



In [1]:

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

## Stack Overflow: Tag Prediction

# 1. Business Problem

## 1.1 Description

### Description

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

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

### Problem Statement

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

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

## 1.2 Source / useful links

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

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

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

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

## 1.3 Real World / Business Objectives and Constraints

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

# 2. Machine Learning problem

## 2.1 Data

### 2.1.1 Data Overview

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

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

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

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

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

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

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

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

### Data Field Explanation

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

**Id** - Unique identifier for each question

**Title** - The question's title

**Body** - The body of the question

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

### 2.1.2 Example Data point

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

**Body :**

```

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

    for(int y=1; y<n+1; y++)\n
    {\n
        cin>>m[y];\n
        cin>>u[y];\n
    }\n
    for(x=1; x<n+1; x++)\n
    {\n
        a[x] = (m[x] + u[x])/2;\n
    }\n
    c=(n*4)-4;\n
    for(int a1=1; a1<n+1; a1++)\n
    {\n\n
        e[a1][0] = m[a1];\n
        e[a1][1] = m[a1]+1;\n
        e[a1][2] = u[a1]-1;\n
        e[a1][3] = u[a1];\n
    }\n
    for(int i=1; i<n+1; i++)\n
    {\n
        for(int l=1; l<=i; l++)\n
        {\n
            if(l!=1)\n
            {\n
                cout<<a[l]<<"\\t";\n
            }\n
        }\n
        for(int j=0; j<4; j++)\n
        {\n
            cout<<e[i][j];\n
            for(int k=0; k<n-(i+1); k++)\n
            {\n
                cout<<a[k]<<"\\t";\n
            }\n
            cout<<"\\n";\n
        }\n
    }\n
    \n\n
    system("PAUSE");\n
    return 0;    \n
}

```

```
} \n
```

```
\n \n
```

The answer should come in the form of a table like

```
\n \n
```

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

```
\n \n
```

if the no of inputs is 3 and their ranges are \n

```
1,100 \n
```

```
1,100 \n
```

```
1,100 \n
```

```
(could be varied too)
```

```
\n \n
```

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

```
\n'
```

```
Tags : 'c++ c'
```

## 2.2 Mapping the real-world problem to a Machine Learning Problem

### 2.2.1 Type of Machine Learning Problem

It is a multi-label classification problem

**Multi-label Classification:** Multilabel classification assigns to each sample a set of target labels. This can be thought as predicting properties of a data-point that are not mutually exclusive, such as topics that are relevant for a document. A question on Stackoverflow might be about any of C, Pointers, FileIO and/or memory-

management at the same time or none of these.

**Credit:** <http://scikit-learn.org/stable/modules/multiclass.html> (<http://scikit-learn.org/stable/modules/multiclass.html>)

## 2.2.2 Performance metric

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

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

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

### 'Micro f1 score':

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

### 'Macro f1 score':

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

<https://www.kaggle.com/wiki/MeanFScore> (<https://www.kaggle.com/wiki/MeanFScore>)  
[http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1\\_score.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html) ([http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1\\_score.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html))

**Hamming loss** : The Hamming loss is the fraction of labels that are incorrectly predicted.

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

## 3. Exploratory Data Analysis

### 3.1 Data Loading and Cleaning

#### 3.1.1 Using Pandas with SQLite to Load the data

In [0]:

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

### 3.1.2 Counting the number of rows

In [0]:

```
if os.path.isfile('train.db'):
    start = datetime.now()
    con = sqlite3.connect('train.db')
    num_rows = pd.read_sql_query("""SELECT count(*) FROM data""", con)
    #Always remember to close the database
    print("Number of rows in the database :", "\n", num_rows['count(*)'].values[0])
    con.close()
    print("Time taken to count the number of rows :", datetime.now() - start)
else:
    print("Please download the train.db file from drive or run the above cell to generate train.db")
```

Number of rows in the database :

6034196

Time taken to count the number of rows : 0:01:15.750352

### 3.1.3 Checking for duplicates

In [0]:

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

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

In [0]:

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

Out[6]:

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

In [0]:

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

number of duplicate questions : 1827881 ( 30.2920389063 % )

In [0]:

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

Out[8]:

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





In [2]:

```

#This method seems more appropriate to work with this much data.
#creating the connection with database file.
if os.path.isfile('train_no_dup.db'):
    start = datetime.now()
    con = sqlite3.connect('train_no_dup.db')
    tag_data = pd.read_sql_query("""SELECT Tags FROM no_dup_train""", con)
    #Always remember to close the database
    con.close()

    # Let's now drop unwanted column.
    tag_data.drop(tag_data.index[0], inplace=True)
    #Printing first 5 columns from our data frame
    tag_data.head()
    print("Time taken to run this cell :", datetime.now() - start)
else:
    print("Please download the train.db file from drive or run the above cells to generate

```

Time taken to run this cell : 0:00:17.142934

## 3.2 Analysis of Tags

### 3.2.1 Total number of unique tags

In [3]:

```

# 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 [4]:

```

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 [5]:

```

#'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', '.bas  
h-profile', '.class-file', '.cs-file', '.doc', '.drv', '.ds-store']

### 3.2.3 Number of times a tag appeared

In [6]:

```
# 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 [7]:

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

Out[7]:

