



In [1]:

```
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/>  
(<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-separated format (all lowercase, should not contain tabs '\t' or ampersands '&')

### 2.1.2 Example Data point

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

**Body :**

```

#include<
iostream>\n
#include<
stdlib.h>\n\n
using namespace std;\n\n
int main()\n
{\n
    int n,a[n],x,c,u[n],m[n],e[n][4];\n
    cout<<"Enter the number of variables";\n          cin>>n;\n
\n
    cout<<"Enter the Lower, and Upper Limits of the variable
s";\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> (<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:

$$F1 = 2 (\text{precision recall}) / (\text{precision} + \text{recall})$$

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) ([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=chu
        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 t

```

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
    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.

```

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

In [0]:

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

Out[6]:

	Title	Body	Tags	cnt_dup
0	Implementing Boundary Value Analysis of S...	<pre>\n#include<iostream>\n#include<...>	c++ c	1
1	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamicall...	c#\nsilverlight\ndata-binding	1
2	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamicall...	c#\nsilverlight\ndata-binding\ncolumns	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>\n<code>...	java jdbc	2

In [0]:

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

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[8]:

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





In [0]:

```

#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

```

Time taken to run this cell : 0:00:52.992676

## 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 [0]:

```

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 [0]:

```

#'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 tages we have : ['.a', '.app', '.asp.net-mvc', '.aspxauth', '.ba  
sh-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 [0]:

```
# 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]:

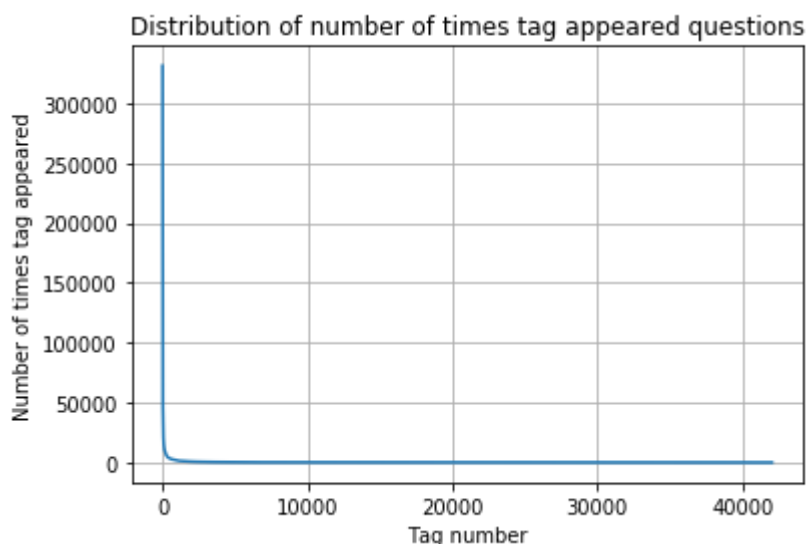
	Tags	Counts
0	.a	18
1	.app	37
2	.asp.net-mvc	1
3	.aspxauth	21
4	.bash-profile	138

In [0]:

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

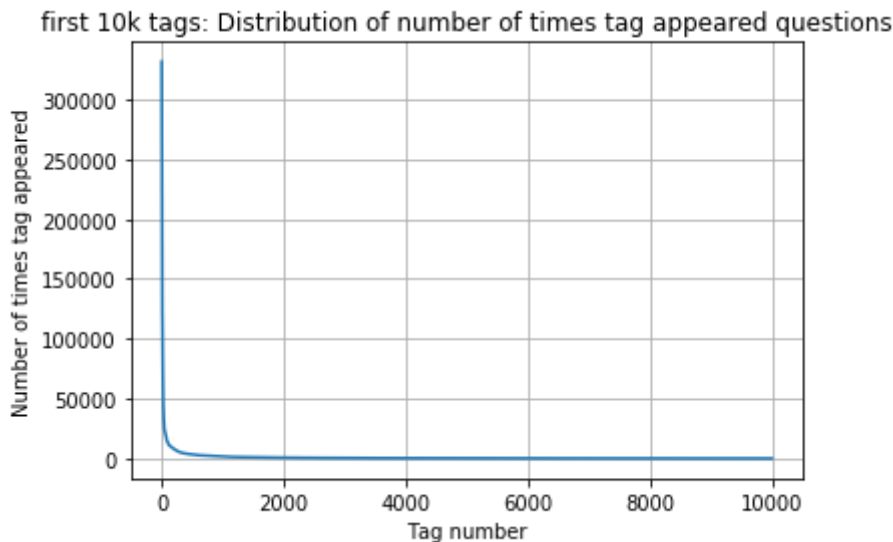
In [0]:

```
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 [0]:

```
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])
```

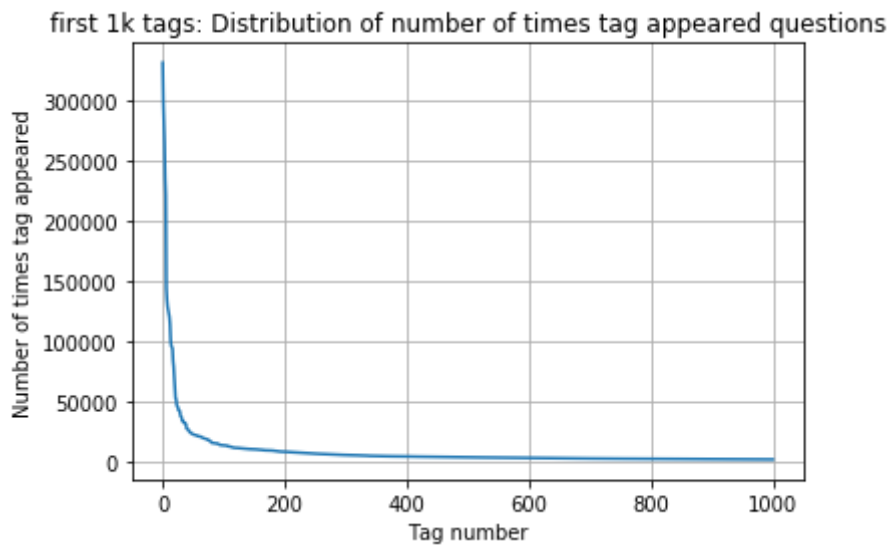


400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7151
6466	5865	5370	4983	4526	4281	4144	3929	3750	3593	
3453	3299	3123	2989	2891	2738	2647	2527	2431	2331	
2259	2186	2097	2020	1959	1900	1828	1770	1723	1673	
1631	1574	1532	1479	1448	1406	1365	1328	1300	1266	
1245	1222	1197	1181	1158	1139	1121	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 [0]:

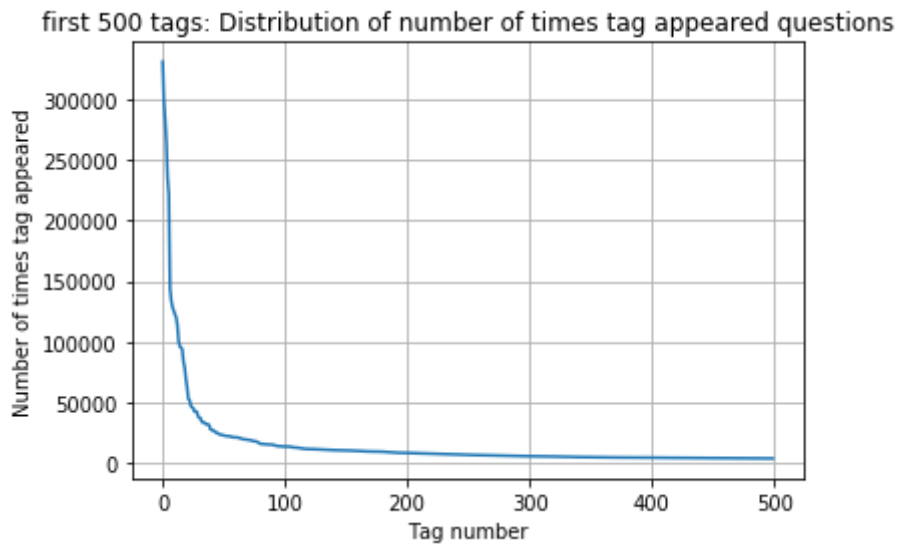
```
plt.plot(tag_counts[0:1000])
plt.title('first 1k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



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

In [0]:

