

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>

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

Research paper : <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tagging-1.pdf>

Research paper : <https://dl.acm.org/citation.cfm?id=2660970&dl=ACM&coll=DL>

1.3 Real World / Business Objectives and Constraints

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

2. Machine Learning problem

2.1 Data

2.1.1 Data Overview

Refer: <https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data>

All of the data is in 2 files: Train and Test.

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

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

Size of Train.csv - 6.75GB

Size of Test.csv - 2GB

Number of rows in Train.csv = 6034195

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

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

Id - Unique identifier for each question

Title - The question's title

Body - The body of the question

Tags - The tags associated with the question in a space-separated format (all lowercase, should not contain tabs '\t' or ampersands '&')

2.1.2 Example Data point

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

Body :

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

    for(int y=1; y<n+1; y++)\n
    {\n
        cin>>m[y];\n
        cin>>u[y];\n
    }\n
    for(x=1; x<n+1; x++)\n
    {\n
        a[x] = (m[x] + u[x])/2;\n
    }\n
    c=(n*4)-4;\n
    for(int al=1; al<n+1; al++)\n
    {\n\n
        e[al][0] = m[al];\n
        e[al][1] = m[al]+1;\n
        e[al][2] = u[al]-1;\n
        e[al][3] = u[al];\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.

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

http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html

<https://www.kaggle.com/wiki/HammingLoss>

3.1 Data Loading and Cleaning

3.1.1 Using Pandas with SQLite to Load the data

```
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
from tqdm import tqdm

from sklearn.model_selection import GridSearchCV
```

```
#Creating db file from csv
#Learn SQL: https://www.w3schools.com/sql/default.asp
if not os.path.isfile('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Train/train.db'):
    start = datetime.now()
    disk_engine = create_engine('sqlite:///E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Train/train.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in tqdm(pd.read_csv('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Train/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 [3]:

```
if os.path.isfile('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Train/train.db'):
    start = datetime.now()
    con = sqlite3.connect('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/Train/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 : 0:01:39.595380

3.1.3 Checking for duplicates

In [4]:

```
#Learn SQL: https://www.w3schools.com/sql/default.asp
if os.path.isfile('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Train/train.db'):
    start = datetime.now()
    con = sqlite3.connect('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/Train/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 : 0:09:01.888052

In [5]:

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

Out[5]:

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

In [6]:

```
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 (30.292038906260256 %)

In [7]:

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

Out[7]:

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

In [8]:

```
df_no_dup['Tags']
```

Out[8]:

```
0          c++ c
1    c# silverlight data-binding
2    c# silverlight data-binding columns
3          jsp jstl
4          java jdbc
...
4206310  wordpress wordpress-plugin
4206311          php mysql text
4206312  php codeigniter character-encoding
4206313          php email outlook mime
4206314          html
Name: Tags, Length: 4206315, dtype: object
```

Checking for NaN values

In [9]:

```
# just to make sure that all Nan containing rows are deleted..
print("No of Nan values in our dataframe : ", sum(df_no_dup.isnull().any()))
```

No of Nan values in our dataframe : 1

Removing Nan Values

In [10]:

```
df_no_dup.dropna()
```

Out[10]:

	Title	Body	Tags	cnt_dup
0	Implementing Boundary Value Analysis of S...	<pre> <code>#include<stream> </code>	c++ c	1
1	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical...	c# silverlight data-binding	1
2	Dynamic Datagrid Binding in Silverlight?	<p>I should do binding for datagrid dynamical...	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> <pre> <code>...	java jdbc	2
...
4206310	(wordpress) Plugin Error:You do not have suffic...	<p>My plugin which adds a menu in admin page h...	wordpress wordpress-plugin	2
4206311	(question mark) getting displayed in place	<p>Everytime i get some text to display on a	php mysql text	1

	Title	Body	Tags	cnt_dup
4206312	⚡ appears using character_limiter() with strip...	<p>I'm getting ⚡ characters when I combine Cod...	php codeigniter character-encoding	1
4206313	⚡ in base64 encoded emails	<p>I have a problem with Swedish language + MS...	php email outlook mime	2
4206314	⚡ odd character	<p>⚡ Odd Character when request file via ajax ...	html	2

4206308 rows × 4 columns

Removing Duplicates

In [11]:

```
dup_bool = df_no_dup.duplicated(['Title', 'Body', 'Tags', 'cnt_dup'])
dups = sum(dup_bool) # by considering all columns..( including timestamp)
print("There are {} duplicate rating entries in the data..".format(dups))
```

There are 0 duplicate rating entries in the data..

In [12]:

```
#https://github.com/priyagunjate/Stack-Overflow-Tag-Prediction/blob/master/Assignment-19.ipynb
```

In [13]:

```
aa_count=[]
hh=[]
for j in range(len(df_no_dup)):
    tex=df_no_dup['Tags'][j]
    #print(tex)
    if tex is not None:
        #print("heyram")
        #start=datetime.now()
        hh.append(tex)
        text=len(tex.split(" "))
        #print(text)
        aa_count.append(text)

print(len(aa_count))
aaa=pd.DataFrame(aa_count,columns=['tag_count'])
hhh=pd.DataFrame(hh,columns=['Tags'])
df_no_dup=pd.concat([hhh,aaa],axis=1)
# adding a new feature number of tags per question
print("Time taken to run this cell :", datetime.now() - start)
df_no_dup.head()
np.where(pd.isnull(df_no_dup))
```

4206308

Time taken to run this cell : 0:16:07.891721

Out[13]:

```
(array([], dtype=int64), array([], dtype=int64))
```

In [14]:

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

Out[14]:

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

In [15]:

```
#Creating a new database with no duplicates
if not os.path.isfile('./BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/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.head()
    no_dup.to_sql('no_dup_train', disk_dup)
```

In [16]:

```
#This method seems more appropriate to work with this much data.
#creating the connection with database file.
if os.path.isfile('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/train_no_dup.db'):
    start = datetime.now()
    con = sqlite3.connect('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/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 : 0:01:39.384676

3.2 Analysis of Tags

3.2.1 Total number of unique tags

In [17]:

```
# 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 [18]:

```
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 [19]:

```
#'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 [20]:

```
# 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 [21]:

```
# Saving this dictionary to csv files.
if not os.path.isfile('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/tag_counts_dict_dtm.csv'):
    with open('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/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("E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/tag_counts_dict_dtm.csv", names=['Tags', 'Counts'])
tag_df.head()
```

Out[21]:

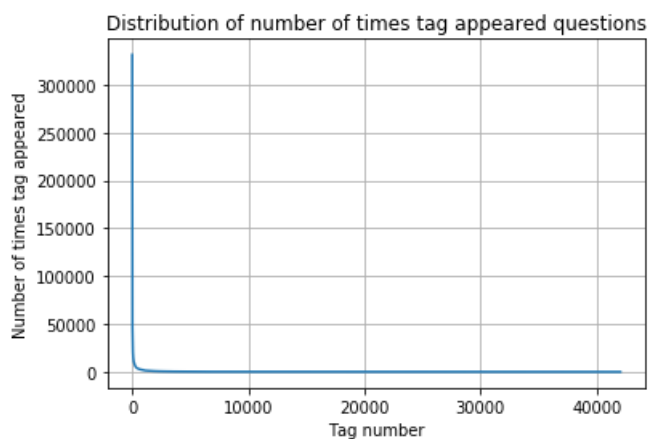
	Tags	Counts
0	.a	18
1	.app	37
2	.asp.net-mvc	1
3	.aspxauth	21
4	.bash-profile	138

In [22]:

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

In [23]:

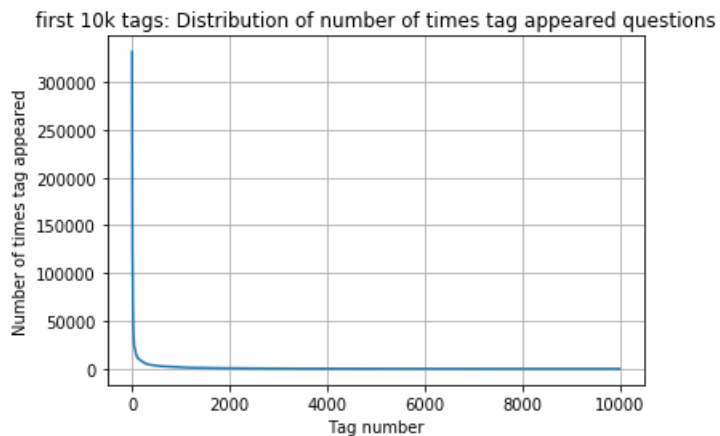
```
plt.plot(tag_counts)
plt.title("Distribution of number of times tag appeared questions")
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
```



In [24]:

```
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])
```

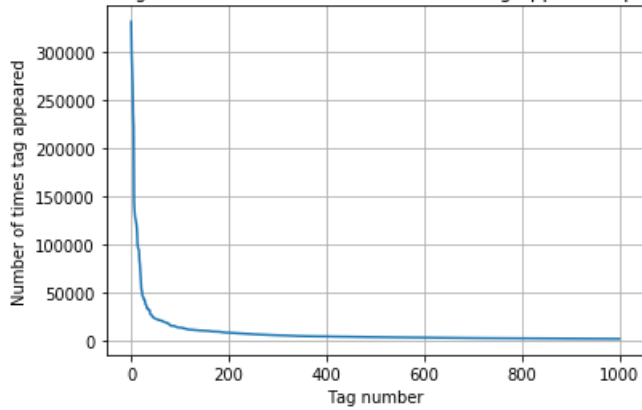


```
400 [331505  44829  22429  17728  13364  11162  10029   9148   8054   7151
 6466   5865   5370   4983   4526   4281   4144   3929   3750   3593
 3453   3299   3123   2986   2891   2738   2647   2527   2431   2331
 2259   2186   2097   2020   1959   1900   1828   1770   1723   1673
 1631   1574   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 [25]:

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

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



```

200 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291 6183 6093 5971 5865 5760 5577 5490 5411
5370 5283 5207 5107 5066 4983 4891 4785 4658 4549
4526 4487 4429 4335 4310 4281 4239 4228 4195 4159
4144 4088 4050 4002 3957 3929 3874 3849 3818 3797
3750 3703 3685 3658 3615 3593 3564 3521 3505 3483
3453 3427 3396 3363 3326 3299 3272 3232 3196 3168
3123 3094 3073 3050 3012 2986 2983 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]

```

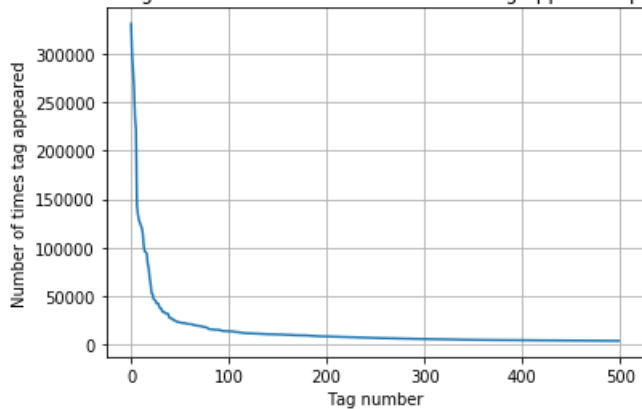
In [26]:

```

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

```

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



```

100 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291 6183 6093 5971 5865 5760 5577 5490 5411
5370 5283 5207 5107 5066 4983 4891 4785 4658 4549
4526 4487 4429 4335 4310 4281 4239 4228 4195 4159
4144 4088 4050 4002 3957 3929 3874 3849 3818 3797
3750 3703 3685 3658 3615 3593 3564 3521 3505 3483
3453 3427 3396 3363 3326 3299 3272 3232 3196 3168
3123 3094 3073 3050 3012 2986 2983 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]

```

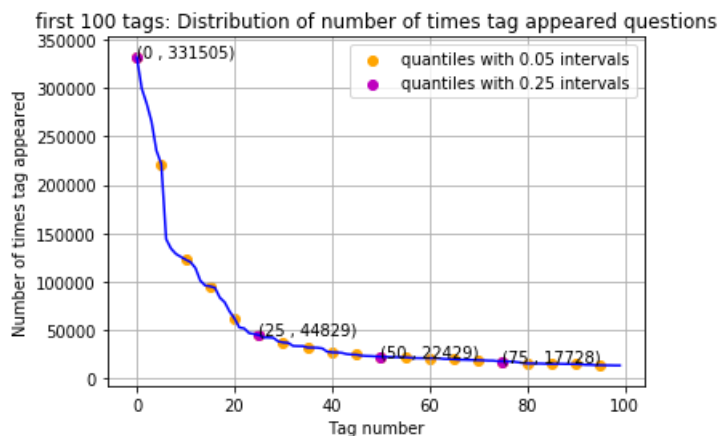
4526	4481	4429	4335	4310	4281	4239	4228	4195	4159
4144	4088	4050	4002	3957	3929	3874	3849	3818	3797
3750	3703	3685	3658	3615	3593	3564	3521	3505	3483]

In [27]:

```
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 i
ntervals")
# 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 in
tervals")

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 [28]:

```
# 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 [29]:

```
#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]

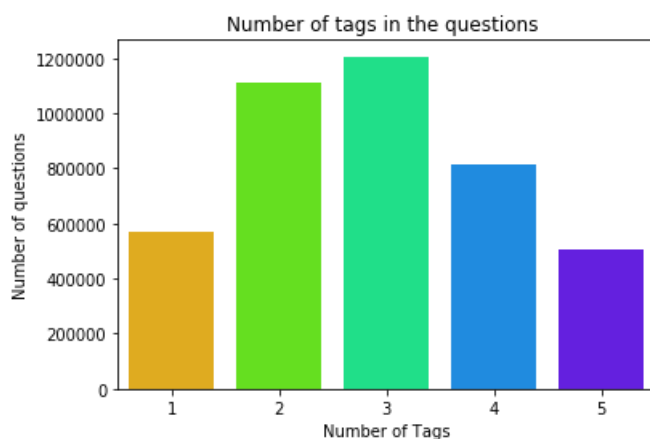
In [30]:

```
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 [31]:

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



Observations:

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

3.2.5 Most Frequent Tags

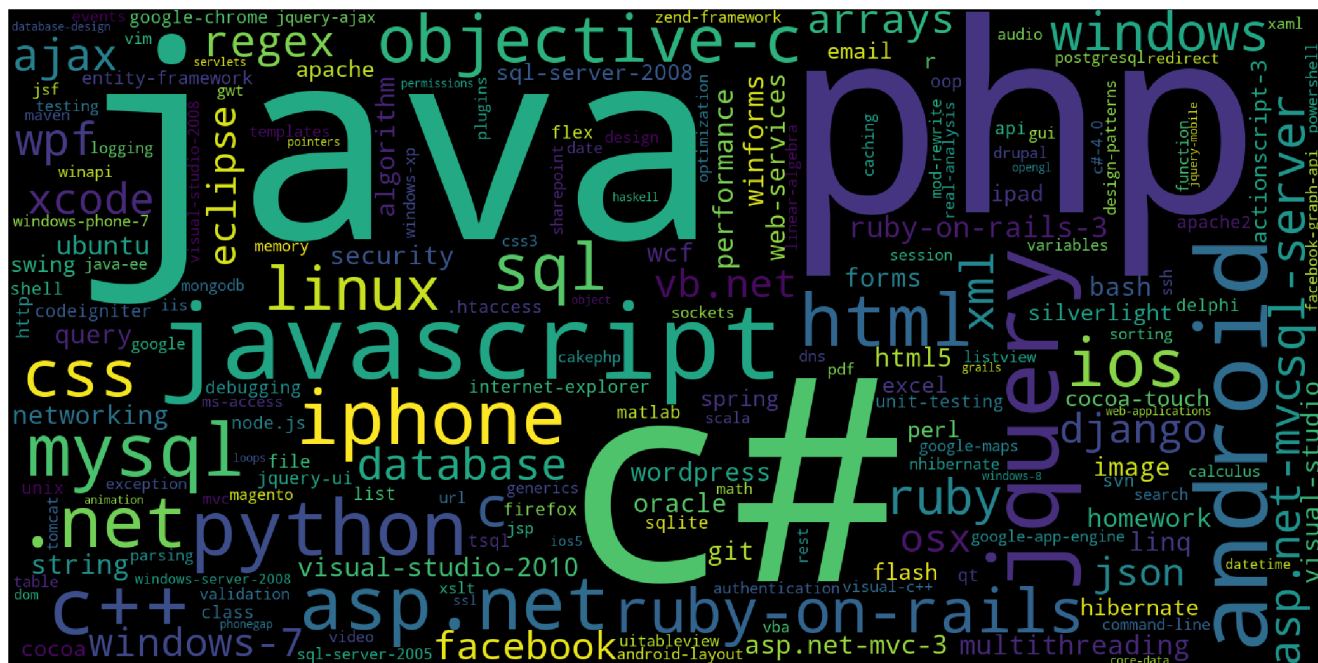
In [32]:

```
# Plotting word cloud
start = datetime.now()

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

fig = plt.figure(figsize=(30,20))
```

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



Time taken to run this cell : 0:00:11.389758

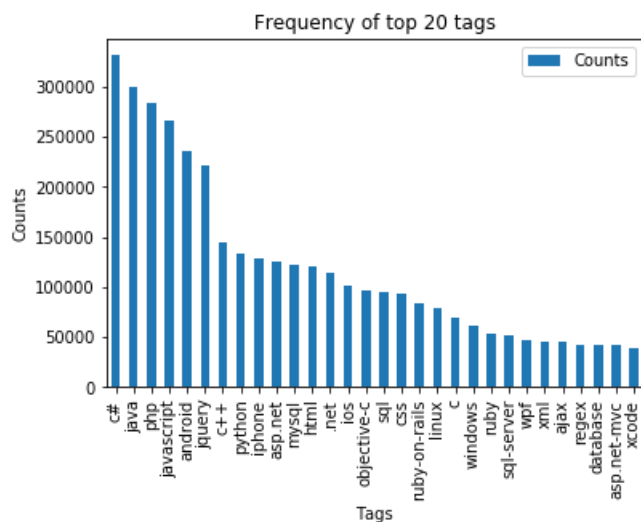
Observations:

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 [33]:

```
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 0.5M 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 [34]:

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

[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\sasha\AppData\Roaming\nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

Out[34]:

True

In [35]:

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

In [36]:

```
#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'"
    cursr.execute(str)
```

```

table_names = cursr.execute(str)
print("Tables in the database:")
tables =table_names.fetchall()
print(tables[0][0])
return(len(tables))

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

sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question text NOT NULL, code
text, tags text, words_pre integer, words_post integer, is_code integer);"""
create_database_table("E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/Processed.db",
sql_create_table)

```

Tables in the database:
QuestionsProcessed

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

In [37]:

```
nltk.download('punkt')
```

```

[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\sesha\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!

```

Out[37]:

True

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

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

In [38]:

```

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("E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/Titlemoreweight.db", sql_create_table)

```

Tables in the database:
QuestionsProcessed

In [39]:

```

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

```



```

read_db = 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/train_no_dup.db'
write_db = 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/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

In [40]:

```
reader
```

Out[40]:

```
<sqlite3.Cursor at 0x262015f91f0>
```

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 [41]:

```

#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_processed = 0
for row in tqdm(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)

```

100111it [03:17, 564.54it/s]

number of questions completed= 100000

200184it [06:27, 623.61it/s]

number of questions completed= 200000

300067it [09:47, 500.62it/s]

number of questions completed= 300000

400086it [13:01, 478.98it/s]

number of questions completed= 400000

500000it [16:15, 512.45it/s]

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 : 0:16:15.722813

In [42]:

```

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

```

```
conn_r.close()
conn_w.close()
```

In [43]:

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

Questions after preprocessed

=====

('dynam datagrid bind silverlight dynam datagrid bind silverlight dynam datagrid bind silverlight
bind datagrid dynam code wrote code debug code block seem bind correct grid come column form come
grid column although necessari bind nthank repli advance..',)

('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid
java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid
java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid follow guid link instal js
tl got follow error tri launch jsp page java.lang.noclassdeffounderror javax servlet jsp tagext ta
glibraryvalid taglib declar instal jstl 1.1 tomcat webapp tri project work also tri version 1.2 js
tl still messag caus solv',)

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

('better way updat feed fb php sdk better way updat feed fb php sdk better way updat feed fb php s
dk novic facebook api read mani tutori still confused.i find post feed api method like correct sec
ond way use curl someth like way better',)

('btnadd click event open two window record ad btnadd click event open two window record ad btnadd
click event open two window record ad open window search.aspx use code hav add button search.aspx
nwhen insert record btnadd click event open anoth window nafter insert record close window',)

('sql inject issu prevent correct form submiss php sql inject issu prevent correct form submiss ph
p sql inject issu prevent correct form submiss php check everyth think make sure input field safe
type sql inject good news safe bad news one tag mess form submiss place even touch life figur exac
t html use templat file forgiv okay entir php script get execut see data post none forum field pos
t problem use someth titl field none data get post current use print post see submit noth work fla
wless statement though also mention script work flawless local machin use host come across problem
state list input test mess',)

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

('hql equival sql queri hql equival sql queri hql equival sql queri hql queri replac name class pr
operti name error occur hql error',)

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

Saving Preprocessed data to a Database

In [44]:

```
#Taking 0.5 Million entries to a dataframe.
write_db = 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/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 [45]:

```
preprocessed_data.head()
```

Out[45]:

	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 [46]:

```
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 [47]:

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

Selecting 500 Tags

In [48]:

```
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

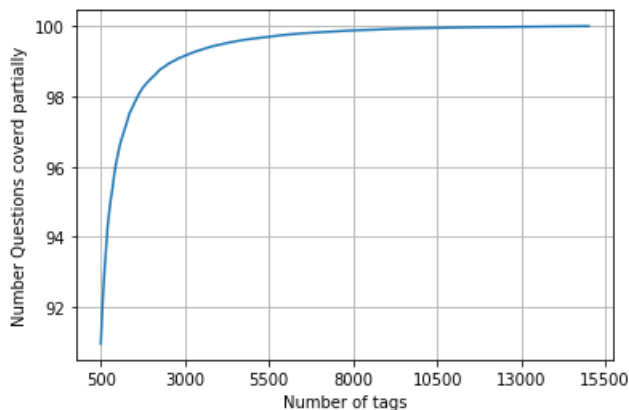
def questions_explained_fn(n):
    multilabel_yn = tags_to_choose(n)
    x= multilabel_yn.sum(axis=1)
    return (np.count_nonzero(x==0))
```

In [49]:

```
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 [50]:

```
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 [51]:

```
# 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 [52]:

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

Out[52]:

```
['preprocessed_data.pkl']
```

In [53]:

```
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 [54]:

```
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 [55]:

```
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=(1,4))
x_train_multilabel = vectorizer.fit_transform(x_train['question'])
x_test_multilabel = vectorizer.transform(x_test['question'])
print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:12:09.257285

In [56]:

```
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, 95585) Y : (400000, 500)
Dimensions of test data X: (100000, 95585) Y: (100000, 500)

4.5.3 OneVsRest Classifier with SGDClassifier using TFIDF

In [57]:

```
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 : 0.23624
Hamming loss 0.00278182
Micro-average quality numbers
Precision: 0.7214, Recall: 0.3255, F1-measure: 0.4486
Macro-average quality numbers
Precision: 0.5473, Recall: 0.2577, F1-measure: 0.3346

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

0	0.94	0.64	0.77	5519
1	0.68	0.27	0.38	8190
2	0.81	0.38	0.51	6529
3	0.81	0.43	0.56	3231
4	0.81	0.41	0.54	6430
5	0.82	0.34	0.48	2879
6	0.88	0.49	0.63	5086
7	0.88	0.54	0.67	4533

7	0.60	0.13	0.21	3000
8	0.60	0.13	0.21	3000
9	0.81	0.52	0.64	2765
10	0.60	0.16	0.26	3051
11	0.69	0.33	0.45	3009
12	0.64	0.24	0.35	2630
13	0.70	0.23	0.35	1426
14	0.90	0.53	0.67	2548
15	0.67	0.18	0.28	2371
16	0.65	0.23	0.34	873
17	0.89	0.61	0.72	2151
18	0.63	0.22	0.33	2204
19	0.72	0.40	0.51	831
20	0.77	0.42	0.54	1860
21	0.27	0.08	0.12	2023
22	0.49	0.22	0.31	1513
23	0.91	0.48	0.63	1207
24	0.57	0.29	0.38	506
25	0.68	0.30	0.42	425
26	0.65	0.39	0.49	793
27	0.59	0.32	0.41	1291
28	0.74	0.36	0.48	1208
29	0.42	0.09	0.14	406
30	0.77	0.18	0.29	504
31	0.29	0.11	0.15	732
32	0.59	0.24	0.34	441
33	0.55	0.17	0.26	1645
34	0.71	0.25	0.37	1058
35	0.83	0.54	0.66	946
36	0.67	0.19	0.29	644
37	0.98	0.67	0.79	136
38	0.64	0.36	0.46	570
39	0.84	0.28	0.43	766
40	0.61	0.28	0.38	1132
41	0.44	0.18	0.25	174
42	0.80	0.53	0.64	210
43	0.80	0.41	0.54	433
44	0.66	0.49	0.56	626
45	0.73	0.32	0.44	852
46	0.75	0.42	0.53	534
47	0.33	0.14	0.20	350
48	0.74	0.50	0.59	496
49	0.80	0.61	0.69	785
50	0.17	0.04	0.06	475
51	0.34	0.11	0.16	305
52	0.44	0.03	0.06	251
53	0.68	0.40	0.50	914
54	0.46	0.15	0.23	728
55	0.25	0.01	0.02	258
56	0.46	0.19	0.27	821
57	0.48	0.09	0.15	541
58	0.78	0.28	0.41	748
59	0.94	0.62	0.75	724
60	0.33	0.07	0.11	660
61	0.84	0.18	0.30	235
62	0.91	0.71	0.80	718
63	0.83	0.63	0.72	468
64	0.53	0.31	0.39	191
65	0.38	0.13	0.19	429
66	0.30	0.05	0.09	415
67	0.75	0.48	0.58	274
68	0.82	0.52	0.64	510
69	0.68	0.45	0.54	466
70	0.29	0.06	0.10	305
71	0.50	0.15	0.24	247
72	0.78	0.48	0.59	401
73	0.98	0.73	0.84	86
74	0.75	0.38	0.50	120
75	0.89	0.68	0.77	129
76	0.67	0.01	0.02	473
77	0.35	0.25	0.29	143
78	0.79	0.44	0.57	347
79	0.73	0.24	0.36	479
80	0.55	0.33	0.41	279
81	0.80	0.18	0.30	461
82	0.13	0.01	0.01	298
83	0.78	0.46	0.58	396
84	0.55	0.33	0.41	184

85	0.69	0.21	0.32	573
86	0.49	0.05	0.09	325
87	0.50	0.26	0.35	273
88	0.42	0.19	0.26	135
89	0.30	0.07	0.12	232
90	0.57	0.31	0.40	409
91	0.63	0.25	0.36	420
92	0.76	0.53	0.63	408
93	0.68	0.48	0.56	241
94	0.31	0.04	0.08	211
95	0.32	0.08	0.12	277
96	0.27	0.03	0.06	410
97	0.89	0.32	0.47	501
98	0.76	0.60	0.67	136
99	0.54	0.31	0.39	239
100	0.52	0.13	0.21	324
101	0.93	0.60	0.73	277
102	0.92	0.71	0.80	613
103	0.52	0.16	0.24	157
104	0.21	0.05	0.08	295
105	0.85	0.34	0.49	334
106	0.80	0.13	0.22	335
107	0.75	0.48	0.58	389
108	0.55	0.24	0.33	251
109	0.53	0.41	0.46	317
110	0.76	0.09	0.15	187
111	0.38	0.06	0.10	140
112	0.60	0.27	0.37	154
113	0.65	0.18	0.28	332
114	0.45	0.28	0.34	323
115	0.48	0.21	0.29	344
116	0.76	0.50	0.60	370
117	0.57	0.23	0.32	313
118	0.78	0.68	0.72	874
119	0.44	0.19	0.27	293
120	0.00	0.00	0.00	200
121	0.76	0.47	0.58	463
122	0.37	0.09	0.15	119
123	0.75	0.01	0.02	256
124	0.91	0.70	0.79	195
125	0.40	0.12	0.18	138
126	0.79	0.48	0.60	376
127	0.14	0.03	0.05	122
128	0.13	0.03	0.05	252
129	0.49	0.13	0.21	144
130	0.39	0.08	0.13	150
131	0.33	0.02	0.04	210
132	0.67	0.25	0.37	361
133	0.93	0.54	0.68	453
134	0.88	0.73	0.80	124
135	0.21	0.03	0.06	91
136	0.69	0.26	0.38	128
137	0.56	0.33	0.42	218
138	0.77	0.15	0.25	243
139	0.36	0.17	0.24	149
140	0.76	0.43	0.55	318
141	0.30	0.11	0.16	159
142	0.66	0.34	0.45	274
143	0.86	0.72	0.78	362
144	0.59	0.16	0.25	118
145	0.65	0.37	0.47	164
146	0.59	0.27	0.37	461
147	0.67	0.40	0.50	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	0.88	0.65	0.75	55
152	0.79	0.45	0.57	387
153	0.47	0.09	0.16	150
154	0.58	0.11	0.19	281
155	0.26	0.05	0.08	202
156	0.75	0.63	0.69	130
157	0.28	0.07	0.12	245
158	0.88	0.58	0.70	177
159	0.49	0.26	0.34	130
160	0.50	0.12	0.20	336
161	0.93	0.57	0.70	220