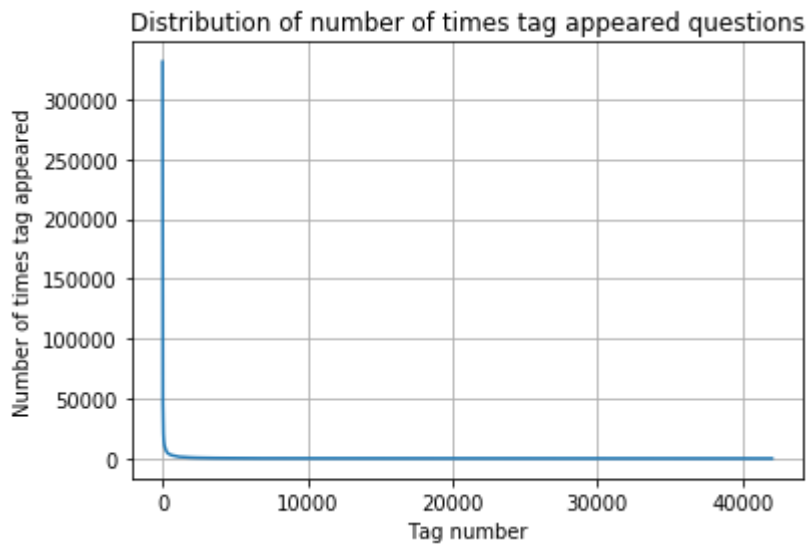
	Tags	Counts
0	.a	18
1	.app	37
2	.asp.net-mvc	1
3	.aspxauth	21
4	.bash-profile	138

In [8]:

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

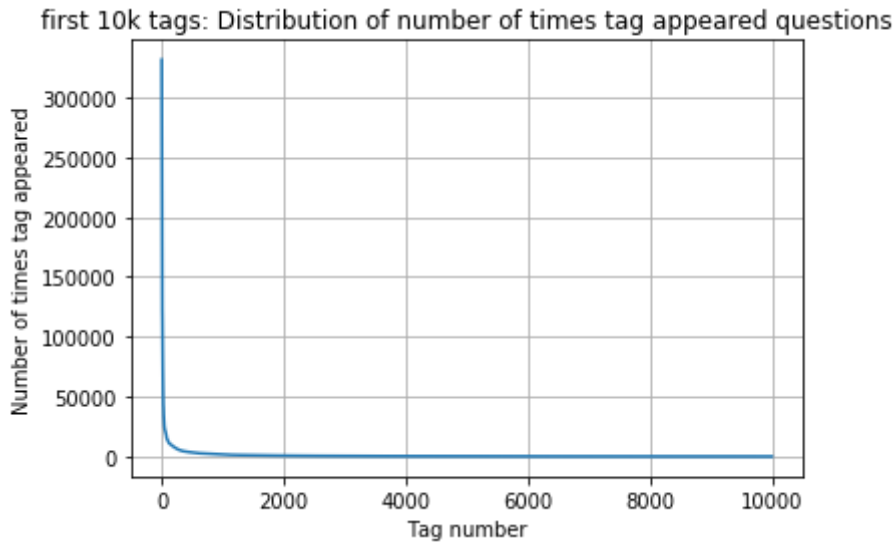
In [9]:

```
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 [10]:

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

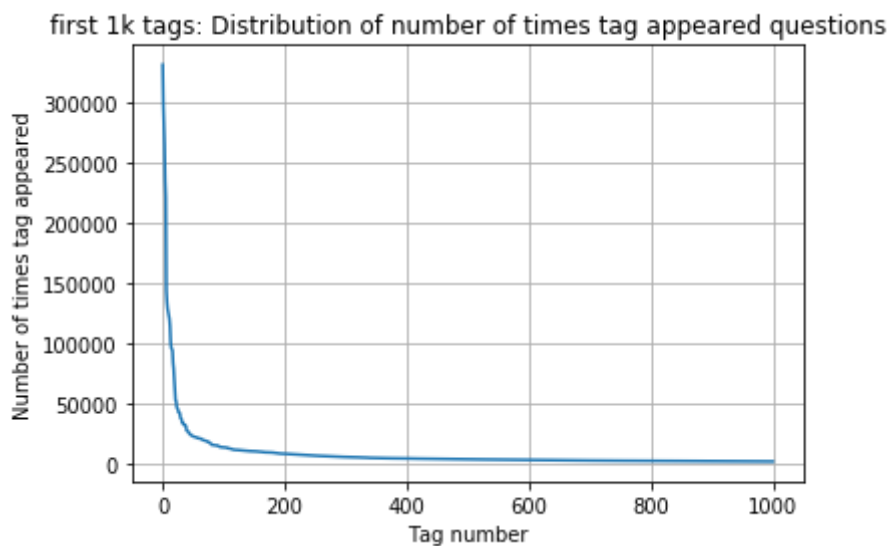


400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7151
6466	5865	5370	4983	4526	4281	4144	3929	3750	3593	
3453	3299	3123	2989	2891	2738	2647	2527	2431	2331	
2259	2186	2097	2020	1959	1900	1828	1770	1723	1673	
1631	1574	1532	1479	1448	1406	1365	1328	1300	1266	
1245	1222	1197	1181	1158	1139	1121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

100	100	99	99	98	98	97	97	96	96
95	95	94	94	93	93	93	92	92	91
91	90	90	89	89	88	88	87	87	86
86	86	85	85	84	84	83	83	83	82
82	82	81	81	80	80	80	79	79	78
78	78	78	77	77	76	76	76	75	75
75	74	74	74	73	73	73	73	72	72]

In [11]:

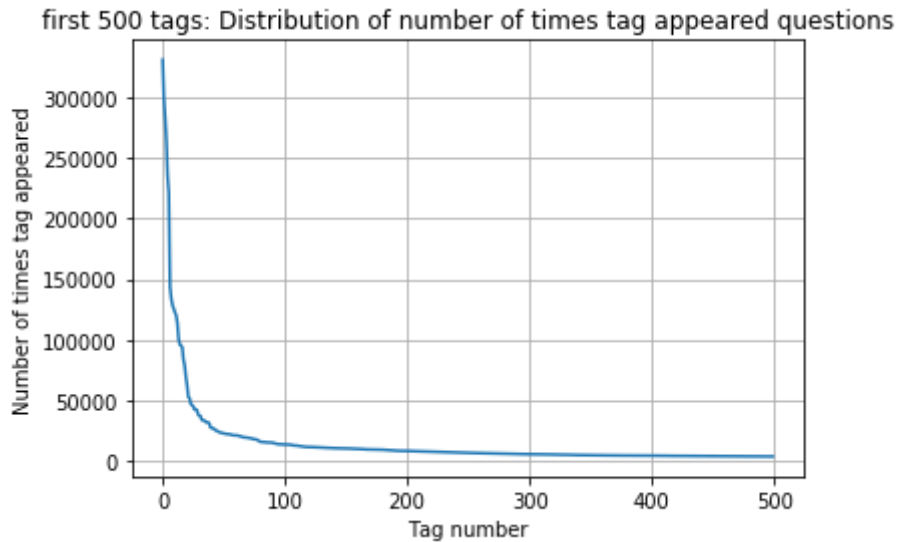
```
plt.plot(tag_counts[0:1000])
plt.title('first 1k tags: Distribution of number of times tag appeared questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