```
plt.plot(tag_counts[0:500])
plt.title('first 500 tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:500:5]), tag_counts[0:500:5])
```



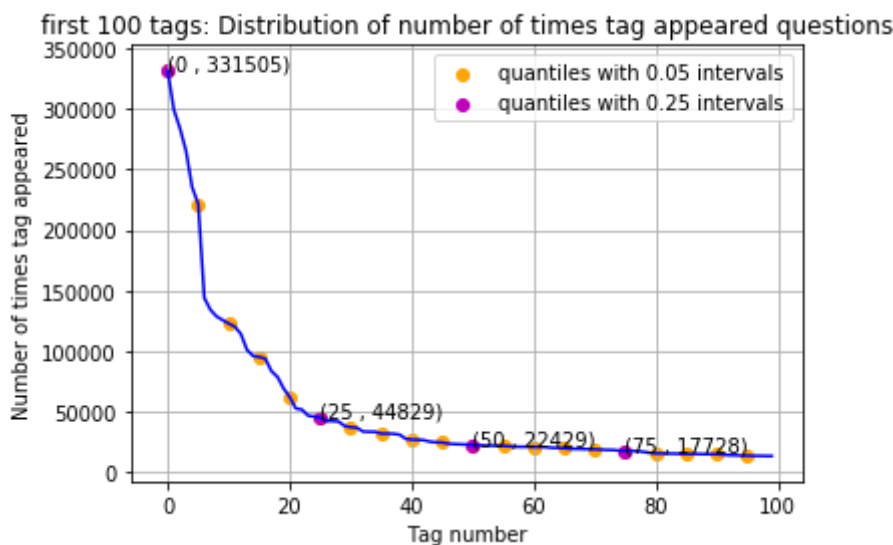
```
100 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291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

In [0]:

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

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 [0]:

```
# 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 [0]:

```
#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
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]

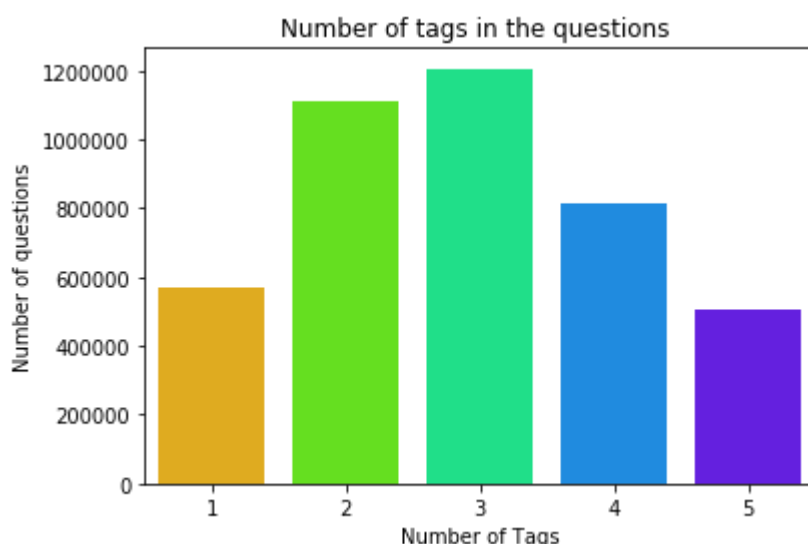
In [0]:

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

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

In [0]:

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



#### Observations:

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



## In [0]:

[illegible]

**Observations:**

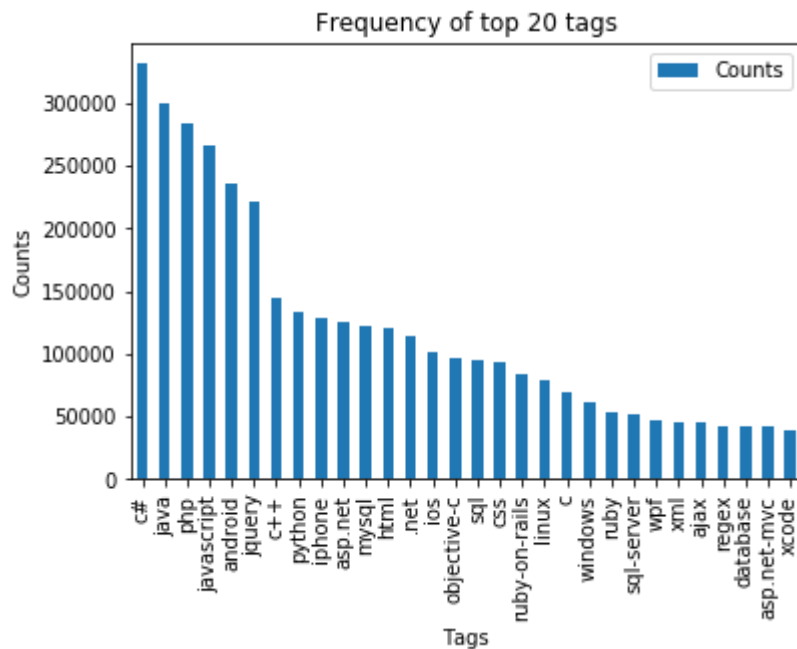
### 3.2.6 The top 20 tags

In [0]:

```

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

In [2]:

```

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

```

In [16]:

```
#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)"""
create_database_table("Processed.db", sql_create_table)
```

Tables in the databse:  
QuestionsProcessed

## 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

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

In [15]:

```
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]:

```
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 [4]:

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

Tables in the database:  
QuestionsProcessed

In [5]:

```
# 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 [7]:

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

start = datetime.now()
preprocessed_data_list=[]
reader.fetchone()
questions_with_code=0
len_pre=0
len_post=0
questions_proccesed = 0
for row in reader:

    is_code = 0

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

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

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

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

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

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

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

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

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

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

    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,
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_pr
```

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

```
[nltk_data] Downloading package punkt to  
[nltk_data] C:\Users\hp\AppData\Roaming\nltk_data...  
[nltk_data] Unzipping tokenizers\punkt.zip.  
number of questions completed= 100000  
number of questions completed= 200000  
number of questions completed= 300000  
number of questions completed= 400000  
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 : 0:19:23.878200
```

In [8]:

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

### Sample quesitons after preprocessing of data

In [9]:

```

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

```

=====
=====
('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index ja
va.sql.sqlexcept microsoft odbc driver manag invalid descriptor index java.s
ql.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 w
ay updat feed fb php sdk novic facebook api read mani tutori still confused.
i find post feed api method like correct second way use curl someth like way
better',)
-----
-----
('btnadd click event open two window record ad btnadd click event open two w
indow record ad btnadd click event open two window record ad open window sea
rch.aspx use code hav add button search.aspx nwhen insert record btnadd clic
k event open anoth window nafter insert record close window',)
-----
-----
('sql inject issu prevent correct form submiss php sql inject issu prevent c
orrect form submiss php sql inject issu prevent correct form submiss php che
ck everyth think make sure input field safe type sql inject good news safe b
ad news one tag mess form submiss place even touch life figur exact html use
templat file forgiv okay entir php script get execut see data post none foru
m field post problem use someth titl field none data get post current use pr
int post see submit noth work flawless statement though also mention script
work flawless local machin use host come across problem state list input tes
t mess',)
-----
-----
('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl
subaddit lebesgu measur let lbrace rbrace sequenc set sigma -algebra mathcal
want show left bigcup right leq sum left right countabl addit measur defin s
et sigma algebra mathcal think use monoton properti somewher proof start app
reci littl help nthank ad han answer make follow addit construct given han a
nswer clear bigcup bigcup cap emptyset neq left bigcup right left bigcup rig
ht sum left right also construct subset monoton left right leq left right fi
nal would sum leq sum result follow',)
-----
-----
('hql equival sql queri hql equival sql queri hql equival sql queri hql quer
i replac name class properti name error occur hql error',)
-----

```



```

-----
('undefin symbol architectur i386 objc class skpsmtpmessag referenc error un
defin symbol architectur i386 objc class skpsmtpmessag referenc error undefi
n symbol architectur i386 objc class skpsmtpmessag referenc error import fra
mework send email applic background import framework i.e skpsmtpmessag someb
odi suggest get error collect2 ld return exit status import framework correc
t sorc taken framework follow mfmailcomposeviewcontrol question lock field u
pdat answer drag drop folder project click copi nthat',)
-----

```

```

-----
('java.lang.nosuchmethoderror javax.servlet.servletcontext.geteffectivesessi
ontrackingmod ljava util set java.lang.nosuchmethoderror javax.servlet.servl
etcontext.geteffectivesessiontrackingmod ljava util set java.lang.nosuchmeth
oderror javax.servlet.servletcontext.geteffectivesessiontrackingmod ljava ut
il set want servlet process input standalon java program deploy servlet jbos
s put servlet.class file web-inf class web.xml gave servlet url map .do java
client program open connect servlet use url object use localhost 8080 .do ge
t folow error error org.apache.catalina.connector.coyoteadapt except error o
ccur contain request process java.lang.nosuchmethoderror javax.servlet.servl
etcontext.geteffectivesessiontrackingmod ljava util set org.apache.catalina.
connector.coyoteadapter.postparserequest coyoteadapter.java 567 org.apache.c
atalina.connector.coyoteadapter.servic coyoteadapter.java 359 org.apache.coy
ote.http11.http11processor.process http11processor.java 877 org.apache.coyot
e.http11.http11protocol http11connectionhandler.process http11protocol.java
654 org.apache.tomcat.util.net.jioendpoint worker.run jioendpoint.java 951 w
eb.xml file content',)
-----

