	0.51	82
123	0.12	0.02	0.03	100
124	0.50	1.00	0.67	4
125	0.75	0.67	0.71	9
126	0.30	0.13	0.18	46
127	0.17	0.04	0.06	54
128	0.90	0.67	0.76	195
129	0.53	0.39	0.45	54
130	0.07	0.01	0.02	96
131	0.65	0.63	0.64	35
132	0.25	0.07	0.11	58
133	0.17	0.06	0.08	36
134	0.43	0.28	0.34	36
135	0.67	0.51	0.58	39
136	0.00	0.00	0.00	97
137	0.33	0.51	0.40	70
138	0.42	0.47	0.44	17
139	0.25	0.13	0.17	119
140	0.90	0.52	0.66	101
141	0.39	0.30	0.33	115
142	0.62	0.16	0.25	94
143	0.60	0.52	0.56	84
144	0.47	0.52	0.49	64
145	0.08	0.05	0.06	61
146	0.13	0.02	0.03	132
147	0.43	0.28	0.34	119
148	0.57	0.45	0.50	62
149	0.31	0.20	0.25	83
150	0.22	0.15	0.18	72
151	0.27	0.26	0.27	23
152	0.15	0.16	0.15	76
153	0.50	0.33	0.40	18
154	0.25	0.18	0.21	17
155	0.11	0.04	0.06	24
156	0.44	0.13	0.20	136
157	0.34	0.16	0.21	129
158	0.35	0.19	0.24	143
159	0.66	0.40	0.50	107
160	0.48	0.31	0.38	78
161	0.46	0.15	0.23	73
162	0.13	0.02	0.03	106
163	0.10	0.03	0.05	126
164	0.51	0.29	0.37	63
165	0.00	0.00	0.00	229
166	0.54	0.17	0.26	115
167	0.31	0.11	0.16	46

168	0.61	0.16	0.25	69
169	0.50	0.33	0.40	70
170	0.78	0.26	0.39	54
171	0.00	0.00	0.00	43
172	0.41	0.18	0.25	76
173	0.27	0.25	0.26	12
174	0.26	0.11	0.15	76
175	0.66	0.53	0.59	91
176	0.76	0.52	0.61	157
177	0.41	0.27	0.32	41
178	0.00	0.00	0.00	0
179	0.50	1.00	0.67	1
180	0.53	0.33	0.40	55
181	0.17	0.06	0.09	62
182	0.00	0.00	0.00	2
183	0.47	0.33	0.39	80
184	0.00	0.00	0.00	206
185	0.30	0.08	0.13	86
186	0.35	0.20	0.25	66
187	0.94	0.51	0.66	59
188	0.76	0.47	0.58	68
189	0.32	0.11	0.16	108
190	0.21	0.14	0.17	85
191	0.78	0.16	0.27	86
192	0.58	0.57	0.57	46
193	0.50	0.28	0.36	18
194	0.46	0.28	0.35	74
195	0.69	0.40	0.51	55
196	0.67	0.47	0.55	38
197	0.55	0.24	0.34	95
198	0.36	0.25	0.30	16
199	0.19	0.08	0.11	39
200	0.56	0.09	0.15	58
201	0.32	0.13	0.18	55
202	0.21	0.10	0.14	58
203	0.07	0.02	0.02	66
204	0.89	0.61	0.72	64
205	0.00	0.00	0.00	10
206	0.21	0.09	0.13	66
207	0.33	0.05	0.09	73
208	0.33	0.04	0.07	54
209	0.16	0.10	0.12	61
210	0.26	0.42	0.32	12
211	0.18	0.05	0.08	59
212	0.69	0.35	0.46	26
213	0.26	0.18	0.21	105
214	0.54	0.38	0.45	50
215	0.50	0.12	0.20	65
216	0.50	0.22	0.30	79
217	0.48	0.25	0.33	55
218	0.00	0.00	0.00	3
219	0.17	0.05	0.08	62
220	0.08	0.01	0.02	81
221	0.09	0.03	0.04	34
222	0.12	0.02	0.03	64
223	0.82	0.38	0.52	61
224	0.31	0.28	0.29	18
225	0.75	0.60	0.67	10
226	0.77	0.51	0.61	99
227	0.78	0.54	0.64	13
228	0.11	0.05	0.07	74
229	0.85	0.66	0.74	50
230	0.23	0.12	0.16	74
231	0.00	0.00	0.00	4
232	0.57	0.15	0.24	26
233	0.03	0.01	0.01	146
234	0.67	0.46	0.54	61
235	0.12	0.15	0.13	13
236	0.30	0.12	0.17	49
237	0.70	0.21	0.32	90
238	0.17	0.02	0.03	58
239	0.22	0.25	0.24	24
240	0.89	0.48	0.63	64
241	0.83	0.52	0.64	75
242	0.58	0.51	0.54	63
243	0.58	0.46	0.51	76
244	0.45	0.32	0.37	63
245	0.23	0.07	0.11	41
246	1.00	0.02	0.05	162
247	0.25	0.09	0.13	22
248	0.76	0.50	0.60	52
249	0.38	0.42	0.40	19
250	0.55	0.48	0.51	23
251	0.57	0.21	0.31	57
252	0.43	0.36	0.39	36
253	0.09	0.05	0.06	41
254	0.00	0.00	0.00	10
255	0.20	0.14	0.16	22

256	0.23	0.38	0.29	8
257	0.00	0.00	0.00	62
258	0.36	0.09	0.15	43
259	0.67	0.43	0.52	87
260	0.00	0.00	0.00	56
261	0.00	0.00	0.00	3
262	0.46	0.30	0.36	20
263	0.00	0.00	0.00	15
264	0.09	0.02	0.03	50
265	0.23	0.24	0.24	25
266	0.26	0.15	0.19	47
267	0.54	0.39	0.45	97
268	0.80	0.56	0.66	36
269	0.64	0.25	0.36	56
270	0.62	0.42	0.50	38
271	0.08	0.05	0.06	58
272	0.43	0.38	0.40	8
273	0.14	0.04	0.06	27
274	0.42	0.04	0.07	123
275	0.35	0.12	0.17	69
276	0.58	0.10	0.17	112
277	0.17	0.10	0.12	31
278	0.22	0.07	0.11	29
279	0.46	0.16	0.24	38
280	0.55	0.34	0.42	50
281	0.91	0.50	0.65	20
282	0.97	0.64	0.77	45
283	0.44	0.27	0.33	15
284	0.42	0.22	0.29	74
285	0.23	0.07	0.10	46
286	0.22	0.07	0.11	29
287	0.06	0.02	0.03	54
288	0.93	0.42	0.58	33
289	0.25	0.04	0.07	26
290	0.85	0.54	0.66	41
291	0.12	0.04	0.06	24
292	0.50	0.07	0.13	40
293	0.37	0.48	0.42	33
294	0.15	0.06	0.09	31
295	0.00	0.00	0.00	47
296	0.25	0.09	0.13	33
297	0.40	0.04	0.08	45
298	0.00	0.00	0.00	59
299	0.00	0.00	0.00	51
300	0.27	0.06	0.10	49
301	0.52	0.39	0.45	38
302	0.80	0.43	0.56	28
303	0.28	0.31	0.29	16
304	0.41	0.22	0.29	32
305	0.40	0.17	0.24	24
306	0.08	0.09	0.08	44
307	0.40	0.33	0.36	6
308	0.00	0.00	0.00	48
309	0.60	0.31	0.41	49
310	0.09	0.05	0.07	38
311	0.36	0.13	0.19	62
312	0.00	0.00	0.00	27
313	0.38	0.06	0.11	49
314	0.22	0.08	0.12	24
315	0.20	0.02	0.03	59
316	0.50	0.10	0.17	10
317	0.39	0.25	0.31	67
318	0.50	0.42	0.45	12
319	0.00	0.00	0.00	14
320	0.50	0.17	0.25	12
321	0.50	0.33	0.40	9
322	0.50	0.22	0.30	23
323	0.70	0.42	0.53	33
324	0.65	0.49	0.56	57
325	0.33	0.16	0.22	25
326	0.25	0.02	0.04	44
327	0.29	0.07	0.12	27
328	0.11	0.03	0.05	34
329	0.29	0.29	0.29	7
330	0.10	0.05	0.06	22
331	0.31	0.20	0.24	25
332	1.00	0.23	0.37	106
333	0.55	0.42	0.47	84
334	0.00	0.00	0.00	36
335	0.57	0.31	0.40	13
336	0.25	0.03	0.05	37
337	0.00	0.00	0.00	38
338	0.91	0.66	0.76	44
339	0.14	0.06	0.08	34
340	0.40	0.35	0.37	40
341	0.87	0.57	0.68	23
342	0.50	0.36	0.42	11
343	0.60	0.50	0.55	12

344	0.36	0.16	0.22	25
345	0.00	0.00	0.00	1
346	0.33	0.10	0.15	41
347	0.16	0.07	0.09	46
348	0.50	0.05	0.10	19
349	0.39	0.32	0.35	38
350	0.40	0.24	0.30	33
351	0.35	0.13	0.19	53
352	0.00	0.00	0.00	49
353	0.30	0.22	0.26	27
354	0.50	0.03	0.06	31
355	0.71	0.42	0.53	12
356	0.56	0.15	0.24	33
357	0.67	0.58	0.62	24
358	0.46	0.35	0.40	34
359	0.70	0.42	0.53	33
360	0.12	0.04	0.06	47
361	0.53	0.46	0.49	39
362	0.84	0.55	0.67	38
363	0.20	0.18	0.19	17
364	0.17	0.03	0.05	33
365	0.25	0.08	0.12	26
366	0.40	0.11	0.17	19
367	0.12	0.01	0.02	98
368	0.52	0.37	0.43	38
369	1.00	0.29	0.44	28
370	0.57	0.27	0.36	15
371	0.13	0.14	0.13	22
372	0.29	0.17	0.21	12
373	0.25	0.33	0.29	6
374	0.57	0.26	0.36	31
375	0.00	0.00	0.00	38
376	0.11	0.02	0.04	42
377	0.12	0.09	0.10	23
378	0.20	0.50	0.29	4
379	0.20	0.05	0.09	37
380	0.29	0.33	0.31	6
381	0.31	0.28	0.29	18
382	0.60	0.15	0.24	40
383	0.00	0.00	0.00	53
384	0.50	0.24	0.32	25
385	0.36	0.08	0.12	53
386	0.88	0.50	0.64	14
387	0.46	0.39	0.42	88
388	0.33	0.06	0.11	16
389	0.00	0.00	0.00	8
390	0.00	0.00	0.00	37
391	0.92	0.46	0.62	52
392	0.00	0.00	0.00	17
393	0.50	0.16	0.24	37
394	0.00	0.00	0.00	19
395	0.40	0.22	0.29	9
396	0.00	0.00	0.00	14
397	1.00	0.48	0.65	29
398	0.44	0.45	0.44	38
399	0.96	0.66	0.78	38
400	0.25	0.08	0.12	36
401	0.67	0.07	0.13	56
402	0.89	0.40	0.55	20
403	0.33	0.18	0.24	11
404	0.79	0.41	0.54	27
405	0.97	0.58	0.73	57
406	0.00	0.00	0.00	95
407	0.12	0.12	0.12	25
408	0.33	0.36	0.35	11
409	0.00	0.00	0.00	27
410	0.20	0.18	0.19	11
411	0.14	0.04	0.06	53
412	0.75	0.48	0.59	31
413	0.60	0.10	0.18	29
414	0.00	0.00	0.00	27
415	0.38	0.10	0.16	30
416	0.12	0.03	0.05	31
417	0.33	0.20	0.25	10
418	0.00	0.00	0.00	23
419	0.29	0.33	0.31	6
420	0.53	0.36	0.43	22
421	0.00	0.00	0.00	1
422	0.14	0.05	0.07	59
423	0.00	0.00	0.00	38
424	1.00	0.01	0.03	76
425	1.00	0.21	0.35	19
426	0.00	0.00	0.00	15
427	0.63	0.75	0.69	48
428	0.60	0.21	0.32	28
429	0.58	0.38	0.45	40
430	0.67	0.07	0.12	29
431	0.00	0.00	0.00	43

432	0.50	0.11	0.17	19
433	0.33	0.03	0.05	34
434	0.00	0.00	0.00	0
435	0.00	0.00	0.00	2
436	0.33	0.10	0.15	40
437	0.47	0.18	0.26	38
438	0.74	0.54	0.62	26
439	0.12	0.03	0.05	36
440	0.20	0.04	0.06	27
441	0.30	0.32	0.31	19
442	0.92	0.57	0.71	21
443	0.15	0.06	0.08	35
444	0.00	0.00	0.00	18
445	0.20	0.08	0.11	25
446	0.73	0.45	0.56	49
447	0.12	0.03	0.05	71
448	0.44	0.21	0.29	19
449	0.55	0.20	0.29	55
450	0.10	0.02	0.03	52
451	0.00	0.00	0.00	25
452	0.69	0.28	0.39	40
453	0.00	0.00	0.00	14
454	0.33	0.13	0.19	15
455	0.00	0.00	0.00	18
456	0.14	0.17	0.15	6
457	0.43	0.14	0.21	22
458	0.14	0.06	0.08	18
459	0.85	0.38	0.52	29
460	0.00	0.00	0.00	24
461	0.27	0.21	0.24	14
462	0.40	0.08	0.13	26
463	0.00	0.00	0.00	22
464	0.89	0.40	0.55	40
465	0.44	0.17	0.25	41
466	0.25	0.07	0.11	42
467	0.44	0.37	0.40	51
468	0.25	0.08	0.12	37
469	0.00	0.00	0.00	5
470	0.20	0.11	0.14	19
471	0.69	0.21	0.32	43
472	0.00	0.00	0.00	55
473	0.35	0.24	0.29	29
474	0.82	0.58	0.68	24
475	0.55	0.44	0.49	68
476	0.29	0.05	0.09	38
477	0.50	0.18	0.27	22
478	0.33	0.06	0.10	53
479	0.00	0.00	0.00	26
480	0.67	0.03	0.06	64
481	0.50	0.08	0.13	26
482	0.33	0.29	0.31	7
483	0.40	0.15	0.22	13
484	0.67	0.35	0.46	23
485	0.70	0.48	0.57	29
486	0.43	0.26	0.32	23
487	0.77	0.32	0.45	31
488	0.57	0.13	0.22	30
489	0.30	0.17	0.21	36
490	0.20	0.12	0.15	16
491	0.00	0.00	0.00	39
492	0.22	0.18	0.20	11
493	0.41	0.36	0.38	25
494	0.00	0.00	0.00	15
495	0.71	0.56	0.63	9
496	0.50	0.58	0.54	19
497	0.00	0.00	0.00	72
498	0.31	0.21	0.25	19
499	0.24	0.12	0.16	32
micro avg	0.53	0.35	0.42	60294
macro avg	0.39	0.23	0.27	60294
weighted avg	0.50	0.35	0.39	60294
samples avg	0.41	0.34	0.34	60294

5. Conclusions

In [37]: `# http://zetcode.com/python/prettytable/`

```
from prettytable import PrettyTable
```

```
#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable
```

```
x = PrettyTable()
```

```
x.field_names = ["Model", "No. of Data Points", "Tags", "Vectorizer", "MicroF1 score", "Macro F1 score"]
```

```
x.add_row(["Logistic Regression", 500000, 500, "Count(1 to 3 grams)", 0.4996, 0.3855])
```

```
x.add_row(["Linear SVM(L1)", 500000, 500, "Count(1 to 3 grams)", 0.4850, 0.3264])
```

```
x.add_row(["-----", "-----", "-----", "-----", "-----", "-----"])
```

```
x.add_row(["Logistic Regression", 100000, 500, "Count(1 to 4 grams)", 0.3863, 0.2588])
```

```
x.add_row(["Linear SVM(L1)", 100000, 500, "Count(1 to 4 grams)", 0.3776, 0.2475])
```

```
x.add_row(["Linear SVM(L2)", 100000, 500, "Count(1 to 4 grams)", 0.4175, 0.2749])
```

```
print(x)
```

Model	No. of Data Points	Tags	Vectorizer	MicroF1 score	Macro F1 score
Logistic Regression	500000	500	Count(1 to 3 grams)	0.4996	0.3855
Linear SVM(L1)	500000	500	Count(1 to 3 grams)	0.485	0.3264
-----	-----	-----	-----	-----	-----
Logistic Regression	100000	500	Count(1 to 4 grams)	0.3863	0.2588
Linear SVM(L1)	100000	500	Count(1 to 4 grams)	0.3776	0.2475
Linear SVM(L2)	100000	500	Count(1 to 4 grams)	0.4175	0.2749