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

In [12]:

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



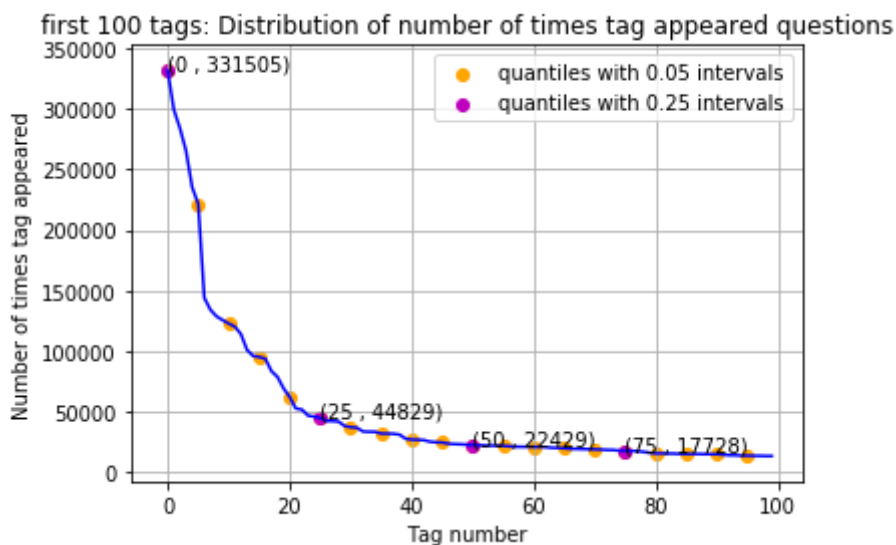
```
100 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291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

In [13]:

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

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

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



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

In [14]:

```
# 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 [15]:

```
#Storing the count of tag in each question in list 'tag_count'
tag_quest_count = tag_dtm.sum(axis=1).tolist()
#Converting each value in the 'tag_quest_count' to integer.
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]

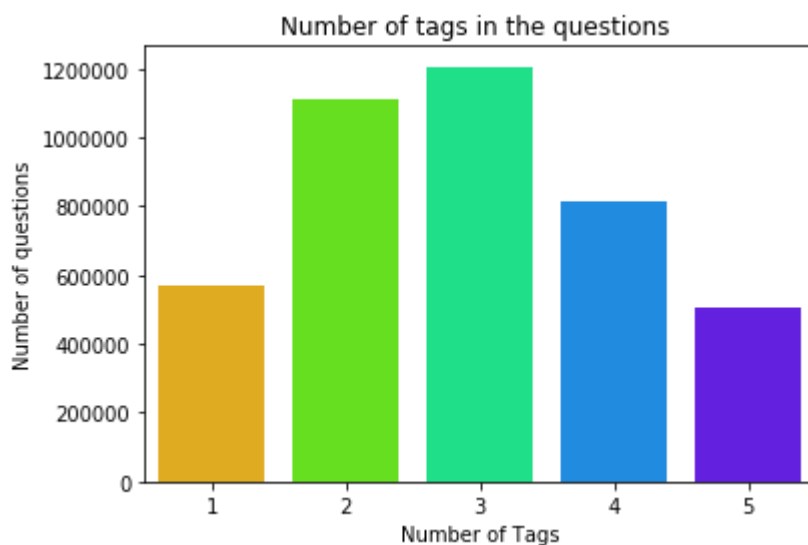
In [16]:

```
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 [17]:

```
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()
```



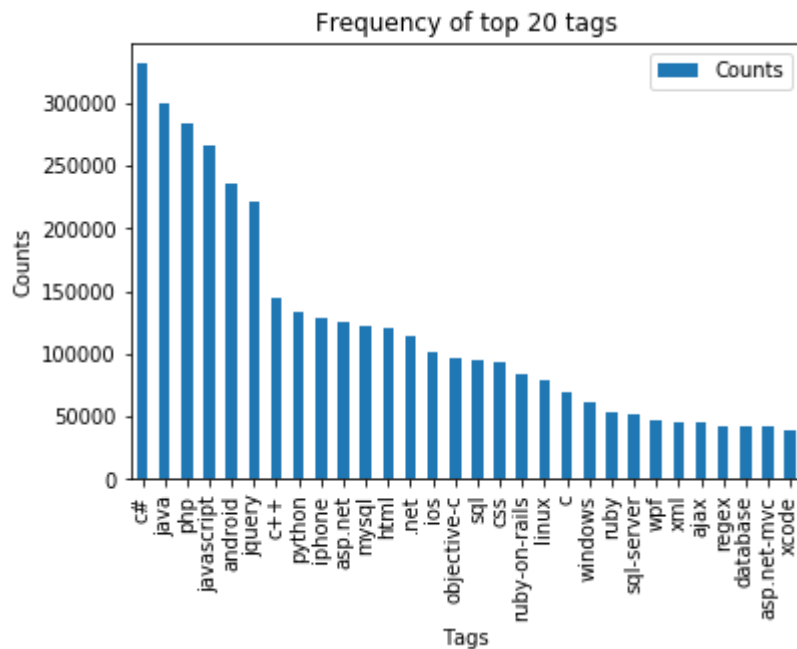
#### Observations:

1. Maximum number of tags per question: 5



In [19]:

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



#### Observations:

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

## 3.3 Cleaning and preprocessing of Questions

### 3.3.1 Preprocessing

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

In [20]:

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

In [3]:

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

    return None

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

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

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

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

Tables in the database:  
QuestionsProcessed

In [22]:

```
# http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/select-random-row-from-a-sqlite-table
start = datetime.now()
read_db = 'train_no_dup.db'
write_db = 'Processed.db'
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        reader = conn_r.cursor()
        reader.execute("SELECT Title, Body, Tags From no_dup_train ORDER BY RANDOM() LIMIT 100")

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

Tables in the databse:

QuestionsProcessed

Cleared All the rows

Time taken to run this cell : 0:11:48.555731

In [24]:

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

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\nisha\AppData\Roaming\nltk_data...
[nltk_data] Unzipping tokenizers\punkt.zip.
```

Out[24]:

True

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

In [25]:

```
#http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/
```

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

    is_code = 0

    title, question, tags = row[0], row[1], row[2]

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

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

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

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

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

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

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

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

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

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

```
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: 1171  
Avg. length of questions(Title+Body) after processing: 327  
Percent of questions containing code: 57  
Time taken to run this cell : 0:39:48.186001
```

In [26]:

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



In [27]:

```

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

```

Questions after preprocessed

```

=====
=====
('knapsack algorithm alway output thing code assum c capac amount item weigh
t item valu item thing knapsack algorithm tri code data set seem case reason
wonder knapsack algorithm taught dimension wherea dimension implement knapsa
ck algorithm variabl',)
-----
-----
('iad error ad inventori unavail ad iad app simul work well load devic adban
nerviewdeleg call bannerview didfailtoreceiveadwitherror descript error user
info think iad network setup correct add ad iad network automat send test ad
simul edit think latenc problem iad network server dispatch ad suppos networ
k ad app began work fine',)
-----
-----
('matlab time stamp read easi way read column matlab format date time ncurre
nt read code c pathc uigetfil txt select data c data file pathc c data dlmre
ad file needless say skip time stamp wonder aeasi way read time stamp plot c
olum agianst time hour file relat question format csv xls file mean except o
fcours xlsread',)
-----
-----
('get rid whitespac string php tri use trim str replac figur googl noth seem
work code thank edit apolog post output expect output output euro equal doll
ar nexpect output euro equal dollar get rid whitespac use money format funct
ion display proper',)
-----
-----
('import packag issu rmi two packag client server client packag contain inte
rfac object class main client server packag contain class impl object implem
ent interfac object final main server notic implement impl object put import
client object nokay precis put import client object client packag server lin
k client problem separ client server put server remot comput client home nca
n someone explain wrong',)
-----
-----
('remov duplic date given variabl one data point data given server need remo
v duplic date pick date cpu highest exampl date new row would go remov dupli
c date cpu check highest',)
-----
-----
('equival hold plot map want add differ color differ state us nhowev run com
mand sequenti get overlay plot updat map keep color expect plot grap acc wor

```

```
k map nso',)
```

```
( 'access file creat separ thread problem aith multithread copi access file s
ervic download unpack zip archiv copi file unzip folder right locat start se
par thread need access copi file code fragment processcopiedfil xdocument loa
d call fail except seem like file copi keep result file lock work synchronu
work without error nhave thought nthx',)
```

```
( 'multilin figur caption center center figur caption long singl line least m
ake text line line text slika current look like',)
```

In [28]:

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

In [29]:

```
preprocessed_data.head()
```

Out[29]:

	question	tags
0	aggreg function claus sqlite simpli put tabl a...	sql sqlite where-clause aggregate-functions
1	knapsack algorithm alway output thing code ass...	algorithm knapsack
2	iad error ad inventori unavail ad iad app simu...	iphone ios iad
3	matlab time stamp read easi way read column ma...	matlab timestamps
4	get rid whitespac string php tri use trim str ...	php

In [30]:

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

```
number of data points in sample : 999997
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 [31]:

```

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

```

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

In [7]:

```

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 [33]:

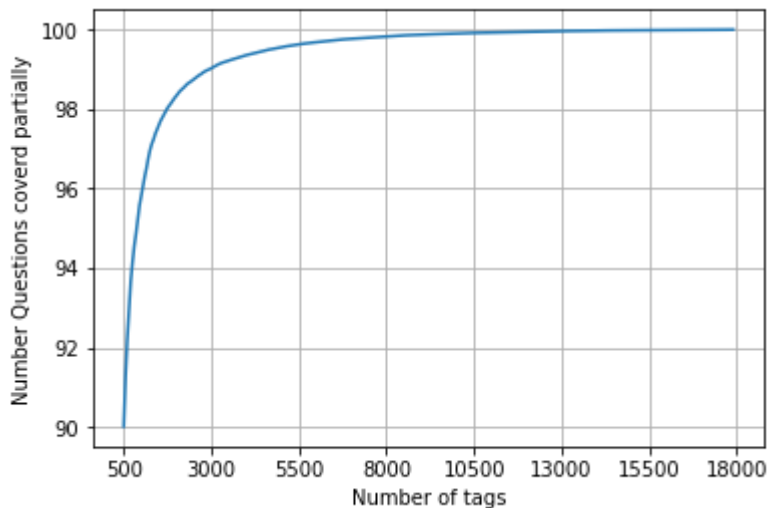
```

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

```

In [35]:

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



with 5500 tags we are covering 99.03 % of questions

In [36]:

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

number of questions that are not covered : 9704 out of 999997

In [37]:

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

Number of tags in sample : 35365

number of tags taken : 5500 ( 15.552099533437014 %)

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

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

In [38]:

```
total_size=preprocessed_data.shape[0]
train_size=int(0.80*total_size)

x_train=preprocessed_data.head(train_size)
x_test=preprocessed_data.tail(total_size - train_size)

y_train = multilabel_yx[0:train_size,:]
y_test = multilabel_yx[train_size:total_size,:]
```

In [39]:

```
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 : (799997, 5500)

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

## 4.3 Featurizing data

In [40]:

```
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=(1,1))
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 : 0:08:42.177370

In [41]:

```
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: (799997, 88195) Y : (799997, 5500)

Dimensions of test data X: (200000, 88195) Y: (200000, 5500)

In [0]:

```
# https://www.analyticsvidhya.com/blog/2017/08/introduction-to-multi-label-classification/
#https://stats.stackexchange.com/questions/117796/scikit-multi-label-classification
# classifier = LabelPowerset(GaussianNB())
"""
from skmultilearn.adapt import MLkNN
classifier = MLkNN(k=21)

# train
classifier.fit(x_train_multilabel, y_train)

# predict
predictions = classifier.predict(x_test_multilabel)
print(accuracy_score(y_test,predictions))
print(metrics.f1_score(y_test, predictions, average = 'macro'))
print(metrics.f1_score(y_test, predictions, average = 'micro'))
print(metrics.hamming_loss(y_test,predictions))

"""
# we are getting memory error because the multilearn package
# is trying to convert the data into dense matrix
# -----
#MemoryError                                Traceback (most recent call last)
#<ipython-input-170-f0e7c7f3e0be> in <module>()
#----> classifier.fit(x_train_multilabel, y_train)
```

Out[92]:

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

## 4.4 Applying Logistic Regression with OneVsRest Classifier

In [0]:

```
# this will be taking so much time try not to run it, download the lr_with_equal_weight.pkl
# This takes about 6-7 hours to run.
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

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

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	0.19	5534

In [0]:

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

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

In [42]:

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

Tables in the database:

QuestionsProcessed

In [43]:

```

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

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

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

```

Tables in the database:

QuestionsProcessed

Cleared All the rows

## 4.5.1 Preprocessing of questions

1. Separate Code from Body
2. Remove Special characters from Question title and description (not in code)
3. **Give more weightage to title : Add title three times to the question**
4. Remove stop words (Except 'C')
5. Remove HTML Tags
6. Convert all the characters into small letters
7. Use SnowballStemmer to stem the words



In [44]:

```

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

    is_code = 0

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

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

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

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

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

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

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

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

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

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

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

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

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

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

```

```
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 : 0:23:53.605409
```

In [45]:

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

**Sample quesitons after preprocessing of data**

In [46]:

```

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 data
grid bind silverlight bind datagrid dynam code wrote code debug code block s
eem bind correct grid come column form come grid column although necessari b
ind nthank repli advance..',)
-----
-----
('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid ja
va.lang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid java.l
ang.noclassdeffounderror javax servlet jsp tagext taglibraryvalid follow gui
d link instal jstl got follow error tri launch jsp page java.lang.noclassdef
founderror javax servlet jsp tagext taglibraryvalid taglib declar instal jst
l 1.1 tomcat webapp tri project work also tri version 1.2 jstl still messag
caus solv',)
-----
-----
('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor index ja
va.sql.sqlexcept microsoft odbc driver manag invalid descriptor index java.s
ql.sqlexcept microsoft odbc driver manag invalid descriptor index use follow
code display caus solv',)
-----
-----
('better way updat feed fb php sdk better way updat feed fb php sdk better w
ay updat feed fb php sdk novic facebook api read mani tutori still confused.
i find post feed api method like correct second way use curl someth like way
better',)
-----
-----
('btnadd click event open two window record ad btnadd click event open two w
indow record ad btnadd click event open two window record ad open window sea
rch.aspx use code hav add button search.aspx nwhen insert record btnadd clic
k event open anoth window nafter insert record close window',)
-----
-----
('sql inject issu prevent correct form submiss php sql inject issu prevent c
orrect form submiss php sql inject issu prevent correct form submiss php che
ck everyth think make sure input field safe type sql inject good news safe b
ad news one tag mess form submiss place even touch life figur exact html use
templat file forgiv okay entir php script get execut see data post none foru
m field post problem use someth titl field none data get post current use pr
int post see submit noth work flawless statement though also mention script
work flawless local machin use host come across problem state list input tes

```

```
t mess',)
```

```
-----
('countabl subaddit lebesgu measur countabl subaddit lebesgu measur countabl
subaddit lebesgu measur let lbrace rbrace sequenc set sigma -algebra mathcal
want show left bigcup right leq sum left right countabl addit measur defin s
et sigma algebra mathcal think use monoton properti somewher proof start app
reci littl help nthank ad han answer make follow addit construct given han a
nswer clear bigcup bigcup cap emptyset neq left bigcup right left bigcup rig
ht sum left right also construct subset monoton left right leq left right fi
nal would sum leq sum result follow',)
-----
```

```
-----
('hql equival sql queri hql equival sql queri hql equival sql queri hql quer
i replac name class properti name error occur hql error',)
-----
```

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

## Saving Preprocessed data to a Database

In [4]:

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

In [5]:

```
preprocessed_data.head()
```

Out[5]:

	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 [6]:

```
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 [50]:

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

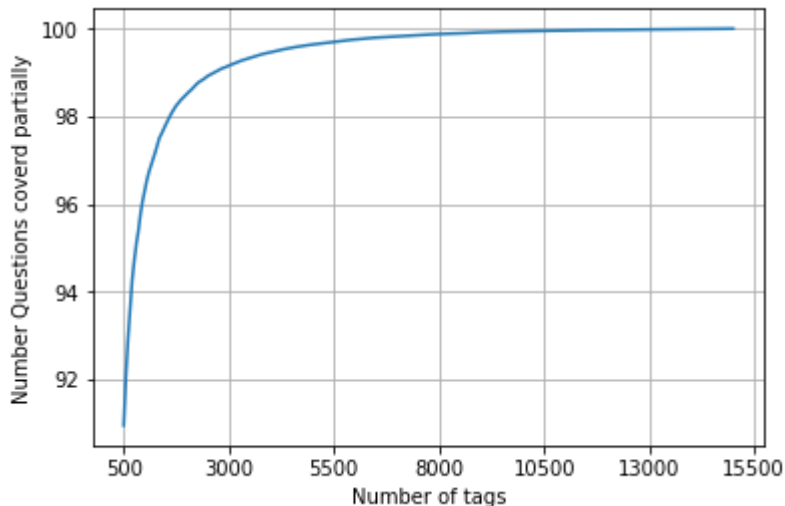
## Selecting 500 Tags

In [51]:

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

In [52]:

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



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

In [55]:

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

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

In [56]:

```
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 [57]:

```
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 [0]:

```
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=(1,1))
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 : 0:03:52.522389

In [0]:

```
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)
```

Diamensions of train data X: (400000, 94927) Y : (400000, 500)

Diamensions of test data X: (100000, 94927) Y: (100000, 500)

## 4.5.3 Applying Logistic Regression with OneVsRest Classifier

In [0]:

```

start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.00001, penalty='l1'), n_
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 : 0.23623  
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	0.33	0.45	3000

In [0]:

```

joblib.dump(classifier, 'lr_with_more_title_weight.pkl')

```

Out[113]:

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



In [0]:

```

start = datetime.now()
classifier_2 = OneVsRestClassifier(LogisticRegression(penalty='l1'), n_jobs=-1)
classifier_2.fit(x_train_multilabel, y_train)
predictions_2 = classifier_2.predict(x_test_multilabel)
print("Accuracy :",metrics.accuracy_score(y_test, predictions_2))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions_2))

precision = precision_score(y_test, predictions_2, average='micro')
recall = recall_score(y_test, predictions_2, average='micro')
f1 = f1_score(y_test, predictions_2, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))

precision = precision_score(y_test, predictions_2, average='macro')
recall = recall_score(y_test, predictions_2, average='macro')
f1 = f1_score(y_test, predictions_2, average='macro')

print("Macro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))

print(metrics.classification_report(y_test, predictions_2))
print("Time taken to run this cell :", datetime.now() - start)

```

```

Accuracy : 0.25108
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.

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)

**Note : Its taking more than 2 days to run on 0.5 million data points so im taking**

## only 100k data points

In [8]:

```
preprocessed_data.head()
```

Out[8]:

	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 [13]:

```
preprocessed_data_100k_data = preprocessed_data[:100000]
print("number of data points in sample :", preprocessed_data_100k_data.shape[0])
print("number of dimensions :", preprocessed_data_100k_data.shape[1])
```

number of data points in sample : 100000

number of dimensions : 2

## Converting string Tags to multilable output variables

In [14]:

```
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data_100k_data['tags'])
```

In [15]:

```
# we will be taking 500 tags
multilabel_yx = tags_to_choose(500)
```

In [17]:

```
multilabel_yx.shape
```

Out[17]:

(100000, 500)

## Splitting data into train and test into 70:30

In [21]:

```
train_datasize = 70000
x_train=preprocessed_data.head(train_datasize)
x_test=preprocessed_data.tail(preprocessed_data_100k_data.shape[0] - train_datasize)

y_train = multilabel_yx[0:train_datasize,:]
y_test = multilabel_yx[train_datasize:preprocessed_data_100k_data.shape[0],:]
```

In [22]:

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

Number of data points in train data : (70000, 500)

Number of data points in test data : (30000, 500)

## 5.1 Featurizing data BOW(upto 4 gram)

In [23]:

```
start = datetime.now()
vectorizer = CountVectorizer(min_df=0.00009,tokenizer = lambda x: x.split(), ngram_range=(1,4))
x_train_bow = vectorizer.fit_transform(x_train['question'])
x_test_bow = vectorizer.transform(x_test['question'])
print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:01:54.977393

### 5.1.1 Saving Bow vectorized data

In [24]:

```
from sklearn.externals import joblib
# data points with 0.5 million data
joblib.dump(x_train_bow, 'x_train_BOW_100k.pkl')
joblib.dump(x_test_bow, 'x_test_BOW_100k.pkl')

# target class i.e multilabel classes with 0.5 million
joblib.dump(y_train, 'y_train_100k.pkl')
joblib.dump(y_test, 'y_test_100k.pkl')
```

Out[24]:

```
['y_test_100k.pkl']
```

### 5.1.2 Loading saved Bow vectorized data

In [3]:

```
from sklearn.externals import joblib
x_train_bow = joblib.load('x_train_BOW_100k.pkl')
x_test_bow = joblib.load('x_test_BOW_100k.pkl')
y_train = joblib.load('y_train_100k.pkl')
y_test = joblib.load('y_test_100k.pkl')
```

In [27]:

```
print(x_train_bow.shape)
print(x_test_bow.shape)
print(y_train.shape)
print(y_test.shape)
```

(70000, 25000)

(30000, 25000)

(70000, 500)

(30000, 500)

## 5.2 Logistic Regression with OneVsRest Classifier Optimized using GridSearchcv

When you use nested estimators with grid search you can scope the parameters with `__` as a separator. In this case the logistic reg model is stored as an attribute named estimator inside the OneVsRestClassifier model:

In [28]:

```
# this will be taking so much time try not to run it
import warnings
warnings.filterwarnings("ignore")

from sklearn.model_selection import GridSearchCV
tuned_parameters = [{'estimator__C': [100, 10, 1, 0.1, 0.01, 0.001, 0.0001]}]
tuned_parameters

#=====#
# Find Optimal C by grid search

log_reg_clf = OneVsRestClassifier(LogisticRegression())
logistic_gs = GridSearchCV(log_reg_clf, tuned_parameters, scoring = 'f1_micro', cv=2, n_jobs=

start = datetime.now()

logistic_gs.fit(x_train_bow, y_train)
print('Time to train',datetime.now()-start)
```

Time to train 3:03:20.956959

In [31]:

```
# predictions =logistic_gs.predict(x_test_bow)
logistic_gs
```

Out[31]:

```
GridSearchCV(cv=2, error_score='raise',
             estimator=OneVsRestClassifier(estimator=LogisticRegression(C=1.0, cla
ss_weight=None, dual=False, fit_intercept=True,
             intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,
             penalty='l2', random_state=None, solver='liblinear', tol=0.0001,
             verbose=0, warm_start=False),
             n_jobs=1),
             fit_params=None, iid=True, n_jobs=-1,
             param_grid=[{'estimator__C': [100, 10, 1, 0.1, 0.01, 0.001, 0.000
1]}],
             pre_dispatch='2*n_jobs', refit=True, return_train_score='warn',
             scoring='f1_micro', verbose=0)
```

### 5.2.1 OneVsRestClassifier with Logistic regression(OvR)

In [33]:

```

logistic_reg_optimal_c = 1.0
start = datetime.now()

classifier_2 = OneVsRestClassifier(LogisticRegression(C=logistic_reg_optimal_c ,penalty='l1'))
classifier_2.fit(x_train_bow, y_train)

predictions_2 = classifier_2.predict(x_test_bow)
print("Accuracy :",metrics.accuracy_score(y_test, predictions_2))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions_2))

precision = precision_score(y_test, predictions_2, average='micro')
recall = recall_score(y_test, predictions_2, average='micro')
f1 = f1_score(y_test, predictions_2, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))

precision = precision_score(y_test, predictions_2, average='macro')
recall = recall_score(y_test, predictions_2, average='macro')
f1 = f1_score(y_test, predictions_2, average='macro')

print("Macro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))

print(metrics.classification_report(y_test, predictions_2))
print("Time taken to run this cell :", datetime.now() - start)

```

Accuracy : 0.02466666666666667

Hamming loss 0.006170666666666666

Micro-average quality numbers

Precision: 0.0252, Recall: 0.0142, F1-measure: 0.0182

Macro-average quality numbers

Precision: 0.0035, Recall: 0.0022, F1-measure: 0.0022

	precision	recall	f1-score	support
0	0.21	0.02	0.03	6668
1	0.12	0.05	0.07	3659
2	0.03	0.01	0.02	971
3	0.08	0.00	0.01	1506
4	0.06	0.05	0.05	1649
5	0.03	0.03	0.03	1113
6	0.04	0.03	0.03	1482
7	0.04	0.06	0.05	980
8	0.06	0.05	0.05	1520
9	0.08	0.01	0.01	1041
10	0.03	0.02	0.02	861
11	0.00	0.00	0.00	336

## 5.3 Linear-SVM (SGDClassifier with loss-hinge) with OneVsRest Classifier Optimized using GridSearchcv

In [34]:

```
# this will be taking so much time try not to run it
# This takes about 20 hours to run.

from sklearn.model_selection import GridSearchCV
tuned_parameters = [{'estimator__alpha': [100, 10, 1, 0.1, 0.01, 0.001, 0.0001]]}
tuned_parameters

#=====#
# Find Optimal C by grid search
lr_svm_clf = OneVsRestClassifier(SGDClassifier(loss='hinge',n_jobs=-1))
linear_svm_gs = GridSearchCV(lr_svm_clf, tuned_parameters,scoring = 'f1_micro', cv=2)

start = datetime.now()

linear_svm_gs.fit(x_train_bow, y_train)
print(linear_svm_gs.best_estimator_)

# predictions =linear_svm_gs.predict(x_test_bow)

print('Time to train',datetime.now()-start)
```

```
OneVsRestClassifier(estimator=SGDClassifier(alpha=0.001, average=False, clas
s_weight=None, epsilon=0.1,
    eta0=0.0, fit_intercept=True, l1_ratio=0.15,
    learning_rate='optimal', loss='hinge', max_iter=None, n_iter=None,
    n_jobs=-1, penalty='l2', power_t=0.5, random_state=None,
    shuffle=True, tol=None, verbose=0, warm_start=False),
    n_jobs=1)
Time to train 0:15:36.758155
```

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

In [36]:

```

lr_svm_optimal_alpha = 0.001

start = datetime.now()
classifier2 = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=lr_svm_optimal_alpha, p
classifier2.fit(x_train_bow, y_train)
predictions2 = classifier2.predict (x_test_bow)

print("Accuracy :",metrics.accuracy_score(y_test, predictions2))
print("Hamming loss ",metrics.hamming_loss(y_test,predictions2))

precision = precision_score(y_test, predictions2, average='micro')
recall = recall_score(y_test, predictions2, average='micro')
f1 = f1_score(y_test, predictions2, average='micro')

print("Micro-average quality numbers")
print("Precision: {:.4f}, Recall: {:.4f}, F1-measure: {:.4f}".format(precision, recall, f1))

precision = precision_score(y_test, predictions2, average='macro')
recall = recall_score(y_test, predictions2, average='macro')
f1 = f1_score(y_test, predictions2, 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, predictions2))
print("Time taken to run this cell :", datetime.now() - start)

```

Accuracy : 0.0316

Hamming loss 0.005414733333333334

Micro-average quality numbers

Precision: 0.0266, Recall: 0.0098, F1-measure: 0.0143

Macro-average quality numbers

Precision: 0.0034, Recall: 0.0017, F1-measure: 0.0017

	precision	recall	f1-score	support
0	0.21	0.01	0.02	6668
1	0.12	0.02	0.04	3659
2	0.03	0.01	0.01	971
3	0.13	0.01	0.02	1506
4	0.06	0.04	0.05	1649
5	0.04	0.03	0.03	1113
6	0.04	0.02	0.02	1482
7	0.04	0.05	0.05	980
8	0.05	0.05	0.05	1520
9	0.10	0.01	0.01	1041
10	0.02	0.01	0.01	861
11	0.00	0.00	0.00	335



In [11]:

```

from prettytable import PrettyTable
print("TF-IDF with 0.5 million dataset")
x = PrettyTable()
x.field_names = ["Model", "Vectorizer", "Accuracy", "Hamming loss", "Precision", "Recall", "Mi

x.add_row(["SDG with loss - log ", 'TF-IDF ', 0.23623, 0.0027, 0.7216, 0.3256, 0.4488])
x.add_row(["LogisticRegression", 'TF-IDF ', 0.25108, 0.00270302, 0.7172, 0.3672, 0.4858])
print(x)

print("\nBOW with tuned hyperparameter with 100k dataset")
x = PrettyTable()
x.field_names = ["Model", "Vectorizer", "Accuracy", "Hamming loss", "Precision", "Recall", "Mi

x.add_row(["Logistic Regression", 'BOW', 0.02466, 0.0061, 0.0252, 0.0142, 0.0182])
x.add_row(["Linear SVM", 'BOW', 0.0316, 0.0054, 0.0266, 0.0098, 0.0143])
print(x)

```

TF-IDF with 0.5 million dataset

Model	Vectorizer	Accuracy	Hamming loss	Precision	Recall	Micro f1
SDG with loss - log	TF-IDF	0.23623	0.0027	0.7216	0.3256	0.4488
LogisticRegression	TF-IDF	0.25108	0.00270302	0.7172	0.3672	0.4858

BOW with tuned hyperparameter with 100k dataset

Model	Vectorizer	Accuracy	Hamming loss	Precision	Recall	Micro f1
Logistic Regression	BOW	0.02466	0.0061	0.0252	0.0142	0.0182
Linear SVM	BOW	0.0316	0.0054	0.0266	0.0098	0.0143

## Conclusion

- According to our problem statement we have to Suggest the tags based on the content that was there in the question posted on Stackoverflow. lets start with step-by-step process to solve this problem.
- lets start -->

- As we know we have dataset which contains 6,034,195 rows. The columns in the table are: ID, TITLE, BODY and TAGS(class of our dataset) which basically decides that question belongs to which class as in this we have multiple tags related/belongs to one question we are treating this as multilable classes for detailed please go above i.e 3.2 how we are doing this.

2. Now lets start with EDA of given dataset so that we will able to analyse deeply about the data so that we will know which features to use which to not and some preprocessing and cleanning of data before that like ckecking for the presence of duplicate rows now of unique questions no of tags present and lots more. And after that we will try to analyse our class lables that is with column name tags as we can see each question has one or more than tags so and this is not the 2 class classification and each question can be tagged to multiple classes at a time so we decide to treat as multilable classes means one question can belongs to one or more than one tags at a time, but before that we will do some EDA on our tag like total no of unique tags present, avg no of tags present to each questions max no of tags and min no of tags belongs to each questions and which tags appears how much time for the questions so that we will able to know which tags are most imp and which are usefull for us.
3. As during EDA we can see that total no if tags present is to much i.e is roughly 42048 see in (3.2.1) that is to many so we will try to plot the Distribution of number of times tag appeared questions and try to take the top Most Frequent Tags. observation see
4. Now after doing lots of EDA on our data set we will clean and preprocess our data set to remove duplicates and unwanted elements like html tags and lots more
5. Before start with our Machine learning models we will Converting tags for multilabel problems using CountVectorizer and w will sample the number of tags instead considering all of them (due to limitation of computing power) and as you can see in with 5500 tags we are covering 99.03 % of questions so we will choose top 5500 tags.
6. And most important we will Featurize our dataset with more weight to tile in our datssset and then we are ready to apply Machine learning models on it i.e Logistic and linear models as we know our dataset is high dim data set and we know linear models woks better on high dim dataset so we use logistic and linear svn( and insted of linear svm we will try sgd and as we know sgd is much more faster and efficient than linear svm) and then try to train our model on it.And most imp our performance matrix we will use Micro f1 score
7. In next step we will try our models with other vectorizer i.e bag of words upto 4 grams and try to do some hyperparameter tuning in order to improve the model performance

In [ ]: