



In [1]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
3 import pandas as pd
4 import sqlite3
5 import csv
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import numpy as np
9 from wordcloud import WordCloud
10 import re
11 import os
12 from sqlalchemy import create_engine # database connection
13 import datetime as dt
14 from nltk.corpus import stopwords
15 from nltk.tokenize import word_tokenize
16 from nltk.stem.snowball import SnowballStemmer
17 from sklearn.feature_extraction.text import CountVectorizer
18 from sklearn.feature_extraction.text import TfidfVectorizer
19 from sklearn.multiclass import OneVsRestClassifier
20 from sklearn.linear_model import SGDClassifier
21 from sklearn import metrics
22 from sklearn.metrics import f1_score,precision_score,recall_score
23 from sklearn import svm
24 from sklearn.linear_model import LogisticRegression
25 from sklearn.naive_bayes import GaussianNB
26 from datetime import datetime
```

# Stack Overflow: Tag Prediction

## 1. Business Problem

### 1.1 Description

#### Description

Stack Overflow is the largest, most trusted online community for developers to learn, share their programming knowledge, and build their careers.

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

#### Problem Statemtent

Suggest the tags based on the content that was there in the question posted on Stackoverflow.

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

## 1.2 Source / useful links

Data Source : <https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data> (<https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data>)  
Youtube : <https://youtu.be/nNDqbUhtIRg> (<https://youtu.be/nNDqbUhtIRg>)  
Research paper : <https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tagging-1.pdf> (<https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tagging-1.pdf>)  
Research paper : <https://dl.acm.org/citation.cfm?id=2660970&dl=ACM&coll=DL> (<https://dl.acm.org/citation.cfm?id=2660970&dl=ACM&coll=DL>)

## 1.3 Real World / Business Objectives and Constraints

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

# 2. Machine Learning problem

## 2.1 Data

### 2.1.1 Data Overview

Refer: <https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data> (<https://www.kaggle.com/c/facebook-recruiting-iii-keyword-extraction/data>)  
All of the data is in 2 files: Train and Test.

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

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

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

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

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

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

### Data Field Explanation

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

**Id** - Unique identifier for each question

**Title** - The question's title

**Body** - The body of the question

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

### 2.1.2 Example Data point

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

```
#include<
iostream>\n
#include<
stdlib.h>\n\n
using namespace std;\n\n
int main()\n
{\n
    int n,a[n],x,c,u[n],m[n],e[n][4];\n
    cout<<"Enter the number of variables";\n          cin>>n;\n\n
    cout<<"Enter the Lower, and Upper Limits of the variables";\n
    for(int y=1; y<n+1; y++)\n
    {\n
        cin>>m[y];\n
        cin>>u[y];\n
    }\n
    for(x=1; x<n+1; x++)\n
    {\n
        a[x] = (m[x] + u[x])/2;\n
    }\n
    c=(n*4)-4;\n
    for(int a1=1; a1<n+1; a1++)\n
    {\n\n
        e[a1][0] = m[a1];\n
        e[a1][1] = m[a1]+1;\n
        e[a1][2] = u[a1]-1;\n
        e[a1][3] = u[a1];\n
    }\n
    for(int i=1; i<n+1; i++)\n
    {\n
        for(int l=1; l<=i; l++)\n
        {\n
            if(l!=1)\n
            {\n
                cout<<a[l]<<"\\t";\n
            }\n
        }\n
        for(int j=0; j<4; j++)\n
        {\n
            cout<<e[i][j];\n
            for(int k=0; k<n-(i+1); k++)\n
            {\n
                cout<<a[k]<<"\\t";\n
            }\n
            cout<<"\\n";\n
        }\n
    }\n
    }\n\n
    system("PAUSE");\n
    return 0;    \n
}\n
```

\n\n

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

1	50	50\n
2	50	50\n
99	50	50\n
100	50	50\n
50	1	50\n
50	2	50\n
50	99	50\n
50	100	50\n
50	50	1\n
50	50	2\n
50	50	99\n
50	50	100\n

\n\n





```
In [0]: 1 # distribution of number of tags per question
        2 df_no_dup.tag_count.value_counts()
```

```
Out[10]: 3      1206157
        2      1111706
        4       814996
        1       568298
        5       505158
        Name: tag_count, dtype: int64
```

```
In [0]: 1 #Creating a new database with no duplicates
        2 if not os.path.isfile('train_no_dup.db'):
        3     disk_dup = create_engine("sqlite:///train_no_dup.db")
        4     no_dup = pd.DataFrame(df_no_dup, columns=['Title', 'Body', 'Tags'])
        5     no_dup.to_sql('no_dup_train',disk_dup)
```

```
In [0]: 1 #This method seems more appropriate to work with this much data.
        2 #creating the connection with database file.
        3 if os.path.isfile('train_no_dup.db'):
        4     start = datetime.now()
        5     con = sqlite3.connect('train_no_dup.db')
        6     tag_data = pd.read_sql_query("""SELECT Tags FROM no_dup_train""", con)
        7     #Always remember to close the database
        8     con.close()
        9
       10     # Let's now drop unwanted column.
       11     tag_data.drop(tag_data.index[0], inplace=True)
       12     #Printing first 5 columns from our data frame
       13     tag_data.head()
       14     print("Time taken to run this cell :", datetime.now() - start)
       15 else:
       16     print("Please download the train.db file from drive or run the above cells to genarate train.db file")
```

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

## 3.2 Analysis of Tags

### 3.2.1 Total number of unique tags

```
In [0]: 1 # Importing & Initializing the "CountVectorizer" object, which
        2 #is scikit-learn's bag of words tool.
        3
        4 #by default 'split()' will tokenize each tag using space.
        5 vectorizer = CountVectorizer(tokenizer = lambda x: x.split())
        6 # fit_transform() does two functions: First, it fits the model
        7 # and learns the vocabulary; second, it transforms our training data
        8 # into feature vectors. The input to fit_transform should be a list of strings.
        9 tag_dtm = vectorizer.fit_transform(tag_data['Tags'])
```

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

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

```
In [0]: 1 #'get_feature_name()' gives us the vocabulary.
        2 tags = vectorizer.get_feature_names()
        3 #Lets look at the tags we have.
        4 print("Some of the tags we have :", tags[:10])
```

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

### 3.2.3 Number of times a tag appeared

```
In [0]: 1 # https://stackoverflow.com/questions/15115765/how-to-access-sparse-matrix-elements
        2 #Lets now store the document term matrix in a dictionary.
        3 freqs = tag_dtm.sum(axis=0).A1
        4 result = dict(zip(tags, freqs))
```

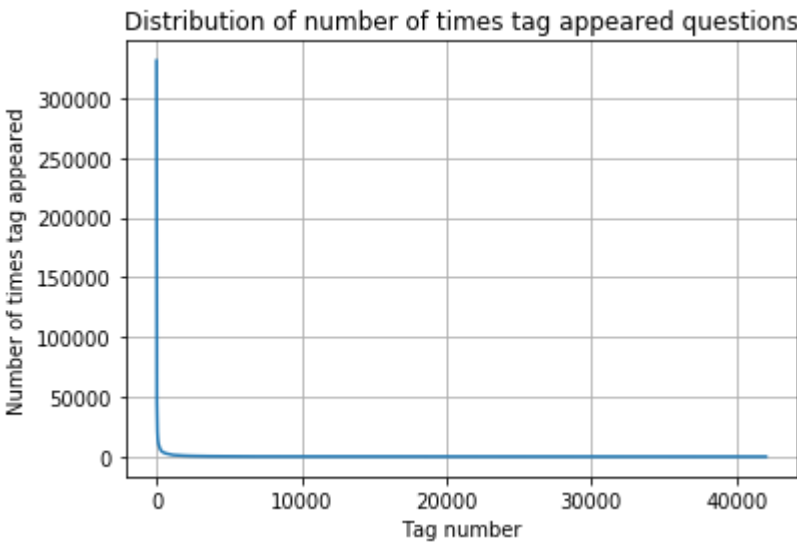
```
In [0]: 1 #Saving this dictionary to csv files.
2 if not os.path.isfile('tag_counts_dict_dtm.csv'):
3     with open('tag_counts_dict_dtm.csv', 'w') as csv_file:
4         writer = csv.writer(csv_file)
5         for key, value in result.items():
6             writer.writerow([key, value])
7 tag_df = pd.read_csv("tag_counts_dict_dtm.csv", names=['Tags', 'Counts'])
8 tag_df.head()
```

Out[17]:

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

```
In [0]: 1 tag_df_sorted = tag_df.sort_values(['Counts'], ascending=False)
2 tag_counts = tag_df_sorted['Counts'].values
```

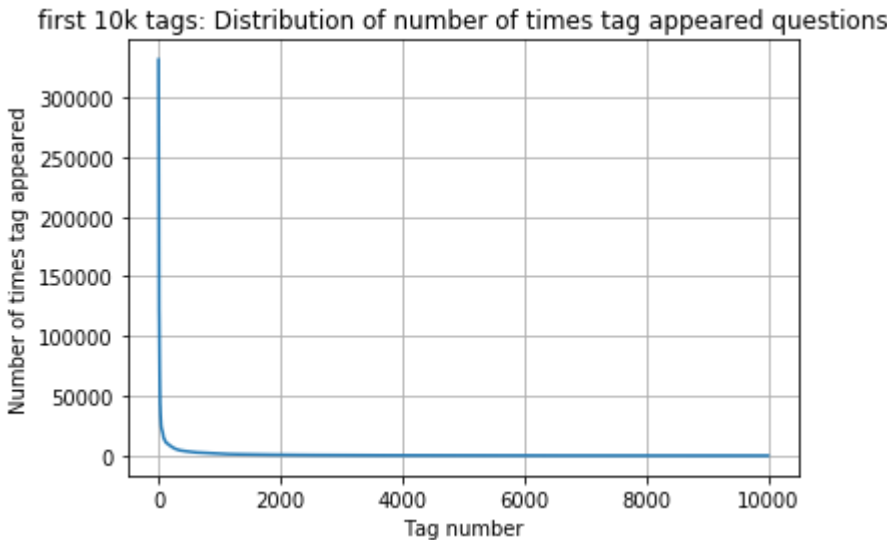
```
In [0]: 1 plt.plot(tag_counts)
2 plt.title("Distribution of number of times tag appeared questions")
3 plt.grid()
4 plt.xlabel("Tag number")
5 plt.ylabel("Number of times tag appeared")
6 plt.show()
```





In [0]:

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

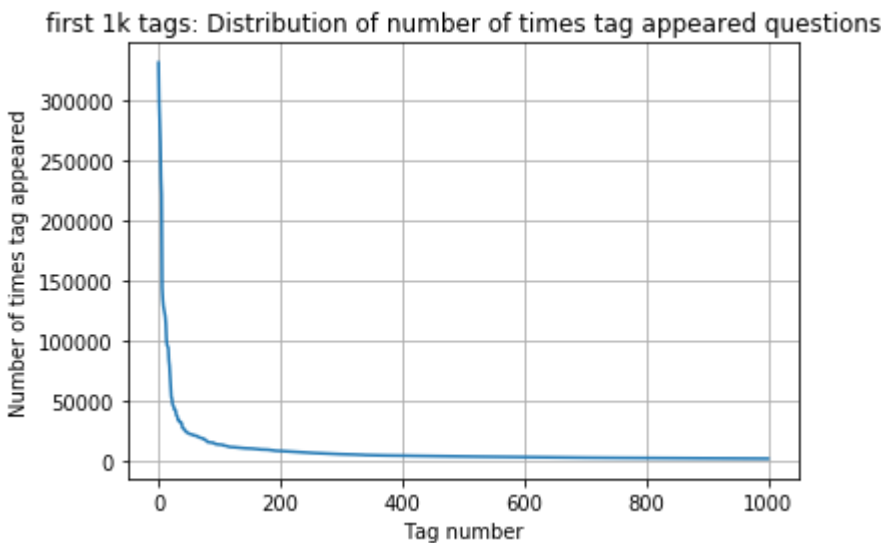


400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7151
6466	5865	5370	4983	4526	4281	4144	3929	3750	3593	
3453	3299	3123	2989	2891	2738	2647	2527	2431	2331	
2259	2186	2097	2020	1959	1900	1828	1770	1723	1673	
1631	1574	1532	1479	1448	1406	1365	1328	1300	1266	
1245	1222	1197	1181	1158	1139	1121	1101	1076	1056	
1038	1023	1006	983	966	952	938	926	911	891	
882	869	856	841	830	816	804	789	779	770	
752	743	733	725	712	702	688	678	671	658	
650	643	634	627	616	607	598	589	583	577	
568	559	552	545	540	533	526	518	512	506	
500	495	490	485	480	477	469	465	457	450	
447	442	437	432	426	422	418	413	408	403	
398	393	388	385	381	378	374	370	367	365	
361	357	354	350	347	344	342	339	336	332	
330	326	323	319	315	312	309	307	304	301	
299	296	293	291	289	286	284	281	278	276	
275	272	270	268	265	262	260	258	256	254	
252	250	249	247	245	243	241	239	238	236	
234	233	232	230	228	226	224	222	220	219	
217	215	214	212	210	209	207	205	204	203	
201	200	199	198	196	194	193	192	191	189	
188	186	185	183	182	181	180	179	178	177	
175	174	172	171	170	169	168	167	166	165	
164	162	161	160	159	158	157	156	156	155	
154	153	152	151	150	149	149	148	147	146	
145	144	143	142	142	141	140	139	138	137	
137	136	135	134	134	133	132	131	130	130	
129	128	128	127	126	126	125	124	124	123	
123	122	122	121	120	120	119	118	118	117	
117	116	116	115	115	114	113	113	112	111	
111	110	109	109	108	108	107	106	106	106	
105	105	104	104	103	103	102	102	101	101	
100	100	99	99	98	98	97	97	96	96	
95	95	94	94	93	93	93	92	92	91	
91	90	90	89	89	88	88	87	87	86	
86	86	85	85	84	84	83	83	83	82	
82	82	81	81	80	80	80	79	79	78	
78	78	78	77	77	76	76	76	75	75	
75	74	74	74	73	73	73	73	72	72]	



In [0]:

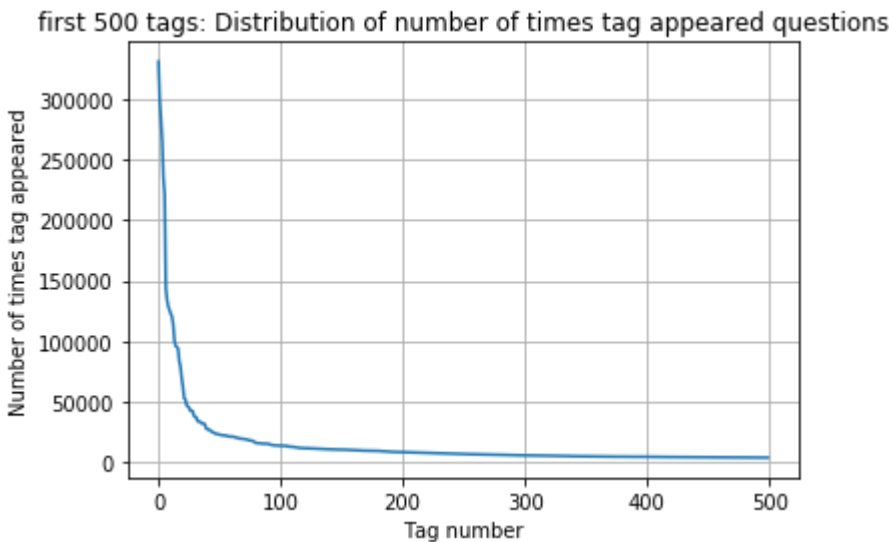
```
1 plt.plot(tag_counts[0:1000])
2 plt.title('first 1k tags: Distribution of number of times tag appeared questions')
3 plt.grid()
4 plt.xlabel("Tag number")
5 plt.ylabel("Number of times tag appeared")
6 plt.show()
7 print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



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

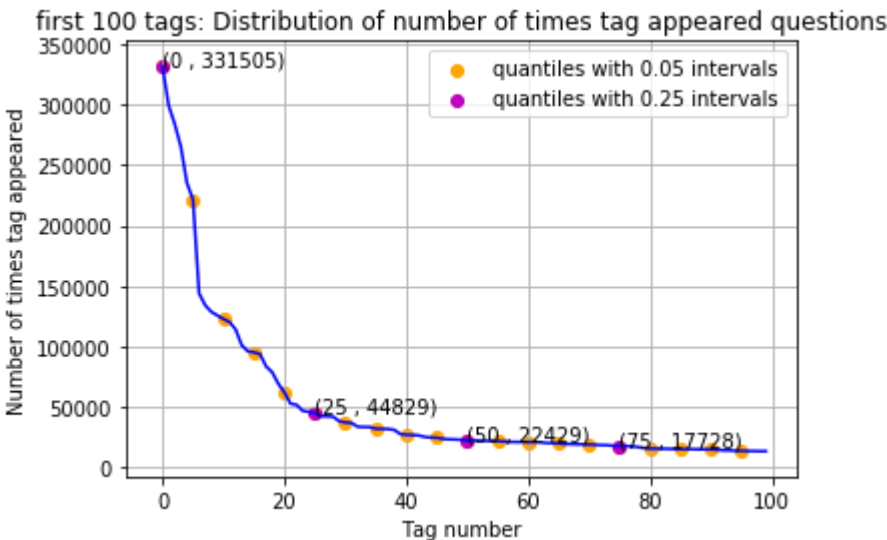
In [0]:

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



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

```
In [0]: 1 plt.plot(tag_counts[0:100], c='b')
2 plt.scatter(x=list(range(0,100,5)), y=tag_counts[0:100:5], c='orange', label="quantiles with 0.05 intervals"
3 # quantiles with 0.25 difference
4 plt.scatter(x=list(range(0,100,25)), y=tag_counts[0:100:25], c='m', label = "quantiles with 0.25 intervals")
5
6 for x,y in zip(list(range(0,100,25)), tag_counts[0:100:25]):
7     plt.annotate(s="({} , {})".format(x,y), xy=(x,y), xytext=(x-0.05, y+500))
8
9 plt.title('first 100 tags: Distribution of number of times tag appeared questions')
10 plt.grid()
11 plt.xlabel("Tag number")
12 plt.ylabel("Number of times tag appeared")
13 plt.legend()
14 plt.show()
15 print(len(tag_counts[0:100:5]), tag_counts[0:100:5])
```



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

```
In [0]: 1 # Store tags greater than 10K in one List
2 lst_tags_gt_10k = tag_df[tag_df.Counts>10000].Tags
3 #Print the Length of the List
4 print ('{} Tags are used more than 10000 times'.format(len(lst_tags_gt_10k)))
5 # Store tags greater than 100K in one List
6 lst_tags_gt_100k = tag_df[tag_df.Counts>100000].Tags
7 #Print the Length of the List.
8 print ('{} Tags are used more than 100000 times'.format(len(lst_tags_gt_100k)))
```

```
153 Tags are used more than 10000 times
14 Tags are used more than 100000 times
```

Observations:

- 1. There are total 153 tags which are used more than 10000 times.
- 2. 14 tags are used more than 100000 times.
- 3. Most frequent tag (i.e. c#) is used 331505 times.
- 4. Since some tags occur much more frequently than others, Micro-averaged F1-score is the appropriate metric for this problem.

3.2.4 Tags Per Question

```
In [0]: 1 #Storing the count of tag in each question in list 'tag_count'
2 tag_quest_count = tag_dtm.sum(axis=1).tolist()
3 #Converting each value in the 'tag_quest_count' to integer.
4 tag_quest_count=[int(j) for i in tag_quest_count for j in i]
5 print ('We have total {} datapoints.'.format(len(tag_quest_count)))
6
7 print(tag_quest_count[:5])
```

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

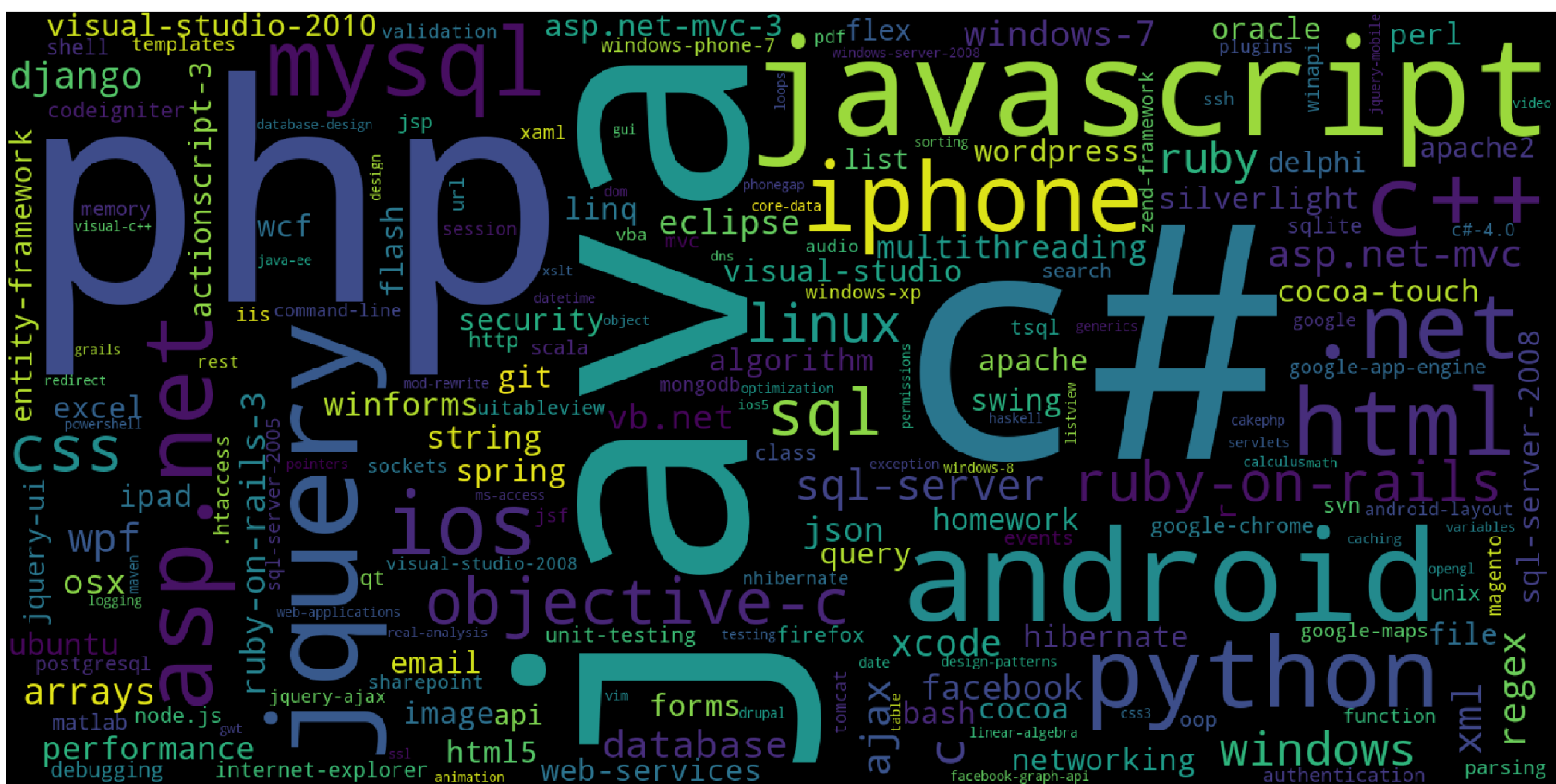
```
In [0]: 1 print( "Maximum number of tags per question: %d"%max(tag_quest_count))
2 print( "Minimum number of tags per question: %d"%min(tag_quest_count))
3 print( "Avg. number of tags per question: %f"% ((sum(tag_quest_count)*1.0)/len(tag_quest_count)))
```

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

Number of Tags	Number of questions
1	570000
2	1120000
3	1210000
4	810000
5	510000

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

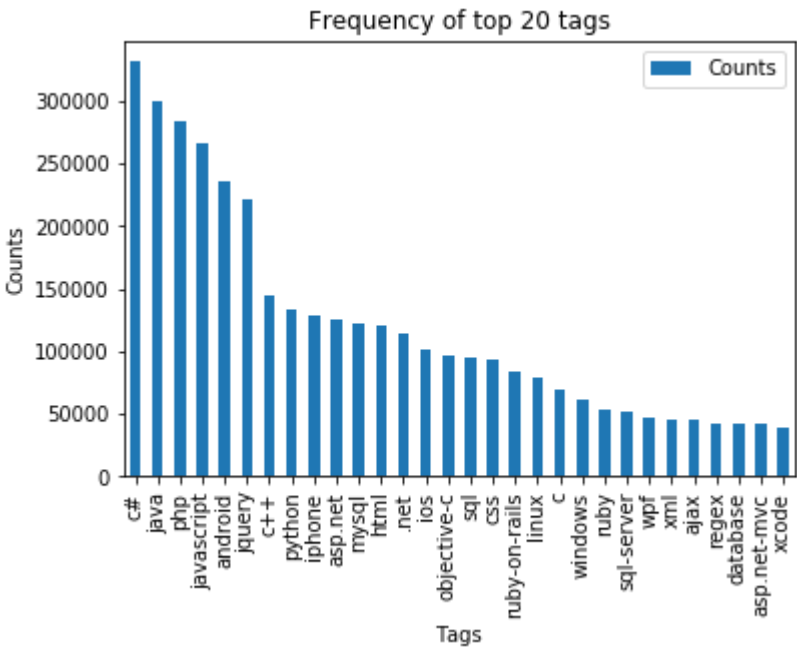
```
In [0]: 1 # Plotting word cloud
2 start = datetime.now()
3
4 # Lets first convert the 'result' dictionary to 'list of tuples'
5 tup = dict(result.items())
6 #Initializing WordCloud using frequencies of tags.
7 wordcloud = WordCloud(    background_color='black',
8                            width=1600,
9                            height=800,
10                           ).generate_from_frequencies(tup)
11
12 fig = plt.figure(figsize=(30,20))
13 plt.imshow(wordcloud)
14 plt.axis('off')
15 plt.tight_layout(pad=0)
16 fig.savefig("tag.png")
17 plt.show()
18 print("Time taken to run this cell :", datetime.now() - start)
```



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

## 11/28

```
In [0]: 1 i=np.arange(30)
2 tag_df_sorted.head(30).plot(kind='bar')
3 plt.title('Frequency of top 20 tags')
4 plt.xticks(i, tag_df_sorted['Tags'])
5 plt.xlabel('Tags')
6 plt.ylabel('Counts')
7 plt.show()
```



Observations:

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

3.3 Cleaning and preprocessing of Questions

3.3.1 Preprocessing

- 1. Sample 1M data points
- 2. Separate out code-snippets from Body
- 3. Remove Spcial 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 [8]: 1 def striphtml(data):
2     cleanr = re.compile('<.*?>')
3     cleantext = re.sub(cleanr, ' ', str(data))
4     return cleantext
5 stop_words = set(stopwords.words('english'))
6 stemmer = SnowballStemmer("english")
```

In [3]:

```
1  #http://www.sqlitetutorial.net/sqlite-python/create-tables/
2  def create_connection(db_file):
3      """ create a database connection to the SQLite database
4          specified by db_file
5      :param db_file: database file
6      :return: Connection object or None
7      """
8      try:
9          conn = sqlite3.connect(db_file)
10         return conn
11     except Error as e:
12         print(e)
13
14     return None
15
16 def create_table(conn, create_table_sql):
17     """ create a table from the create_table_sql statement
18     :param conn: Connection object
19     :param create_table_sql: a CREATE TABLE statement
20     :return:
21     """
22     try:
23         c = conn.cursor()
24         c.execute(create_table_sql)
25     except Error as e:
26         print(e)
27
28 def checkTableExists(dbcon):
29     cursr = dbcon.cursor()
30     str = "select name from sqlite_master where type='table'"
31     table_names = cursr.execute(str)
32     print("Tables in the databse:")
33     tables =table_names.fetchall()
34     print(tables[0][0])
35     return(len(tables))
36
37 def create_database_table(database, query):
38     conn = create_connection(database)
39     if conn is not None:
40         create_table(conn, query)
41         checkTableExists(conn)
42     else:
43         print("Error! cannot create the database connection.")
44     conn.close()
45
46 sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question text NOT NULL, code text, tags
47 create_database_table("Processed.db", sql_create_table)
```

Tables in the databse:  
QuestionsProcessed

In [0]:

```
1  # http://www.sqlitetutorial.net/sqlite-delete/
2  # https://stackoverflow.com/questions/2279706/select-random-row-from-a-sqlite-table
3  start = datetime.now()
4  read_db = 'train_no_dup.db'
5  write_db = 'Processed.db'
6  if os.path.isfile(read_db):
7      conn_r = create_connection(read_db)
8      if conn_r is not None:
9          reader =conn_r.cursor()
10         reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY RANDOM() LIMIT 1000000;")
11
12 if os.path.isfile(write_db):
13     conn_w = create_connection(write_db)
14     if conn_w is not None:
15         tables = checkTableExists(conn_w)
16         writer =conn_w.cursor()
17         if tables != 0:
18             writer.execute("DELETE FROM QuestionsProcessed WHERE 1")
19             print("Cleared All the rows")
20 print("Time taken to run this cell :", datetime.now() - start)
```

Tables in the databse:  
QuestionsProcessed  
Cleared All the rows  
Time taken to run this cell : 0:06:32.806567

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



In [0]:

```
1  #http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/
2
3  start = datetime.now()
4  preprocessed_data_list=[]
5  reader.fetchone()
6  questions_with_code=0
7  len_pre=0
8  len_post=0
9  questions_proccesed = 0
10 for row in reader:
11
12     is_code = 0
13
14     title, question, tags = row[0], row[1], row[2]
15
16     if '<code>' in question:
17         questions_with_code+=1
18         is_code = 1
19     x = len(question)+len(title)
20     len_pre+=x
21
22     code = str(re.findall(r'<code>(.*?)</code>', question, flags=re.DOTALL))
23
24     question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTILINE|re.DOTALL)
25     question=stripthtml(question.encode('utf-8'))
26
27     title=title.encode('utf-8')
28
29     question=str(title)+" "+str(question)
30     question=re.sub(r'^A-Za-z+', ' ',question)
31     words=word_tokenize(str(question.lower()))
32
33     #Removing all single letter and and stopwords from question exceptt for the letter 'c'
34     question=' '.join(str(stemmer.stem(j)) for j in words if j not in stop_words and (len(j)!=1 or j=='c'))
35
36     len_post+=len(question)
37     tup = (question,code,tags,x,len(question),is_code)
38     questions_proccesed += 1
39     writer.execute("insert into QuestionsProcessed(question,code,tags,words_pre,words_post,is_code) values ("
40     if (questions_proccesed%100000==0):
41         print("number of questions completed=",questions_proccesed)
42
43 no_dup_avg_len_pre=(len_pre*1.0)/questions_proccesed
44 no_dup_avg_len_post=(len_post*1.0)/questions_proccesed
45
46 print( "Avg. length of questions(Title+Body) before processing: %d"%no_dup_avg_len_pre)
47 print( "Avg. length of questions(Title+Body) after processing: %d"%no_dup_avg_len_post)
48 print( "Percent of questions containing code: %d"%((questions_with_code*100.0)/questions_proccesed))
49
50 print("Time taken to run this cell :", datetime.now() - start)
```

number of questions completed= 100000  
number of questions completed= 200000  
number of questions completed= 300000  
number of questions completed= 400000  
number of questions completed= 500000  
number of questions completed= 600000  
number of questions completed= 700000  
number of questions completed= 800000  
number of questions completed= 900000  
Avg. length of questions(Title+Body) before processing: 1169  
Avg. length of questions(Title+Body) after processing: 327  
Percent of questions containing code: 57  
Time taken to run this cell : 0:47:05.946582

In [0]:

```
1  # dont forget to close the connections, or else you will end up with locks
2  conn_r.commit()
3  conn_w.commit()
4  conn_r.close()
5  conn_w.close()
```

In [0]:

```
1 if os.path.isfile(write_db):
2     conn_r = create_connection(write_db)
3     if conn_r is not None:
4         reader =conn_r.cursor()
5         reader.execute("SELECT question From QuestionsProcessed LIMIT 10")
6         print("Questions after preprocessed")
7         print('='*100)
8         reader.fetchone()
9         for row in reader:
10             print(row)
11             print('-'*100)
12 conn_r.commit()
13 conn_r.close()
```

Questions after preprocessed

=====

('ef code first defin one mani relationship differ key troubl defin one zero mani relationship entiti ef object model look like use fluent api object composi pk defin batch id batch detail id use fluent api object composi pk defin batch detail id compani id map exist databas tpt basic idea submittedtransact zero mani submittedsplit transact associ navig realli need one way submittedtransact submittedsplittransact need dbcontext class onmodel cr overrid map class lazi load occur submittedtransact submittedsplittransact help would much appreci edit take n advic made follow chang dbcontext class ad follow onmodelcr overrid must miss someth get follow except thrown submittedtransact key batch id batch detail id zero one mani submittedsplittransact key batch detail id compani id rather assum convent creat relationship two object configur requir sinc obvious wrong',)

-----

('explan new statement review section c code came accross statement block come accross new oper use way someon explain new call way',)

-----

('error function notat function solv logic riddl iloczyni list structur list possibl candid solut list possibl coordin matrix wan na choos one candid compar possibl candid element equal wan na delet coordin call function s kasuj look like ni knowledg haskel cant see what wrong',)

-----

('step plan move one isp anoth one work busi plan switch isp realli soon need chang lot inform dns wan wan wifi question guy help mayb peopl plan correct chang current isp new one first dns know receiv new ip isp major chan g need take consider exchang server owa vpn two site link wireless connect km away citrix server vmware exchang domain control link place import server crucial step inform need know avoid downtim busi regard ndavid',)

-----

('use ef migrat creat databas googl migrat tutori af first run applic creat databas ef enabl migrat way creat d atabas migrat rune applic tri',)

-----

('magento unit test problem magento site recent look way check integr magento site given point unit test jump o ne method would assum would big job write whole lot test check everyth site work anyon involv unit test magento advis follow possibl test whole site custom modul nis exampl test would amaz given site heavili link databas wo uld nbe possibl fulli test site without disturb databas better way automaticlli check integr magento site say i ntegr realli mean fault site ship payment etc work correct',)

-----

('find network devic without bonjour write mac applic need discov mac pcs iphon ipad connect wifi network bonjo ur seem reason choic turn problem mani type router mine exampl work block bonjour servic need find ip devic tri connect applic specif port determin process run best approach accomplish task without violat app store sandbo x',)

-----

('send multipl row mysql databas want send user mysql databas column user skill time nnow want abl add one row user differ time etc would code send databas nthen use help schema',)

-----

('insert data mysql php powerpoint event powerpoint present run continu way updat slide present automat data my sql databas websit',)

-----

In [0]:

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

In [0]:

```
1 preprocessed_data.head()
```

Out[47]:

	question	tags
0	resiz root window tkinter resizable root window re...	python tkinter
1	ef code first defin one mani relationship diff...	entity-framework-4.1
2	explan new statement review section c code cam...	c++
3	error function notat function solv logic riddl...	haskell logic
4	step plan move one isp anoth one work busi pla...	dns isp

In [0]:

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

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



## 4. Machine Learning Models

### 4.1 Converting tags for multilabel problems

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

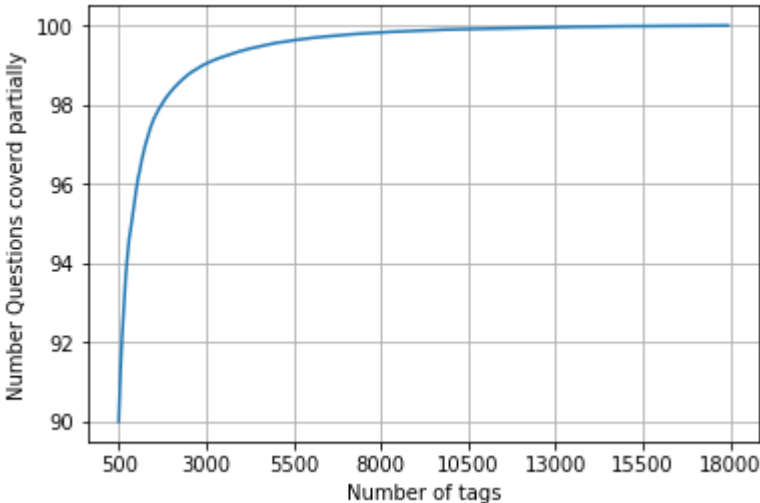
```
In [0]: 1 # binary='true' will give a binary vectorizer
2 vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
3 multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

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

```
In [12]: 1 def tags_to_choose(n):
2     t = multilabel_y.sum(axis=0).tolist()[0]
3     sorted_tags_i = sorted(range(len(t)), key=lambda i: t[i], reverse=True)
4     multilabel_yn=multilabel_y[:,sorted_tags_i[:n]]
5     return multilabel_yn
6
7 def questions_explained_fn(n):
8     multilabel_yn = tags_to_choose(n)
9     x= multilabel_yn.sum(axis=1)
10    return (np.count_nonzero(x==0))
```

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

```
In [0]: 1 fig, ax = plt.subplots()
2 ax.plot(questions_explained)
3 xlabel = list(500+np.array(range(-50,450,50))*50)
4 ax.set_xticklabels(xlabel)
5 plt.xlabel("Number of tags")
6 plt.ylabel("Number Questions covered partially")
7 plt.grid()
8 plt.show()
9 # you can choose any number of tags based on your computing power, minimun is 50(it covers 90% of the tags)
10 print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
```



with 5500 tags we are covering 99.04 % of questions

```
In [0]: 1 multilabel_yx = tags_to_choose(5500)
2 print("number of questions that are not covered :", questions_explained_fn(5500),"out of ", total_qs)
```

number of questions that are not covered : 9599 out of 999999

```
In [0]: 1 print("Number of tags in sample :", multilabel_y.shape[1])
2 print("number of tags taken :", multilabel_yx.shape[1], "(",(multilabel_yx.shape[1]/multilabel_y.shape[1])*100
```

Number of tags in sample : 35422  
number of tags taken : 5500 ( 15.527073570097679 %)

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

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

```
In [0]: 1 total_size=preprocessed_data.shape[0]
2 train_size=int(0.80*total_size)
3
4 x_train=preprocessed_data.head(train_size)
5 x_test=preprocessed_data.tail(total_size - train_size)
6
7 y_train = multilabel_yx[0:train_size,:]
8 y_test = multilabel_yx[train_size:total_size,:]
```

```
In [0]: 1 print("Number of data points in train data :", y_train.shape)
2 print("Number of data points in test data :", y_test.shape)
```

Number of data points in train data : (799999, 5500)  
Number of data points in test data : (200000, 5500)

### 4.3 Featurizing data

```
In [0]: 1 start = datetime.now()
2 vectorizer = TfidfVectorizer(min_df=0.00009, max_features=200000, smooth_idf=True, norm="l2", \
3                               tokenizer = lambda x: x.split(), sublinear_tf=False, ngram_range=(1,3))
4 x_train_multilabel = vectorizer.fit_transform(x_train['question'])
5 x_test_multilabel = vectorizer.transform(x_test['question'])
6 print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:09:50.460431

```
In [0]: 1 print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.shape)
2 print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

Dimensions of train data X: (799999, 88244) Y : (799999, 5500)  
Dimensions of test data X: (200000, 88244) Y: (200000, 5500)

```
In [0]: 1 # https://www.analyticsvidhya.com/blog/2017/08/introduction-to-multi-label-classification/
2 #https://stats.stackexchange.com/questions/117796/scikit-multi-label-classification
3 # classifier = LabelPowerset(GaussianNB())
4 """
5 from skmultilearn.adapt import MLkNN
6 classifier = MLkNN(k=21)
7
8 # train
9 classifier.fit(x_train_multilabel, y_train)
10
11 # predict
12 predictions = classifier.predict(x_test_multilabel)
13 print(accuracy_score(y_test,predictions))
14 print(metrics.f1_score(y_test, predictions, average = 'macro'))
15 print(metrics.f1_score(y_test, predictions, average = 'micro'))
16 print(metrics.hamming_loss(y_test,predictions))
17
18 """
19 # we are getting memory error because the multilearn package
20 # is trying to convert the data into dense matrix
21 # -----
22 #MemoryError                                Traceback (most recent call last)
23 #<ipython-input-170-f0e7c7f3e0be> in <module>()
24 #----> classifier.fit(x_train_multilabel, y_train)
```

Out[92]: "\nfrom skmultilearn.adapt import MLkNN\nclassifier = MLkNN(k=21)\n\n# train\nclassifier.fit(x\_train\_multilabel, y\_train)\n\n# predict\npredictions = classifier.predict(x\_test\_multilabel)\nprint(accuracy\_score(y\_test,predictions))\nprint(metrics.f1\_score(y\_test, predictions, average = 'macro'))\nprint(metrics.f1\_score(y\_test, predictions, average = 'micro'))\nprint(metrics.hamming\_loss(y\_test,predictions))\n\n"

### 4.4 Applying Logistic Regression with OneVsRest Classifier

```
In [0]: 1 # this will be taking so much time try not to run it, download the lr_with_equal_weight.pkl file and use to
2 # This takes about 6-7 hours to run.
3 classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_jobs=-1)
4 classifier.fit(x_train_multilabel, y_train)
5 predictions = classifier.predict(x_test_multilabel)
6
7 print("accuracy :",metrics.accuracy_score(y_test,predictions))
8 print("macro f1 score :",metrics.f1_score(y_test, predictions, average = 'macro'))
9 print("micro f1 scoore :",metrics.f1_score(y_test, predictions, average = 'micro'))
10 print("hamming loss :",metrics.hamming_loss(y_test,predictions))
11 print("Precision recall report :\n",metrics.classification_report(y_test, predictions))
12
```

```
accuracy : 0.081965
macro f1 score : 0.0963020140154
micro f1 scoore : 0.374270748817
hamming loss : 0.00041225090909090907
Precision recall report :

```

	precision	recall	f1-score	support
0	0.62	0.23	0.33	15760
1	0.79	0.43	0.56	14039
2	0.82	0.55	0.66	13446
3	0.76	0.42	0.54	12730
4	0.94	0.76	0.84	11229
5	0.85	0.64	0.73	10561
6	0.70	0.30	0.42	6958
7	0.87	0.61	0.72	6309
8	0.70	0.40	0.50	6032
9	0.78	0.43	0.55	6020
10	0.86	0.62	0.72	5707
11	0.52	0.17	0.25	5723
12	0.55	0.10	0.10	5531

```
In [0]: 1 from sklearn.externals import joblib
2 joblib.dump(classifier, 'lr_with_equal_weight.pkl')
```

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

```
In [0]: 1 sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (question text NOT NULL, code text, tags
2 create_database_table("Titlemoreweight.db", sql_create_table)
```

Tables in the databse:  
QuestionsProcessed

```
In [0]: 1 # http://www.sqlitetutorial.net/sqlite-delete/
2 # https://stackoverflow.com/questions/2279706/select-random-row-from-a-sqlite-table
3
4 read_db = 'train_no_dup.db'
5 write_db = 'Titlemoreweight.db'
6 train_datasize = 400000
7 if os.path.isfile(read_db):
8     conn_r = create_connection(read_db)
9     if conn_r is not None:
10         reader =conn_r.cursor()
11         # for selecting first 0.5M rows
12         reader.execute("SELECT Title, Body, Tags From no_dup_train LIMIT 500001;")
13         # for selecting random points
14         #reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY RANDOM() LIMIT 500001;")
15
16 if os.path.isfile(write_db):
17     conn_w = create_connection(write_db)
18     if conn_w is not None:
19         tables = checkTableExists(conn_w)
20         writer =conn_w.cursor()
21         if tables != 0:
22             writer.execute("DELETE FROM QuestionsProcessed WHERE 1")
23             print("Cleared All the rows")
```

Tables in the databse:  
QuestionsProcessed  
Cleared All the rows

### 4.5.1 Preprocessing of questions

1. Separate Code from Body
2. Remove Spcial 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 [0]:

```
1 #http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/
2 start = datetime.now()
3 preprocessed_data_list=[]
4 reader.fetchone()
5 questions_with_code=0
6 len_pre=0
7 len_post=0
8 questions_proccesed = 0
9 for row in reader:
10
11     is_code = 0
12
13     title, question, tags = row[0], row[1], str(row[2])
14
15     if '<code>' in question:
16         questions_with_code+=1
17         is_code = 1
18     x = len(question)+len(title)
19     len_pre+=x
20
21     code = str(re.findall(r'<code>(.*?)</code>', question, flags=re.DOTALL))
22
23     question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTILINE|re.DOTALL)
24     question=stripthtml(question.encode('utf-8'))
25
26     title=title.encode('utf-8')
27
28     # adding title three time to the data to increase its weight
29     # add tags string to the training data
30
31     question=str(title)+" "+str(title)+" "+str(title)+" "+question
32
33     # if questions_proccesed<=train_datasize:
34     #     question=str(title)+" "+str(title)+" "+str(title)+" "+question+" "+str(tags)
35     # else:
36     #     question=str(title)+" "+str(title)+" "+str(title)+" "+question
37
38     question=re.sub(r'^A-Za-z0-9#+.\-]+', ' ',question)
39     words=word_tokenize(str(question.lower()))
40
41     #Removing all single letter and and stopwords from question exceptt for the letter 'c'
42     question=' '.join(str(stemmer.stem(j)) for j in words if j not in stop_words and (len(j)!=1 or j=='c'))
43
44     len_post+=len(question)
45     tup = (question,code,tags,x,len(question),is_code)
46     questions_proccesed += 1
47     writer.execute("insert into QuestionsProcessed(question,code,tags,words_pre,words_post,is_code) values (
48     if (questions_proccesed%100000==0):
49         print("number of questions completed=",questions_proccesed)
50
51 no_dup_avg_len_pre=(len_pre*1.0)/questions_proccesed
52 no_dup_avg_len_post=(len_post*1.0)/questions_proccesed
53
54 print( "Avg. length of questions(Title+Body) before processing: %d"%no_dup_avg_len_pre)
55 print( "Avg. length of questions(Title+Body) after processing: %d"%no_dup_avg_len_post)
56 print( "Percent of questions containing code: %d"%((questions_with_code*100.0)/questions_proccesed))
57
58 print("Time taken to run this cell :", datetime.now() - start)
```

number of questions completed= 100000  
number of questions completed= 200000  
number of questions completed= 300000  
number of questions completed= 400000  
number of questions completed= 500000  
Avg. length of questions(Title+Body) before processing: 1239  
Avg. length of questions(Title+Body) after processing: 424  
Percent of questions containing code: 57  
Time taken to run this cell : 0:23:12.329039

In [0]:

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

Sample quesitons after preprocessing of data

In [0]:

```
1  if os.path.isfile(write_db):
2      conn_r = create_connection(write_db)
3      if conn_r is not None:
4          reader =conn_r.cursor()
5          reader.execute("SELECT question From QuestionsProcessed LIMIT 10")
6          print("Questions after preprocessed")
7          print('='*100)
8          reader.fetchone()
9          for row in reader:
10             print(row)
11             print('- '*100)
12 conn_r.commit()
13 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 necess ari 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 foll ow guid link instal jstl got follow error tri launch jsp page java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid taglib declar instal jstl 1.1 tomcat webapp tri project work also tri version 1.2 jstl s till messag caus solv',)

-----

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

-----

('better way updat feed fb php sdk better way updat feed fb php sdk better way updat feed fb php sdk novic face book api read mani tutori still confused.i find post feed api method like correct second way use curl someth li ke 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 btnad d click event open anoth window nafter insert record close window',)

-----

('sql inject issu prevent correct form submiss php sql inject issu prevent correct form submiss php sql inject issu prevent correct form submiss php check everyth think make sure input field safe type sql inject good news safe bad news one tag mess form submiss place even touch life figur exact html use templat file forgiv okay ent ir php script get execut see data post none forum field post problem use someth titl field none data get post c urrent use print post see submit noth work flawless statement though also mention script work flawless local ma chin use host come across problem state list input test mess',)

-----

('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl subaddit lebesgu measur let lbrace rbrace sequenc set sigma -algebra mathcal want show left bigcup right leq sum left right countabl addit measur defin set sigma algebra mathcal think use monoton properti somewher proof start appreci littl help nthank ad ha n answer make follow addit construct given han answer clear bigcup bigcup cap emptyset neq left bigcup right le ft 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 properti name e rror 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 impo rt framework send email applic background import framework i.e skpsmtpmessag somebodi suggest get error collect 2 ld return exit status import framework correct sorc taken framework follow mfmcomposeviewcontrol question lock field updat answer drag drop folder project click copi nthat',)

-----

Saving Preprocessed data to a Database

In [0]:

```
1  #Taking 0.5 Million entries to a dataframe.
2  write_db = 'Titlemoreweight.db'
3  if os.path.isfile(write_db):
4      conn_r = create_connection(write_db)
5      if conn_r is not None:
6          preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM QuestionsProcessed""", conn_r)
7  conn_r.commit()
8  conn_r.close()
```

In [0]:

```
1  preprocessed_data.head()
```

Out[100]:

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

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

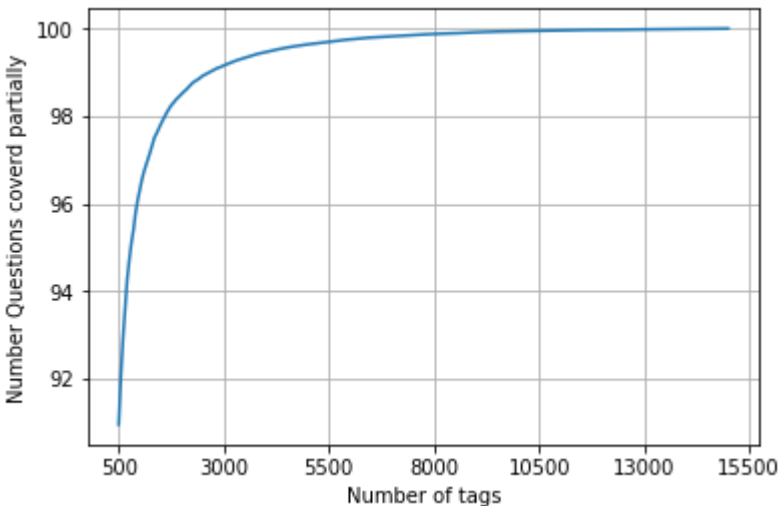
Converting string Tags to multilable output variables

```
In [0]: 1 vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
2 multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

Selecting 500 Tags

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

```
In [0]: 1 fig, ax = plt.subplots()
2 ax.plot(questions_explained)
3 xlabel = list(500+np.array(range(-50,450,50))*50)
4 ax.set_xticklabels(xlabel)
5 plt.xlabel("Number of tags")
6 plt.ylabel("Number Questions covered partially")
7 plt.grid()
8 plt.show()
9 # you can choose any number of tags based on your computing power, minimun is 500(it covers 90% of the tags)
10 print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
11 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 [0]: 1 # we will be taking 500 tags
2 multilabel_yx = tags_to_choose(500)
3 print("number of questions that are not covered :", questions_explained_fn(500),"out of ", total_qs)
```

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

```
In [0]: 1 x_train=preprocessed_data.head(train_datasize)
2 x_test=preprocessed_data.tail(preprocessed_data.shape[0] - 400000)
3
4 y_train = multilabel_yx[0:train_datasize,:]
5 y_test = multilabel_yx[train_datasize:preprocessed_data.shape[0],:]
```

```
In [0]: 1 print("Number of data points in train data :", y_train.shape)
2 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 [0]: 1 start = datetime.now()
2 vectorizer = TfidfVectorizer(min_df=0.00009, max_features=200000, smooth_idf=True, norm="l2", \
3                             tokenizer = lambda x: x.split(), sublinear_tf=False, ngram_range=(1,3))
4 x_train_multilabel = vectorizer.fit_transform(x_train['question'])
5 x_test_multilabel = vectorizer.transform(x_test['question'])
6 print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:03:52.522389

```
In [0]: 1 print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.shape)
2 print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

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

4.5.3 Applying Logistic Regression with OneVsRest Classifier

```
In [0]: 1 start = datetime.now()
2 classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_jobs=-1)
3 classifier.fit(x_train_multilabel, y_train)
4 predictions = classifier.predict (x_test_multilabel)
5
6
7 print("Accuracy :",metrics.accuracy_score(y_test, predictions))
8 print("Hamming loss ",metrics.hamming_loss(y_test,predictions))
9
10
11 precision = precision_score(y_test, predictions, average='micro')
12 recall = recall_score(y_test, predictions, average='micro')
13 f1 = f1_score(y_test, predictions, average='micro')
14
15 print("Micro-average quality numbers")
16 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
17
18 precision = precision_score(y_test, predictions, average='macro')
19 recall = recall_score(y_test, predictions, average='macro')
20 f1 = f1_score(y_test, predictions, average='macro')
21
22 print("Macro-average quality numbers")
23 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
24
25 print (metrics.classification_report(y_test, predictions))
26 print("Time taken to run this cell :", datetime.now() - start)
```

Accuracy : 0.23623  
Hamming loss 0.00278088  
Micro-average quality numbers  
Precision: 0.7216, Recall: 0.3256, F1-measure: 0.4488  
Macro-average quality numbers  
Precision: 0.5473, Recall: 0.2572, F1-measure: 0.3339

	precision	recall	f1-score	support
0	0.94	0.64	0.76	5519
1	0.69	0.26	0.38	8190
2	0.81	0.37	0.51	6529
3	0.81	0.43	0.56	3231
4	0.81	0.40	0.54	6430
5	0.82	0.33	0.47	2879
6	0.87	0.50	0.63	5086
7	0.87	0.54	0.67	4533
8	0.60	0.13	0.22	3000
9	0.81	0.53	0.64	2765
10	0.59	0.17	0.26	3051
11	0.70	0.22	0.45	2000

```
In [0]: 1 joblib.dump(classifier, 'lr_with_more_title_weight.pkl')
```

Out[113]: ['lr\_with\_more\_title\_weight.pkl']



```
In [0]: 1 start = datetime.now()
2 classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1'), n_jobs=-1)
3 classifier_2.fit(x_train_multilabel, y_train)
4 predictions_2 = classifier_2.predict(x_test_multilabel)
5 print("Accuracy :",metrics.accuracy_score(y_test, predictions_2))
6 print("Hamming loss ",metrics.hamming_loss(y_test,predictions_2))
7
8
9 precision = precision_score(y_test, predictions_2, average='micro')
10 recall = recall_score(y_test, predictions_2, average='micro')
11 f1 = f1_score(y_test, predictions_2, average='micro')
12
13 print("Micro-average quality numbers")
14 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
15
16 precision = precision_score(y_test, predictions_2, average='macro')
17 recall = recall_score(y_test, predictions_2, average='macro')
18 f1 = f1_score(y_test, predictions_2, average='macro')
19
20 print("Macro-average quality numbers")
21 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
22
23 print (metrics.classification_report(y_test, predictions_2))
24 print("Time taken to run this cell :", datetime.now() - start)
```

Accuracy : 0.25108  
Hamming loss 0.00270302  
Micro-average quality numbers  
Precision: 0.7172, Recall: 0.3672, F1-measure: 0.4858  
Macro-average quality numbers  
Precision: 0.5570, Recall: 0.2950, F1-measure: 0.3710

	precision	recall	f1-score	support
0	0.94	0.72	0.82	5519
1	0.70	0.34	0.45	8190
2	0.80	0.42	0.55	6529
3	0.82	0.49	0.61	3231
4	0.80	0.44	0.57	6430
5	0.82	0.38	0.52	2879
6	0.86	0.53	0.66	5086
7	0.87	0.58	0.70	4533
8	0.60	0.13	0.22	3000
9	0.82	0.57	0.67	2765
10	0.60	0.20	0.30	3051
11	0.60	0.20	0.30	3000

## 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)

```
In [41]: 1 #Taking 100k Million entries to a dataframe due to memory constraints.
2 write_db = 'Titlemoreweight.db'
3 if os.path.isfile(write_db):
4     conn_r = create_connection(write_db)
5     if conn_r is not None:
6         preprocessed_data = pd.read_sql_query("""SELECT question, Tags FROM QuestionsProcessed LIMIT 100000""")
7     conn_r.commit()
8     conn_r.close()
```

```
In [42]: 1 preprocessed_data.head()
```

Out[42]:

	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 [43]: 1 print("number of data points in sample :", preprocessed_data.shape[0])
2 print("number of dimensions :", preprocessed_data.shape[1])
```

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

### Converting string Tags to multilable output variables

In [44]:

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

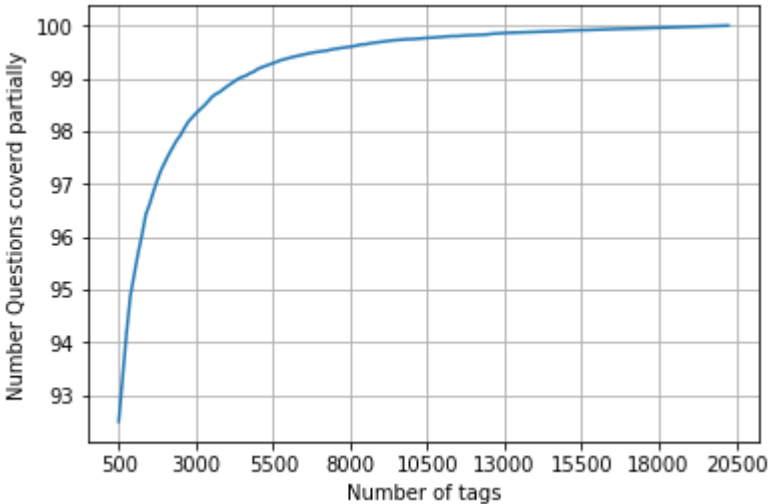
Selecting 500 Tags

In [45]:

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

In [46]:

```
1 fig, ax = plt.subplots()
2 ax.plot(questions_explained)
3 xlabel = list(500+np.array(range(-50,450,50))*50)
4 ax.set_xticklabels(xlabel)
5 plt.xlabel("Number of tags")
6 plt.ylabel("Number Questions covered partially")
7 plt.grid()
8 plt.show()
9 # you can choose any number of tags based on your computing power, minimun is 500(it covers 90% of the tags)
10 print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
11 print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



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

In [47]:

```
1 # we will be taking 500 tags
2 multilabel_yx = tags_to_choose(500)
3 print("number of questions that are not covered :", questions_explained_fn(500),"out of ", total_qs)
```

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

In [54]:

```
1 train_datasize = 80000
2 x_train=preprocessed_data.head(train_datasize)
3 x_test=preprocessed_data.tail(preprocessed_data.shape[0] - train_datasize)
4
5 y_train = multilabel_yx[0:train_datasize,:]
6 y_test = multilabel_yx[train_datasize:preprocessed_data.shape[0],:]
```

In [55]:

```
1 print("Number of data points in train data :", y_train.shape)
2 print("Number of data points in test data :", y_test.shape)
```

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

In [56]:

```
1 start = datetime.now()
2 vectorizer = TfidfVectorizer(min_df=0.00009, max_features=200000, smooth_idf=True, norm="l2", \
3                             tokenizer = lambda x: x.split(), sublinear_tf=False, ngram_range=(1,3))
4 x_train_multilabel = vectorizer.fit_transform(x_train['question'])
5 x_test_multilabel = vectorizer.transform(x_test['question'])
6 print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:01:26.459812

In [57]:

```
1 print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.shape)
2 print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

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

4.5.3 Applying Linear SVM with OneVsRest Classifier

```
In [58]: 1 start = datetime.now()
2 classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.00001, penalty='l1'), n_jobs=1)
3 classifier.fit(x_train_multilabel, y_train)
4 predictions = classifier.predict (x_test_multilabel)
5
6
7 print("Accuracy :",metrics.accuracy_score(y_test, predictions))
8 print("Hamming loss ",metrics.hamming_loss(y_test,predictions))
9
10
11 precision = precision_score(y_test, predictions, average='micro')
12 recall = recall_score(y_test, predictions, average='micro')
13 f1 = f1_score(y_test, predictions, average='micro')
14
15 print("Micro-average quality numbers")
16 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
17
18 precision = precision_score(y_test, predictions, average='macro')
19 recall = recall_score(y_test, predictions, average='macro')
20 f1 = f1_score(y_test, predictions, average='macro')
21
22 print("Macro-average quality numbers")
23 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
24
25 print (metrics.classification_report(y_test, predictions))
26 print("Time taken to run this cell :", datetime.now() - start)
```

Accuracy : 0.20745  
Hamming loss 0.0029409  
Micro-average quality numbers  
Precision: 0.7766, Recall: 0.3021, F1-measure: 0.4350  
Macro-average quality numbers  
Precision: 0.4148, Recall: 0.2165, F1-measure: 0.2621

	precision	recall	f1-score	support
0	0.80	0.37	0.51	820
1	0.74	0.16	0.26	1931
2	0.64	0.14	0.23	544
3	0.68	0.28	0.40	222
4	0.82	0.45	0.58	1311
5	0.90	0.49	0.64	1014
6	0.81	0.39	0.53	1374
7	0.90	0.57	0.70	702
8	0.95	0.65	0.77	1424
9	0.72	0.73	0.73	1037
10	0.81	0.50	0.62	797
11	0.60	0.40	0.50	155

Notes:

- 1. Hamming Loss of SVM OVR (0.0029409)
- 2. Micro-average quality numbers Precision: 0.7766 || Recall: 0.3021 || F1-measure: 0.4350
- 3. Macro-average quality numbers Precision: 0.4148 || Recall: 0.2165 || F1-measure: 0.2621

CountVectorizer for LogRegression One Vs Rest with HyperParameter Tuning

```
In [60]: 1 start = datetime.now()
2 vectorizer = CountVectorizer(min_df=0.00009, max_features=10000, ngram_range=(1,4))
3 x_train_multilabel = vectorizer.fit_transform(x_train['question'])
4 x_test_multilabel = vectorizer.transform(x_test['question'])
5 print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:01:57.833166

```
In [61]: 1 print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",y_train.shape)
2 print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",y_test.shape)
```

Dimensions of train data X: (80000, 10000) Y : (80000, 500)  
Dimensions of test data X: (20000, 10000) Y: (20000, 500)

In [62]:

```
1 from sklearn.model_selection import GridSearchCV
2 from sklearn.linear_model import LogisticRegression
3
4 param = [{"estimator__C": [10**-4, 10**-2, 10**0, 10**2, 10**4]}]
5 classifier = OneVsRestClassifier(LogisticRegression(penalty='l1'), n_jobs=1)
6 gsearch_cv = GridSearchCV(estimator=classifier, param_grid=param, cv=2, verbose=1, scoring='f1_micro', n_jobs
7 gsearch_cv.fit(x_train_multilabel, y_train)
```

Fitting 2 folds for each of 5 candidates, totalling 10 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1143: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.  
'precision', 'predicted', average, warn\_for)  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1143: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.  
'precision', 'predicted', average, warn\_for)  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\metrics\classification.py:1143: UndefinedMetricWarning: F-score is ill-defined and being set to 0.0 due to no predicted samples.  
'precision', 'predicted', average, warn\_for)  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\svm\base.py:922: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.  
"the number of iterations.", ConvergenceWarning)  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\svm\base.py:922: ConvergenceWarning: Liblinear failed to converge, increase the number of iterations.  
"the number of iterations.", ConvergenceWarning)  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 165 is present in all training examples.  
str(classes[c]))  
C:\Users\asus\Anaconda3\lib\site-packages\sklearn\multiclass.py:76: UserWarning: Label not 435 is present in all training examples.  
str(classes[c]))  
[Parallel(n\_jobs=1)]: Done 10 out of 10 | elapsed: 111.0min finished

Out[62]:

```
GridSearchCV(cv=2, error_score='raise-deprecating',
             estimator=OneVsRestClassifier(estimator=LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
             intercept_scaling=1, max_iter=100, multi_class='warn',
             n_jobs=None, penalty='l1', random_state=None, solver='warn',
             tol=0.0001, verbose=0, warm_start=False),
             n_jobs=1),
             fit_params=None, iid='warn', n_jobs=1,
             param_grid=[{'estimator__C': [0.0001, 0.01, 1, 100, 10000]}],
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='f1_micro', verbose=1)
```

In [63]:

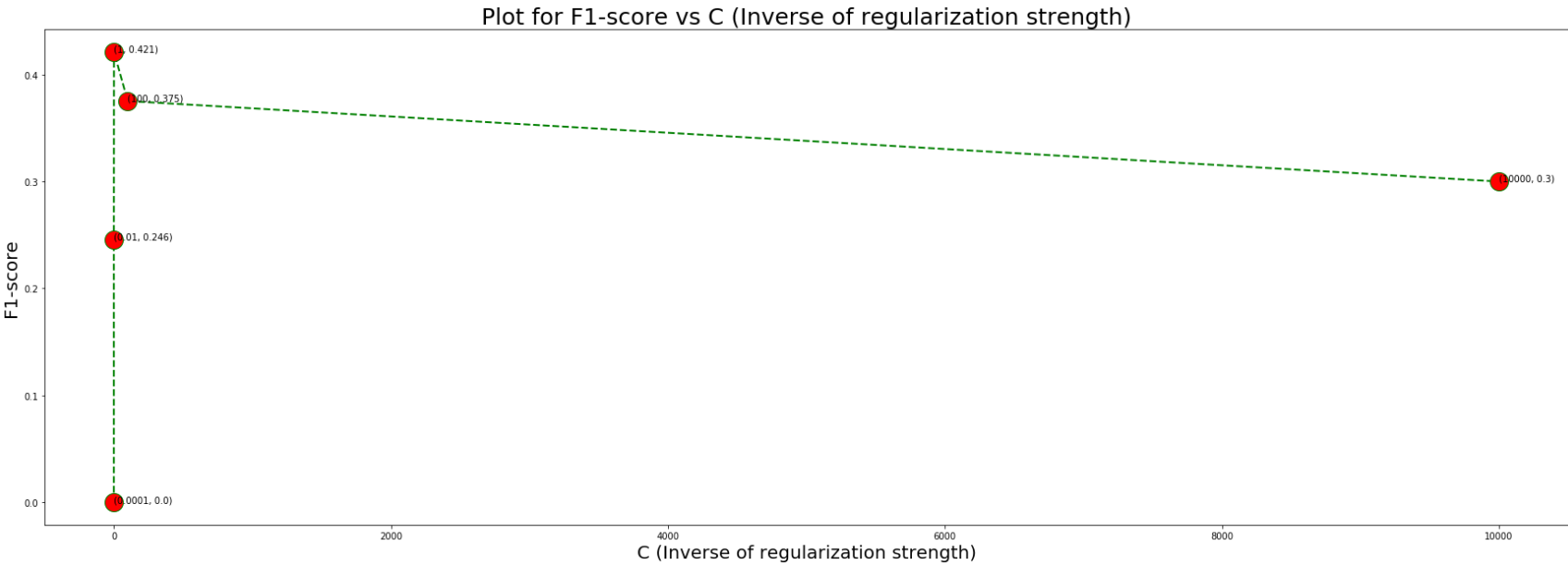
```
1 print("Best estimator for the model :\n ",gsearch_cv.best_estimator_)
2 print("Best Score for the model : ",gsearch_cv.best_score_)
```

Best estimator for the model :  
OneVsRestClassifier(estimator=LogisticRegression(C=1, class\_weight=None, dual=False, fit\_intercept=True, intercept\_scaling=1, max\_iter=100, multi\_class='warn', n\_jobs=None, penalty='l1', random\_state=None, solver='warn', tol=0.0001, verbose=0, warm\_start=False), n\_jobs=1)  
Best Score for the model : 0.4210240550627034

Best estimator C = 1

In [64]:

```
1 #Source:Github.com
2 # Here we obtain the c values and their corresponding mean test scores.
3 cv_result = gsearch_cv.cv_results_
4 mts = cv_result["mean_test_score"] #list that will hold the mean of cross validation scores for each c
5 c = cv_result["params"]
6
7 c_values = [] #list that will hold all the c values that the grid search cross validator tried
8 for i in range(0,len(c)):
9     c_values.append(c[i]["estimator__C"])
10
11 #Plot F1-score vs C values
12 plt.figure(figsize=(30,10))
13 plt.plot(c_values , mts, color='green', linestyle='dashed', linewidth=2, marker='o', markerfacecolor='red',
14 for xy in zip(c_values, np.round(mts,3)):
15     plt.annotate('(%s, %s)' % xy, xy=xy, textcoords='data')
16 plt.title('Plot for F1-score vs C (Inverse of regularization strength)',fontsize=25)
17 plt.xlabel('C (Inverse of regularization strength)',fontsize=20)
18 plt.ylabel('F1-score',fontsize=20)
19 plt.show()
```



In [65]:

```
1 classifier = gsearch_cv.best_estimator_
2 classifier.fit(x_train_multilabel, y_train)
3
4 predictions = classifier.predict(x_test_multilabel)
5 print("Accuracy :",metrics.accuracy_score(y_test, predictions))
6 print("Hamming loss ",metrics.hamming_loss(y_test,predictions))
7
8
9 precision = precision_score(y_test, predictions, average='micro')
10 recall = recall_score(y_test, predictions, average='micro')
11 f1 = f1_score(y_test, predictions, average='micro')
12
13 print("Micro-average quality numbers")
14 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
15
16 precision = precision_score(y_test, predictions, average='macro')
17 recall = recall_score(y_test, predictions, average='macro')
18 f1 = f1_score(y_test, predictions, average='macro')
19
20 print("Macro-average quality numbers")
21 print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))
22
23 print (metrics.classification_report(y_test, predictions))
```

Accuracy : 0.1586  
Hamming loss 0.0035862  
Micro-average quality numbers  
Precision: 0.5323, Recall: 0.3543, F1-measure: 0.4254  
Macro-average quality numbers  
Precision: 0.3882, Recall: 0.2883, F1-measure: 0.3186

	precision	recall	f1-score	support
0	0.73	0.37	0.49	820
1	0.50	0.19	0.28	1931
2	0.31	0.18	0.23	544
3	0.47	0.20	0.28	222
4	0.67	0.50	0.58	1311
5	0.73	0.50	0.59	1014
6	0.62	0.41	0.49	1374
7	0.69	0.61	0.65	702
8	0.82	0.62	0.71	1424
9	0.73	0.65	0.69	1037
10	0.53	0.40	0.45	797
11	0.53	0.40	0.45	156

Notes:

- 1. Hamming Loss for LogRegression OVR = 0.0035862
- 2. Micro-average quality numbers Precision: 0.5323 || Recall: 0.3543 || F1-measure: 0.4254
- 3. Macro-average quality numbers Precision: 0.3882 || Recall: 0.2883 || F1-measure: 0.3186



## Procedure for solving CaseStudy:-

- Firstly, the Business Problem is that  
**Suggest the tags based on the content that was there in the question posted on Stackoverflow.**
- The major objective is Predict as many tags as possible with high precision and recall And the Performance metric used = Mean Average F1Score because it takes Frequency of Tags into Account.
- We loaded Data using pandas.
- Did Ananalysis of Tag and found most frequent tags are Programming Languages like:

C#,Java,php,Javascript,Android

- Did Analysis of Titles by finding titles similar to above programming languages.
- Did Preprocessing and cleaning of Data and took 1M points for Analysis and 5500 tags.
- Due to compute resources took only 0.5M points with 500 tags .
- Furthher reduced to 100K points and 500 tags.
- performed TFIDF Vectorzer and did SVM One Vs Rest .
- Performd CountVect with n\_gram=(1,4), took top 10K featureess.
- Applied GridSearch for hypertuning with Logistic Regression One Vs Rest.
- Found Best HyperParam of C = 1
- Applied Logistic Regression One Vs Rest

In [66]:

```
1 from prettytable import PrettyTable
2 table = PrettyTable()
3 table.field_names = [ "Model","hyperparameter", "Accuracy","HammingLoss" , "MicroAvg F1-Score","MacroAvg F1-
4 table.add_row(["Logistic Regression OneVsRestClassifier","C = 1",'0.1586'," 0.0035862","0.4254","0.3186"])
5 table.add_row(["Linear-SVM OneVsRestClassifier","alpha = 0.00001", "0.20745" , "0.0029409", "0.4350","0.2621
6 print(table)
```

Model	hyperparameter	Accuracy	HammingLoss	MicroAvg F1-Score	MacroAvg F1-Score
Logistic Regression OneVsRestClassifier	C = 1	0.1586	0.0035862	0.4254	0.3186
Linear-SVM OneVsRestClassifier	alpha = 0.00001	0.20745	0.0029409	0.4350	0.2621