161	0.95	0.57	0.70	220
162	0.13	0.02	0.04	229
163	0.90	0.41	0.56	316
164	0.74	0.36	0.48	283
165	0.64	0.32	0.43	197
166	0.47	0.21	0.29	101
167	0.47	0.18	0.26	231
168	0.58	0.23	0.33	370
169	0.41	0.18	0.25	258
170	0.30	0.06	0.10	101
171	0.40	0.22	0.29	89
172	0.51	0.34	0.40	193
173	0.42	0.21	0.28	309
174	0.49	0.13	0.21	172
175	0.94	0.71	0.81	95
176	0.94	0.58	0.72	346
177	0.93	0.43	0.58	322
178	0.65	0.47	0.55	232
179	0.35	0.06	0.11	125
180	0.54	0.26	0.35	145
181	0.40	0.10	0.16	77
182	0.16	0.02	0.04	182
183	0.62	0.32	0.42	257
184	0.08	0.01	0.02	216
185	0.33	0.07	0.11	242
186	0.38	0.15	0.22	165
187	0.76	0.57	0.65	263
188	0.31	0.09	0.14	174
189	0.75	0.30	0.43	136
190	0.88	0.50	0.63	202
191	0.43	0.13	0.20	134
192	0.72	0.40	0.51	230
193	0.43	0.18	0.25	90
194	0.58	0.47	0.52	185
195	0.19	0.04	0.06	156
196	0.36	0.07	0.12	160
197	0.64	0.07	0.12	266
198	0.38	0.05	0.09	284
199	0.39	0.06	0.11	145
200	0.94	0.69	0.79	212
201	0.68	0.21	0.32	317
202	0.78	0.52	0.63	427
203	0.32	0.09	0.14	232
204	0.49	0.22	0.30	217
205	0.49	0.44	0.47	527
206	0.14	0.02	0.03	124
207	0.48	0.10	0.16	103
208	0.90	0.48	0.63	287
209	0.33	0.08	0.13	193
210	0.72	0.32	0.45	220
211	0.82	0.19	0.31	140
212	0.14	0.02	0.03	161
213	0.50	0.22	0.31	72
214	0.61	0.45	0.52	396
215	0.86	0.31	0.46	134
216	0.50	0.06	0.11	400
217	0.54	0.25	0.35	75
218	0.96	0.75	0.85	219
219	0.75	0.35	0.48	210
220	0.90	0.59	0.71	298
221	0.97	0.60	0.74	266
222	0.78	0.43	0.55	290
223	0.08	0.01	0.01	128
224	0.79	0.40	0.53	159
225	0.60	0.29	0.39	164
226	0.64	0.36	0.46	144
227	0.58	0.31	0.40	276
228	0.15	0.02	0.03	235
229	0.36	0.02	0.04	216
230	0.35	0.18	0.23	228
231	0.71	0.47	0.57	64
232	0.38	0.06	0.10	103
233	0.71	0.29	0.41	216
234	0.82	0.08	0.14	116
235	0.55	0.36	0.44	77
236	0.96	0.64	0.77	67
237	0.56	0.06	0.12	218
238	0.27	0.06	0.10	120

238	0.27	0.00	0.10	139
239	0.17	0.01	0.02	94
240	0.54	0.27	0.36	77
241	0.52	0.09	0.15	167
242	0.84	0.30	0.44	86
243	0.48	0.17	0.25	58
244	0.63	0.17	0.27	269
245	0.17	0.05	0.08	112
246	0.95	0.73	0.83	255
247	0.47	0.26	0.33	58
248	0.22	0.02	0.04	81
249	0.00	0.00	0.00	131
250	0.42	0.20	0.28	93
251	0.66	0.28	0.39	154
252	0.35	0.05	0.08	129
253	0.57	0.33	0.42	83
254	0.38	0.09	0.14	191
255	0.15	0.02	0.04	219
256	0.25	0.04	0.07	130
257	0.47	0.28	0.35	93
258	0.71	0.45	0.55	217
259	0.31	0.10	0.15	141
260	0.95	0.13	0.23	143
261	0.55	0.12	0.20	219
262	0.56	0.29	0.38	107
263	0.39	0.23	0.29	236
264	0.26	0.17	0.20	119
265	0.38	0.15	0.22	72
266	0.00	0.00	0.00	70
267	0.30	0.13	0.18	107
268	0.67	0.43	0.53	169
269	0.30	0.10	0.15	129
270	0.74	0.53	0.62	159
271	0.81	0.34	0.48	190
272	0.62	0.23	0.33	248
273	0.92	0.70	0.79	264
274	0.89	0.64	0.74	105
275	0.53	0.08	0.13	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	0.57	0.44	0.50	80
281	0.67	0.54	0.60	217
282	0.74	0.46	0.57	175
283	0.31	0.05	0.09	269
284	0.67	0.27	0.38	74
285	0.86	0.50	0.63	206
286	0.90	0.59	0.71	227
287	0.81	0.30	0.44	130
288	0.29	0.05	0.09	129
289	0.50	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	0.84	0.42	0.56	109
294	0.40	0.26	0.31	140
295	0.25	0.08	0.12	241
296	0.23	0.10	0.14	72
297	0.24	0.04	0.06	107
298	0.79	0.38	0.51	61
299	0.94	0.38	0.54	77
300	0.15	0.05	0.08	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.56	0.35	0.43	176
304	0.96	0.71	0.82	230
305	0.96	0.58	0.73	156
306	0.50	0.36	0.41	146
307	0.23	0.06	0.10	98
308	0.00	0.00	0.00	78
309	0.78	0.07	0.14	94
310	0.79	0.35	0.49	162
311	0.81	0.52	0.63	116
312	0.50	0.28	0.36	57
313	0.75	0.05	0.09	65
314	0.51	0.36	0.42	138
315	0.52	0.21	0.22	105

315	0.53	0.21	0.30	195
316	0.44	0.25	0.31	69
317	0.35	0.10	0.16	134
318	0.49	0.33	0.40	148
319	0.84	0.43	0.57	161
320	0.23	0.14	0.18	104
321	0.85	0.54	0.66	156
322	0.57	0.31	0.40	134
323	0.57	0.39	0.46	232
324	0.43	0.16	0.24	92
325	0.46	0.29	0.36	197
326	0.12	0.02	0.04	126
327	0.50	0.04	0.08	115
328	0.98	0.64	0.78	198
329	0.60	0.31	0.41	125
330	0.84	0.20	0.32	81
331	0.53	0.09	0.15	94
332	0.50	0.02	0.03	56
333	0.17	0.03	0.06	260
334	0.20	0.03	0.06	60
335	0.28	0.07	0.12	110
336	0.64	0.42	0.51	71
337	0.13	0.03	0.05	66
338	0.46	0.32	0.38	150
339	0.00	0.00	0.00	54
340	0.86	0.55	0.67	195
341	0.88	0.19	0.31	79
342	0.47	0.21	0.29	38
343	0.67	0.37	0.48	43
344	0.47	0.22	0.30	68
345	0.68	0.38	0.49	73
346	0.30	0.03	0.05	116
347	0.88	0.33	0.48	111
348	0.27	0.10	0.14	63
349	0.82	0.59	0.69	104
350	0.64	0.48	0.55	44
351	0.67	0.20	0.31	40
352	0.95	0.40	0.56	136
353	0.40	0.19	0.25	54
354	0.45	0.04	0.07	134
355	0.53	0.29	0.38	120
356	0.53	0.21	0.31	228
357	0.66	0.26	0.37	269
358	0.71	0.36	0.48	80
359	0.87	0.44	0.58	140
360	0.33	0.11	0.17	125
361	0.89	0.63	0.74	169
362	0.11	0.04	0.05	56
363	0.94	0.65	0.77	154
364	0.33	0.05	0.09	58
365	0.25	0.13	0.17	71
366	1.00	0.65	0.79	54
367	0.29	0.03	0.06	116
368	0.00	0.00	0.00	54
369	0.00	0.00	0.00	71
370	0.20	0.03	0.06	61
371	0.46	0.08	0.14	71
372	0.66	0.48	0.56	52
373	0.79	0.37	0.50	150
374	0.37	0.14	0.20	93
375	0.15	0.03	0.05	67
376	0.00	0.00	0.00	76
377	0.75	0.20	0.31	106
378	0.27	0.03	0.06	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	0.51	0.29	0.37	110
384	0.00	0.00	0.00	155
385	0.40	0.08	0.13	50
386	0.21	0.09	0.13	64
387	0.36	0.05	0.09	93
388	0.60	0.29	0.39	102
389	0.07	0.01	0.02	108
390	0.96	0.64	0.77	178
391	0.62	0.17	0.27	115
392	0.81	0.48	0.54	40

392	0.81	0.40	0.54	42
393	0.00	0.00	0.00	134
394	0.40	0.02	0.03	112
395	0.43	0.13	0.20	176
396	0.46	0.09	0.15	125
397	0.72	0.24	0.36	224
398	0.89	0.52	0.66	63
399	0.00	0.00	0.00	59
400	0.49	0.35	0.41	63
401	0.46	0.16	0.24	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	0.30	0.07	0.11	92
406	0.80	0.20	0.31	41
407	0.67	0.37	0.48	43
408	0.81	0.32	0.46	160
409	0.17	0.10	0.13	50
410	0.00	0.00	0.00	19
411	0.37	0.10	0.15	175
412	0.29	0.06	0.09	72
413	0.50	0.05	0.10	95
414	0.18	0.03	0.05	97
415	0.33	0.17	0.22	48
416	0.48	0.29	0.36	83
417	0.50	0.07	0.13	40
418	0.37	0.08	0.13	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	0.61	0.34	0.44	73
423	0.47	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	0.69	0.80	29
429	0.99	0.66	0.79	389
430	0.63	0.36	0.46	167
431	0.52	0.09	0.15	123
432	0.48	0.36	0.41	39
433	0.27	0.15	0.19	82
434	1.00	0.67	0.80	66
435	0.66	0.45	0.54	93
436	0.52	0.25	0.34	87
437	0.25	0.06	0.09	86
438	0.75	0.46	0.57	104
439	0.62	0.13	0.21	100
440	0.25	0.01	0.01	141
441	0.42	0.25	0.31	110
442	0.38	0.12	0.19	123
443	0.47	0.11	0.18	71
444	0.41	0.06	0.11	109
445	0.35	0.17	0.23	48
446	0.44	0.26	0.33	76
447	0.24	0.11	0.15	38
448	0.68	0.53	0.60	81
449	0.55	0.14	0.22	132
450	0.46	0.26	0.33	81
451	0.88	0.29	0.44	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.94	0.43	0.59	70
455	0.48	0.06	0.11	155
456	0.54	0.16	0.25	43
457	0.50	0.18	0.27	72
458	0.28	0.08	0.12	62
459	0.82	0.13	0.23	69
460	0.07	0.01	0.02	119
461	0.75	0.11	0.20	79
462	0.69	0.23	0.35	47
463	0.17	0.03	0.05	104
464	0.66	0.35	0.46	106
465	0.50	0.11	0.18	64
466	0.58	0.29	0.39	173
467	0.80	0.34	0.47	107
468	0.80	0.13	0.22	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	0.39	0.27	0.32	143
473	0.70	0.32	0.44	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	0.68	176
478	0.94	0.71	0.81	24
479	0.09	0.01	0.02	92
480	0.81	0.48	0.60	100
481	0.50	0.17	0.26	103
482	0.49	0.23	0.31	74
483	0.84	0.58	0.69	105
484	0.25	0.02	0.04	83
485	0.25	0.02	0.04	82
486	0.38	0.11	0.17	71
487	0.44	0.19	0.27	120
488	0.33	0.02	0.04	105
489	0.74	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	0.52	0.18	0.27	61
495	0.98	0.58	0.73	344
496	0.37	0.19	0.25	52
497	0.63	0.20	0.30	137
498	0.25	0.03	0.05	98
499	0.68	0.16	0.27	79
micro avg	0.72	0.33	0.45	173812
macro avg	0.55	0.26	0.33	173812
weighted avg	0.67	0.33	0.42	173812
samples avg	0.41	0.31	0.33	173812

Time taken to run this cell : 0:25:37.386384

In [58]:

```
joblib.dump(classifier, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/lr_with_more_title_weight.pkl')
```

Out[58]:

```
['E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/lr_with_more_title_weight.pkl']
```

5. Assignments

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

TASK 1: BOW and (1,2,3,4) n-grams with Lgistic Regression(OvR)

In [59]:

```
alpha=[10**-3,10**-2,10**-1]
```

In [60]:

```
start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009, max_features=200000, \
                             tokenizer = lambda x: x.split(), ngram_range=(1,4))
```

In [61]:

```
x_train_multilabel = vectorizer.fit_transform(x_train['question'])
```

Time taken to run this cell : 0:16:12.674519

```
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)
```

Dump and load train and test data into joblib

```
joblib.dump(x_train_multilabel, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag  
Predictor/x_train_BOW.pkl')  
joblib.dump(x_test_multilabel, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag  
Predictor/x_test_BOW.pkl')  
joblib.dump(y_train, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/y_train.pkl')  
joblib.dump(y_test, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/y_test.pkl')
```

```
['E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/y test.pkl']
```

```
x_train_multilabel = joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/x_train_BOW.pkl')
y_train = joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/y_train.pkl')
```

```
x_test_multilabel = joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/x_test_BOW.pkl')
y_test = joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/y_test.pkl')
```

OneVsRestClassifier with Logistic regression (alpha tuning using Gridsearch)

OneVsRestClassifier with SGDClassifier(penalty=l2, loss=log)==> {Logistic regression}

[illegible]

```

print("Gridsearchcv")
best_model1=model11.fit(x_train_multilabel, y_train)
print('fit model')
Train_model_score=best_model1.score(x_train_multilabel,
                                     y_train)

#print("best_model1")
cv_scores.append(Train_model_score.mean())

fscore = [x for x in cv_scores]

# determining best alpha
optimal_alpha21 = alpha[fscore.index(max(fscore))]
print('\n The optimal value of alpha with penalty=l2 and loss= log is %d.' % optimal_alpha21)

# Plots
fig4 = plt.figure( facecolor='c', edgecolor='k')
plt.plot(alpha, fscore,color='green', marker='o', linestyle='dashed',
linewidth=2, markersize=12)

for xy in zip(alpha, np.round(fscore,3)):
    plt.annotate('(%s, %s)' % xy, xy=xy, textcoords='data')

plt.xlabel('Hyper parameter Alpha')
plt.ylabel('F1_Score value ')
plt.show()

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

```

0%| | 1/3 [21:36
[00:00<?, ?it/s]

0.001
{'estimator__alpha': [0.001], 'estimator__loss': ['log'], 'estimator__penalty': ['l2']}

Gridsearchcv
fit model

33%| | 1/3 [21:36
3:12, 1296.06s/it]

0.01
{'estimator__alpha': [0.01], 'estimator__loss': ['log'], 'estimator__penalty': ['l2']}

Gridsearchcv
fit model

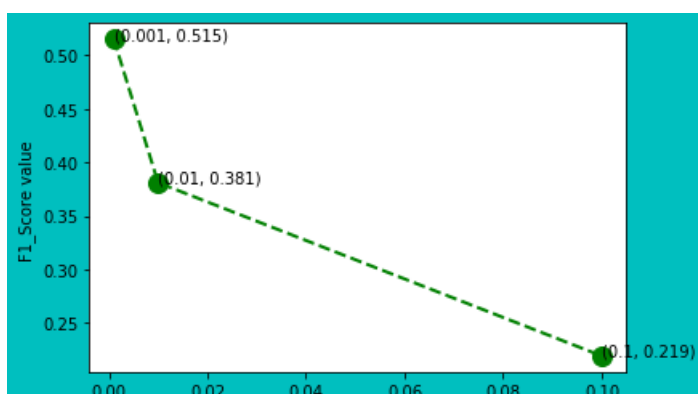
67%| | 2/3
[41:07<20:58, 1258.71s/it]

0.1
{'estimator__alpha': [0.1], 'estimator__loss': ['log'], 'estimator__penalty': ['l2']}

Gridsearchcv
fit model

100%| | 3/3
[1:06:19<00:00, 1326.47s/it]

The optimal value of alpha with penalty=l2 and loss= log is 0.



Time taken to run this cell : 1:06:20.853640

In [70]:

```
print(optimal_alpha21)
```

0.001

In [71]:

```
start = datetime.now()
best_model1 = OneVsRestClassifier(SGDClassifier(loss='log', alpha=optimal_alpha21,
                                              penalty='l2'), n_jobs=-1)
best_model1.fit(x_train_multilabel, y_train)
```

Out[71]:

```
OneVsRestClassifier(estimator=SGDClassifier(alpha=0.001, average=False,
                                           class_weight=None,
                                           early_stopping=False, epsilon=0.1,
                                           eta0=0.0, fit_intercept=True,
                                           l1_ratio=0.15,
                                           learning_rate='optimal', loss='log',
                                           max_iter=1000, n_iter_no_change=5,
                                           n_jobs=None, penalty='l2',
                                           power_t=0.5, random_state=None,
                                           shuffle=True, tol=0.001,
                                           validation_fraction=0.1, verbose=0,
                                           warm_start=False),
                   n_jobs=-1)
```

In [72]:

```
joblib.dump(best_model1, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/best_model1_LR.pkl')
```

Out[72]:

```
['E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/best_model1_LR.pkl']
```

In [73]:

```
best_model1=joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/best_model1_LR.pkl')
```

In [74]:

```
predictions = best_model1.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-averasge 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)) #printing classification report for all
500 labels
```



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

Accuracy : 0.21171

Hamming loss 0.00296114

Micro-average quality numbers

Precision: 0.6532, Recall: 0.3159, F1-measure: 0.4259

Macro-average quality numbers

Precision: 0.4968, Recall: 0.2351, F1-measure: 0.3059

	precision	recall	f1-score	support
0	0.95	0.65	0.77	5519
1	0.68	0.25	0.37	8190
2	0.81	0.35	0.49	6529
3	0.82	0.42	0.55	3231
4	0.81	0.42	0.55	6430
5	0.79	0.35	0.48	2879
6	0.88	0.49	0.63	5086
7	0.87	0.56	0.68	4533
8	0.59	0.14	0.23	3000
9	0.81	0.57	0.67	2765
10	0.59	0.21	0.31	3051
11	0.70	0.34	0.46	3009
12	0.64	0.25	0.36	2630
13	0.74	0.27	0.39	1426
14	0.89	0.50	0.64	2548
15	0.65	0.13	0.22	2371
16	0.63	0.26	0.37	873
17	0.86	0.61	0.72	2151
18	0.65	0.24	0.35	2204
19	0.71	0.43	0.54	831
20	0.78	0.39	0.52	1860
21	0.28	0.11	0.16	2023
22	0.47	0.28	0.35	1513
23	0.90	0.49	0.63	1207
24	0.50	0.33	0.40	506
25	0.61	0.29	0.39	425
26	0.59	0.44	0.50	793
27	0.56	0.40	0.47	1291
28	0.71	0.31	0.44	1208
29	0.35	0.10	0.16	406
30	0.58	0.15	0.24	504
31	0.27	0.17	0.21	732
32	0.57	0.28	0.38	441
33	0.52	0.26	0.35	1645
34	0.70	0.23	0.35	1058
35	0.83	0.57	0.68	946
36	0.61	0.22	0.32	644
37	0.98	0.65	0.78	136
38	0.60	0.44	0.51	570
39	0.85	0.22	0.35	766
40	0.60	0.30	0.40	1132
41	0.48	0.23	0.31	174
42	0.69	0.43	0.53	210
43	0.76	0.40	0.53	433
44	0.63	0.49	0.55	626
45	0.64	0.31	0.42	852
46	0.70	0.45	0.55	534
47	0.28	0.23	0.25	350
48	0.72	0.50	0.59	496
49	0.79	0.63	0.70	785
50	0.21	0.12	0.15	475
51	0.29	0.14	0.19	305
52	0.39	0.06	0.10	251
53	0.67	0.39	0.49	914
54	0.45	0.21	0.28	728
55	0.00	0.00	0.00	258
56	0.39	0.25	0.31	821
57	0.42	0.11	0.18	541
58	0.80	0.24	0.37	748
59	0.95	0.58	0.72	724
60	0.26	0.07	0.11	660
61	0.85	0.19	0.31	235
62	0.87	0.69	0.77	718
63	0.84	0.54	0.66	468
64	0.50	0.45	0.47	191
65	0.26	0.15	0.19	429
66	0.26	0.14	0.18	415

66	0.20	0.14	0.10	410
67	0.66	0.47	0.55	274
68	0.84	0.48	0.61	510
69	0.65	0.43	0.52	466
70	0.28	0.13	0.18	305
71	0.34	0.18	0.24	247
72	0.75	0.42	0.53	401
73	0.92	0.65	0.76	86
74	0.71	0.33	0.45	120
75	0.90	0.60	0.72	129
76	0.40	0.01	0.02	473
77	0.36	0.35	0.35	143
78	0.76	0.38	0.50	347
79	0.69	0.21	0.32	479
80	0.47	0.41	0.44	279
81	0.77	0.11	0.19	461
82	0.21	0.07	0.11	298
83	0.71	0.42	0.53	396
84	0.46	0.40	0.42	184
85	0.49	0.25	0.34	573
86	0.24	0.07	0.11	325
87	0.48	0.25	0.33	273
88	0.32	0.25	0.28	135
89	0.22	0.18	0.20	232
90	0.49	0.40	0.44	409
91	0.62	0.32	0.43	420
92	0.75	0.45	0.57	408
93	0.50	0.49	0.49	241
94	0.30	0.10	0.15	211
95	0.27	0.18	0.22	277
96	0.23	0.08	0.12	410
97	0.87	0.16	0.27	501
98	0.77	0.57	0.66	136
99	0.49	0.30	0.37	239
100	0.46	0.18	0.26	324
101	0.91	0.50	0.65	277
102	0.90	0.64	0.75	613
103	0.45	0.20	0.27	157
104	0.20	0.14	0.17	295
105	0.65	0.37	0.48	334
106	0.76	0.06	0.11	335
107	0.75	0.49	0.59	389
108	0.51	0.34	0.41	251
109	0.47	0.36	0.41	317
110	0.45	0.09	0.15	187
111	0.36	0.07	0.12	140
112	0.48	0.26	0.34	154
113	0.61	0.14	0.23	332
114	0.43	0.33	0.38	323
115	0.42	0.17	0.24	344
116	0.71	0.46	0.56	370
117	0.55	0.21	0.31	313
118	0.80	0.46	0.58	874
119	0.35	0.23	0.28	293
120	0.14	0.04	0.07	200
121	0.78	0.42	0.54	463
122	0.37	0.22	0.28	119
123	0.50	0.00	0.01	256
124	0.91	0.64	0.75	195
125	0.39	0.24	0.30	138
126	0.79	0.49	0.61	376
127	0.16	0.06	0.08	122
128	0.21	0.09	0.13	252
129	0.39	0.10	0.16	144
130	0.42	0.09	0.14	150
131	0.19	0.03	0.06	210
132	0.58	0.22	0.32	361
133	0.94	0.39	0.55	453
134	0.88	0.66	0.76	124
135	0.12	0.01	0.02	91
136	0.52	0.28	0.37	128
137	0.44	0.34	0.38	218
138	0.40	0.08	0.13	243
139	0.33	0.24	0.28	149
140	0.68	0.32	0.44	318
141	0.18	0.14	0.16	159
142	0.65	0.42	0.51	274
143	0.84	0.60	0.70	260

143	0.04	0.03	0.12	302
144	0.50	0.23	0.31	118
145	0.59	0.41	0.48	164
146	0.57	0.28	0.37	461
147	0.65	0.44	0.53	159
148	0.34	0.16	0.21	166
149	0.97	0.31	0.47	346
150	0.56	0.07	0.12	350
151	0.88	0.42	0.57	55
152	0.76	0.45	0.56	387
153	0.40	0.04	0.07	150
154	0.55	0.06	0.11	281
155	0.30	0.15	0.20	202
156	0.73	0.55	0.63	130
157	0.29	0.11	0.16	245
158	0.89	0.47	0.62	177
159	0.43	0.28	0.34	130
160	0.49	0.24	0.33	336
161	0.85	0.51	0.64	220
162	0.18	0.10	0.13	229
163	0.90	0.29	0.44	316
164	0.70	0.28	0.40	283
165	0.55	0.28	0.37	197
166	0.32	0.20	0.24	101
167	0.41	0.23	0.30	231
168	0.43	0.23	0.30	370
169	0.43	0.31	0.36	258
170	0.21	0.09	0.13	101
171	0.53	0.22	0.31	89
172	0.39	0.34	0.36	193
173	0.42	0.28	0.34	309
174	0.50	0.12	0.19	172
175	0.91	0.74	0.81	95
176	0.93	0.43	0.59	346
177	0.95	0.24	0.38	322
178	0.57	0.44	0.50	232
179	0.50	0.06	0.10	125
180	0.43	0.21	0.28	145
181	0.48	0.21	0.29	77
182	0.13	0.07	0.09	182
183	0.56	0.34	0.42	257
184	0.13	0.06	0.08	216
185	0.31	0.14	0.19	242
186	0.27	0.18	0.21	165
187	0.77	0.46	0.57	263
188	0.32	0.17	0.22	174
189	0.77	0.32	0.45	136
190	0.94	0.36	0.52	202
191	0.36	0.13	0.19	134
192	0.65	0.30	0.41	230
193	0.32	0.19	0.24	90
194	0.59	0.52	0.55	185
195	0.08	0.04	0.05	156
196	0.23	0.07	0.11	160
197	0.09	0.01	0.02	266
198	0.42	0.09	0.15	284
199	0.14	0.03	0.05	145
200	0.93	0.51	0.66	212
201	0.54	0.21	0.31	317
202	0.72	0.42	0.53	427
203	0.25	0.13	0.17	232
204	0.43	0.24	0.31	217
205	0.48	0.38	0.42	527
206	0.09	0.04	0.06	124
207	0.35	0.15	0.21	103
208	0.81	0.34	0.48	287
209	0.25	0.11	0.15	193
210	0.68	0.25	0.37	220
211	0.65	0.08	0.14	140
212	0.08	0.06	0.07	161
213	0.54	0.29	0.38	72
214	0.60	0.42	0.49	396
215	0.79	0.17	0.28	134
216	0.42	0.07	0.12	400
217	0.44	0.24	0.31	75
218	0.97	0.51	0.67	219
219	0.76	0.29	0.42	210
220	0.04	0.24	0.50	200

220	0.94	0.34	0.50	298
221	0.96	0.41	0.58	266
222	0.71	0.28	0.40	290
223	0.24	0.04	0.07	128
224	0.75	0.36	0.49	159
225	0.35	0.22	0.27	164
226	0.55	0.35	0.43	144
227	0.51	0.41	0.46	276
228	0.07	0.02	0.03	235
229	0.23	0.02	0.04	216
230	0.35	0.27	0.30	228
231	0.67	0.45	0.54	64
232	0.16	0.08	0.10	103
233	0.72	0.20	0.32	216
234	0.56	0.13	0.21	116
235	0.57	0.43	0.49	77
236	0.93	0.60	0.73	67
237	0.62	0.05	0.09	218
238	0.16	0.11	0.13	139
239	0.23	0.03	0.06	94
240	0.41	0.16	0.23	77
241	0.45	0.09	0.15	167
242	0.77	0.23	0.36	86
243	0.46	0.21	0.29	58
244	0.43	0.24	0.31	269
245	0.17	0.06	0.09	112
246	0.96	0.54	0.69	255
247	0.39	0.21	0.27	58
248	0.33	0.06	0.10	81
249	0.03	0.01	0.01	131
250	0.29	0.22	0.25	93
251	0.60	0.28	0.38	154
252	0.21	0.05	0.09	129
253	0.56	0.36	0.44	83
254	0.22	0.10	0.14	191
255	0.17	0.07	0.10	219
256	0.13	0.03	0.05	130
257	0.41	0.32	0.36	93
258	0.66	0.35	0.46	217
259	0.24	0.11	0.15	141
260	0.85	0.12	0.21	143
261	0.55	0.11	0.18	219
262	0.43	0.26	0.33	107
263	0.33	0.28	0.31	236
264	0.21	0.18	0.19	119
265	0.33	0.21	0.25	72
266	0.17	0.07	0.10	70
267	0.25	0.13	0.17	107
268	0.61	0.33	0.43	169
269	0.25	0.19	0.22	129
270	0.69	0.50	0.58	159
271	0.49	0.17	0.26	190
272	0.58	0.21	0.30	248
273	0.93	0.43	0.59	264
274	0.87	0.50	0.64	105
275	0.13	0.03	0.05	104
276	0.10	0.02	0.03	115
277	0.86	0.51	0.64	170
278	0.63	0.19	0.29	145
279	0.89	0.31	0.46	230
280	0.54	0.34	0.42	80
281	0.67	0.48	0.56	217
282	0.72	0.39	0.51	175
283	0.36	0.10	0.16	269
284	0.67	0.27	0.38	74
285	0.86	0.36	0.51	206
286	0.92	0.43	0.58	227
287	0.76	0.25	0.37	130
288	0.26	0.07	0.11	129
289	0.16	0.07	0.10	80
290	0.15	0.12	0.14	99
291	0.83	0.21	0.33	208
292	0.38	0.12	0.18	67
293	0.80	0.33	0.47	109
294	0.33	0.36	0.34	140
295	0.18	0.14	0.16	241
296	0.21	0.15	0.18	72
297	0.07	0.11	0.16	107

297	0.27	0.11	0.16	107
298	0.67	0.43	0.52	61
299	0.86	0.39	0.54	77
300	0.17	0.09	0.12	111
301	0.00	0.00	0.00	126
302	0.33	0.01	0.03	73
303	0.53	0.40	0.46	176
304	0.96	0.44	0.61	230
305	0.94	0.40	0.57	156
306	0.43	0.38	0.41	146
307	0.28	0.11	0.16	98
308	0.08	0.04	0.05	78
309	0.33	0.02	0.04	94
310	0.57	0.28	0.38	162
311	0.68	0.38	0.49	116
312	0.48	0.26	0.34	57
313	0.67	0.03	0.06	65
314	0.47	0.31	0.37	138
315	0.48	0.24	0.32	195
316	0.41	0.33	0.37	69
317	0.17	0.07	0.10	134
318	0.41	0.30	0.34	148
319	0.70	0.29	0.41	161
320	0.17	0.21	0.19	104
321	0.81	0.42	0.56	156
322	0.55	0.32	0.41	134
323	0.50	0.41	0.45	232
324	0.35	0.21	0.26	92
325	0.33	0.28	0.30	197
326	0.06	0.02	0.03	126
327	0.28	0.04	0.08	115
328	0.97	0.31	0.47	198
329	0.53	0.32	0.40	125
330	0.54	0.09	0.15	81
331	0.19	0.04	0.07	94
332	0.33	0.02	0.03	56
333	0.13	0.08	0.10	260
334	0.50	0.03	0.06	60
335	0.25	0.12	0.16	110
336	0.65	0.42	0.51	71
337	0.12	0.06	0.08	66
338	0.46	0.34	0.39	150
339	0.00	0.00	0.00	54
340	0.88	0.33	0.48	195
341	0.75	0.19	0.30	79
342	0.36	0.32	0.34	38
343	0.57	0.30	0.39	43
344	0.50	0.21	0.29	68
345	0.60	0.38	0.47	73
346	0.08	0.03	0.04	116
347	0.92	0.22	0.35	111
348	0.22	0.08	0.12	63
349	0.91	0.39	0.55	104
350	0.52	0.30	0.38	44
351	0.50	0.15	0.23	40
352	1.00	0.18	0.31	136
353	0.50	0.30	0.37	54
354	0.26	0.04	0.08	134
355	0.47	0.23	0.31	120
356	0.44	0.22	0.30	228
357	0.54	0.22	0.31	269
358	0.68	0.31	0.43	80
359	0.66	0.26	0.38	140
360	0.35	0.19	0.25	125
361	0.89	0.33	0.48	169
362	0.12	0.05	0.07	56
363	0.95	0.47	0.63	154
364	0.30	0.05	0.09	58
365	0.22	0.21	0.22	71
366	1.00	0.37	0.54	54
367	0.20	0.05	0.08	116
368	0.20	0.02	0.03	54
369	0.12	0.03	0.05	71
370	0.11	0.03	0.05	61
371	0.33	0.06	0.10	71
372	0.62	0.35	0.44	52
373	0.60	0.17	0.26	150
374	0.00	0.00	0.00	00

374	0.39	0.23	0.29	93
375	0.38	0.07	0.12	67
376	0.00	0.00	0.00	76
377	0.66	0.18	0.28	106
378	0.17	0.01	0.02	86
379	0.25	0.07	0.11	14
380	0.94	0.14	0.24	122
381	0.11	0.05	0.07	104
382	0.21	0.09	0.13	66
383	0.51	0.26	0.35	110
384	0.20	0.01	0.02	155
385	0.22	0.04	0.07	50
386	0.21	0.14	0.17	64
387	0.19	0.03	0.06	93
388	0.54	0.20	0.29	102
389	0.09	0.02	0.03	108
390	0.95	0.33	0.49	178
391	0.56	0.16	0.24	115
392	0.47	0.21	0.30	42
393	0.00	0.00	0.00	134
394	0.06	0.01	0.02	112
395	0.41	0.21	0.28	176
396	0.19	0.02	0.04	125
397	0.69	0.21	0.32	224
398	0.86	0.30	0.45	63
399	0.20	0.02	0.03	59
400	0.41	0.29	0.34	63
401	0.26	0.16	0.20	98
402	0.35	0.07	0.12	162
403	0.39	0.19	0.26	83
404	0.76	0.68	0.72	19
405	0.20	0.12	0.15	92
406	0.75	0.22	0.34	41
407	0.72	0.30	0.43	43
408	0.64	0.18	0.28	160
409	0.28	0.20	0.23	50
410	0.00	0.00	0.00	19
411	0.28	0.14	0.18	175
412	0.31	0.06	0.09	72
413	0.40	0.04	0.08	95
414	0.18	0.10	0.13	97
415	0.21	0.12	0.16	48
416	0.43	0.29	0.35	83
417	0.14	0.03	0.04	40
418	0.26	0.11	0.16	91
419	0.43	0.28	0.34	90
420	0.15	0.08	0.11	37
421	0.10	0.06	0.08	66
422	0.55	0.38	0.45	73
423	0.42	0.20	0.27	56
424	0.95	0.58	0.72	33
425	0.05	0.01	0.02	76
426	0.19	0.06	0.09	81
427	1.00	0.32	0.48	150
428	0.94	0.55	0.70	29
429	1.00	0.07	0.13	389
430	0.63	0.20	0.31	167
431	0.30	0.06	0.10	123
432	0.39	0.28	0.33	39
433	0.43	0.32	0.36	82
434	1.00	0.42	0.60	66
435	0.60	0.39	0.47	93
436	0.55	0.21	0.30	87
437	0.24	0.05	0.08	86
438	0.77	0.35	0.48	104
439	0.52	0.11	0.18	100
440	0.36	0.04	0.06	141
441	0.36	0.32	0.34	110
442	0.26	0.16	0.20	123
443	0.33	0.01	0.03	71
444	0.21	0.03	0.05	109
445	0.21	0.12	0.16	48
446	0.33	0.20	0.25	76
447	0.17	0.13	0.15	38
448	0.68	0.48	0.57	81
449	0.51	0.18	0.27	132
450	0.48	0.28	0.36	81

451	0.80	0.16	0.26	76
452	0.00	0.00	0.00	44
453	0.09	0.02	0.04	44
454	0.73	0.31	0.44	70
455	0.24	0.10	0.14	155
456	0.33	0.21	0.26	43
457	0.38	0.21	0.27	72
458	0.17	0.06	0.09	62
459	0.57	0.12	0.19	69
460	0.04	0.03	0.03	119
461	0.74	0.25	0.38	79
462	0.33	0.11	0.16	47
463	0.24	0.09	0.13	104
464	0.57	0.30	0.40	106
465	0.53	0.12	0.20	64
466	0.58	0.25	0.35	173
467	0.62	0.22	0.33	107
468	0.48	0.08	0.14	126
469	0.00	0.00	0.00	114
470	0.95	0.51	0.66	140
471	0.62	0.06	0.11	79
472	0.29	0.22	0.25	143
473	0.49	0.16	0.24	158
474	0.28	0.05	0.09	138
475	0.07	0.03	0.04	59
476	0.62	0.28	0.39	88
477	0.85	0.42	0.56	176
478	0.93	0.54	0.68	24
479	0.17	0.04	0.07	92
480	0.86	0.31	0.46	100
481	0.37	0.21	0.27	103
482	0.29	0.27	0.28	74
483	0.82	0.30	0.44	105
484	0.06	0.02	0.03	83
485	0.11	0.02	0.04	82
486	0.38	0.15	0.22	71
487	0.37	0.21	0.27	120
488	0.25	0.02	0.04	105
489	0.59	0.20	0.29	87
490	0.95	0.56	0.71	32
491	0.10	0.01	0.03	69
492	0.25	0.02	0.04	49
493	0.07	0.01	0.02	117
494	0.43	0.05	0.09	61
495	1.00	0.08	0.15	344
496	0.31	0.15	0.21	52
497	0.57	0.12	0.19	137
498	0.42	0.05	0.09	98
499	0.45	0.06	0.11	79

micro avg	0.65	0.32	0.43	173812
macro avg	0.50	0.24	0.31	173812
weighted avg	0.64	0.32	0.41	173812
samples avg	0.39	0.30	0.31	173812

Time taken to run this cell : 0:06:49.603184

OneVsRestClassifier with Logistic regression(penalty=l1)¶

In [76]:

```
start = datetime.now()
import warnings
warnings.filterwarnings('ignore')

# hp1={'estimator__C':alpha}

cv_scores = []
for i in tqdm(alpha):
    print(i)
    hp1={'estimator__alpha':[i],
        'estimator__loss':['log'],
        'estimator__penalty':['l1']}
    print(hp1)
    classifier = OneVsRestClassifier(SGDClassifier())
```

```

classifier = OneVsRestClassifier(svm_classifier)

model11 = GridSearchCV(classifier, hp1,
                        cv=3, scoring='f1_micro', n_jobs=-1)
print("Gridsearchcv")
best_model11 = model11.fit(x_train_multilabel, y_train)
print('fit model')
Train_model_score = best_model11.score(x_train_multilabel,
                                       y_train)

#print("best_model11")
cv_scores.append(Train_model_score.mean())

fscore = [x for x in cv_scores]

# determining best alpha
optimal_alpha22 = alpha[fscore.index(max(fscore))]
print('\n The optimal value of alpha with penalty=l1 and loss= log is %d.' % optimal_alpha22)

# Plots
fig4 = plt.figure( facecolor='c', edgecolor='k')
plt.plot(alpha, fscore, color='green', marker='o', linestyle='dashed',
         linewidth=2, markersize=12)

for xy in zip(alpha, np.round(fscore,3)):
    plt.annotate('(%s, %s)' % xy, xy=xy, textcoords='data')

plt.xlabel('Hyper parameter Alpha')
plt.ylabel('F1_Score value ')
plt.show()

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

```

```

0%|
[00:00<?, ?it/s]

```

```

0.001
{'estimator__alpha': [0.001], 'estimator__loss': ['log'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

```

```

33%|
0:08, 3304.41s/it]

```

```

0.01
{'estimator__alpha': [0.01], 'estimator__loss': ['log'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

```

```

67%|
[1:35:24<50:38, 3038.99s/it]

```

```

0.1
{'estimator__alpha': [0.1], 'estimator__loss': ['log'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

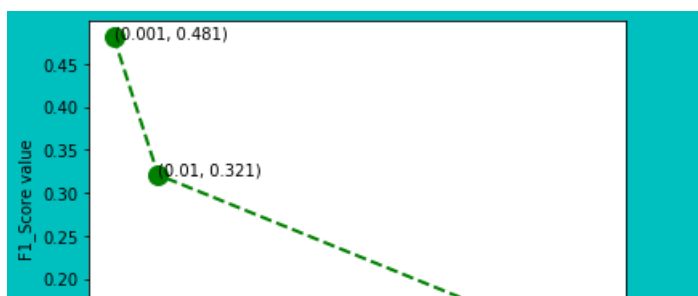
```

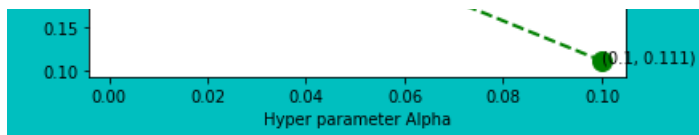
```

100%|
[2:21:10<00:00, 2823.36s/it]

```

The optimal value of alpha with penalty=l1 and loss= log is 0.





Time taken to run this cell : 2:21:11.178399

In [77]:

```
start = datetime.now()
best_model2 = OneVsRestClassifier(SGDClassifier(loss='log', alpha=optimal_alpha22,
                                              penalty='l1'), n_jobs=-1)
best_model2.fit(x_train_multilabel, y_train)
```

Out[77]:

```
OneVsRestClassifier(estimator=SGDClassifier(alpha=0.001, average=False,
                                           class_weight=None,
                                           early_stopping=False, epsilon=0.1,
                                           eta0=0.0, fit_intercept=True,
                                           l1_ratio=0.15,
                                           learning_rate='optimal', loss='log',
                                           max_iter=1000, n_iter_no_change=5,
                                           n_jobs=None, penalty='l1',
                                           power_t=0.5, random_state=None,
                                           shuffle=True, tol=0.001,
                                           validation_fraction=0.1, verbose=0,
                                           warm_start=False),
                   n_jobs=-1)
```

In [78]:

```
joblib.dump(best_model2, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/best_model2_LR.pkl')
```

Out[78]:

```
['E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/best_model2_LR.pkl']
```

In [79]:

```
best_model2=joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/best_model2_LR.pkl')
```

Logistic regression with l1 penalty

In [80]:

```
start = datetime.now()
#classifier = OneVsRestClassifier(LogisticRegression(penalty='l1'), n_jobs=-1)
#classifier.fit(x_train_multilabel, y_train)
predictions = best_model2.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)
```

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

Accuracy : 0.19363

Hamming loss 0.00313556

Micro-average quality numbers

Precision: 0.5955, Recall: 0.3055, F1-measure: 0.4039

Macro-average quality numbers

Precision: 0.4181, Recall: 0.2335, F1-measure: 0.2828

	precision	recall	f1-score	support
0	0.85	0.62	0.71	5519
1	0.52	0.19	0.28	8190
2	0.80	0.30	0.44	6529
3	0.74	0.38	0.51	3231
4	0.80	0.36	0.49	6430
5	0.56	0.34	0.42	2879
6	0.83	0.49	0.62	5086
7	0.87	0.53	0.65	4533
8	0.53	0.14	0.22	3000
9	0.60	0.55	0.58	2765
10	0.47	0.17	0.25	3051
11	0.75	0.29	0.42	3009
12	0.67	0.22	0.33	2630
13	0.61	0.13	0.22	1426
14	0.79	0.61	0.69	2548
15	0.44	0.14	0.21	2371
16	0.58	0.25	0.35	873
17	0.75	0.66	0.70	2151
18	0.68	0.21	0.32	2204
19	0.63	0.39	0.49	831
20	0.75	0.43	0.54	1860
21	0.25	0.09	0.14	2023
22	0.37	0.23	0.29	1513
23	0.77	0.58	0.66	1207
24	0.46	0.35	0.40	506
25	0.66	0.32	0.43	425
26	0.57	0.38	0.46	793
27	0.55	0.31	0.40	1291
28	0.65	0.37	0.47	1208
29	0.26	0.12	0.16	406
30	0.59	0.21	0.31	504
31	0.26	0.16	0.20	732
32	0.54	0.31	0.39	441
33	0.31	0.13	0.18	1645
34	0.73	0.22	0.34	1058
35	0.80	0.58	0.67	946
36	0.64	0.18	0.28	644
37	0.83	0.82	0.83	136
38	0.56	0.36	0.44	570
39	0.80	0.33	0.47	766
40	0.48	0.29	0.36	1132
41	0.38	0.18	0.25	174
42	0.67	0.45	0.54	210
43	0.51	0.41	0.45	433
44	0.63	0.46	0.53	626
45	0.61	0.25	0.35	852
46	0.54	0.41	0.46	534
47	0.20	0.15	0.17	350
48	0.70	0.46	0.55	496
49	0.80	0.56	0.66	785
50	0.15	0.17	0.16	475
51	0.28	0.12	0.17	305
52	0.27	0.06	0.10	251
53	0.52	0.52	0.52	914
54	0.39	0.20	0.26	728
55	0.13	0.02	0.03	258
56	0.36	0.16	0.22	821
57	0.44	0.13	0.20	541
58	0.70	0.33	0.45	748
59	0.85	0.73	0.79	724
60	0.28	0.07	0.11	660
61	0.89	0.17	0.29	235
62	0.91	0.67	0.77	718
63	0.77	0.65	0.70	468
64	0.34	0.50	0.41	191
65	0.22	0.10	0.13	429
66	0.15	0.08	0.11	415

67	0.65	0.63	0.64	274
68	0.84	0.49	0.62	510
69	0.62	0.45	0.52	466
70	0.24	0.09	0.13	305
71	0.27	0.25	0.26	247
72	0.69	0.49	0.57	401
73	0.93	0.80	0.86	86
74	0.67	0.36	0.47	120
75	0.87	0.65	0.74	129
76	0.00	0.00	0.00	473
77	0.22	0.34	0.27	143
78	0.78	0.43	0.55	347
79	0.73	0.21	0.32	479
80	0.40	0.34	0.37	279
81	0.60	0.17	0.27	461
82	0.10	0.01	0.02	298
83	0.71	0.38	0.49	396
84	0.38	0.36	0.37	184
85	0.38	0.18	0.25	573
86	0.26	0.03	0.06	325
87	0.50	0.23	0.32	273
88	0.31	0.25	0.28	135
89	0.18	0.08	0.11	232
90	0.46	0.32	0.38	409
91	0.41	0.34	0.37	420
92	0.71	0.42	0.53	408
93	0.51	0.56	0.53	241
94	0.21	0.13	0.16	211
95	0.23	0.12	0.15	277
96	0.23	0.02	0.04	410
97	0.94	0.12	0.21	501
98	0.11	0.60	0.19	136
99	0.51	0.29	0.37	239
100	0.46	0.06	0.11	324
101	0.86	0.52	0.65	277
102	0.92	0.64	0.75	613
103	0.48	0.22	0.30	157
104	0.18	0.11	0.13	295
105	0.63	0.34	0.44	334
106	0.41	0.03	0.05	335
107	0.57	0.53	0.55	389
108	0.26	0.33	0.29	251
109	0.46	0.34	0.39	317
110	0.25	0.07	0.11	187
111	0.27	0.14	0.19	140
112	0.28	0.14	0.19	154
113	0.59	0.22	0.32	332
114	0.38	0.24	0.30	323
115	0.38	0.10	0.16	344
116	0.69	0.42	0.52	370
117	0.48	0.22	0.31	313
118	0.77	0.71	0.74	874
119	0.39	0.19	0.26	293
120	0.04	0.01	0.02	200
121	0.67	0.55	0.61	463
122	0.26	0.24	0.25	119
123	0.00	0.00	0.00	256
124	0.87	0.75	0.81	195
125	0.33	0.23	0.27	138
126	0.70	0.38	0.50	376
127	0.16	0.06	0.08	122
128	0.17	0.07	0.10	252
129	0.12	0.01	0.02	144
130	0.10	0.02	0.03	150
131	0.14	0.01	0.02	210
132	0.47	0.07	0.12	361
133	0.92	0.46	0.61	453
134	0.76	0.86	0.81	124
135	0.00	0.00	0.00	91
136	0.42	0.22	0.29	128
137	0.40	0.27	0.32	218
138	0.00	0.00	0.00	243
139	0.33	0.20	0.25	149
140	0.70	0.45	0.55	318
141	0.13	0.12	0.13	159
142	0.68	0.35	0.46	274
143	0.76	0.83	0.79	362

143	0.70	0.00	0.70	302
144	0.48	0.20	0.29	118
145	0.40	0.42	0.41	164
146	0.58	0.25	0.35	461
147	0.62	0.41	0.49	159
148	0.32	0.13	0.18	166
149	0.94	0.55	0.70	346
150	0.48	0.03	0.06	350
151	0.87	0.47	0.61	55
152	0.66	0.51	0.57	387
153	0.30	0.21	0.25	150
154	0.18	0.10	0.13	281
155	0.28	0.14	0.18	202
156	0.70	0.61	0.65	130
157	0.36	0.13	0.19	245
158	0.60	0.67	0.63	177
159	0.51	0.36	0.42	130
160	0.43	0.14	0.21	336
161	0.84	0.57	0.68	220
162	0.11	0.05	0.07	229
163	0.84	0.40	0.54	316
164	0.66	0.18	0.28	283
165	0.54	0.26	0.35	197
166	0.14	0.11	0.12	101
167	0.25	0.29	0.27	231
168	0.31	0.11	0.16	370
169	0.40	0.27	0.32	258
170	0.14	0.09	0.11	101
171	0.38	0.15	0.21	89
172	0.33	0.32	0.32	193
173	0.41	0.30	0.35	309
174	0.42	0.10	0.17	172
175	0.93	0.74	0.82	95
176	0.91	0.54	0.68	346
177	0.95	0.35	0.51	322
178	0.55	0.41	0.47	232
179	0.50	0.03	0.06	125
180	0.41	0.32	0.36	145
181	0.36	0.13	0.19	77
182	0.09	0.04	0.05	182
183	0.37	0.44	0.40	257
184	0.21	0.02	0.04	216
185	0.27	0.08	0.12	242
186	0.28	0.15	0.19	165
187	0.74	0.53	0.62	263
188	0.25	0.12	0.16	174
189	0.65	0.23	0.34	136
190	0.66	0.53	0.59	202
191	0.28	0.11	0.16	134
192	0.73	0.37	0.49	230
193	0.29	0.17	0.21	90
194	0.54	0.37	0.44	185
195	0.06	0.08	0.06	156
196	0.00	0.00	0.00	160
197	0.00	0.00	0.00	266
198	0.48	0.05	0.09	284
199	0.17	0.03	0.05	145
200	0.88	0.78	0.83	212
201	0.47	0.09	0.15	317
202	0.59	0.46	0.52	427
203	0.21	0.09	0.13	232
204	0.28	0.14	0.19	217
205	0.45	0.33	0.38	527
206	0.04	0.01	0.01	124
207	0.13	0.07	0.09	103
208	0.74	0.62	0.68	287
209	0.20	0.08	0.12	193
210	0.52	0.29	0.37	220
211	0.23	0.02	0.04	140
212	0.08	0.07	0.08	161
213	0.30	0.18	0.23	72
214	0.62	0.43	0.51	396
215	0.79	0.37	0.51	134
216	0.50	0.01	0.02	400
217	0.49	0.25	0.33	75
218	0.94	0.73	0.82	219
219	0.79	0.33	0.47	210
220	0.86	0.42	0.56	298

220	0.00	0.72	0.00	230
221	0.97	0.56	0.71	266
222	0.71	0.44	0.54	290
223	0.29	0.04	0.07	128
224	0.79	0.31	0.45	159
225	0.45	0.25	0.32	164
226	0.49	0.48	0.49	144
227	0.34	0.46	0.39	276
228	0.06	0.01	0.02	235
229	0.00	0.00	0.00	216
230	0.32	0.19	0.24	228
231	0.52	0.73	0.61	64
232	0.11	0.05	0.07	103
233	0.74	0.25	0.37	216
234	0.00	0.00	0.00	116
235	0.45	0.52	0.48	77
236	0.94	0.67	0.78	67
237	0.03	0.01	0.01	218
238	0.05	0.01	0.02	139
239	0.11	0.05	0.07	94
240	0.35	0.10	0.16	77
241	0.33	0.02	0.03	167
242	0.85	0.26	0.39	86
243	0.47	0.12	0.19	58
244	0.21	0.07	0.11	269
245	0.18	0.09	0.12	112
246	0.79	0.75	0.77	255
247	0.42	0.28	0.33	58
248	0.36	0.05	0.09	81
249	0.00	0.00	0.00	131
250	0.28	0.24	0.26	93
251	0.52	0.21	0.30	154
252	0.09	0.02	0.03	129
253	0.46	0.33	0.38	83
254	0.24	0.08	0.12	191
255	0.15	0.04	0.06	219
256	0.07	0.02	0.03	130
257	0.40	0.33	0.36	93
258	0.67	0.33	0.44	217
259	0.25	0.07	0.11	141
260	0.94	0.10	0.19	143
261	0.44	0.08	0.13	219
262	0.41	0.29	0.34	107
263	0.32	0.28	0.30	236
264	0.15	0.12	0.13	119
265	0.18	0.26	0.21	72
266	0.14	0.13	0.13	70
267	0.30	0.10	0.15	107
268	0.66	0.35	0.46	169
269	0.20	0.12	0.15	129
270	0.72	0.55	0.62	159
271	0.40	0.25	0.31	190
272	0.43	0.11	0.18	248
273	0.91	0.61	0.73	264
274	0.82	0.58	0.68	105
275	0.00	0.00	0.00	104
276	0.05	0.01	0.01	115
277	0.84	0.52	0.64	170
278	0.48	0.14	0.21	145
279	0.93	0.40	0.56	230
280	0.54	0.40	0.46	80
281	0.61	0.66	0.64	217
282	0.76	0.44	0.56	175
283	0.42	0.04	0.07	269
284	0.60	0.28	0.39	74
285	0.86	0.41	0.56	206
286	0.90	0.52	0.66	227
287	0.83	0.23	0.36	130
288	0.26	0.07	0.11	129
289	0.14	0.01	0.02	80
290	0.15	0.12	0.14	99
291	0.79	0.24	0.36	208
292	0.33	0.12	0.18	67
293	0.60	0.26	0.36	109
294	0.26	0.23	0.24	140
295	0.16	0.15	0.15	241
296	0.14	0.12	0.13	72
297	0.31	0.11	0.16	107

297	0.31	0.11	0.10	107
298	0.32	0.16	0.22	61
299	0.74	0.30	0.43	77
300	0.06	0.08	0.07	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.50	0.43	0.46	176
304	0.97	0.57	0.71	230
305	0.96	0.54	0.69	156
306	0.34	0.33	0.33	146
307	0.14	0.05	0.08	98
308	0.33	0.03	0.05	78
309	0.40	0.02	0.04	94
310	0.68	0.26	0.38	162
311	0.72	0.41	0.52	116
312	0.48	0.28	0.36	57
313	0.00	0.00	0.00	65
314	0.39	0.29	0.33	138
315	0.44	0.22	0.29	195
316	0.33	0.39	0.36	69
317	0.00	0.00	0.00	134
318	0.36	0.30	0.33	148
319	0.84	0.32	0.46	161
320	0.17	0.19	0.18	104
321	0.74	0.42	0.53	156
322	0.46	0.25	0.32	134
323	0.28	0.26	0.27	232
324	0.13	0.17	0.15	92
325	0.32	0.08	0.13	197
326	0.07	0.03	0.04	126
327	0.17	0.01	0.02	115
328	0.96	0.69	0.80	198
329	0.52	0.26	0.35	125
330	0.33	0.01	0.02	81
331	0.19	0.03	0.05	94
332	0.00	0.00	0.00	56
333	0.04	0.01	0.01	260
334	0.00	0.00	0.00	60
335	0.21	0.13	0.16	110
336	0.49	0.45	0.47	71
337	0.14	0.14	0.14	66
338	0.46	0.32	0.38	150
339	0.00	0.00	0.00	54
340	0.85	0.51	0.64	195
341	0.73	0.10	0.18	79
342	0.27	0.32	0.29	38
343	0.16	0.44	0.23	43
344	0.00	0.00	0.00	68
345	0.43	0.40	0.41	73
346	0.19	0.09	0.13	116
347	0.80	0.47	0.59	111
348	0.10	0.03	0.05	63
349	0.88	0.47	0.61	104
350	0.70	0.32	0.44	44
351	0.00	0.00	0.00	40
352	1.00	0.22	0.36	136
353	0.39	0.28	0.33	54
354	0.00	0.00	0.00	134
355	0.25	0.09	0.13	120
356	0.30	0.07	0.12	228
357	0.75	0.06	0.10	269
358	0.57	0.35	0.43	80
359	0.71	0.29	0.41	140
360	0.21	0.05	0.08	125
361	0.92	0.46	0.61	169
362	0.09	0.05	0.07	56
363	0.83	0.75	0.79	154
364	0.00	0.00	0.00	58
365	0.15	0.13	0.14	71
366	0.88	0.78	0.82	54
367	0.17	0.08	0.11	116
368	0.00	0.00	0.00	54
369	0.00	0.00	0.00	71
370	0.06	0.03	0.04	61
371	0.29	0.08	0.13	71
372	0.62	0.48	0.54	52
373	0.79	0.22	0.34	150
374	0.44	0.10	0.26	02

374	0.44	0.10	0.20	33
375	0.40	0.03	0.06	67
376	0.00	0.00	0.00	76
377	0.42	0.21	0.28	106
378	0.20	0.01	0.02	86
379	0.10	0.07	0.08	14
380	0.97	0.25	0.40	122
381	0.12	0.05	0.07	104
382	0.24	0.15	0.19	66
383	0.46	0.24	0.31	110
384	0.00	0.00	0.00	155
385	0.09	0.02	0.03	50
386	0.22	0.20	0.21	64
387	0.00	0.00	0.00	93
388	0.50	0.21	0.29	102
389	0.00	0.00	0.00	108
390	0.96	0.42	0.59	178
391	0.62	0.16	0.25	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	0.02	0.03	176
396	0.00	0.00	0.00	125
397	0.60	0.12	0.20	224
398	0.83	0.30	0.44	63
399	0.00	0.00	0.00	59
400	0.35	0.29	0.31	63
401	0.12	0.02	0.03	98
402	0.37	0.04	0.08	162
403	0.36	0.14	0.21	83
404	0.81	0.68	0.74	19
405	0.11	0.07	0.08	92
406	0.33	0.17	0.23	41
407	0.57	0.19	0.28	43
408	0.00	0.00	0.00	160
409	0.21	0.16	0.18	50
410	0.00	0.00	0.00	19
411	0.26	0.16	0.20	175
412	0.09	0.01	0.02	72
413	0.00	0.00	0.00	95
414	0.12	0.07	0.09	97
415	0.25	0.10	0.15	48
416	0.36	0.24	0.29	83
417	0.00	0.00	0.00	40
418	0.19	0.07	0.10	91
419	0.39	0.33	0.36	90
420	0.14	0.11	0.12	37
421	0.06	0.05	0.05	66
422	0.59	0.27	0.37	73
423	0.34	0.20	0.25	56
424	0.93	0.82	0.87	33
425	0.09	0.01	0.02	76
426	0.25	0.01	0.02	81
427	1.00	0.53	0.69	150
428	0.80	0.69	0.74	29
429	0.00	0.00	0.00	389
430	0.62	0.20	0.31	167
431	0.00	0.00	0.00	123
432	0.44	0.36	0.39	39
433	0.34	0.22	0.27	82
434	1.00	0.53	0.69	66
435	0.56	0.38	0.45	93
436	0.30	0.03	0.06	87
437	0.40	0.07	0.12	86
438	0.52	0.36	0.42	104
439	0.05	0.01	0.02	100
440	0.33	0.01	0.01	141
441	0.31	0.30	0.30	110
442	0.15	0.09	0.11	123
443	0.00	0.00	0.00	71
444	0.00	0.00	0.00	109
445	0.25	0.15	0.18	48
446	0.33	0.17	0.22	76
447	0.11	0.08	0.09	38
448	0.64	0.42	0.51	81
449	0.30	0.06	0.10	132
450	0.44	0.23	0.31	81
451	0.00	0.11	0.10	76

451	0.80	0.11	0.19	76
452	0.00	0.00	0.00	44
453	0.00	0.00	0.00	44
454	0.71	0.21	0.33	70
455	0.10	0.01	0.01	155
456	0.21	0.14	0.17	43
457	0.38	0.14	0.20	72
458	0.19	0.10	0.13	62
459	0.00	0.00	0.00	69
460	0.25	0.03	0.06	119
461	0.62	0.16	0.26	79
462	0.25	0.04	0.07	47
463	0.12	0.01	0.02	104
464	0.49	0.33	0.40	106
465	0.00	0.00	0.00	64
466	0.55	0.20	0.29	173
467	0.52	0.36	0.43	107
468	0.00	0.00	0.00	126
469	0.00	0.00	0.00	114
470	0.93	0.63	0.75	140
471	0.00	0.00	0.00	79
472	0.31	0.34	0.33	143
473	0.33	0.01	0.01	158
474	0.00	0.00	0.00	138
475	0.14	0.07	0.09	59
476	0.61	0.49	0.54	88
477	0.84	0.45	0.59	176
478	0.88	0.62	0.73	24
479	0.00	0.00	0.00	92
480	0.81	0.38	0.52	100
481	0.00	0.00	0.00	103
482	0.27	0.22	0.24	74
483	0.78	0.43	0.55	105
484	0.15	0.05	0.07	83
485	0.06	0.07	0.07	82
486	0.17	0.07	0.10	71
487	0.36	0.22	0.27	120
488	0.00	0.00	0.00	105
489	0.69	0.25	0.37	87
490	1.00	0.62	0.77	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	0.75	0.05	0.09	61
495	0.00	0.00	0.00	344
496	0.08	0.02	0.03	52
497	0.24	0.18	0.21	137
498	0.33	0.01	0.02	98
499	0.67	0.03	0.05	79
micro avg	0.60	0.31	0.40	173812
macro avg	0.42	0.23	0.28	173812
weighted avg	0.56	0.31	0.38	173812
samples avg	0.37	0.29	0.30	173812

Time taken to run this cell : 0:00:15.414688

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

In [81]:

```
start = datetime.now()
import warnings
warnings.filterwarnings('ignore')

# hp1={'estimator__C':alpha}

cv_scores = []
for i in tqdm(alpha):
    print(i)
    hp1={'estimator__alpha':[i],
          'estimator__loss':['hinge'],
          'estimator__penalty':['l1']}
    print(hp1)
    classifier = OneVsRestClassifier(SGDClassifier())
```



```

model11 =GridSearchCV(classifier,hp1,
                      cv=3, scoring='f1_micro',n_jobs=-1)
print("Gridsearchcv")
best_model11=model11.fit(x_train_multilabel, y_train)
print('fit model')
Train_model_score=best_model11.score(x_train_multilabel,
                                     y_train)

#print("best_model1")
cv_scores.append(Train_model_score.mean())

fscore = [x for x in cv_scores]

# determining best alpha
optimal_alpha23 = alpha[fscore.index(max(fscore))]
print('\n The optimal value of alpha with penalty=l1 and loss= log is %d.' % optimal_alpha23)

# Plots
fig4 = plt.figure( facecolor='c', edgecolor='k')
plt.plot(alpha, fscore,color='green', marker='o', linestyle='dashed',
linewidth=2, markersize=12)

for xy in zip(alpha, np.round(fscore,3)):
    plt.annotate('(%s, %s)' % xy, xy=xy, textcoords='data')

plt.xlabel('Hyper parameter Alpha')
plt.ylabel('F1_Score value ')
plt.show()

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

```

```

0%|
[00:00<?, ?it/s]

```

```

0.001
{'estimator__alpha': [0.001], 'estimator__loss': ['hinge'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

```

```

33%|
8:50, 2665.11s/it]

```

```

0.01
{'estimator__alpha': [0.01], 'estimator__loss': ['hinge'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

```

```

67%|
[1:20:26<41:53, 2513.93s/it]

```

```

0.1
{'estimator__alpha': [0.1], 'estimator__loss': ['hinge'], 'estimator__penalty': ['l1']}
Gridsearchcv
fit model

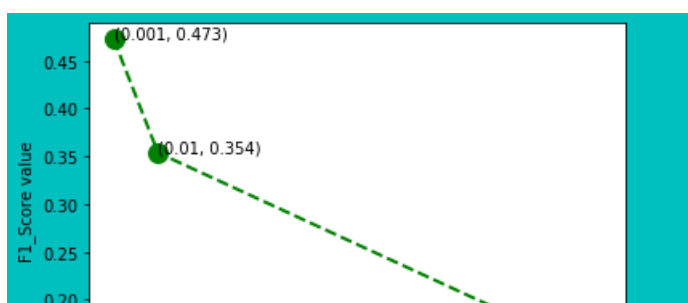
```

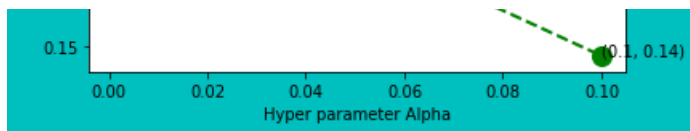
```

100%|
[2:23:44<00:00, 2874.80s/it]

```

The optimal value of alpha with penalty=l1 and loss= log is 0.





Time taken to run this cell : 2:23:44.664128

OneVsRestClassifier with SGDClassifier for optimal alpha with hinge loss

In [82]:

```
start = datetime.now()
classifier2 = OneVsRestClassifier(SGDClassifier(loss='hinge',
                                              alpha=optimal_alpha23,
                                              penalty='l1'))
classifier2=classifier2.fit(x_train_multilabel, y_train)
```

In [83]:

```
joblib.dump(classifier2, 'E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/classifier2.pkl')
```

Out[83]:

```
['E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag Predictor/classifier2.pkl']
```

In [84]:

```
classifier2=joblib.load('E:/BOOKS NEW/Cases datasets/5. Stack Overflow Tag
Predictor/classifier2.pkl')
```

In [85]:

```
predictions = classifier2.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-averasge 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)) #printing classification report for all
500 labels
print("Time taken to run this cell :", datetime.now() - start)
```

```
Accuracy : 0.19161
Hamming loss  0.00313352
Micro-averasge quality numbers
Precision: 0.6005, Recall: 0.2946, F1-measure: 0.3952
Macro-average quality numbers
Precision: 0.3264, Recall: 0.2217, F1-measure: 0.2469
```

	precision	recall	f1-score	support
0	0.90	0.62	0.73	5519
1	0.54	0.19	0.28	8190
2	0.81	0.29	0.42	6529
3	0.60	0.44	0.51	2021

	precision	recall	f1-score	support
0	0.90	0.62	0.73	5519
1	0.54	0.19	0.28	8190
2	0.81	0.29	0.42	6529
3	0.60	0.44	0.51	2021

3	0.68	0.44	0.54	3231
4	0.81	0.31	0.45	6430
5	0.59	0.30	0.40	2879
6	0.76	0.57	0.65	5086
7	0.71	0.58	0.64	4533
8	0.47	0.15	0.22	3000
9	0.71	0.52	0.60	2765
10	0.18	0.01	0.02	3051
11	0.69	0.32	0.44	3009
12	0.65	0.25	0.36	2630
13	0.33	0.12	0.18	1426
14	0.78	0.63	0.70	2548
15	0.76	0.11	0.19	2371
16	0.52	0.35	0.42	873
17	0.75	0.68	0.72	2151
18	0.64	0.19	0.29	2204
19	0.50	0.46	0.48	831
20	0.72	0.48	0.58	1860
21	0.00	0.00	0.00	2023
22	0.00	0.00	0.00	1513
23	0.85	0.54	0.66	1207
24	0.00	0.00	0.00	506
25	0.70	0.34	0.46	425
26	0.56	0.38	0.45	793
27	0.59	0.31	0.40	1291
28	0.60	0.37	0.46	1208
29	0.29	0.16	0.20	406
30	0.69	0.26	0.38	504
31	0.22	0.07	0.11	732
32	0.68	0.09	0.16	441
33	0.28	0.22	0.25	1645
34	0.68	0.23	0.34	1058
35	0.61	0.56	0.58	946
36	0.45	0.31	0.36	644
37	0.83	0.82	0.83	136
38	0.55	0.39	0.46	570
39	0.76	0.32	0.45	766
40	0.47	0.23	0.31	1132
41	0.37	0.25	0.30	174
42	0.58	0.50	0.54	210
43	0.59	0.52	0.56	433
44	0.61	0.45	0.52	626
45	0.54	0.28	0.37	852
46	0.63	0.40	0.49	534
47	0.00	0.00	0.00	350
48	0.50	0.55	0.52	496
49	0.71	0.61	0.66	785
50	0.14	0.17	0.15	475
51	0.09	0.16	0.12	305
52	0.00	0.00	0.00	251
53	0.67	0.32	0.43	914
54	0.00	0.00	0.00	728
55	0.00	0.00	0.00	258
56	0.29	0.08	0.12	821
57	0.00	0.00	0.00	541
58	0.73	0.29	0.42	748
59	0.82	0.75	0.79	724
60	0.20	0.00	0.01	660
61	0.72	0.26	0.38	235
62	0.83	0.80	0.82	718
63	0.76	0.60	0.67	468
64	0.33	0.50	0.39	191
65	0.00	0.00	0.00	429
66	0.05	0.00	0.00	415
67	0.64	0.55	0.59	274
68	0.75	0.60	0.66	510
69	0.47	0.52	0.49	466
70	0.11	0.19	0.14	305
71	0.00	0.00	0.00	247
72	0.66	0.55	0.60	401
73	0.84	0.84	0.84	86
74	0.61	0.45	0.52	120
75	0.72	0.74	0.73	129
76	0.00	0.00	0.00	473
77	0.30	0.43	0.35	143
78	0.71	0.52	0.60	347
79	0.53	0.28	0.37	479
80	0.41	0.10	0.00	070

80	0.41	0.19	0.26	279
81	0.71	0.18	0.29	461
82	0.00	0.00	0.00	298
83	0.70	0.44	0.54	396
84	0.34	0.33	0.33	184
85	0.28	0.16	0.21	573
86	0.60	0.02	0.04	325
87	0.39	0.18	0.25	273
88	0.17	0.32	0.23	135
89	0.00	0.00	0.00	232
90	0.00	0.00	0.00	409
91	0.60	0.28	0.38	420
92	0.67	0.67	0.67	408
93	0.46	0.56	0.51	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	0.48	0.23	0.31	501
98	0.58	0.66	0.62	136
99	0.47	0.24	0.32	239
100	0.00	0.00	0.00	324
101	0.91	0.66	0.76	277
102	0.68	0.82	0.74	613
103	0.47	0.22	0.30	157
104	0.00	0.00	0.00	295
105	0.45	0.42	0.44	334
106	0.00	0.00	0.00	335
107	0.75	0.41	0.53	389
108	0.13	0.01	0.02	251
109	0.33	0.49	0.40	317
110	0.00	0.00	0.00	187
111	0.38	0.14	0.20	140
112	0.44	0.09	0.15	154
113	0.43	0.32	0.37	332
114	0.00	0.00	0.00	323
115	0.00	0.00	0.00	344
116	0.68	0.35	0.46	370
117	0.66	0.08	0.14	313
118	0.75	0.51	0.61	874
119	0.40	0.01	0.03	293
120	0.00	0.00	0.00	200
121	0.71	0.44	0.55	463
122	0.00	0.00	0.00	119
123	0.00	0.00	0.00	256
124	0.79	0.81	0.80	195
125	0.20	0.18	0.19	138
126	0.71	0.45	0.55	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	0.67	0.01	0.03	150
131	0.00	0.00	0.00	210
132	1.00	0.00	0.01	361
133	0.76	0.65	0.70	453
134	0.68	0.81	0.74	124
135	0.00	0.00	0.00	91
136	0.85	0.09	0.16	128
137	0.30	0.40	0.34	218
138	0.00	0.00	0.00	243
139	0.00	0.00	0.00	149
140	0.68	0.38	0.49	318
141	0.00	0.00	0.00	159
142	0.61	0.39	0.48	274
143	0.74	0.63	0.68	362
144	0.53	0.19	0.29	118
145	0.38	0.50	0.43	164
146	0.41	0.34	0.37	461
147	0.54	0.55	0.55	159
148	0.00	0.00	0.00	166
149	0.92	0.59	0.72	346
150	0.34	0.12	0.18	350
151	0.80	0.58	0.67	55
152	0.68	0.44	0.53	387
153	0.41	0.13	0.20	150
154	0.50	0.08	0.14	281
155	0.00	0.00	0.00	202
156	0.60	0.68	0.64	130
157	0.00	0.00	0.00	215

157	0.22	0.05	0.08	245
158	0.69	0.63	0.66	177
159	0.65	0.25	0.36	130
160	0.24	0.20	0.22	336
161	0.69	0.67	0.68	220
162	0.00	0.00	0.00	229
163	0.77	0.59	0.67	316
164	0.51	0.21	0.30	283
165	0.35	0.08	0.13	197
166	0.00	0.00	0.00	101
167	0.00	0.00	0.00	231
168	0.00	0.00	0.00	370
169	0.39	0.21	0.27	258
170	0.00	0.00	0.00	101
171	0.20	0.29	0.24	89
172	0.49	0.16	0.24	193
173	0.32	0.30	0.31	309
174	0.00	0.00	0.00	172
175	0.69	0.68	0.69	95
176	0.79	0.45	0.57	346
177	0.84	0.43	0.57	322
178	0.51	0.49	0.50	232
179	0.00	0.00	0.00	125
180	0.39	0.39	0.39	145
181	0.16	0.21	0.18	77
182	0.00	0.00	0.00	182
183	0.48	0.30	0.37	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	0.56	0.67	0.61	263
188	0.00	0.00	0.00	174
189	0.47	0.05	0.09	136
190	0.91	0.47	0.62	202
191	0.34	0.15	0.21	134
192	0.68	0.57	0.62	230
193	0.00	0.00	0.00	90
194	0.35	0.46	0.40	185
195	0.00	0.00	0.00	156
196	0.00	0.00	0.00	160
197	0.00	0.00	0.00	266
198	0.00	0.00	0.00	284
199	0.00	0.00	0.00	145
200	0.85	0.77	0.81	212
201	0.00	0.00	0.00	317
202	0.53	0.45	0.49	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	0.31	0.05	0.08	103
208	0.71	0.54	0.61	287
209	0.00	0.00	0.00	193
210	0.00	0.00	0.00	220
211	0.65	0.08	0.14	140
212	0.00	0.00	0.00	161
213	0.56	0.12	0.20	72
214	0.56	0.37	0.45	396
215	0.72	0.34	0.46	134
216	0.00	0.00	0.00	400
217	0.48	0.31	0.37	75
218	0.89	0.82	0.85	219
219	0.42	0.31	0.36	210
220	0.47	0.62	0.54	298
221	0.95	0.60	0.74	266
222	0.70	0.40	0.51	290
223	0.00	0.00	0.00	128
224	0.47	0.47	0.47	159
225	0.74	0.20	0.31	164
226	0.29	0.38	0.33	144
227	0.39	0.38	0.39	276
228	0.00	0.00	0.00	235
229	0.00	0.00	0.00	216
230	0.15	0.30	0.20	228
231	0.56	0.47	0.51	64
232	0.00	0.00	0.00	103
233	0.63	0.33	0.44	216

234	0.00	0.00	0.00	116
235	0.40	0.43	0.41	77
236	0.87	0.70	0.78	67
237	0.08	0.01	0.02	218
238	0.05	0.01	0.01	139
239	0.00	0.00	0.00	94
240	0.33	0.18	0.24	77
241	0.00	0.00	0.00	167
242	0.17	0.31	0.22	86
243	0.00	0.00	0.00	58
244	0.28	0.19	0.23	269
245	0.00	0.00	0.00	112
246	0.49	0.64	0.56	255
247	0.18	0.26	0.21	58
248	0.00	0.00	0.00	81
249	0.10	0.04	0.05	131
250	0.31	0.20	0.25	93
251	0.41	0.21	0.28	154
252	0.00	0.00	0.00	129
253	0.27	0.39	0.32	83
254	0.19	0.16	0.17	191
255	0.05	0.00	0.01	219
256	0.00	0.00	0.00	130
257	0.39	0.30	0.34	93
258	0.58	0.54	0.56	217
259	0.00	0.00	0.00	141
260	0.77	0.24	0.36	143
261	0.00	0.00	0.00	219
262	0.00	0.00	0.00	107
263	0.00	0.00	0.00	236
264	0.29	0.03	0.06	119
265	0.20	0.32	0.25	72
266	0.00	0.00	0.00	70
267	0.16	0.13	0.14	107
268	0.49	0.52	0.50	169
269	0.00	0.00	0.00	129
270	0.60	0.50	0.55	159
271	0.91	0.11	0.20	190
272	0.00	0.00	0.00	248
273	0.85	0.66	0.74	264
274	0.68	0.77	0.72	105
275	0.00	0.00	0.00	104
276	0.00	0.00	0.00	115
277	0.70	0.65	0.67	170
278	0.48	0.21	0.29	145
279	0.81	0.45	0.58	230
280	0.38	0.46	0.42	80
281	0.54	0.56	0.55	217
282	0.69	0.58	0.63	175
283	0.00	0.00	0.00	269
284	0.63	0.45	0.52	74
285	0.62	0.56	0.59	206
286	0.80	0.55	0.65	227
287	0.75	0.33	0.46	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	0.45	0.30	0.36	208
292	0.00	0.00	0.00	67
293	0.33	0.27	0.30	109
294	0.00	0.00	0.00	140
295	0.11	0.26	0.16	241
296	0.02	0.03	0.02	72
297	0.00	0.00	0.00	107
298	0.61	0.23	0.33	61
299	0.95	0.23	0.38	77
300	0.07	0.12	0.09	111
301	0.00	0.00	0.00	126
302	0.00	0.00	0.00	73
303	0.48	0.45	0.46	176
304	0.85	0.54	0.66	230
305	0.84	0.45	0.59	156
306	0.23	0.35	0.28	146
307	0.00	0.00	0.00	98
308	0.00	0.00	0.00	78
309	0.26	0.10	0.14	94
310	0.53	0.17	0.25	162

311	0.67	0.47	0.56	116
312	0.41	0.51	0.46	57
313	0.00	0.00	0.00	65
314	0.27	0.28	0.27	138
315	0.32	0.13	0.19	195
316	0.35	0.39	0.37	69
317	0.20	0.13	0.16	134
318	0.00	0.00	0.00	148
319	0.79	0.38	0.51	161
320	0.25	0.03	0.05	104
321	0.79	0.60	0.68	156
322	0.60	0.19	0.29	134
323	0.57	0.17	0.26	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	0.90	0.53	0.66	198
329	0.31	0.47	0.38	125
330	0.81	0.21	0.33	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	0.35	0.55	0.43	71
337	0.00	0.00	0.00	66
338	0.56	0.13	0.22	150
339	0.00	0.00	0.00	54
340	0.55	0.49	0.52	195
341	0.00	0.00	0.00	79
342	0.00	0.00	0.00	38
343	0.32	0.30	0.31	43
344	0.15	0.22	0.18	68
345	0.28	0.40	0.33	73
346	0.00	0.00	0.00	116
347	0.82	0.37	0.51	111
348	0.09	0.02	0.03	63
349	0.60	0.67	0.63	104
350	0.47	0.43	0.45	44
351	0.00	0.00	0.00	40
352	0.64	0.41	0.50	136
353	0.00	0.00	0.00	54
354	0.00	0.00	0.00	134
355	0.18	0.11	0.14	120
356	0.50	0.02	0.03	228
357	0.57	0.13	0.21	269
358	0.39	0.44	0.41	80
359	0.76	0.39	0.52	140
360	0.00	0.00	0.00	125
361	0.85	0.61	0.71	169
362	0.15	0.12	0.14	56
363	0.74	0.61	0.67	154
364	0.00	0.00	0.00	58
365	0.23	0.23	0.23	71
366	0.44	0.85	0.58	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	0.47	0.52	0.49	52
373	0.78	0.17	0.27	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	0.00	0.00	0.00	14
380	0.91	0.34	0.50	122
381	0.00	0.00	0.00	104
382	0.21	0.09	0.13	66
383	0.33	0.38	0.36	110
384	0.00	0.00	0.00	155
385	0.00	0.00	0.00	50
386	0.00	0.00	0.00	64
387	0.00	0.00	0.00	93

388	0.56	0.20	0.29	102
389	0.00	0.00	0.00	108
390	0.90	0.61	0.73	178
391	0.50	0.21	0.29	115
392	0.77	0.40	0.53	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	0.47	0.47	0.47	224
398	0.70	0.41	0.52	63
399	0.00	0.00	0.00	59
400	0.24	0.49	0.33	63
401	0.00	0.00	0.00	98
402	0.28	0.13	0.18	162
403	0.00	0.00	0.00	83
404	0.61	0.89	0.72	19
405	0.00	0.00	0.00	92
406	0.34	0.46	0.39	41
407	0.47	0.21	0.29	43
408	0.00	0.00	0.00	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	0.25	0.03	0.06	95
414	0.00	0.00	0.00	97
415	0.00	0.00	0.00	48
416	0.29	0.24	0.26	83
417	0.00	0.00	0.00	40
418	0.11	0.10	0.10	91
419	0.25	0.28	0.26	90
420	0.00	0.00	0.00	37
421	0.05	0.02	0.02	66
422	0.56	0.33	0.41	73
423	0.19	0.36	0.25	56
424	0.81	0.79	0.80	33
425	0.00	0.00	0.00	76
426	0.00	0.00	0.00	81
427	0.93	0.75	0.83	150
428	0.73	0.66	0.69	29
429	0.00	0.00	0.00	389
430	0.43	0.31	0.36	167
431	0.00	0.00	0.00	123
432	0.35	0.28	0.31	39
433	0.00	0.00	0.00	82
434	0.95	0.56	0.70	66
435	0.51	0.45	0.48	93
436	0.56	0.23	0.33	87
437	0.00	0.00	0.00	86
438	0.42	0.40	0.41	104
439	0.29	0.13	0.18	100
440	0.00	0.00	0.00	141
441	0.35	0.31	0.33	110
442	0.00	0.00	0.00	123
443	0.00	0.00	0.00	71
444	0.00	0.00	0.00	109
445	0.00	0.00	0.00	48
446	0.48	0.30	0.37	76
447	0.00	0.00	0.00	38
448	0.62	0.56	0.58	81
449	0.26	0.11	0.16	132
450	0.26	0.21	0.23	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	0.76	0.53	0.62	70
455	0.00	0.00	0.00	155
456	0.00	0.00	0.00	43
457	0.00	0.00	0.00	72
458	0.00	0.00	0.00	62
459	0.00	0.00	0.00	69
460	0.00	0.00	0.00	119
461	0.56	0.06	0.11	79
462	0.18	0.04	0.07	47
463	0.21	0.11	0.14	104
464	0.67	0.23	0.34	106

465	0.00	0.00	0.00	64
466	0.35	0.26	0.30	173
467	0.52	0.24	0.33	107
468	0.02	0.03	0.03	126
469	0.00	0.00	0.00	114
470	0.85	0.89	0.87	140
471	0.00	0.00	0.00	79
472	0.35	0.33	0.34	143
473	0.70	0.19	0.30	158
474	0.00	0.00	0.00	138
475	0.00	0.00	0.00	59
476	0.57	0.30	0.39	88
477	0.70	0.68	0.69	176
478	0.90	0.75	0.82	24
479	0.00	0.00	0.00	92
480	0.55	0.61	0.58	100
481	0.13	0.08	0.10	103
482	0.00	0.00	0.00	74
483	0.69	0.53	0.60	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	0.56	0.10	0.16	105
489	0.52	0.31	0.39	87
490	1.00	0.78	0.88	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	1.00	0.02	0.03	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	0.38	0.15	0.22	79
micro avg	0.60	0.29	0.40	173812
macro avg	0.33	0.22	0.25	173812
weighted avg	0.50	0.29	0.35	173812
samples avg	0.37	0.28	0.30	173812

Time taken to run this cell : 0:50:43.137304

Observation¶

In [87]:

```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Sr.No", "MODEL", "FEATURIZATION", "PENALTY", "ALPHA", "LOSS", "MICRO_F1_SCORE"]

x.add_row(["1", 'OneVsRest+SGD Classifier', "Tf-idf", "l1", 0.0001, "log", 0.4486])
x.add_row(["2", 'OneVsRest+SGD(log)=LR', "Bag-of-words", "l2", 0.001, "log", 0.4259])
x.add_row(["3", 'OneVsRest+SGD(log)=LR', "Bag-of-words", "l1", 0.001, "log", 0.4039])
x.add_row(["4", 'OneVsRest+SGD Classifier', "Bag-of-words", "l1", 0.001, "Hinge", 0.3952])

print(x)
```

Sr.No	MODEL	FEATURIZATION	PENALTY	ALPHA	LOSS	MICRO_F1_SCORE
1	OneVsRest+SGD Classifier	Tf-idf	l1	0.0001	log	0.4486
2	OneVsRest+SGD(log)=LR	Bag-of-words	l2	0.001	log	0.4259
3	OneVsRest+SGD(log)=LR	Bag-of-words	l1	0.001	log	0.4039
4	OneVsRest+SGD Classifier	Bag-of-words	l1	0.001	Hinge	0.3952