```

```

-----
('obtain updat locat use gps servic obtain updat locat use gps servic obtain
updat locat use gps servic app two button start track stop track strart trac
k button click gps start listen locat stop listen use besid toast everi new
updat locat want thing use background servic alway updat locat even activ cl
osed.a toast appear everi new updat location.pleas hint link would apprec
i',)
-----

```

## Saving Preprocessed data to a Database

In [10]:

```

#Taking 0.5 Million entries to a dataframe.
write_db = 'Titlmoreweight.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 QuestionsProces
conn_r.commit()
conn_r.close()

```

In [11]:

```
preprocessed_data.head()
```

Out[11]:

	question	tags
0	java.lang.noclassdeffounderror javax servlet j...	jsp jstl
1	java.sql.sqlexcept microsoft odbc driver manag...	java jdbc
2	better way updat feed fb php sdk better way up...	facebook api facebook-php-sdk
3	btnadd click event open two window record ad b...	javascript asp.net web
4	sql inject issu prevent correct form submiss p...	php forms

In [12]:

```
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 : 499998

number of dimensions : 2

## Converting string Tags to multilable output variables

In [13]:

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

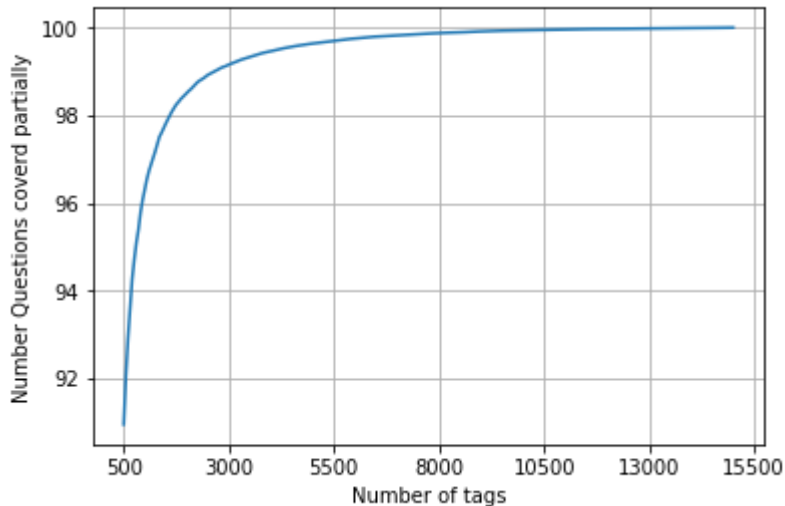
## Selecting 500 Tags

In [17]:

```
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
```

In [19]:

```
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
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 [20]:

```
# 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 ",
number of questions that are not covered : 45221 out of 499998
```

In [21]:

```
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 [22]:

```
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 : (99998, 500)

## 5. Assignments

# USE BoW WITH UPTO 4 GRAMS AND COMPUTE micro F1 score with L.R(1 Vs Rest)

## 5.1 Featurizing data with Count vectorizer

In [28]:

```
start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009, max_features=200000, tokenizer = lambda x: x.split())

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 : 0:06:31.754828

In [30]:

```
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, 95586) Y : (400000, 500)

Dimensions of test data X: (99998, 95586) Y: (99998, 500)

## Logistic Regression with OneVsRestClassifier using count vectorizer

Hyperparameter tuning(Using grid search) :-

In [32]:

```

from sklearn.model_selection import GridSearchCV
param={'estimator__alpha': [10**-9, 10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1]}
classifier = OneVsRestClassifier(SGDClassifier(loss='log', penalty='l1'))
gsv = GridSearchCV(estimator = classifier, param_grid=param, verbose=0, scoring='f1_micro',
gsv.fit(x_train_multilabel, y_train)

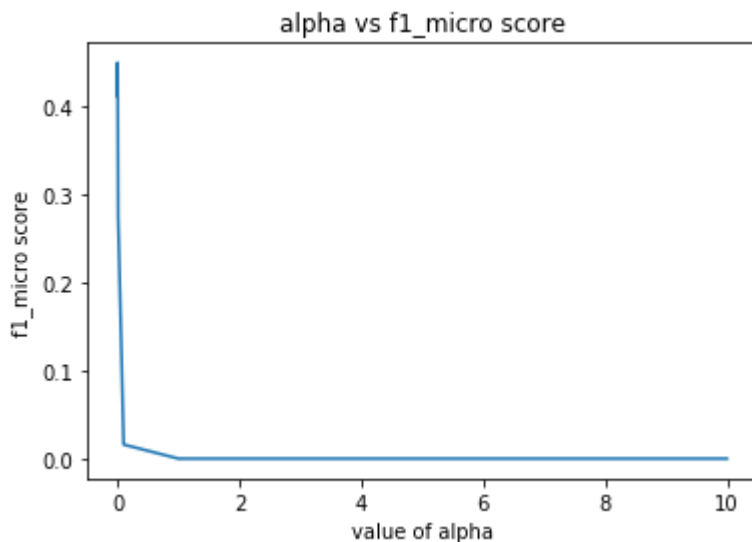
best_alpha = gsv.best_estimator_.get_params()['estimator__alpha']
print('value of alpha after hyperparameter tuning : ',best_alpha)
print('-----')
# plotting C vs f1_micro_score
x_1=[]
y_1=[]
for x in gsv.grid_scores_:
    x_1.append(x[0]['estimator__alpha'])
    y_1.append(x[1])
plt.plot(x_1,y_1)
plt.xlabel('value of alpha')
plt.ylabel('f1_micro score')
plt.title('alpha vs f1_micro score')
plt.show()

```

value of alpha after hyperparameter tuning : 0.001

-----

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\model\_selection\\_search.py:761: DeprecationWarning: The grid\_scores\_ attribute was deprecated in version 0.18 in favor of the more elaborate cv\_results\_ attribute. The grid\_scores\_ attribute will not be available from 0.20  
DeprecationWarning)



**Applying model using best hyperparameter:-**

In [33]:

```

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.001, penalty='l1'), n_jobs=4)
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 : 0.1868137362747255

Hamming loss 0.003220684413688274

Micro-average quality numbers

Precision: 0.5624, Recall: 0.3312, F1-measure: 0.4169

Macro-average quality numbers

Precision: 0.4101, Recall: 0.2382, F1-measure: 0.2832

	precision	recall	f1-score	support
0	0.78	0.66	0.71	5519
1	0.40	0.36	0.38	8189
2	0.77	0.35	0.48	6529
3	0.56	0.48	0.52	3231
4	0.67	0.47	0.55	6430
5	0.71	0.29	0.42	2878
6	0.72	0.56	0.63	5086
7	0.76	0.61	0.68	4533
8	0.57	0.14	0.22	3000
9	0.68	0.57	0.62	2765
10	0.45	0.18	0.26	3051
11	0.67	0.36	0.47	3009
12	0.56	0.25	0.35	2630
13	0.56	0.18	0.27	1425
14	0.80	0.61	0.69	2548
15	0.73	0.09	0.15	2371
16	0.52	0.34	0.41	873
17	0.78	0.65	0.71	2151
18	0.47	0.27	0.34	2204
19	0.47	0.45	0.46	831
20	0.79	0.35	0.48	1860
21	0.24	0.12	0.16	2023
22	0.38	0.21	0.27	1513
23	0.83	0.56	0.67	1207

24	0.44	0.35	0.39	506
25	0.68	0.37	0.48	425
26	0.52	0.41	0.46	793
27	0.54	0.33	0.41	1291
28	0.72	0.30	0.43	1208
29	0.27	0.07	0.11	406
30	0.57	0.15	0.24	504
31	0.26	0.17	0.20	732
32	0.56	0.30	0.39	441
33	0.28	0.19	0.22	1645
34	0.50	0.40	0.44	1058
35	0.69	0.66	0.68	946
36	0.51	0.26	0.35	644
37	0.77	0.82	0.79	136
38	0.56	0.38	0.45	570
39	0.81	0.31	0.45	766
40	0.51	0.27	0.35	1132
41	0.32	0.29	0.30	174
42	0.64	0.48	0.55	210
43	0.55	0.50	0.52	433
44	0.42	0.53	0.47	626
45	0.52	0.35	0.42	852
46	0.66	0.38	0.48	534
47	0.24	0.18	0.21	350
48	0.71	0.46	0.56	496
49	0.78	0.55	0.64	785
50	0.18	0.10	0.13	475
51	0.24	0.12	0.16	305
52	0.29	0.06	0.10	251
53	0.66	0.38	0.49	914
54	0.29	0.18	0.22	728
55	0.00	0.00	0.00	258
56	0.36	0.13	0.19	821
57	0.40	0.08	0.14	541
58	0.80	0.25	0.38	748
59	0.91	0.64	0.75	724
60	0.24	0.06	0.10	660
61	0.84	0.25	0.38	235
62	0.84	0.81	0.82	718
63	0.83	0.55	0.66	468
64	0.53	0.45	0.49	191
65	0.22	0.11	0.14	429
66	0.16	0.10	0.13	415
67	0.73	0.49	0.59	274
68	0.78	0.57	0.66	510
69	0.51	0.53	0.52	466
70	0.27	0.09	0.13	305
71	0.38	0.16	0.23	247
72	0.59	0.58	0.58	401
73	0.92	0.78	0.84	86
74	0.37	0.41	0.39	120
75	0.83	0.78	0.80	129
76	0.07	0.01	0.01	473
77	0.28	0.28	0.28	143
78	0.76	0.45	0.56	347
79	0.64	0.27	0.38	479
80	0.36	0.33	0.34	279
81	0.81	0.12	0.20	461
82	0.09	0.01	0.02	298
83	0.75	0.41	0.53	396
84	0.38	0.35	0.36	184

85	0.48	0.17	0.25	573
86	0.24	0.08	0.12	325
87	0.50	0.23	0.31	273
88	0.28	0.20	0.23	135
89	0.20	0.19	0.19	232
90	0.40	0.39	0.39	409
91	0.46	0.44	0.45	420
92	0.65	0.64	0.64	408
93	0.62	0.45	0.52	241
94	0.30	0.09	0.13	211
95	0.21	0.14	0.17	277
96	0.21	0.04	0.07	410
97	0.88	0.22	0.35	501
98	0.53	0.62	0.57	136
99	0.48	0.25	0.33	239
100	0.40	0.07	0.12	324
101	0.91	0.55	0.68	277
102	0.91	0.67	0.78	613
103	0.43	0.19	0.26	157
104	0.16	0.21	0.18	295
105	0.65	0.36	0.46	334
106	0.51	0.08	0.14	335
107	0.76	0.47	0.58	389
108	0.45	0.21	0.29	251
109	0.51	0.42	0.46	317
110	0.41	0.05	0.09	187
111	0.24	0.11	0.15	140
112	0.11	0.04	0.06	154
113	0.45	0.37	0.41	332
114	0.33	0.23	0.27	323
115	0.32	0.14	0.19	344
116	0.71	0.44	0.54	370
117	0.47	0.17	0.25	313
118	0.79	0.52	0.63	874
119	0.36	0.19	0.25	293
120	0.00	0.00	0.00	200
121	0.66	0.47	0.55	463
122	0.18	0.20	0.19	119
123	0.00	0.00	0.00	256
124	0.89	0.71	0.79	195
125	0.33	0.22	0.26	138
126	0.73	0.47	0.57	376
127	0.13	0.07	0.09	122
128	0.12	0.11	0.11	252
129	0.00	0.00	0.00	144
130	0.10	0.02	0.03	150
131	0.17	0.01	0.02	210
132	0.17	0.14	0.15	361
133	0.89	0.52	0.66	453
134	0.67	0.83	0.74	124
135	0.00	0.00	0.00	91
136	0.32	0.26	0.28	128
137	0.41	0.35	0.38	218
138	0.00	0.00	0.00	243
139	0.32	0.24	0.28	149
140	0.77	0.37	0.50	318
141	0.10	0.18	0.13	159
142	0.58	0.32	0.41	274
143	0.69	0.81	0.75	362
144	0.51	0.21	0.30	118
145	0.52	0.35	0.42	164



146	0.56	0.25	0.35	461
147	0.62	0.43	0.51	159
148	0.33	0.14	0.20	166
149	0.95	0.52	0.67	346
150	0.43	0.09	0.14	350
151	0.80	0.60	0.69	55
152	0.75	0.44	0.55	387
153	0.35	0.06	0.10	150
154	0.69	0.07	0.13	281
155	0.23	0.14	0.18	202
156	0.62	0.69	0.65	130
157	0.26	0.10	0.14	245
158	0.86	0.48	0.62	177
159	0.45	0.25	0.32	130
160	0.46	0.19	0.27	336
161	0.84	0.65	0.73	220
162	0.10	0.03	0.05	229
163	0.84	0.47	0.60	316
164	0.57	0.16	0.25	283
165	0.30	0.38	0.33	197
166	0.12	0.08	0.09	101
167	0.24	0.24	0.24	231
168	0.30	0.14	0.19	370
169	0.39	0.24	0.29	258
170	0.13	0.06	0.08	101
171	0.36	0.20	0.26	89
172	0.31	0.28	0.29	193
173	0.33	0.36	0.35	309
174	0.41	0.09	0.15	172
175	0.92	0.73	0.81	95
176	0.94	0.50	0.65	346
177	0.97	0.29	0.44	322
178	0.47	0.50	0.48	232
179	0.56	0.04	0.07	125
180	0.42	0.19	0.26	145
181	0.08	0.16	0.11	77
182	0.12	0.05	0.07	182
183	0.57	0.30	0.39	257
184	0.13	0.05	0.07	216
185	0.20	0.05	0.07	242
186	0.27	0.18	0.22	165
187	0.76	0.51	0.61	263
188	0.29	0.11	0.16	174
189	0.60	0.09	0.15	136
190	0.93	0.43	0.59	202
191	0.30	0.10	0.15	134
192	0.76	0.32	0.45	230
193	0.28	0.18	0.22	90
194	0.56	0.44	0.50	185
195	0.11	0.07	0.08	156
196	0.14	0.01	0.01	160
197	0.00	0.00	0.00	266
198	0.22	0.14	0.17	284
199	0.14	0.18	0.15	145
200	0.91	0.66	0.76	212
201	0.51	0.25	0.33	317
202	0.72	0.41	0.52	427
203	0.19	0.10	0.13	232
204	0.32	0.19	0.24	217
205	0.46	0.34	0.39	527
206	0.04	0.01	0.01	124

207	0.25	0.01	0.02	103
208	0.73	0.61	0.67	287
209	0.13	0.10	0.11	193
210	0.43	0.32	0.37	220
211	0.67	0.01	0.03	140
212	0.08	0.07	0.08	161
213	0.48	0.18	0.26	72
214	0.60	0.43	0.51	396
215	0.49	0.32	0.39	134
216	0.00	0.00	0.00	400
217	0.51	0.28	0.36	75
218	0.94	0.76	0.84	219
219	0.66	0.35	0.46	210
220	0.30	0.24	0.27	298
221	0.96	0.52	0.68	266
222	0.83	0.25	0.38	290
223	0.10	0.12	0.11	128
224	0.76	0.32	0.45	159
225	0.38	0.19	0.25	164
226	0.51	0.38	0.43	144
227	0.37	0.47	0.42	276
228	0.05	0.02	0.03	235
229	0.29	0.04	0.07	216
230	0.31	0.21	0.25	228
231	0.52	0.56	0.54	64
232	0.09	0.06	0.07	103
233	0.71	0.24	0.35	216
234	0.00	0.00	0.00	116
235	0.66	0.51	0.57	77
236	0.92	0.67	0.78	67
237	0.00	0.00	0.00	218
238	0.09	0.04	0.06	139
239	0.17	0.02	0.04	94
240	0.35	0.18	0.24	77
241	0.40	0.02	0.05	167
242	0.51	0.29	0.37	86
243	0.24	0.19	0.21	58
244	0.22	0.09	0.12	269
245	0.11	0.05	0.07	112
246	0.95	0.62	0.75	255
247	0.37	0.19	0.25	58
248	0.30	0.04	0.07	81
249	0.00	0.00	0.00	131
250	0.28	0.12	0.17	93
251	0.43	0.23	0.30	154
252	0.21	0.05	0.08	129
253	0.45	0.28	0.34	83
254	0.23	0.09	0.13	191
255	0.00	0.00	0.00	219
256	0.11	0.03	0.05	130
257	0.39	0.35	0.37	93
258	0.65	0.46	0.54	217
259	0.15	0.06	0.09	141
260	0.94	0.10	0.19	143
261	0.50	0.10	0.17	219
262	0.47	0.26	0.34	107
263	0.31	0.28	0.29	236
264	0.14	0.11	0.12	119
265	0.05	0.18	0.08	72
266	0.13	0.04	0.06	70
267	0.18	0.04	0.06	107

268	0.65	0.33	0.43	169
269	0.18	0.08	0.11	129
270	0.59	0.63	0.61	159
271	0.63	0.18	0.28	190
272	0.41	0.09	0.15	248
273	0.90	0.64	0.75	264
274	0.63	0.69	0.66	105
275	0.00	0.00	0.00	104
276	0.09	0.02	0.03	115
277	0.63	0.63	0.63	170
278	0.50	0.30	0.37	145
279	0.66	0.51	0.57	230
280	0.58	0.38	0.45	80
281	0.68	0.52	0.59	217
282	0.74	0.56	0.64	175
283	0.48	0.05	0.09	269
284	0.60	0.36	0.45	74
285	0.79	0.47	0.59	206
286	0.89	0.54	0.67	227
287	0.59	0.44	0.50	130
288	0.24	0.07	0.11	129
289	0.08	0.01	0.02	80
290	0.15	0.12	0.13	99
291	0.75	0.27	0.40	208
292	0.28	0.07	0.12	67
293	0.60	0.28	0.39	109
294	0.17	0.26	0.21	140
295	0.17	0.14	0.15	241
296	0.19	0.18	0.18	72
297	0.28	0.10	0.15	107
298	0.83	0.16	0.27	61
299	0.74	0.38	0.50	77
300	0.12	0.05	0.07	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.51	0.45	0.48	176
304	0.96	0.61	0.74	230
305	0.93	0.67	0.78	156
306	0.36	0.56	0.44	146
307	0.19	0.08	0.11	98
308	0.05	0.03	0.03	78
309	0.67	0.02	0.04	94
310	0.69	0.26	0.38	162
311	0.69	0.41	0.51	116
312	0.53	0.30	0.38	57
313	0.00	0.00	0.00	65
314	0.29	0.39	0.34	138
315	0.36	0.30	0.33	195
316	0.45	0.39	0.42	69
317	0.00	0.00	0.00	134
318	0.30	0.19	0.23	148
319	0.80	0.45	0.57	161
320	0.16	0.17	0.17	104
321	0.70	0.63	0.66	156
322	0.49	0.23	0.31	134
323	0.54	0.31	0.39	232
324	0.21	0.11	0.14	92
325	0.30	0.09	0.13	197
326	0.04	0.01	0.01	126
327	0.33	0.01	0.02	115
328	0.97	0.57	0.71	198

329	0.52	0.26	0.35	125
330	0.67	0.10	0.17	81
331	0.18	0.03	0.05	94
332	0.00	0.00	0.00	56
333	0.03	0.00	0.01	260
334	0.00	0.00	0.00	60
335	0.25	0.14	0.18	110
336	0.58	0.39	0.47	71
337	0.20	0.21	0.21	66
338	0.46	0.31	0.37	150
339	0.00	0.00	0.00	54
340	0.82	0.45	0.58	195
341	0.00	0.00	0.00	79
342	0.32	0.32	0.32	38
343	0.36	0.23	0.28	43
344	0.20	0.01	0.03	68
345	0.52	0.40	0.45	73
346	0.15	0.06	0.09	116
347	0.48	0.38	0.42	111
348	0.11	0.05	0.07	63
349	0.89	0.49	0.63	104
350	0.59	0.43	0.50	44
351	0.00	0.00	0.00	40
352	1.00	0.20	0.33	136
353	0.42	0.31	0.36	54
354	0.00	0.00	0.00	134
355	0.30	0.12	0.17	120
356	0.27	0.06	0.10	228
357	0.55	0.09	0.15	269
358	0.64	0.29	0.40	80
359	0.79	0.27	0.40	140
360	0.18	0.04	0.07	125
361	0.92	0.34	0.50	169
362	0.07	0.04	0.05	56
363	0.93	0.56	0.70	154
364	0.33	0.02	0.03	58
365	0.11	0.20	0.14	71
366	0.96	0.50	0.66	54
367	0.06	0.01	0.01	116
368	1.00	0.02	0.04	54
369	0.00	0.00	0.00	71
370	0.11	0.02	0.03	61
371	0.57	0.06	0.10	71
372	0.73	0.42	0.54	52
373	0.60	0.06	0.11	150
374	0.30	0.31	0.30	93
375	0.03	0.03	0.03	67
376	0.00	0.00	0.00	76
377	0.91	0.09	0.17	106
378	0.20	0.01	0.02	86
379	0.11	0.07	0.09	14
380	1.00	0.23	0.37	122
381	0.10	0.06	0.07	104
382	0.22	0.11	0.14	66
383	0.48	0.29	0.36	110
384	0.00	0.00	0.00	155
385	0.07	0.06	0.06	50
386	0.19	0.19	0.19	64
387	0.00	0.00	0.00	93
388	0.31	0.27	0.29	102
389	0.00	0.00	0.00	108

390	0.85	0.51	0.63	178
391	0.58	0.10	0.16	115
392	0.92	0.26	0.41	42
393	0.00	0.00	0.00	134
394	0.00	0.00	0.00	112
395	0.25	0.01	0.01	176
396	0.00	0.00	0.00	125
397	0.37	0.35	0.36	224
398	0.79	0.35	0.48	63
399	0.00	0.00	0.00	59
400	0.39	0.30	0.34	63
401	0.16	0.04	0.07	98
402	0.32	0.05	0.09	162
403	0.29	0.14	0.19	83
404	0.75	0.63	0.69	19
405	0.10	0.05	0.07	92
406	0.33	0.15	0.20	41
407	0.22	0.28	0.25	43
408	0.00	0.00	0.00	160
409	0.16	0.16	0.16	50
410	0.00	0.00	0.00	19
411	0.32	0.12	0.18	175
412	0.12	0.01	0.02	72
413	0.17	0.09	0.12	95
414	0.11	0.07	0.09	97
415	0.27	0.15	0.19	48
416	0.38	0.25	0.30	83
417	0.00	0.00	0.00	40
418	0.15	0.09	0.11	91
419	0.39	0.23	0.29	90
420	0.22	0.16	0.19	37
421	0.05	0.02	0.02	66
422	0.55	0.29	0.38	73
423	0.35	0.20	0.25	56
424	0.88	0.67	0.76	33
425	0.08	0.01	0.02	76
426	0.03	0.02	0.03	81
427	0.99	0.63	0.77	150
428	0.82	0.79	0.81	29
429	0.00	0.00	0.00	389
430	0.59	0.20	0.30	167
431	0.00	0.00	0.00	123
432	0.31	0.49	0.38	39
433	0.34	0.24	0.29	82
434	0.94	0.73	0.82	66
435	0.53	0.32	0.40	93
436	0.27	0.03	0.06	87
437	0.17	0.10	0.13	86
438	0.53	0.34	0.41	104
439	0.00	0.00	0.00	100
440	0.33	0.01	0.01	141
441	0.30	0.27	0.28	110
442	0.21	0.09	0.12	123
443	0.00	0.00	0.00	71
444	0.28	0.05	0.08	109
445	0.20	0.12	0.15	48
446	0.44	0.21	0.29	76
447	0.10	0.05	0.07	38
448	0.60	0.49	0.54	81
449	0.44	0.06	0.11	132
450	0.45	0.25	0.32	81

451	0.00	0.00	0.00	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.73	0.39	0.50	70
455	0.17	0.04	0.06	155
456	0.21	0.16	0.18	43
457	0.31	0.22	0.26	72
458	0.21	0.10	0.13	62
459	0.24	0.07	0.11	69
460	0.21	0.03	0.06	119
461	0.62	0.16	0.26	79
462	0.15	0.04	0.07	47
463	0.18	0.02	0.03	104
464	0.50	0.27	0.35	106
465	0.20	0.02	0.03	64
466	0.51	0.18	0.27	173
467	0.77	0.32	0.45	107
468	0.50	0.01	0.02	126
469	0.00	0.00	0.00	114
470	0.93	0.61	0.74	140
471	0.00	0.00	0.00	79
472	0.27	0.36	0.31	143
473	0.23	0.03	0.06	158
474	0.57	0.03	0.06	138
475	0.11	0.07	0.08	59
476	0.62	0.33	0.43	88
477	0.69	0.74	0.71	176
478	0.11	0.54	0.19	24
479	0.00	0.00	0.00	92
480	0.80	0.33	0.47	100
481	0.50	0.02	0.04	103
482	0.09	0.23	0.13	74
483	0.76	0.48	0.58	105
484	0.07	0.02	0.04	83
485	0.11	0.01	0.02	82
486	0.18	0.06	0.09	71
487	0.38	0.23	0.28	120
488	0.00	0.00	0.00	105
489	0.46	0.53	0.49	87
490	1.00	0.66	0.79	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	0.16	0.05	0.07	61
495	0.00	0.00	0.00	344
496	0.00	0.00	0.00	52
497	0.67	0.12	0.20	137
498	0.00	0.00	0.00	98
499	0.00	0.00	0.00	79

avg / total	0.54	0.33	0.39	173809
-------------	------	------	------	--------

Time taken to run this cell : 0:14:06.786004

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples.

'precision', 'predicted', average, warn\_for)

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 in labels with no predicted samples.

```
'precision', 'predicted', average, warn_for)
C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:11
35: UndefinedMetricWarning: Precision and F-score are ill-defined and being
set to 0.0 in labels with no predicted samples.
'precision', 'predicted', average, warn_for)
```

## **Linear SVM with OneVsRestClassifier (With Hinge loss):-**

### **Hyperparameter tuning:-**

In [34]:

```

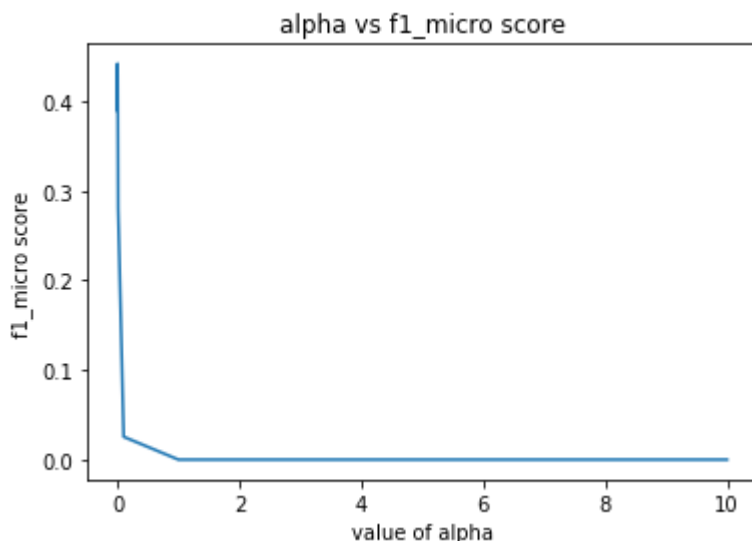
param={'estimator__alpha': [10**-5,10**-4, 10**-3, 10**-2, 10**-1, 10**0, 10**1]}
classifier = OneVsRestClassifier(SGDClassifier( loss='hinge', penalty='l1'))
gsv = GridSearchCV(estimator = classifier, param_grid=param, verbose=0, scoring='f1_micro')
gsv.fit(x_train_multilabel, y_train)

best_alpha = gsv.best_estimator_.get_params()['estimator__alpha']
print('value of alpha after hyperparameter tuning : ',best_alpha)
print('-----')
# plotting C vs f1_micro_score
x_1=[]
y_1=[]
for x in gsv.grid_scores_:
    x_1.append(x[0]['estimator__alpha'])
    y_1.append(x[1])
plt.plot(x_1,y_1)
plt.xlabel('value of alpha')
plt.ylabel('f1_micro score')
plt.title('alpha vs f1_micro score')
plt.show()

```

value of alpha after hyperparameter tuning : 0.001

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\model\_selection\\_search.py:761: DeprecationWarning: The grid\_scores\_ attribute was deprecated in version 0.18 in favor of the more elaborate cv\_results\_ attribute. The grid\_scores\_ attribute will not be available from 0.20  
 DeprecationWarning)



**Applying model using best hyperparameter:-**



In [35]:

```

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.001, penalty='l1'), n_
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 : 0.18101362027240545

Hamming loss 0.0032552451049020982

Micro-average quality numbers

Precision: 0.5552, Recall: 0.3199, F1-measure: 0.4059

Macro-average quality numbers

Precision: 0.3223, Recall: 0.2400, F1-measure: 0.2587

	precision	recall	f1-score	support
0	0.81	0.67	0.73	5519
1	0.45	0.19	0.26	8189
2	0.68	0.38	0.49	6529
3	0.58	0.47	0.52	3231
4	0.71	0.42	0.52	6430
5	0.59	0.39	0.47	2878
6	0.73	0.56	0.63	5086
7	0.82	0.56	0.66	4533
8	0.50	0.17	0.26	3000
9	0.69	0.58	0.63	2765
10	0.42	0.11	0.17	3051
11	0.72	0.32	0.44	3009
12	0.58	0.27	0.37	2630
13	0.24	0.29	0.26	1425
14	0.75	0.62	0.68	2548
15	0.65	0.13	0.21	2371
16	0.51	0.30	0.38	873
17	0.70	0.72	0.71	2151
18	0.57	0.19	0.29	2204
19	0.47	0.49	0.48	831
20	0.73	0.48	0.58	1860
21	0.00	0.00	0.00	2023
22	0.08	0.00	0.00	1513
23	0.65	0.61	0.63	1207

24	0.67	0.07	0.12	506
25	0.70	0.35	0.46	425
26	0.40	0.45	0.42	793
27	0.55	0.31	0.40	1291
28	0.60	0.37	0.46	1208
29	0.26	0.17	0.20	406
30	0.58	0.27	0.36	504
31	0.10	0.04	0.05	732
32	0.44	0.33	0.38	441
33	0.27	0.02	0.04	1645
34	0.52	0.36	0.42	1058
35	0.63	0.59	0.61	946
36	0.49	0.29	0.37	644
37	0.77	0.82	0.80	136
38	0.37	0.46	0.41	570
39	0.68	0.36	0.47	766
40	0.52	0.26	0.34	1132
41	0.21	0.28	0.24	174
42	0.47	0.64	0.55	210
43	0.58	0.52	0.55	433
44	0.48	0.47	0.48	626
45	0.40	0.29	0.34	852
46	0.50	0.46	0.48	534
47	0.04	0.02	0.02	350
48	0.55	0.55	0.55	496
49	0.68	0.59	0.64	785
50	0.00	0.00	0.00	475
51	0.10	0.06	0.08	305
52	0.00	0.00	0.00	251
53	0.65	0.34	0.45	914
54	0.06	0.00	0.00	728
55	0.00	0.00	0.00	258
56	0.24	0.23	0.24	821
57	0.28	0.05	0.08	541
58	0.70	0.29	0.41	748
59	0.85	0.74	0.79	724
60	0.00	0.00	0.00	660
61	0.73	0.26	0.38	235
62	0.81	0.82	0.82	718
63	0.72	0.66	0.69	468
64	0.42	0.46	0.44	191
65	0.18	0.14	0.16	429
66	0.00	0.00	0.00	415
67	0.64	0.66	0.65	274
68	0.73	0.61	0.66	510
69	0.41	0.32	0.36	466
70	0.20	0.14	0.17	305
71	0.22	0.19	0.20	247
72	0.63	0.50	0.56	401
73	0.81	0.74	0.78	86
74	0.46	0.43	0.45	120
75	0.30	0.76	0.43	129
76	0.00	0.00	0.00	473
77	0.23	0.34	0.27	143
78	0.63	0.56	0.59	347
79	0.53	0.39	0.45	479
80	0.32	0.35	0.33	279
81	0.63	0.16	0.26	461
82	0.00	0.00	0.00	298
83	0.67	0.47	0.55	396
84	0.32	0.21	0.25	184

85	0.71	0.09	0.16	573
86	0.33	0.08	0.13	325
87	0.50	0.18	0.27	273
88	0.28	0.31	0.29	135
89	0.00	0.00	0.00	232
90	0.29	0.36	0.32	409
91	0.45	0.39	0.41	420
92	0.65	0.58	0.61	408
93	0.46	0.50	0.48	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	0.82	0.17	0.29	501
98	0.58	0.58	0.58	136
99	0.46	0.22	0.30	239
100	0.09	0.06	0.07	324
101	0.89	0.73	0.80	277
102	0.81	0.74	0.77	613
103	0.53	0.18	0.27	157
104	0.00	0.00	0.00	295
105	0.36	0.36	0.36	334
106	0.00	0.00	0.00	335
107	0.50	0.47	0.49	389
108	0.34	0.14	0.19	251
109	0.43	0.50	0.46	317
110	0.00	0.00	0.00	187
111	0.61	0.10	0.17	140
112	0.36	0.27	0.31	154
113	0.42	0.29	0.34	332
114	0.00	0.00	0.00	323
115	0.10	0.01	0.01	344
116	0.50	0.54	0.52	370
117	0.49	0.26	0.34	313
118	0.71	0.48	0.57	874
119	0.25	0.25	0.25	293
120	0.00	0.00	0.00	200
121	0.58	0.56	0.57	463
122	0.20	0.14	0.17	119
123	0.00	0.00	0.00	256
124	0.58	0.81	0.67	195
125	0.15	0.27	0.19	138
126	0.53	0.39	0.45	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	0.04	0.02	0.03	210
132	0.31	0.06	0.10	361
133	0.76	0.66	0.70	453
134	0.72	0.91	0.80	124
135	0.00	0.00	0.00	91
136	0.68	0.10	0.18	128
137	0.35	0.28	0.31	218
138	0.00	0.00	0.00	243
139	0.11	0.25	0.15	149
140	0.63	0.48	0.55	318
141	0.00	0.00	0.00	159
142	0.58	0.41	0.48	274
143	0.72	0.72	0.72	362
144	0.22	0.34	0.27	118
145	0.41	0.49	0.45	164

146	0.46	0.30	0.37	461
147	0.49	0.40	0.44	159
148	0.00	0.00	0.00	166
149	0.92	0.54	0.68	346
150	0.17	0.11	0.13	350
151	0.29	0.58	0.39	55
152	0.62	0.52	0.56	387
153	0.00	0.00	0.00	150
154	0.50	0.07	0.12	281
155	0.36	0.02	0.05	202
156	0.49	0.70	0.58	130
157	0.16	0.16	0.16	245
158	0.74	0.74	0.74	177
159	0.37	0.33	0.35	130
160	0.28	0.21	0.24	336
161	0.75	0.65	0.70	220
162	0.00	0.00	0.00	229
163	0.78	0.51	0.62	316
164	0.67	0.29	0.40	283
165	0.67	0.03	0.06	197
166	0.19	0.20	0.19	101
167	0.00	0.00	0.00	231
168	0.13	0.10	0.11	370
169	0.32	0.31	0.32	258
170	0.00	0.00	0.00	101
171	0.22	0.22	0.22	89
172	0.22	0.40	0.29	193
173	0.35	0.35	0.35	309
174	0.20	0.18	0.19	172
175	0.71	0.88	0.79	95
176	0.80	0.58	0.67	346
177	0.83	0.32	0.46	322
178	0.52	0.45	0.48	232
179	0.27	0.14	0.19	125
180	0.38	0.28	0.32	145
181	0.26	0.12	0.16	77
182	0.00	0.00	0.00	182
183	0.31	0.35	0.33	257
184	0.00	0.00	0.00	216
185	0.19	0.13	0.16	242
186	0.22	0.21	0.22	165
187	0.56	0.58	0.57	263
188	0.00	0.00	0.00	174
189	0.22	0.06	0.09	136
190	0.78	0.45	0.57	202
191	0.00	0.00	0.00	134
192	0.63	0.45	0.53	230
193	0.40	0.11	0.17	90
194	0.43	0.55	0.48	185
195	0.08	0.04	0.05	156
196	0.00	0.00	0.00	160
197	0.00	0.00	0.00	266
198	0.19	0.11	0.14	284
199	0.23	0.08	0.11	145
200	0.80	0.79	0.79	212
201	0.10	0.02	0.03	317
202	0.62	0.57	0.59	427
203	0.00	0.00	0.00	232
204	0.00	0.00	0.00	217
205	0.45	0.44	0.45	527
206	0.00	0.00	0.00	124

207	0.26	0.17	0.21	103
208	0.70	0.56	0.62	287
209	0.00	0.00	0.00	193
210	0.41	0.25	0.32	220
211	0.90	0.06	0.12	140
212	0.00	0.00	0.00	161
213	0.09	0.21	0.13	72
214	0.55	0.34	0.42	396
215	0.81	0.29	0.43	134
216	0.00	0.00	0.00	400
217	0.31	0.28	0.30	75
218	0.89	0.67	0.76	219
219	0.05	0.05	0.05	210
220	0.83	0.61	0.70	298
221	0.92	0.58	0.71	266
222	0.73	0.46	0.56	290
223	0.00	0.00	0.00	128
224	0.54	0.47	0.50	159
225	0.47	0.17	0.25	164
226	0.40	0.48	0.43	144
227	0.52	0.19	0.28	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	0.60	0.64	0.62	64
232	0.00	0.00	0.00	103
233	0.37	0.35	0.36	216
234	0.00	0.00	0.00	116
235	0.53	0.48	0.50	77
236	0.84	0.73	0.78	67
237	0.00	0.00	0.00	218
238	0.00	0.00	0.00	139
239	0.00	0.00	0.00	94
240	0.36	0.35	0.36	77
241	0.00	0.00	0.00	167
242	0.75	0.31	0.44	86
243	0.16	0.21	0.18	58
244	0.17	0.07	0.10	269
245	0.00	0.00	0.00	112
246	0.91	0.73	0.81	255
247	0.34	0.22	0.27	58
248	0.00	0.00	0.00	81
249	0.00	0.00	0.00	131
250	0.38	0.20	0.27	93
251	0.32	0.29	0.30	154
252	0.00	0.00	0.00	129
253	0.25	0.31	0.28	83
254	0.11	0.23	0.15	191
255	0.00	0.00	0.00	219
256	0.00	0.00	0.00	130
257	0.19	0.40	0.26	93
258	0.61	0.62	0.61	217
259	0.25	0.02	0.04	141
260	0.66	0.17	0.28	143
261	0.60	0.11	0.19	219
262	0.30	0.36	0.33	107
263	0.25	0.34	0.29	236
264	0.17	0.12	0.14	119
265	0.16	0.25	0.19	72
266	0.14	0.14	0.14	70
267	0.30	0.03	0.05	107

268	0.50	0.49	0.50	169
269	0.00	0.00	0.00	129
270	0.48	0.67	0.56	159
271	0.00	0.00	0.00	190
272	0.25	0.19	0.21	248
273	0.83	0.69	0.76	264
274	0.59	0.69	0.63	105
275	0.00	0.00	0.00	104
276	0.00	0.00	0.00	115
277	0.72	0.62	0.67	170
278	0.50	0.41	0.45	145
279	0.83	0.54	0.66	230
280	0.38	0.28	0.32	80
281	0.56	0.63	0.59	217
282	0.65	0.51	0.57	175
283	0.00	0.00	0.00	269
284	0.52	0.46	0.49	74
285	0.61	0.49	0.54	206
286	0.85	0.72	0.78	227
287	0.59	0.44	0.50	130
288	0.38	0.02	0.04	129
289	0.00	0.00	0.00	80
290	0.00	0.00	0.00	99
291	0.73	0.31	0.44	208
292	0.00	0.00	0.00	67
293	0.43	0.20	0.28	109
294	0.22	0.08	0.12	140
295	0.14	0.16	0.15	241
296	0.13	0.17	0.14	72
297	0.00	0.00	0.00	107
298	0.73	0.18	0.29	61
299	0.83	0.19	0.32	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	0.32	0.59	0.42	176
304	0.89	0.58	0.71	230
305	0.90	0.66	0.76	156
306	0.37	0.27	0.32	146
307	0.09	0.15	0.11	98
308	0.00	0.00	0.00	78
309	0.23	0.20	0.22	94
310	0.00	0.00	0.00	162
311	0.63	0.61	0.62	116
312	0.37	0.60	0.46	57
313	0.00	0.00	0.00	65
314	0.32	0.33	0.32	138
315	0.34	0.26	0.30	195
316	0.33	0.41	0.36	69
317	0.33	0.09	0.14	134
318	0.35	0.09	0.15	148
319	0.78	0.52	0.62	161
320	0.00	0.00	0.00	104
321	0.47	0.55	0.51	156
322	0.40	0.48	0.43	134
323	0.39	0.34	0.36	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	0.97	0.45	0.62	198

329	0.29	0.41	0.34	125
330	0.38	0.19	0.25	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	0.08	0.02	0.03	60
335	0.00	0.00	0.00	110
336	0.36	0.42	0.39	71
337	0.08	0.06	0.07	66
338	0.24	0.40	0.30	150
339	0.00	0.00	0.00	54
340	0.75	0.47	0.58	195
341	0.00	0.00	0.00	79
342	0.20	0.24	0.21	38
343	0.58	0.33	0.42	43
344	0.00	0.00	0.00	68
345	0.44	0.45	0.45	73
346	0.00	0.00	0.00	116
347	0.77	0.53	0.63	111
348	0.03	0.05	0.04	63
349	0.76	0.50	0.60	104
350	0.57	0.48	0.52	44
351	0.00	0.00	0.00	40
352	0.82	0.36	0.50	136
353	0.27	0.37	0.31	54
354	0.11	0.07	0.09	134
355	0.23	0.17	0.20	120
356	0.00	0.00	0.00	228
357	0.00	0.00	0.00	269
358	0.35	0.36	0.36	80
359	0.69	0.26	0.37	140
360	0.00	0.00	0.00	125
361	0.79	0.71	0.75	169
362	0.15	0.14	0.15	56
363	0.76	0.74	0.75	154
364	0.00	0.00	0.00	58
365	0.00	0.00	0.00	71
366	0.93	0.46	0.62	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	0.04	0.07	0.05	61
371	0.36	0.13	0.19	71
372	0.40	0.40	0.40	52
373	0.49	0.31	0.38	150
374	0.29	0.28	0.28	93
375	0.00	0.00	0.00	67
376	0.00	0.00	0.00	76
377	0.83	0.09	0.17	106
378	0.00	0.00	0.00	86
379	0.00	0.00	0.00	14
380	0.82	0.15	0.25	122
381	0.00	0.00	0.00	104
382	0.24	0.15	0.19	66
383	0.26	0.26	0.26	110
384	0.00	0.00	0.00	155
385	0.00	0.00	0.00	50
386	0.17	0.27	0.20	64
387	0.00	0.00	0.00	93
388	0.26	0.13	0.17	102
389	0.00	0.00	0.00	108

390	0.84	0.66	0.74	178
391	0.41	0.19	0.26	115
392	0.74	0.33	0.46	42
393	0.00	0.00	0.00	134
394	0.02	0.02	0.02	112
395	0.00	0.00	0.00	176
396	0.00	0.00	0.00	125
397	0.54	0.37	0.44	224
398	0.60	0.43	0.50	63
399	0.00	0.00	0.00	59
400	0.25	0.43	0.32	63
401	0.00	0.00	0.00	98
402	0.00	0.00	0.00	162
403	0.24	0.34	0.28	83
404	0.70	0.84	0.76	19
405	0.00	0.00	0.00	92
406	0.38	0.51	0.44	41
407	0.25	0.51	0.33	43
408	0.00	0.00	0.00	160
409	0.19	0.26	0.22	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	0.28	0.42	0.34	83
417	0.00	0.00	0.00	40
418	0.00	0.00	0.00	91
419	0.31	0.12	0.18	90
420	0.14	0.19	0.16	37
421	0.00	0.00	0.00	66
422	0.44	0.29	0.35	73
423	0.25	0.27	0.26	56
424	0.86	0.91	0.88	33
425	0.00	0.00	0.00	76
426	0.00	0.00	0.00	81
427	0.99	0.57	0.72	150
428	0.79	0.76	0.77	29
429	0.00	0.00	0.00	389
430	0.44	0.32	0.37	167
431	0.00	0.00	0.00	123
432	0.32	0.41	0.36	39
433	0.25	0.32	0.28	82
434	0.94	0.70	0.80	66
435	0.46	0.40	0.43	93
436	0.00	0.00	0.00	87
437	0.00	0.00	0.00	86
438	0.46	0.44	0.45	104
439	0.06	0.09	0.07	100
440	0.00	0.00	0.00	141
441	0.00	0.00	0.00	110
442	0.13	0.11	0.12	123
443	0.40	0.03	0.05	71
444	0.00	0.00	0.00	109
445	0.00	0.00	0.00	48
446	0.38	0.33	0.35	76
447	0.00	0.00	0.00	38
448	0.52	0.60	0.56	81
449	0.32	0.10	0.15	132
450	0.00	0.00	0.00	81



451	0.58	0.14	0.23	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.49	0.39	0.43	70
455	0.57	0.03	0.05	155
456	0.00	0.00	0.00	43
457	0.37	0.22	0.28	72
458	0.10	0.16	0.12	62
459	0.00	0.00	0.00	69
460	0.00	0.00	0.00	119
461	0.00	0.00	0.00	79
462	0.08	0.21	0.12	47
463	0.00	0.00	0.00	104
464	0.41	0.32	0.36	106
465	0.00	0.00	0.00	64
466	0.37	0.15	0.21	173
467	0.64	0.23	0.34	107
468	0.00	0.00	0.00	126
469	0.00	0.00	0.00	114
470	0.65	0.80	0.72	140
471	0.00	0.00	0.00	79
472	0.41	0.13	0.19	143
473	0.71	0.16	0.26	158
474	0.00	0.00	0.00	138
475	0.00	0.00	0.00	59
476	0.56	0.10	0.17	88
477	0.72	0.53	0.61	176
478	0.94	0.67	0.78	24
479	0.00	0.00	0.00	92
480	0.59	0.43	0.50	100
481	0.17	0.02	0.03	103
482	0.08	0.20	0.12	74
483	0.64	0.58	0.61	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	0.32	0.17	0.23	120
488	0.00	0.00	0.00	105
489	0.65	0.30	0.41	87
490	0.95	0.62	0.75	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	0.16	0.26	61
495	0.00	0.00	0.00	344
496	0.00	0.00	0.00	52
497	0.00	0.00	0.00	137
498	0.00	0.00	0.00	98
499	0.43	0.25	0.32	79

avg / total	0.48	0.32	0.37	173809
-------------	------	------	------	--------

Time taken to run this cell : 0:11:26.150570

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in labels with no predicted samples.

'precision', 'predicted', average, warn\_for)

C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1135: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 in labels with no predicted samples.

```
'precision', 'predicted', average, warn_for)
C:\Users\hp\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:11
35: UndefinedMetricWarning: Precision and F-score are ill-defined and being
set to 0.0 in labels with no predicted samples.
'precision', 'predicted', average, warn_for)
```

## -----CONCLUSION-----

----- PrettyTable for model comparisions:-----

-----Micro-average quality numbers(Precision, Recall, F1-  
measure)-----

In [12]:

```

# Creating table using PrettyTable Library
from prettytable import PrettyTable

# Names of models
names = ['Logistic Regression(1VsRestClassifier)', 'Linear SVM(1VsRestClassifier)']

alpha=[0.001, 0.001]

Precision = [0.5624,0.5552]

Recall = [0.3312,0.3199]

F1_measure = [0.4169, 0.4059]

Accuracy=[0.186 ,0.181 ]

Hamming_loss=[0.00322068,0.003255245]

penalty=['L1', 'L1']

# Initializing prettytable
ptable = PrettyTable()

# Adding columns

ptable.add_column("MODEL",names)
ptable.add_column("Penalty",penalty)
ptable.add_column("Hyperparameter",alpha)
ptable.add_column("Precision",Precision)
ptable.add_column("Recall",Recall)
ptable.add_column("F1-measure",F1_measure)
ptable.add_column("Hamming loss",Hamming_loss)
ptable.add_column("Accuracy",Accuracy)

# Printing the Table
print(ptable)

```

```

+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
|          MODEL          | Penalty | Hyperparameter | Precis
ion | Recall | F1-measure | Hamming loss | Accuracy |
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
| Logistic Regression(1VsRestClassifier) | L1 | 0.001 | 0.56
24 | 0.3312 | 0.4169 | 0.00322068 | 0.186 |
| Linear SVM(1VsRestClassifier) | L1 | 0.001 | 0.55
52 | 0.3199 | 0.4059 | 0.003255245 | 0.181 |
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+

```

-----Macro-average quality numbers(Precision, Recall, F1-measure)-----

In [14]:

```

# Creating table using PrettyTable Library
from prettytable import PrettyTable

# Names of models
names = ['Logistic Regression(1VsRestClassifier)', 'Linear SVM(1VsRestClassifier)']

alpha=[0.001, 0.001]

Precision = [ 0.4101,0.3223]

Recall = [0.2382, 0.2400]

F1_measure = [0.2832, 0.2587]

Accuracy=[0.186 ,0.181 ]

Hamming_loss=[0.00322068,0.003255245]

penalty=['L1', 'L1']

# Initializing prettytable
ptable = PrettyTable()

# Adding columns

ptable.add_column("MODEL",names)
ptable.add_column("Penalty",penalty)
ptable.add_column("Hyperparameter",alpha)
ptable.add_column("Precision",Precision)
ptable.add_column("Recall",Recall)
ptable.add_column("F1-measure",F1_measure)
ptable.add_column("Hamming loss",Hamming_loss)
ptable.add_column("Accuracy",Accuracy)

# Printing the Table
print(ptable)

```

```

+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
|          MODEL          | Penalty | Hyperparameter | Precis
ion | Recall | F1-measure | Hamming loss | Accuracy |
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+
| Logistic Regression(1VsRestClassifier) | L1 | 0.001 | 0.41
01 | 0.2382 | 0.2832 | 0.00322068 | 0.186 |
| Linear SVM(1VsRestClassifier) | L1 | 0.001 | 0.32
23 | 0.24 | 0.2587 | 0.003255245 | 0.181 |
+-----+-----+-----+-----+-----+
+-----+-----+-----+-----+-----+

```

## STEPS FOLLOWED:

1----->> Exploratory Data Analysis.

A: Data Loading and Cleaning.

B: Analysis of Tags.

C: Checking for duplicates.

## 2----->> Machine Learning Models.

A: Converting tags for multilabel problems.

B: Preprocessing of questions.

C: Separate Code from Body.

D: Remove Special characters from Question title and description (not in code).

E: Give more weightage to title : Add title three times to the question.

F: Remove stop words (Except 'C').

G: Remove HTML Tags.

H: Convert all the characters into small letters.

I: Use SnowballStemmer to stem the words.

## 3----->> USE BoW WITH UPTO 4 GRAMS AND COMPUTE micro F1 score with Logistic Regression(1 Vs Rest)

## 4----->> USE BoW WITH UPTO 4 GRAMS AND COMPUTE micro F1 score with Linear SVM(1 Vs Rest) (With Hinge loss)

# OBSERVATIONS:

**note: HERE WE HAVE TAKEN ONLY 500 TAGS AND VERY LESS NO OF QUESTION PAIRS i.e .5M because we have low configuration system and limited time.**

1:Here we have taken linear models i.e logistic regression and linear svm because they work very well with high dimension data, so we use tfidf featurization for these models.

2:Here we have taken v.less no.of points and tags If we take more data points and all tags our performance matrix can improve much.

3:As we observed the most important thing is time taken by our simple linear model which are v.fast as compared to other models here, i.e they also take decent much amount of time.

4:If we take time into consideration model(Logistic Regression,Linear SVM) are best for productization in these type of real world problems(i.e multilabel classification).

In [ ]: