



```
In [1]: import warnings
warnings.filterwarnings("ignore")
import pandas as pd
import sqlite3
import csv
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
from wordcloud import WordCloud
import re
import os
from sqlalchemy import create_engine # database connection
import datetime as dt
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from nltk.stem.snowball import SnowballStemmer
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.multiclass import OneVsRestClassifier
from sklearn.linear_model import SGDClassifier
from sklearn import metrics
from sklearn.metrics import f1_score, precision_score, recall_score
from sklearn import svm
from sklearn.linear_model import LogisticRegression
from skmultilearn.adapt import mlknn
from skmultilearn.problem_transform import ClassifierChain
from skmultilearn.problem_transform import BinaryRelevance
from skmultilearn.problem_transform import LabelPowerset
from sklearn.naive_bayes import GaussianNB
from datetime import datetime
```

Stack Overflow: Tag Prediction

1. Business Problem

1.1 Description

Description

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

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

Problem Statement

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

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

1.2 Source / useful links

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

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

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

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

1.3 Real World / Business Objectives and Constraints

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

2. Machine Learning problem

2.1 Data

2.1.1 Data Overview

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

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

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

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

Size of Train.csv - 6.75GB

Size of Test.csv - 2GB

Number of rows in Train.csv = 6034195

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

Data Field Explanation

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

Id - Unique identifier for each question

Title - The question's title

Body - The body of the question

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

2.1.2 Example Data point

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

Body :

```

#include<
iostream>\n
#include<
stdlib.h>\n\n
using namespace std;\n\n
int main()\n
{\n
    int n,a[n],x,c,u[n],m[n],e[n][4];\n

    cout<<"Enter the number of variables";\n
    cin>>n;\n\n
    cout<<"Enter the Lower, and Upper Limits
of the variables";\n
    for(int y=1; y<n+1; y++)\n
    {\n
        cin>>m[y];\n
        cin>>u[y];\n
    }\n
    for(x=1; x<n+1; x++)\n
    {\n
        a[x] = (m[x] + u[x])/2;\n
    }\n
    c=(n*4)-4;\n
    for(int 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
    system("PAUSE");\n
    return 0;    \n
}\n

```

\n\n

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

1

50

50\n

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

\n\n

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

1,100\n

1,100\n

1,100\n

(could be varied too)

\n\n

The output is not coming,can anyone correct the code or tell me

what\'s wrong?

\n'

Tags : 'c++ c'

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

2.2.1 Type of Machine Learning Problem

It is a multi-label classification problem

Multi-label Classification: Multilabel classification assigns to each sample a set of target labels. This can be thought as predicting properties of a data-point that are not mutually exclusive, such as topics that are relevant for a document. A question on Stackoverflow might be about any of C, Pointers, FileIO and/or memory-management at the same time or none of these.

__Credit__: <http://scikit-learn.org/stable/modules/multiclass.html>

2.2.2 Performance metric

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

$$F1 = 2 * (precision * recall) / (precision + recall)$$

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

'Micro f1 score':

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

'Macro f1 score':

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

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

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>

3. Exploratory Data Analysis

3.1 Data Loading and Cleaning

3.1.1 Using Pandas with SQLite to Load the data

```
In [2]: #Creating db file from csv
#Learn SQL: https://www.w3schools.com/sql/default.asp
if not os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\facebook-recruiting-iii-keyword-extraction\\Train\\train_S0.db'):
    start = datetime.now()
    disk_engine = create_engine('sqlite:///G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\facebook-recruiting-iii-keyword-extraction\\Train\\train_S0.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\facebook-recruiting-iii-keyword-extraction\\Train\\Train.csv', names=['Id', 'Title', 'Body', 'Tags'], chunksize=chunksize, iterator=True, encoding='utf-8', ):
        df.index += index_start
        j+=1
        print('{} rows'.format(j*chunksize))
        df.to_sql('data_S0', disk_engine, if_exists='append')
        index_start = df.index[-1] + 1
    print("Time taken to run this cell :", datetime.now() - start)
```

3.1.2 Counting the number of rows

```

In [3]: '''
if os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stack
overflow tag predictor\\assignment\\facebook-recruiting-iii-keyword-ext
raction\\Train\\train_S0.db'):
    start = datetime.now()
    con = sqlite3.connect('G:\\machine_learning\\case_study\\Case study
5 Stackoverflow tag predictor\\assignment\\facebook-recruiting-iii-key
word-extraction\\Train\\train_S0.db')
    num_rows = pd.read_sql_query("""SELECT count(*) FROM data_S0 LIMIT
20000""", 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 abov
e cell to generate train.db file")

'''

conn_r = sqlite3.connect(r'G:\\machine_learning\\case_study\\Case study
5 Stackoverflow tag predictor\\assignment\\facebook-recruiting-iii-key
word-extraction\\Train\\train_S0.db')
num_rows = pd.read_sql_query("""SELECT * FROM data_S0 LIMIT 20000""", c
onn_r)
conn_r.commit()
conn_r.close()

```

3.1.3 Checking for duplicates

```

In [4]: #Learn SQL: https://www.w3schools.com/sql/default.asp
if os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stack
overflow tag predictor\\assignment\\facebook-recruiting-iii-keyword-ext
raction\\Train\\train_S0.db'):
    start = datetime.now()
    con = sqlite3.connect('G:\\machine_learning\\case_study\\Case study

```

Time taken to run this cell : 0:13:04.340158

Out[5]:

```
In [6]: # number of times each question appeared in our database
df.no_dup.cnt.dup.value_counts()
```

PDFCROWD

```
3      1464
Name: cnt_dup, dtype: int64
```

```
In [7]: #Creating a new database with no duplicates
if not os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\train_no_dup.db'):
    disk_dup = create_engine("sqlite:///G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\train_no_dup.db")
    no_dup = pd.DataFrame(df_no_dup, columns=['Title', 'Body', 'Tags'])
    no_dup.to_sql('no_dup_train', disk_dup)
```

```
In [8]: #This method seems more appropriate to work with this much data.
#creating the connection with database file.
if os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\train_no_dup.db'):
    start = datetime.now()
    con = sqlite3.connect('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\train_no_dup.db')
    tag_data = pd.read_sql_query("""SELECT Tags FROM no_dup_train""", con)
    #Always remember to close the database
    con.close()

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

Time taken to run this cell : 0:01:51.718765

3.2 Analysis of Tags

3.2.1 Total number of unique tags

```
In [9]: # 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 [10]: 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 [11]: #'get_feature_name()' gives us the vocabulary.
        tags = vectorizer.get_feature_names()
        #Lets look at the tags we have.
        print("Some of the tags we have :", tags[:10])
```

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

3.2.3 Number of times a tag appeared

```
In [12]: # 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))
```

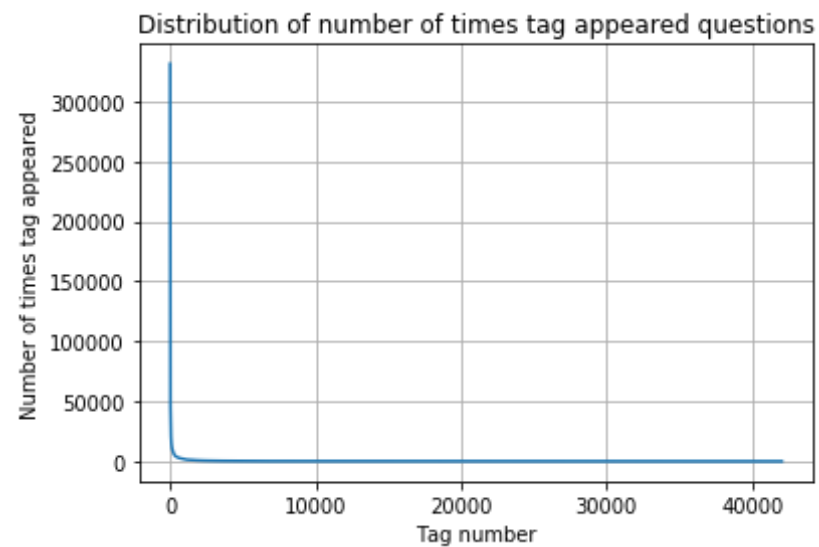
```
In [13]: #Saving this dictionary to csv files.
if not os.path.isfile('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\tag_counts_dict_dtm.csv'):
    with open('G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\tag_counts_dict_dtm.csv', 'w') as csv_file:
        writer = csv.writer(csv_file)
        for key, value in result.items():
            writer.writerow([key, value])
tag_df = pd.read_csv("G:\\machine_learning\\case_study\\Case study 5 Stackoverflow tag predictor\\assignment\\tag_counts_dict_dtm.csv", names=['Tags', 'Counts'])
tag_df.head()
```

Out[13]:

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

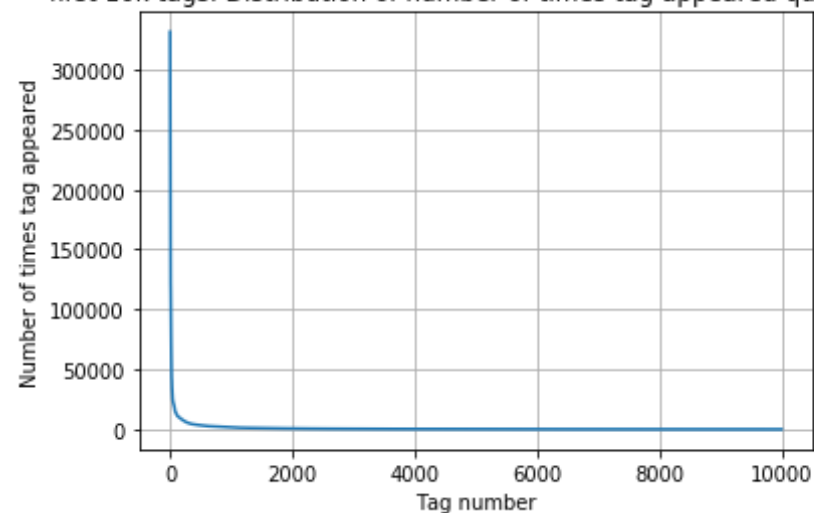
```
In [14]: tag_df_sorted = tag_df.sort_values(['Counts'], ascending=False)
tag_counts = tag_df_sorted['Counts'].values
```

```
In [15]: 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 [16]: 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])
```

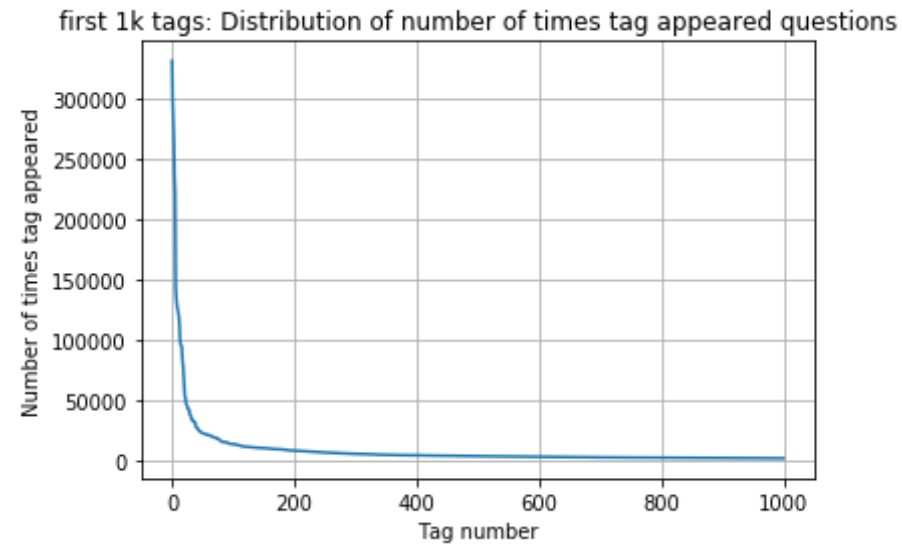
first 10k tags: Distribution of number of times tag appeared questions



400	[331505	44829	22429	17728	13364	11162	10029	9148	8054	7
151										
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]

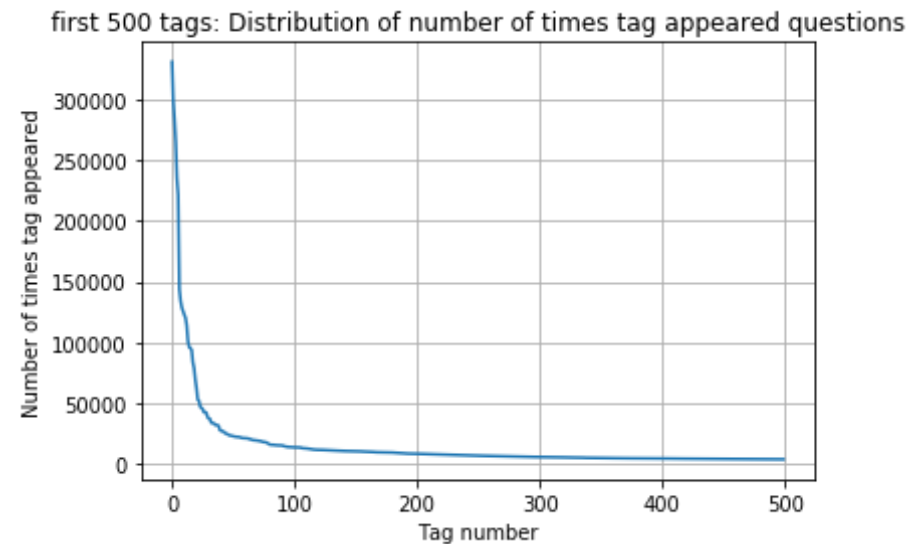
```
In [17]: plt.plot(tag_counts[0:1000])
plt.title('first 1k tags: Distribution of number of times tag appeared
questions')
plt.grid()
plt.xlabel("Tag number")
plt.ylabel("Number of times tag appeared")
plt.show()
print(len(tag_counts[0:1000:5]), tag_counts[0:1000:5])
```



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 [18]: 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])
```

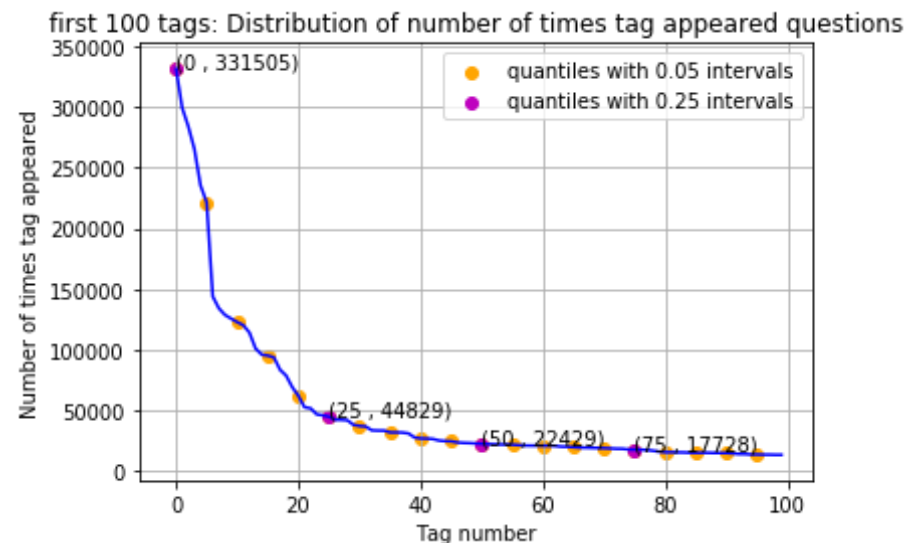


```
100 [331505 221533 122769 95160 62023 44829 37170 31897 26925 24
537
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703
13364 13157 12407 11658 11228 11162 10863 10600 10350 10224
10029 9884 9719 9411 9252 9148 9040 8617 8361 8163
8054 7867 7702 7564 7274 7151 7052 6847 6656 6553
6466 6291 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 [19]: plt.plot(tag_counts[0:100], c='b')
plt.scatter(x=list(range(0,100,5)), y=tag_counts[0:100:5], c='orange',
label="quantiles with 0.05 intervals")
# quantiles with 0.25 difference
plt.scatter(x=list(range(0,100,25)), y=tag_counts[0:100:25], c='m', lab
el = "quantiles with 0.25 intervals")

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

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



```
20 [331505 221533 122769 95160 62023 44829 37170 31897 26925 245
37
22429 21820 20957 19758 18905 17728 15533 15097 14884 13703]
```

```
In [20]: # 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 [21]: #Storing the count of tag in each question in list 'tag_count'
tag_quest_count = tag_dtm.sum(axis=1).tolist()
#Converting list of lists into single list, we will get [[3], [4], [2], [2], [3]] and we are converting this to [3, 4, 2, 2, 3]
tag_quest_count=[int(j) for i in tag_quest_count for j in i]
print ('We have total {} datapoints.'.format(len(tag_quest_count)))

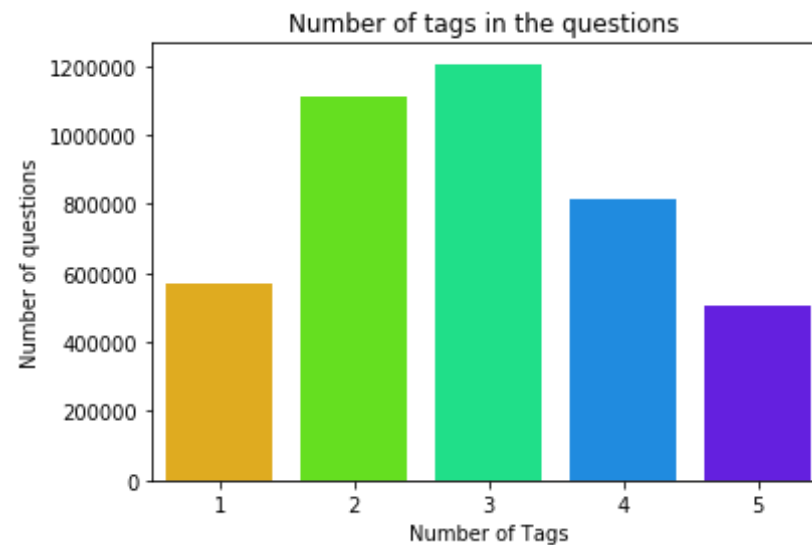
print(tag_quest_count[:5])
```

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

```
In [22]: 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 [23]: sns.countplot(tag_quest_count, palette='gist_rainbow')
plt.title("Number of tags in the questions ")
plt.xlabel("Number of Tags")
plt.ylabel("Number of questions")
plt.show()
```



Observations:

1. Maximum number of tags per question: 5
2. Minimum number of tags per question: 1
3. Avg. number of tags per question: 2.899

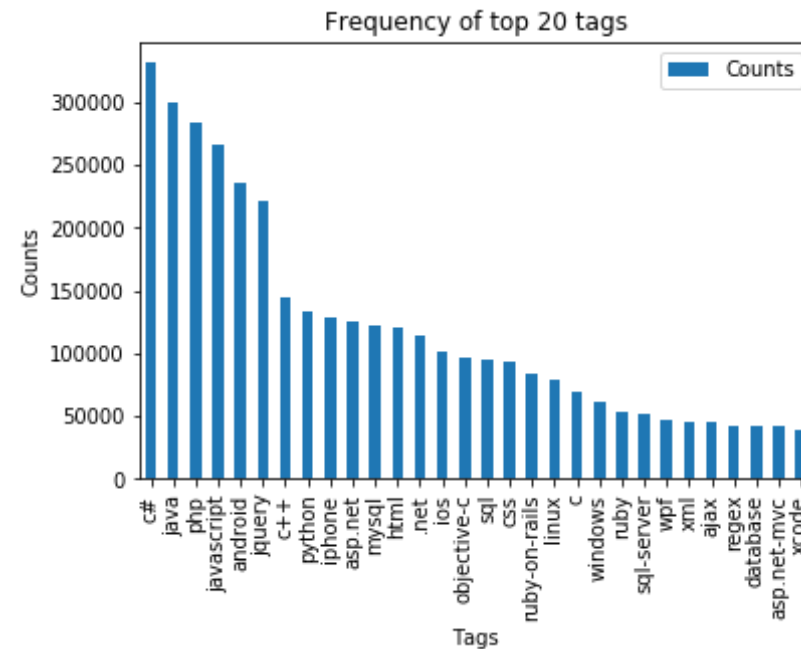
4. Most of the questions are having 2 or 3 tags

3.2.5 Most Frequent Tags

```
In [24]: # Plotting word cloud
start = datetime.now()

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

fig = plt.figure(figsize=(30,20))
plt.imshow(wordcloud)
plt.axis('off')
plt.tight_layout(pad=0)
fig.savefig("tag.png")
plt.show()
print("Time taken to run this cell :", datetime.now() - start)
```

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

```
In [27]: #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 database:")
    tables = table_names.fetchall()
    print(tables[0][0])
    return(len(tables))

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

sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (qu
estion text NOT NULL, code text, tags text, words_pre integer, words_po
st integer, is_code integer);"""
create_database_table("G:\\machine_learning\\case_study\\Case study 5 S
tackoverflow tag predictor\\assignment\\Processed.db", sql_create_table
)

```

Tables in the database:
QuestionsProcessed

```

In [28]: # http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/select-random-row-from-a-
sqlite-table
start = datetime.now()
read_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflow
tag predictor\\assignment\\train_no_dup.db'
write_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflo
w tag predictor\\assignment\\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 ORDE
R BY RANDOM() LIMIT 20000;")

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 database:

QuestionsProcessed

Cleared All the rows

Time taken to run this cell : 0:02:06.989444

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

In [29]: [#http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/](http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/)

```

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

    is_code = 0

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

```

```

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

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

question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTIL
INE|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 s
top_words and (len(j)!=1 or j=='c'))

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

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

print( "Avg. length of questions(Title+Body) before processing: %d"%no_
dup_avg_len_pre)
print( "Avg. length of questions(Title+Body) after processing: %d"%no_d

```

```

up_avg_len_post)
print ("Percent of questions containing code: %d"%((questions_with_code
*100.0)/questions_proccesed))

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

```

Avg. length of questions(Title+Body) before processing: 1177
 Avg. length of questions(Title+Body) after processing: 330
 Percent of questions containing code: 57
 Time taken to run this cell : 0:00:48.267339

In [30]:

```

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

```

In [31]:

```

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 1
0")

        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

=====

('angular access bound element may issu way architect applic keep run n
 eed abl access dom element via item array bound html basic anoth piec c
 ode want run scope item make chang div posit base div height',)

('insert post data form mysql dynam form dynam name field post php file
want take field post php file current want insert data mysql tabl forma
t issu intrest post data valu wp http refer nbut submit use refer point
part stuggl would get part post data dynam would get requir array key i
nsert ntabl dynam mani thank',)

('preload class librari dot net compact framework let say load form sho
w button let access form click button load form see take bit time load
class librari dll make ui look unrespons see form class librari load ma
ke form still shown librari load pretti fast pre load class librari per
hap applic start tri put statement form compact framework good wont loa
d actual requir form talk second delay still look bad anyway overcom pr
oblem cant forc cf load dll file load system dll system window form dll
etc updat could load class librari use still unabl load follow file for
m load even possibl',)

('jqueryi mobil button row controlgroup html page footer want today refr
esh button line control group button left right align controlgroup cent
er easi way jqueryi mobil',)

('track usag column tabl microsoft sql server collect old horribl desig
n databas given green light tear restructur howev databas normal numer
field empti year purpos without document slew legaci applic public webs
it use various piec data one hous idea could determin field tabl use wa
y sql server use third parti tool need see histori usag set addit log d
etermin usag usag ideal mean last updat insert often includ select stat
ement addit migrat sql solut use either server type would work abl brin
g dbs complianc',)

('differ applicationpool ident anonymi ident processmodel ident imperso
n ident thread ident littl bit confus various differ ident asp net ii t
ell exact differ applic pool ident anonymi ident processmodel ident imp
erson ident thread ident pleas also tell asp net use ident',)


```

-----
('pcie interrupt rout current implement pcie endpoint devic xilinx pfga
problem regard interrupt driver init map interrupt irq howev interrupt
fire irq seem rout anoth pin irq use pcie msi interrupt could caus prob
lem',)
-----
-----
('would alter queryset filter handl user input white space django view
search databas name includ text user submit form use search check user
queri field record work problem user enter full name search box whitesp
ac eg john smith instead john smith function would return result quit n
ew sure would go chang function even form could lazi prevent enter spac
e key presum possibl someth like learn actual solut problem form view p
retti simpl complet side question realis use text editor line wrap forg
ot still know safe add thing line line break python indent matter apolo
g scroll code',)
-----
-----
('open sourc android librari reusabl view viewgroup adapt etc realli wa
nt site collect usabl compon android found various small list biggest o
pen intent librari list mark murphi hi mark also list librari project p
ublish none order say cocoa control cocoa object go grab usabl librari
android edit librari project though would ideal flip side publish open
sourc librari project github',)
-----
-----

```

```

In [32]: #Taking 1 Million entries to a dataframe.
write_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflo
w tag predictor\\assignment\\Processed.db'
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        preprocessed_data = pd.read_sql_query("""SELECT question, Tags
FROM QuestionsProcessed LIMIT 20000""", conn_r)
    conn_r.commit()
    conn_r.close()

```

```
In [33]: preprocessed_data.head()
```

Out[33]:

	question	tags
0	menu slug use creat option page follow page de...	plugin-development admin-menu settings
1	angular access bound element may issu way arch...	javascript angularjs
2	insert post data form mysql dynam form dynam n...	php mysql
3	preload class librari dot net compact framewor...	c# dll windows-mobile compact-framework load
4	jqueri mobil button row controlgroup html page...	jquery jquery-mobile

```
In [34]: 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 : 19999
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 [35]: # binary='true' will give a binary vectorizer
vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

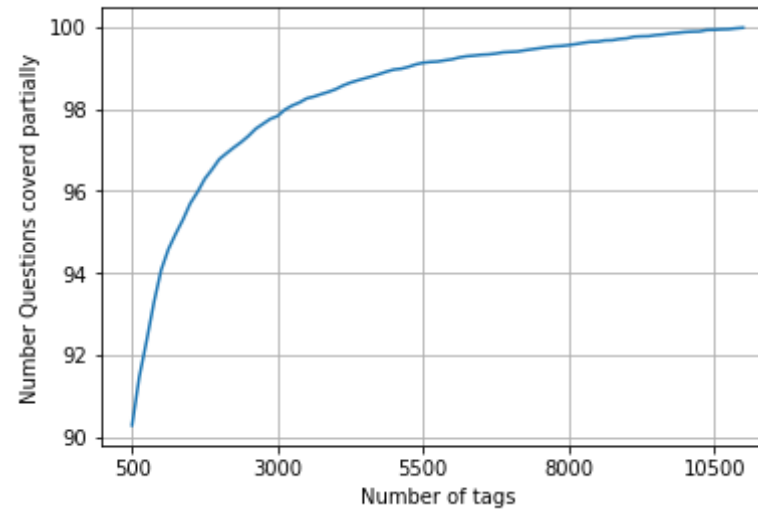
We will sample the number of tags instead considering all of them (due to limitation of computing power)

```
In [36]: 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 [37]: questions_explained = []
        total_tags=multilabel_y.shape[1]
        total_qs=preprocessed_data.shape[0]
        for i in range(500, total_tags, 100):
            questions_explained.append(np.round(((total_qs-questions_explained_fn(i))/total_qs)*100,3))
```

```
In [38]: fig, ax = plt.subplots()
        ax.plot(questions_explained)
        xlabel = list(500+np.array(range(-50,450,50))*50)
        ax.set_xticklabels(xlabel)
        plt.xlabel("Number of tags")
        plt.ylabel("Number Questions covered partially")
        plt.grid()
        plt.show()
        # you can choose any number of tags based on your computing power, minimum is 50(it covers 90% of the tags)
        print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
```



with 5500 tags we are covering 99.355 % of questions

```
In [39]: multilabel_yx = tags_to_choose(5500)
print("number of questions that are not covered :", questions_explained_
_fn(5500),"out of ", total_qs)
```

number of questions that are not covered : 129 out of 19999

```
In [40]: 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 : 8986
number of tags taken : 5500 (61.20632094369019 %)

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

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

```
In [41]: 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 [42]: print(x_train.shape)
print(y_train.shape)

print(x_test.shape)
print(y_test.shape)
```

```
(15999, 2)
(15999, 5500)
(4000, 2)
(4000, 5500)
```

```
In [43]: 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 : (15999, 5500)
Number of data points in test data : (4000, 5500)
```

4.3 Featurizing data

```
In [44]: start = datetime.now()
vectorizer = TfidfVectorizer(min_df=0.00009, max_features=20000, smooth
_idf=True, norm="l2", \
                           tokenizer = lambda x: x.split(), sublinear
_tf=False, ngram_range=(1,3))
x_train_multilabel = vectorizer.fit_transform(x_train['question'])
x_test_multilabel = vectorizer.transform(x_test['question'])
print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:00:09.876975

```
In [45]: 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: (15999, 20000) Y : (15999, 5500)
Dimensions of test data X: (4000, 20000) Y: (4000, 5500)

```
In [46]: # 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[46]: `from skmultilearn.adapt import MLkNN\nclassifier = MLkNN(k=21)\n\n#`

```
Out[40]: \nfrom sklearn.metrics import MLKNNClassifier = MLKNN(K=21)\n\n#
train\nclassifier.fit(x_train_multilabel, y_train)\n\n# predict\npredic
tions = classifier.predict(x_test_multilabel)\nprint(accuracy_score(y_t
est,predictions))\nprint(metrics.f1_score(y_test, predictions, average
= 'macro'))\nprint(metrics.f1_score(y_test, predictions, average = 'mic
ro'))\nprint(metrics.hamming_loss(y_test,predictions))\n\n"
```

4.4 Applying Logistic Regression with OneVsRest Classifier

```
In [47]: # this will be taking so much time try not to run it, download the lr_w
ith_equal_weight.pkl file and use to predict
# This takes about 6-7 hours to run.
classifier = OneVsRestClassifier(SGDClassifier(loss='log', alpha=0.0000
1, penalty='l1'), n_jobs=-1)
classifier.fit(x_train_multilabel, y_train)
predictions = classifier.predict(x_test_multilabel)

print("accuracy :",metrics.accuracy_score(y_test,predictions))
print("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_te
st, predictions))
```

```
accuracy : 0.06625
macro f1 score : 0.04862066134340646
micro f1 scoore : 0.351889852386311
hamming loss : 0.00044504545454545453
```

```
C:\Users\hemant\AnacondaNew\lib\site-packages\sklearn\metrics\_classifi
cation.py:1511: UndefinedMetricWarning: F-score is ill-defined and bein
g set to 0.0 in labels with no true nor predicted samples. Use `zero_di
vision` parameter to control this behavior.
average, "true nor predicted", 'F-score is', len(true_sum)
```

```
In [48]: from sklearn.externals import joblib
joblib.dump(classifier, 'lr_with_equal_weight.pkl')
```

```
Out[48]: ['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 [49]: sql_create_table = """CREATE TABLE IF NOT EXISTS QuestionsProcessed (qu
estion text NOT NULL, code text, tags text, words_pre integer, words_po
st integer, is_code integer);"""
create_database_table("G:\\machine_learning\\case_study\\Case study 5 S
tackoverflow tag predictor\\assignment\\Titlmoreweight.db", sql_create
_table)
```

Tables in the database:
QuestionsProcessed

```
In [50]: # http://www.sqlitetutorial.net/sqlite-delete/
# https://stackoverflow.com/questions/2279706/select-random-row-from-a-
sqlite-table

read_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflow
tag predictor\\assignment\\train_no_dup.db'
write_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflo
w tag predictor\\assignment\\Titlmoreweight.db'
train_datasize = 20000
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
T 20000;")
        # for selecting random points
        #reader.execute("SELECT Title, Body, Tags From no_dup_train ORD
ER 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 [51]: *#<http://www.bernzilla.com/2008/05/13/selecting-a-random-row-from-an-sqlite-table/>*

```

start = datetime.now()
preprocessed_data_list=[]
reader.fetchone()
questions_with_code=0
len_pre=0
len_post=0
questions_processed = 0
for row in 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.DOT
ALL))

question=re.sub('<code>(.*?)</code>', '', question, flags=re.MULTIL
INE|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)+" "+questio
n+" "+str(tags)
# else:
# question=str(title)+" "+str(title)+" "+str(title)+" "+questio
n

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

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

```

```

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

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

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

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

```

Avg. length of questions(Title+Body) before processing: 1279
 Avg. length of questions(Title+Body) after processing: 437
 Percent of questions containing code: 60
 Time taken to run this cell : 0:01:10.834027

In [52]:

```

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

```

Sample quesitons after preprocessing of data

In [53]:

```

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 1
0")
        print("Questions after preprocessed")
        print('='*100)
        reader.fetchone()
        for row in reader:
            print(row)
            print('-'*100)
conn_r.commit()
conn_r.close()

```

Questions after preprocessed

```

=====
=====
('dynam datagrid bind silverlight dynam datagrid bind silverlight dynam
datagrid bind silverlight bind datagrid dynam code wrote code debug cod
e block seem bind correct grid come column form come grid column althou
gh necessari bind nthank repli advance..',)
-----
-----
('java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryval
id java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryva
lid java.lang.noclassdeffounderror javax servlet jsp tagext taglibraryv
alid follow 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 v
ersion 1.2 jstl still messag caus solv',)
-----
-----
('java.sql.sqlexcept microsoft odbc driver manag invalid descriptor ind
ex java.sql.sqlexcept microsoft odbc driver manag invalid descriptor in
dex java.sql.sqlexcept microsoft odbc driver manag invalid descriptor i
ndex use follow code display caus solv',)
-----
-----
('better way updat feed fb php sdk better way updat feed fb php sdk bet
ter way 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 s
ometh like way better',)
-----

```

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

```
-----  
( 'sql inject issu prevent correct form submiss php sql inject issu prev  
ent correct form submiss php sql inject issu prevent correct form submi  
ss php check everyth think make sure input field safe type sql inject g  
ood news safe bad news one tag mess form submiss place even touch life  
figur exact html use templat file forgiv okay entir php script get exec  
ut see data post none forum field post problem use someth titl field no  
ne data get post current use print post see submit noth work flawless s  
tatement though also mention script work flawless local machin use host  
come across problem state list input test mess', )  
-----
```

```
-----  
( 'countabl subaddit lebesgu measur countabl subaddit lebesgu measur cou  
ntabl subaddit lebesgu measur let lbrace rbrace sequenc set sigma -alge  
bra mathcal want show left bigcup right leq sum left right countabl add  
it measur defin set sigma algebra mathcal think use monoton properti so  
mewher proof start appreci littl help nthank ad han answer make follow  
addit construct given han answer clear bigcup bigcup cap emptyset neq l  
eft bigcup right left bigcup right sum left right also construct subset  
monoton left right leq left right final would sum leq sum result follo  
w', )  
-----
```

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

```
-----  
( 'undefin symbol architectur i386 objc class skpsmtpmessag referenc err  
or undefin symbol architectur i386 objc class skpsmtpmessag referenc er  
ror undefin symbol architectur i386 objc class skpsmtpmessag referenc e  
rror import framework send email applic background import framework i.e  
skpsmtpmessag somebodi suggest get error collect2 ld return exit status  
import framework correct core taken framework follow mfmilcomposeview
```

```

import framework correct sort taken framework follow initialcomposeview
ontrol question lock field updat answer drag drop folder project click
copi nthat',)
-----
-----

```

Saving Preprocessed data to a Database

```

In [54]: #Taking 0.5 Million entries to a dataframe.
write_db = 'G:\\machine_learning\\case_study\\Case study 5 Stackoverflo
w tag predictor\\assignment\\Titlmoreweight.db'
if os.path.isfile(write_db):
    conn_r = create_connection(write_db)
    if conn_r is not None:
        preprocessed_data = pd.read_sql_query("""SELECT question, Tags
FROM QuestionsProcessed LIMIT 20000""", conn_r)
    conn_r.commit()
    conn_r.close()

```

```

In [55]: preprocessed_data.head()

```

Out[55]:

	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 [56]: 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 : 19999
number of dimensions : 2

```

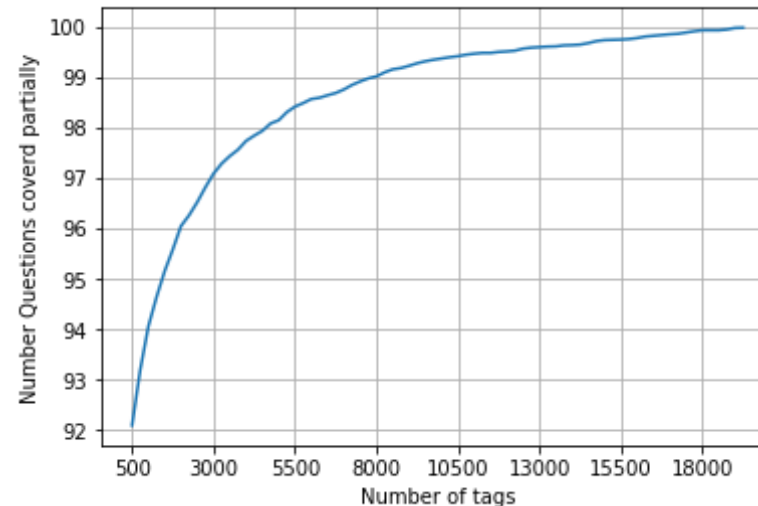
Converting string Tags to multilable output variables

```
In [57]: vectorizer = CountVectorizer(tokenizer = lambda x: x.split(), binary='true')
        multilabel_y = vectorizer.fit_transform(preprocessed_data['tags'])
```

Selecting 500 Tags

```
In [58]: questions_explained = []
        total_tags=multilabel_y.shape[1]
        total_qs=preprocessed_data.shape[0]
        for i in range(500, total_tags, 100):
            questions_explained.append(np.round(((total_qs-questions_explained_
            fn(i))/total_qs)*100,3))
```

```
In [59]: fig, ax = plt.subplots()
        ax.plot(questions_explained)
        xlabel = list(500+np.array(range(-50,450,50))*50)
        ax.set_xticklabels(xlabel)
        plt.xlabel("Number of tags")
        plt.ylabel("Number Questions covered partially")
        plt.grid()
        plt.show()
        # you can choose any number of tags based on your computing power, minimum is 500(it covers 90% of the tags)
        print("with ",5500,"tags we are covering ",questions_explained[50],"% of questions")
        print("with ",500,"tags we are covering ",questions_explained[0],"% of questions")
```



with 5500 tags we are covering 99.6 % of questions
 with 500 tags we are covering 92.1 % of questions

```
In [60]: # we will be taking 500 tags
          multilabel_yx = tags_to_choose(500)
          print("number of questions that are not covered :", questions_explained_
                _fn(500),"out of ", total_qs)
```

number of questions that are not covered : 1580 out of 19999

```
In [61]: """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],:]"""
```

```
Out[61]: 'x_train=preprocessed_data.head(train_datasize)\nx_test=preprocessed_data.tail(preprocessed_data.shape[0] - 400000)\n\nny_train = multilabel_yx[0:train_datasize,:]\nny_test = multilabel_yx[train_datasize:preprocessed_data.shape[0],:]'
```



```
In [62]: 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 : (15999, 5500)
Number of data points in test data : (4000, 5500)

4.5.2 Featurizing data with Tfidf vectorizer

```
In [63]: 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 [64]: start = datetime.now()
vectorizer = CountVectorizer(ngram_range=(4,4))
x_train_multilabel = vectorizer.fit_transform(x_train['question'])
x_test_multilabel = vectorizer.transform(x_test['question'])
print("Time taken to run this cell :", datetime.now() - start)
```

Time taken to run this cell : 0:00:06.783344

```
In [65]: print("Dimensions of train data X:",x_train_multilabel.shape, "Y :",x_train_multilabel.shape)
print("Dimensions of test data X:",x_test_multilabel.shape,"Y:",x_test_multilabel.shape)
```

Dimensions of train data X: (15999, 854897) Y : (15999, 854897)
Dimensions of test data X: (4000, 854897) Y: (4000, 854897)

4.5.3 Applying Logistic Regression with OneVsRest Classifier

```

In [66]: start = datetime.now()
classifier = OneVsRestClassifier(SGDClassifier(loss='hinge', alpha=0.00001))
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='micro')
recall = recall_score(y_test, predictions, average='micro')
f1 = f1_score(y_test, predictions, average='micro')

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 : 0.0755
Hamming loss 0.0038565
Micro-average quality numbers
Precision: 0.3039, Recall: 0.0239, F1-measure: 0.0444
Macro-average quality numbers
Precision: 0.3039, Recall: 0.0239, F1-measure: 0.0444
Time taken to run this cell : 0:01:15.376991

```

```

In [67]: joblib.dump(classifier, 'lr_with_more_title_weight.pkl')

```

Out[67]: ['lr_with_more_title_weight.pkl']

```
In [68]: start = datetime.now()

classifier_2 = OneVsRestClassifier(LogisticRegression())
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='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))

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

Accuracy : 0.07225
Hamming loss  0.0037415
Micro-average quality numbers
Precision: 0.4444, Recall: 0.0005, F1-measure: 0.0011
Micro-average quality numbers
Precision: 0.4444, Recall: 0.0005, F1-measure: 0.0011
Time taken to run this cell : 0:48:31.555129
```

```

In [69]: from sklearn.metrics import roc_auc_score
alpha = [10**-5,10**-4,10**-3,10**-2,10**-1,10**0,10]

#alpha_values = np.arange(7)

#acc = np.empty(len(alpha_values))

cv_auc = []

for i in alpha:
    clf = OneVsRestClassifier(SGDClassifier(loss = 'log',alpha = i, pen
alty = 'l1'))
    clf.fit(x_train_multilabel, y_train)

    pred = clf.predict(x_test_multilabel)

    # evaluate CV accuracy
    #acc[i] = f1_score(y_test, pred, average='macro')
    cv_auc.append(metrics.accuracy_score(y_test, pred))

    #print('\nCV accuracy for k = %d is %d%%' % (i, acc))
    #error[i] = 100-acc[i]

# optimal_k = int(min(error))
# print('\nThe optimal number of neighbors is %d.' % optimal_k)

#Generate plot

d = max(cv_auc)

i = np.where(cv_auc == d)

i = i[0][0]
best_alpha = float(alpha[i])
print("Best alpha is:-",best_alpha)

```

Best alpha is:- 0.0001

```
In [70]: classifier_2 = OneVsRestClassifier(SGDClassifier(loss = 'log',alpha = best_alpha, penalty = 'l1'))
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='micro')
recall = recall_score(y_test, predictions_2, average='micro')
f1 = f1_score(y_test, predictions_2, average='micro')

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))

Accuracy : 0.074
Hamming loss  0.0037845
Micro-average quality numbers
```

```
Precision: 0.3239, Recall: 0.0107, F1-measure: 0.0207  
Macro-average quality numbers  
Precision: 0.3239, Recall: 0.0107, F1-measure: 0.0207
```

In []:

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 []: