

# Quora Question Pairs

## 1. Business Problem

### 1.1 Description

Quora is a place to gain and share knowledge—about anything. It's a platform to ask questions and connect with people who contribute unique insights and quality answers. This empowers people to learn from each other and to better understand the world.

Over 100 million people visit Quora every month, so it's no surprise that many people ask similarly worded questions. Multiple questions with the same intent can cause seekers to spend more time finding the best answer to their question, and make writers feel they need to answer multiple versions of the same question. Quora values canonical questions because they provide a better experience to active seekers and writers, and offer more value to both of these groups in the long term.

Credits: Kaggle

#### Problem Statement

- Identify which questions asked on Quora are duplicates of questions that have already been asked.
- This could be useful to instantly provide answers to questions that have already been answered.
- We are tasked with predicting whether a pair of questions are duplicates or not.

### 1.2 Sources/Useful Links

- Source : <https://www.kaggle.com/c/quora-question-pairs>

#### Useful Links

- Discussions : <https://www.kaggle.com/anokas/data-analysis-xgboost-starter-0-35460-lb/comments>
- Kaggle Winning Solution and other approaches:  
<https://www.dropbox.com/sh/93968nfnrz8bp5/AACZdtsApc1QSTQc7X0H3QZ5a?dl=0>
- Blog 1 : <https://engineering.quora.com/Semantic-Question-Matching-with-Deep-Learning>
- Blog 2 : <https://towardsdatascience.com/identifying-duplicate-questions-on-quora-top-12-on-kaggle-4c1cf93f1c30>

### 1.3 Real world/Business Objectives and Constraints

1. The cost of a mis-classification can be very high.
2. You would want a probability of a pair of questions to be duplicates so that you can choose any threshold of choice.
3. No strict latency concerns.
4. Interpretability is partially important.

## 2. Machine Learning Problem

### 2.1 Data

#### 2.1.1 Data Overview

- Data will be in a file Train.csv
- Train.csv contains 5 columns : qid1, qid2, question1, question2, is\_duplicate

- Size of Train.csv - 60MB
- Number of rows in Train.csv = 404,290

## 2.1.2 Example Data point

```
"id","qid1","qid2","question1","question2","is_duplicate"  
"0","1","2","What is the step by step guide to invest in share market in india?","What is  
the step by step guide to invest in share market?","0"  
"1","3","4","What is the story of Kohinoor (Koh-i-Noor) Diamond?","What would happen if the  
Indian government stole the Kohinoor (Koh-i-Noor) diamond back?","0"  
"7","15","16","How can I be a good geologist?","What should I do to be a great  
geologist?","1"  
"11","23","24","How do I read and find my YouTube comments?","How can I see all my Youtube  
comments?","1"
```

## 2.2 Mapping the real world problem to an ML problem

### 2.2.1 Type of Machine Learning Problem

It is a binary classification problem, for a given pair of questions we need to predict if they are duplicate or not.

### 2.2.2 Performance Metric

Source: <https://www.kaggle.com/c/quora-question-pairs#evaluation>

Metric(s):

- log-loss : <https://www.kaggle.com/wiki/LogarithmicLoss>
- Binary Confusion Matrix

## 2.3 Train and Test Construction

We build train and test by randomly splitting in the ratio of 70:30 or 80:20 whatever we choose as we have sufficient points to work with.

## 3. Exploratory Data Analysis

In [46]:

```
import numpy as np  
import pandas as pd  
import seaborn as sns  
import matplotlib.pyplot as plt  
from subprocess import check_output  
%matplotlib inline  
import plotly.offline as py  
py.init_notebook_mode(connected=True)  
import plotly.graph_objs as go  
import plotly.tools as tls  
import gc  
import sys  
import os  
from tqdm import tqdm  
  
import re  
from nltk.corpus import stopwords  
import distance  
from nltk.stem import PorterStemmer  
from bs4 import BeautifulSoup  
from fuzzywuzzy import fuzz  
from sklearn.manifold import TSNE
```

```

from wordcloud import WordCloud, STOPWORDS
from os import path
from PIL import Image
color = sns.color_palette()

```

## 3.1 Reading data and basic stats

In [2]:

```

df = pd.read_csv("train.csv")

print("Number of data points:",df.shape[0])

```

Number of data points: 404290

In [3]:

```
df.head()
```

Out[3]:

	id	qid1	qid2	question1	question2	is_duplicate
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0
3	3	7	8	Why am I mentally very lonely? How can I solve...	Find the remainder when $23^{24}$ i...	0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

In [4]:

```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 404290 entries, 0 to 404289
Data columns (total 6 columns):
id                404290 non-null int64
qid1              404290 non-null int64
qid2              404290 non-null int64
question1         404289 non-null object
question2         404288 non-null object
is_duplicate      404290 non-null int64
dtypes: int64(4), object(2)
memory usage: 18.5+ MB

```

We are given a minimal number of data fields here, consisting of:

- id: Looks like a simple rowID
- qid{1, 2}: The unique ID of each question in the pair
- question{1, 2}: The actual textual contents of the questions.
- is\_duplicate: The label that we are trying to predict - whether the two questions are duplicates of each other.

### 3.2.1 Distribution of data points among output classes

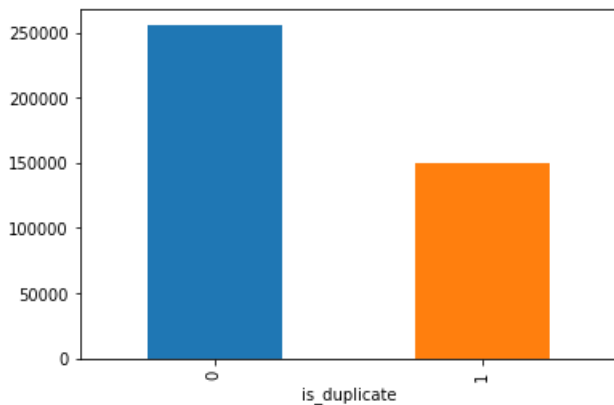
- Number of duplicate(smilar) and non-duplicate(non similar) questions

In [5]:

```
df.groupby("is_duplicate")['id'].count().plot.bar()
```

Out[5]:

<matplotlib.axes.\_subplots.AxesSubplot at 0x152a1b6588>



In [6]:

```
print('~> Total number of question pairs for training:\n {}'.format(len(df)))
```

>> Total number of question pairs for training:  
404290

In [7]:

```
print('~> Question pairs are not Similar (is_duplicate = 0):\n {}'.format(100 -  
round(df['is_duplicate'].mean()*100, 2)))  
print('\n~> Question pairs are Similar (is_duplicate = 1):\n {}'.format(round(df['is_duplicate']  
.mean()*100, 2)))
```

>> Question pairs are not Similar (is\_duplicate = 0):  
63.08%

>> Question pairs are Similar (is\_duplicate = 1):  
36.92%

### 3.2.2 Number of unique questions

In [8]:

```
qids = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())  
unique_qs = len(np.unique(qids))  
qs_morethan_onetime = np.sum(qids.value_counts() > 1)  
print ('Total number of Unique Questions are: {}\n'.format(unique_qs))  
#print len(np.unique(qids))  
  
print ('Number of unique questions that appear more than one time: {}  
({})\n'.format(qs_morethan_onetime, qs_morethan_onetime/unique_qs*100))  
  
print ('Max number of times a single question is repeated: {}\n'.format(max(qids.value_counts())))  
  
q_vals=qids.value_counts()  
q_vals=q_vals.values
```

Total number of Unique Questions are: 537933

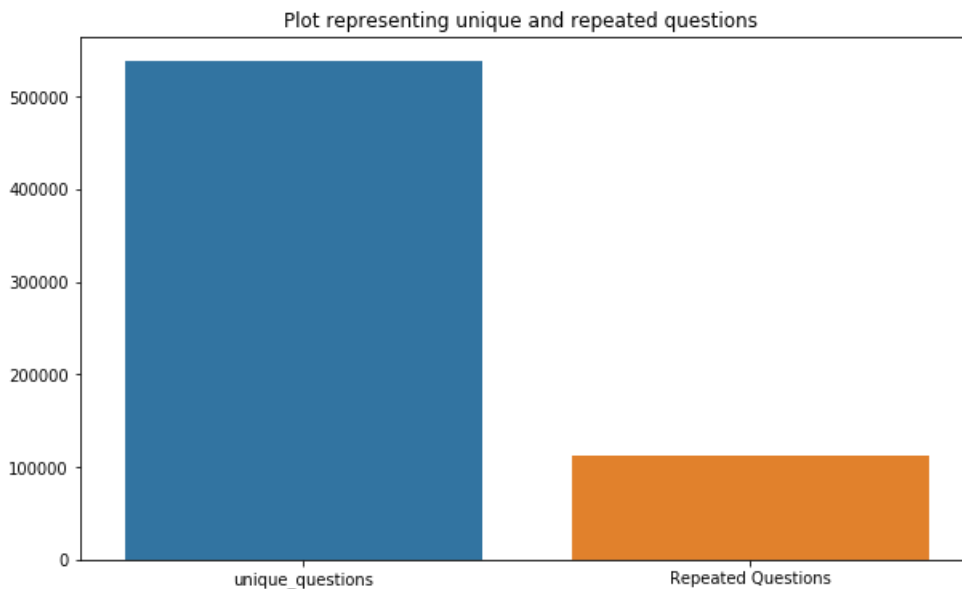
Number of unique questions that appear more than one time: 111780 (20.77953945937505%)

Max number of times a single question is repeated: 157

In [9]:

```
x = ["unique_questions" , "Repeated Questions"]
y = [unique_qs , qs_morethan_onetime]

plt.figure(figsize=(10, 6))
plt.title ("Plot representing unique and repeated questions ")
sns.barplot(x,y)
plt.show()
```



### 3.2.3 Checking for Duplicates

In [10]:

```
#checking whether there are any repeated pair of questions

pair_duplicates =
df[['qid1','qid2','is_duplicate']].groupby(['qid1','qid2']).count().reset_index()

print ("Number of duplicate questions", (pair_duplicates).shape[0] - df.shape[0])
```

Number of duplicate questions 0

### 3.2.4 Number of occurrences of each question

In [11]:

```
plt.figure(figsize=(20, 10))

plt.hist(qids.value_counts(), bins=160)

plt.yscale('log', nonposy='clip')

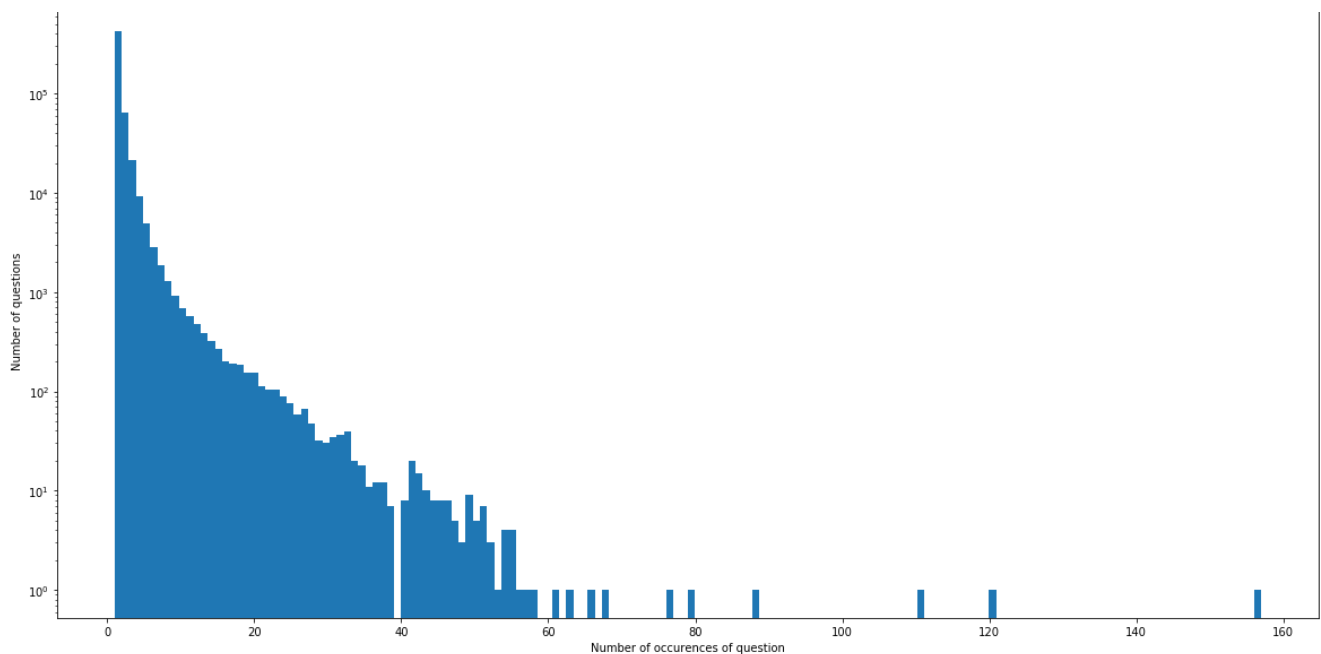
plt.title('Log-Histogram of question appearance counts')

plt.xlabel('Number of occurrences of question')

plt.ylabel('Number of questions')

print ('Maximum number of times a single question is repeated: {}'.format(max(qids.value_counts(
))))
```

Maximum number of times a single question is repeated: 157



### 3.2.5 Checking for NULL values

In [12]:

```
#Checking whether there are any rows with null values
nan_rows = df[df.isnull().any(1)]
print (nan_rows)
```

	id	qid1	qid2	question1 \	question2	is_duplicate
105780	105780	174363	174364	How can I develop android app?	NaN	0
201841	201841	303951	174364	How can I create an Android app?	NaN	0
363362	363362	493340	493341	My Chinese name is Haichao Yu. What English na...	NaN	0

- There are two rows with null values in question2

In [13]:

```
# Filling the null values with ' '
df = df.fillna(' ')
nan_rows = df[df.isnull().any(1)]
print (nan_rows)
```

```
Empty DataFrame
Columns: [id, qid1, qid2, question1, question2, is_duplicate]
Index: []
```

### 3.3 Basic Feature Extraction (before cleaning)

Let us now construct a few features like:

- **freq\_qid1** = Frequency of qid1's
- **freq\_qid2** = Frequency of qid2's
- **q1len** = Length of q1
- **q2len** = Length of q2
- **q1\_n\_words** = Number of words in Question 1
- **q2\_n\_words** = Number of words in Question 2
- **word\_Common** = (Number of common unique words in Question 1 and Question 2)
- **word\_Total** =(Total num of words in Question 1 + Total num of words in Question 2)

- **word\_share** = (word\_common)/(word\_Total)
- **freq\_q1+freq\_q2** = sum total of frequency of qid1 and qid2
- **freq\_q1-freq\_q2** = absolute difference of frequency of qid1 and qid2

In [41]:

```
if os.path.isfile('df_fe_without_preprocessing_train.csv'):
    df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
else:
    df['freq_qid1'] = df.groupby('qid1')['qid1'].transform('count')
    df['freq_qid2'] = df.groupby('qid2')['qid2'].transform('count')
    df['q1len'] = df['question1'].str.len()
    df['q2len'] = df['question2'].str.len()
    df['q1_n_words'] = df['question1'].apply(lambda row: len(row.split(" ")))
    df['q2_n_words'] = df['question2'].apply(lambda row: len(row.split(" ")))

    def normalized_word_Common(row):
        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
        return 1.0 * len(w1 & w2)
    df['word_Common'] = df.apply(normalized_word_Common, axis=1)

    def normalized_word_Total(row):
        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
        return 1.0 * (len(w1) + len(w2))
    df['word_Total'] = df.apply(normalized_word_Total, axis=1)

    def normalized_word_share(row):
        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
        return 1.0 * len(w1 & w2) / (len(w1) + len(w2))
    df['word_share'] = df.apply(normalized_word_share, axis=1)

    df['freq_q1+q2'] = df['freq_qid1']+df['freq_qid2']
    df['freq_q1-q2'] = abs(df['freq_qid1']-df['freq_qid2'])

    df.to_csv("df_fe_without_preprocessing_train.csv", index=False)

df.head()
```

Out[41]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0	1	1	73	59	14	10	4.0
3	3	7	8	Why am I mentally very lonely?	Find the remainder when $123456789$	0	1	1	50	65	11	9	0.0

	id	qid1	qid2	How can I question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0	3	1	76	39	13	7	2.0

## number of words distribution in question

In [44]:

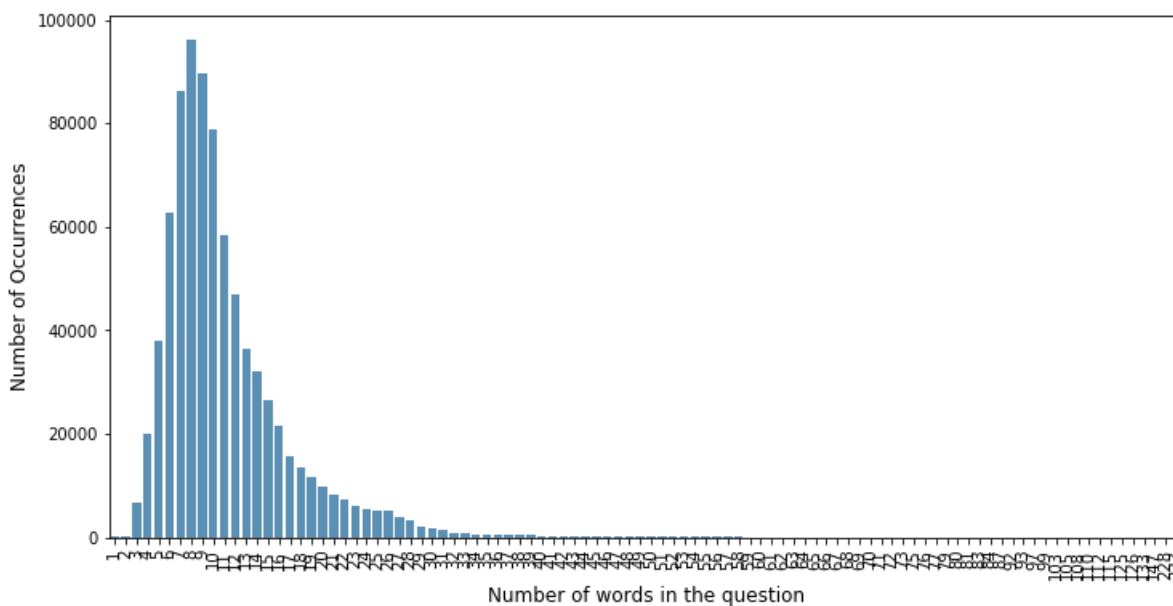
```
all_qes_df = pd.DataFrame(pd.concat([df['question1'], df['question2']]))
all_qes_df.columns = ["questions"]

all_qes_df["num_of_words"] = all_qes_df["questions"].apply(lambda x : len(str(x).split()))
```

In [47]:

```
cnt_srs = all_qes_df['num_of_words'].value_counts()

plt.figure(figsize=(12,6))
sns.barplot(cnt_srs.index, cnt_srs.values, alpha=0.8, color=color[0])
plt.ylabel('Number of Occurrences', fontsize=12)
plt.xlabel('Number of words in the question', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



observation:-The distribution here is right skewed with upto 237 words in question. There are also few questions with one or two words as well.

## exploration of characters

In [48]:

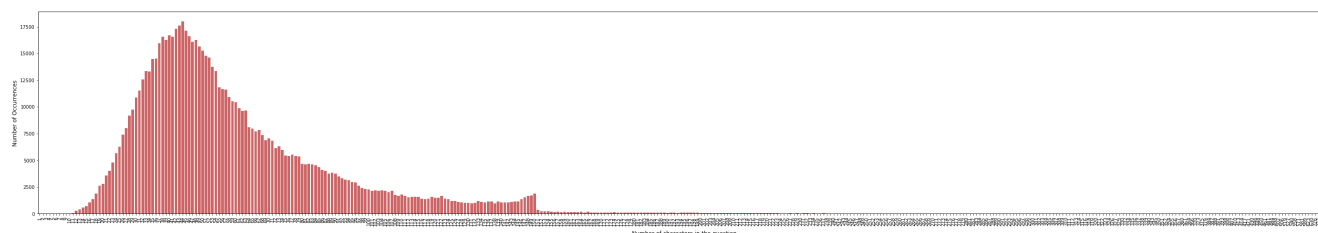
```
all_qes_df["num_of_chars"] = all_qes_df["questions"].apply(lambda x : len(str(x)))
cnt_srs = all_qes_df['num_of_chars'].value_counts()

plt.figure(figsize=(50,8))
sns.barplot(cnt_srs.index, cnt_srs.values, alpha=0.8, color=color[3])
plt.ylabel('Number of Occurrences', fontsize=12)
plt.xlabel('Number of characters in the question', fontsize=12)
```



```
plt.xticks(rotation='vertical')
plt.show()
```

```
del all_ques_df
```



observation:- 1.It is also right skewed distribution.

2.There is sudden dip at 160 character mark

## Common unigrams between given question pairs

In [51]:

```
from nltk import word_tokenize, ngrams
eng_stopwords = set(stopwords.words('english'))
def get_unigrams(que):
    return [word for word in word_tokenize(que.lower()) if word not in eng_stopwords]

def get_common_unigrams(row):
    return len( set(row["unigrams_ques1"]).intersection(set(row["unigrams_ques2"])) )

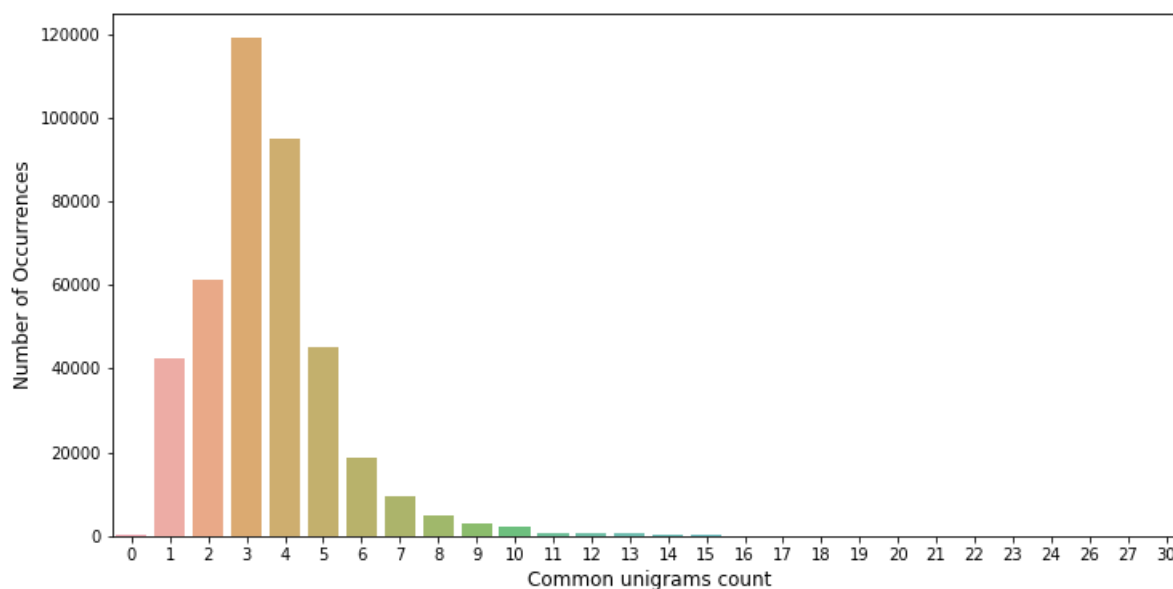
def get_common_unigram_ratio(row):
    return float(row["unigrams_common_count"]) / max(len( set(row["unigrams_ques1"]).union(set(row["unigrams_ques2"])) ),1)

df["unigrams_ques1"] = df['question1'].apply(lambda x: get_unigrams(str(x)))
df["unigrams_ques2"] = df['question2'].apply(lambda x: get_unigrams(str(x)))
df["unigrams_common_count"] = df.apply(lambda row: get_common_unigrams(row),axis=1)
df["unigrams_common_ratio"] = df.apply(lambda row: get_common_unigram_ratio(row), axis=1)
```

In [52]:

```
cnt_srs = df['unigrams_common_count'].value_counts()

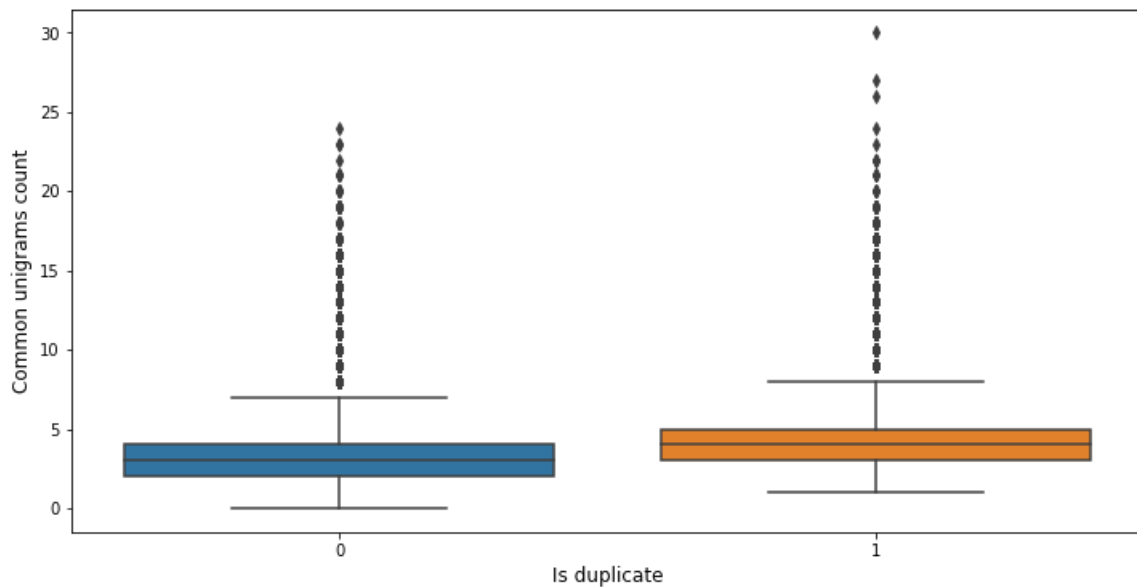
plt.figure(figsize=(12,6))
sns.barplot(cnt_srs.index, cnt_srs.values, alpha=0.8)
plt.ylabel('Number of Occurrences', fontsize=12)
plt.xlabel('Common unigrams count', fontsize=12)
plt.show()
```



from the above graph one can notice that there are few question pairs with no common words

In [54]:

```
plt.figure(figsize=(12,6))
sns.boxplot(x="is_duplicate", y="unigrams_common_count", data=df)
plt.xlabel('Is duplicate', fontsize=12)
plt.ylabel('Common unigrams count', fontsize=12)
plt.show()
```

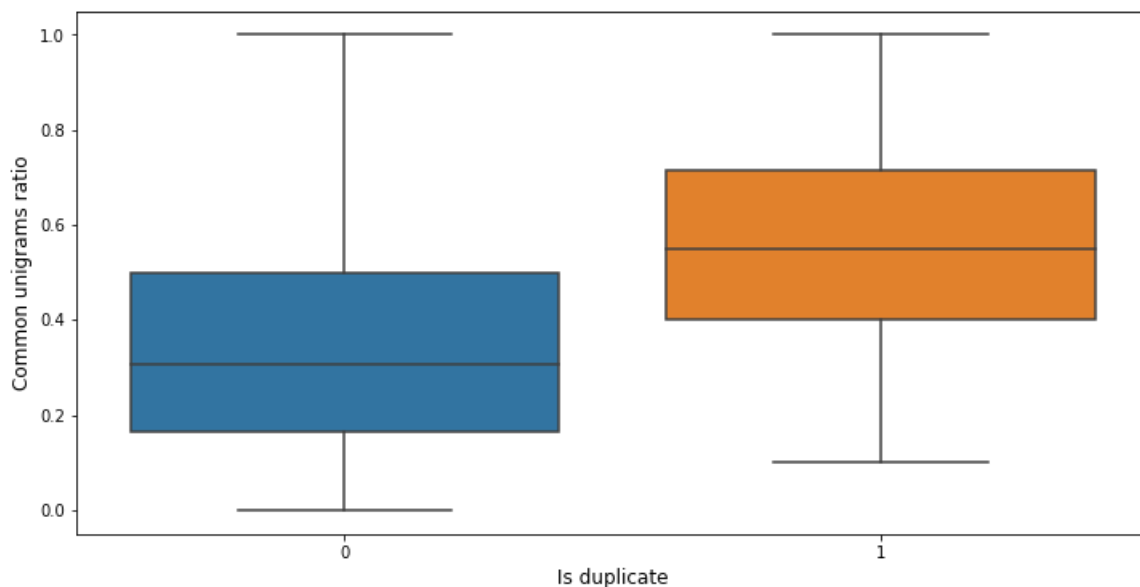


There is some good difference between 0 and 1 class using the common unigram count variable.

## Let us look at the same graph using common unigrams ratio

In [55]:

```
plt.figure(figsize=(12,6))
sns.boxplot(x="is_duplicate", y="unigrams_common_ratio", data=df)
plt.xlabel('Is duplicate', fontsize=12)
plt.ylabel('Common unigrams ratio', fontsize=12)
plt.show()
```



# q1 -q2 neighbour intersection count

In [57]:

```
ques = pd.concat([df[['question1', 'question2']]].reset_index(drop='index')
```

In [58]:

```
from collections import defaultdict
q_dict = defaultdict(set)
for i in range(ques.shape[0]):
    q_dict[ques.question1[i]].add(ques.question2[i])
    q_dict[ques.question2[i]].add(ques.question1[i])
```

In [60]:

```
def q1_freq(row):
    return(len(q_dict[row['question1']]))

def q2_freq(row):
    return(len(q_dict[row['question2']]))

def q1_q2_intersect(row):
    return(len(set(q_dict[row['question1']]).intersection(set(q_dict[row['question2']]))))

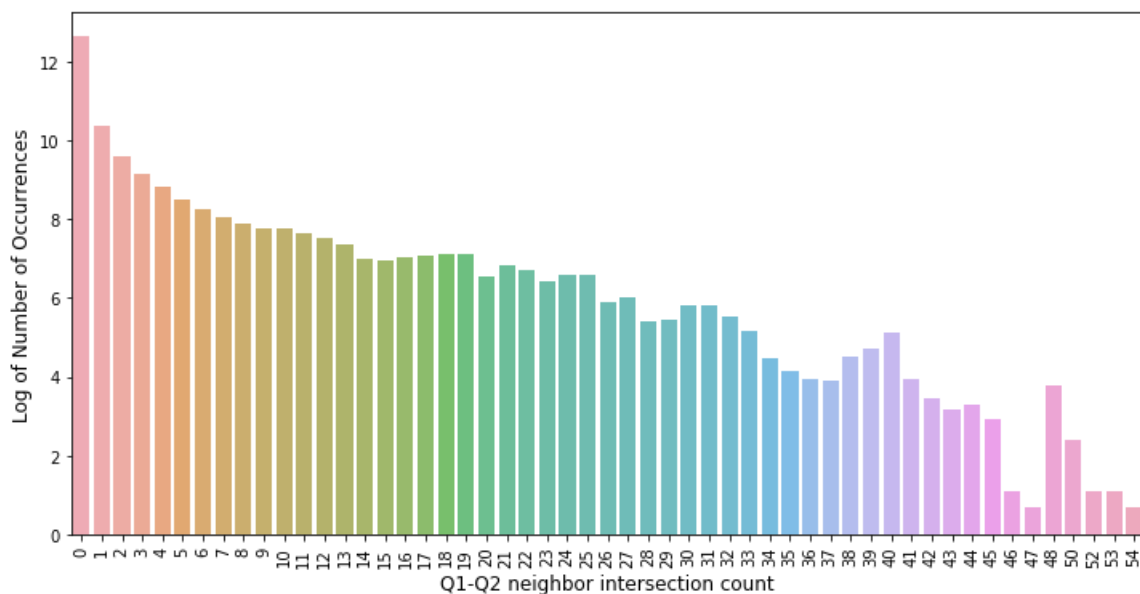
df['q1_q2_intersect'] = df.apply(q1_q2_intersect, axis=1, raw=True)
df['q1_freq'] = df.apply(q1_freq, axis=1, raw=True)
df['q2_freq'] = df.apply(q2_freq, axis=1, raw=True)
```

Let us first do simple count plots and see the distribution

In [61]:

```
cnt_srs = df['q1_q2_intersect'].value_counts()

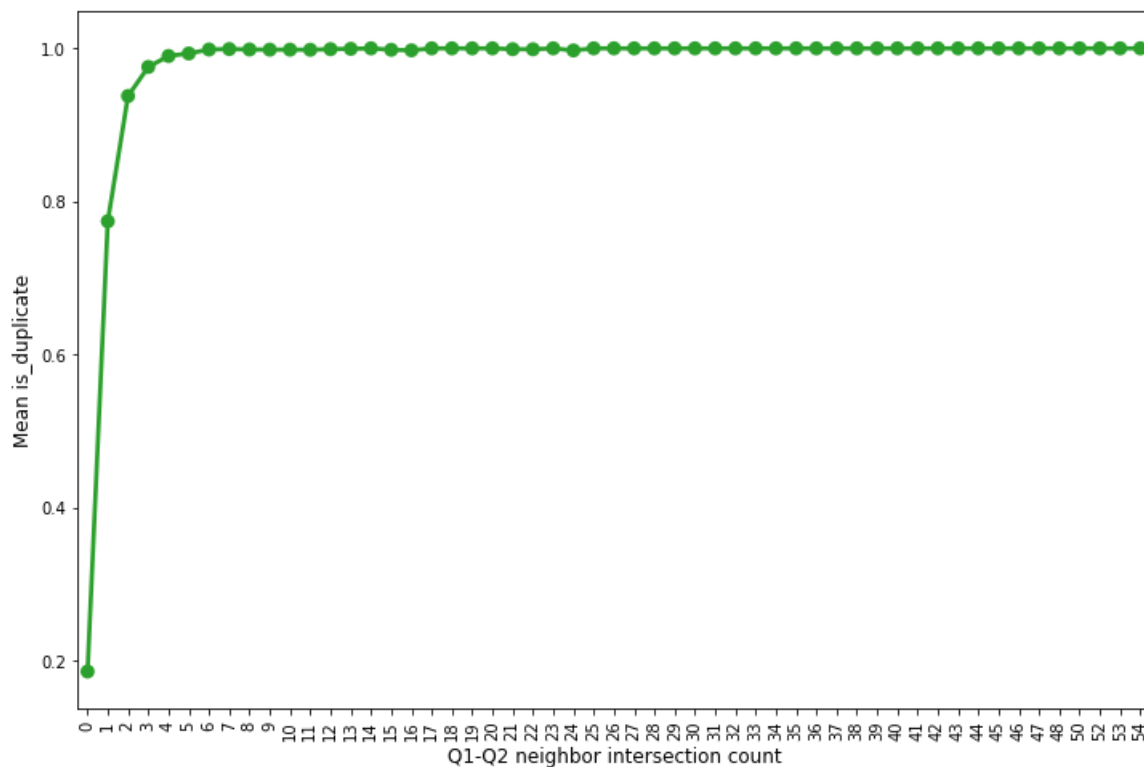
plt.figure(figsize=(12,6))
sns.barplot(cnt_srs.index, np.log1p(cnt_srs.values), alpha=0.8)
plt.xlabel('Q1-Q2 neighbor intersection count', fontsize=12)
plt.ylabel('Log of Number of Occurrences', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



In [62]:

```
grouped_df = df.groupby('q1_q2_intersect')['is_duplicate'].aggregate(np.mean).reset_index()
plt.figure(figsize=(12,8))
```

```
sns.pointplot(grouped_df["q1_q2_intersect"].values, grouped_df["is_duplicate"].values, alpha=0.8, color=color[2])
plt.ylabel('Mean is_duplicate', fontsize=12)
plt.xlabel('Q1-Q2 neighbor intersection count', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



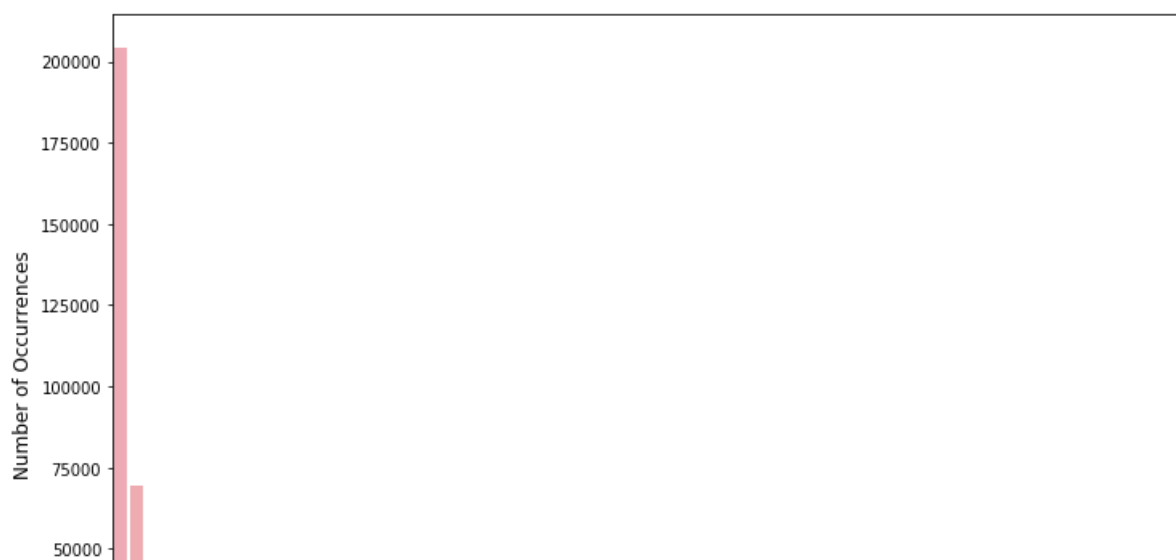
from the above graph we can see that this feature is super useful

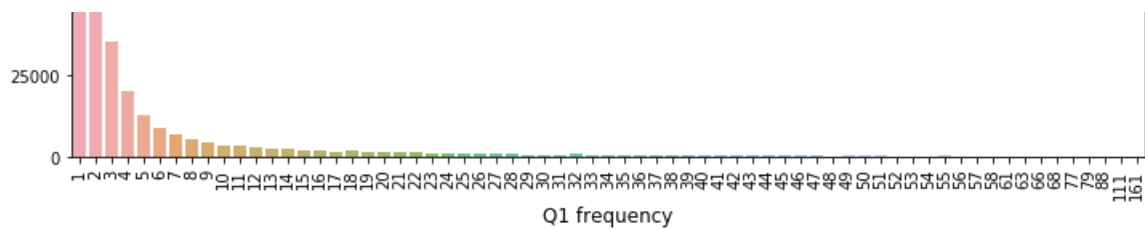
## Question 1 Frequency

In [63]:

```
cnt_srs = df['q1_freq'].value_counts()

plt.figure(figsize=(12,8))
sns.barplot(cnt_srs.index, cnt_srs.values, alpha=0.8)
plt.xlabel('Q1 frequency', fontsize=12)
plt.ylabel('Number of Occurrences', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```

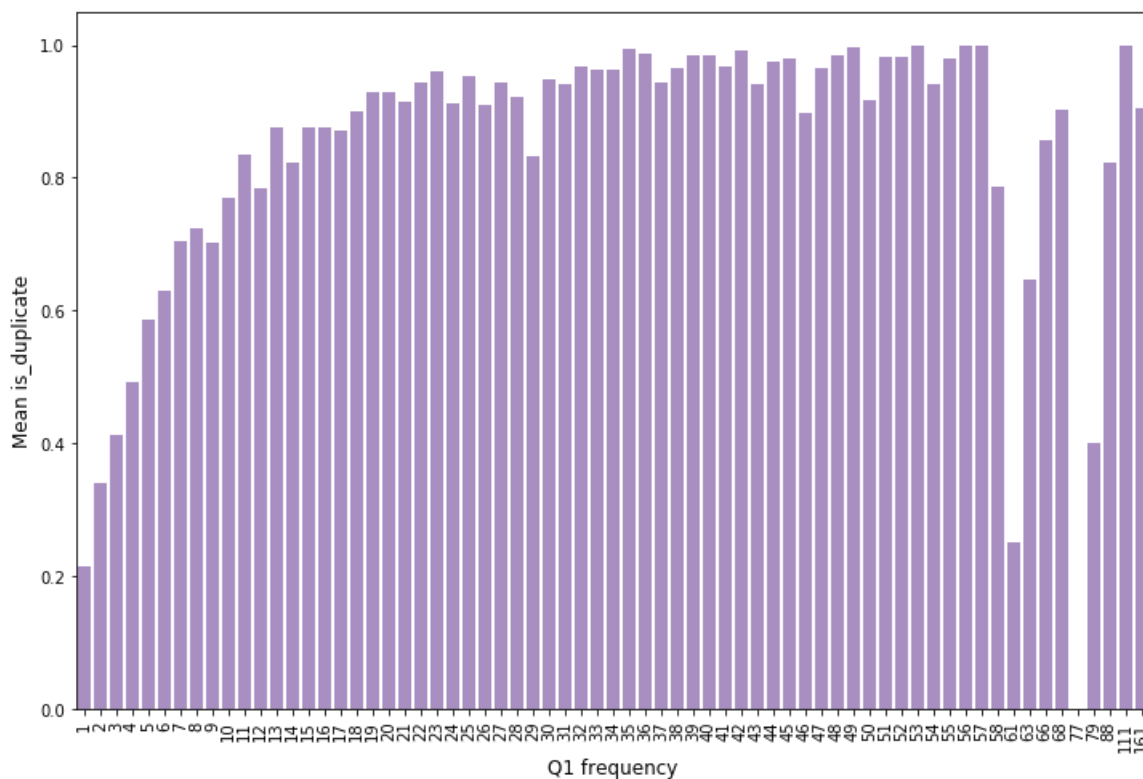




We could see a long tail here as well. Now let us check the target variable distribution.

In [64]:

```
plt.figure(figsize=(12,8))
grouped_df = df.groupby('q1_freq')['is_duplicate'].aggregate(np.mean).reset_index()
sns.barplot(grouped_df["q1_freq"].values, grouped_df["is_duplicate"].values, alpha=0.8, color=color
[4])
plt.ylabel('Mean is_duplicate', fontsize=12)
plt.xlabel('Q1 frequency', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



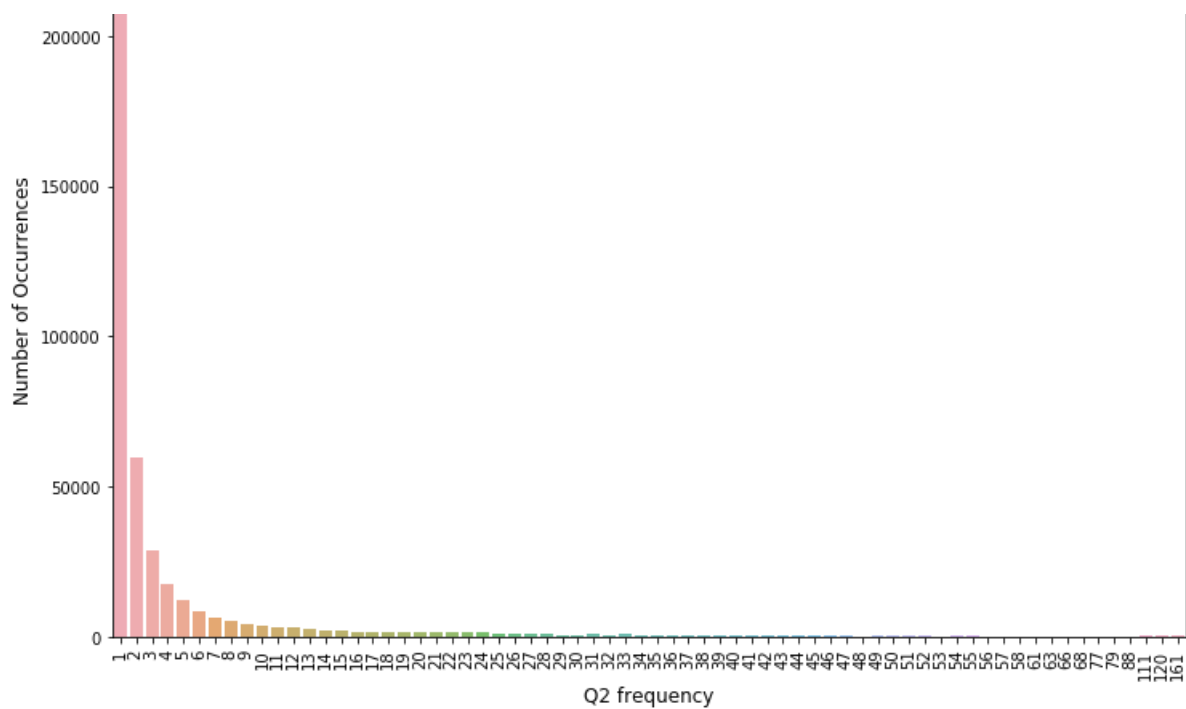
From the above plot there is increase in the mean target rate as the frequency increases.

## Question2 Frequency

In [65]:

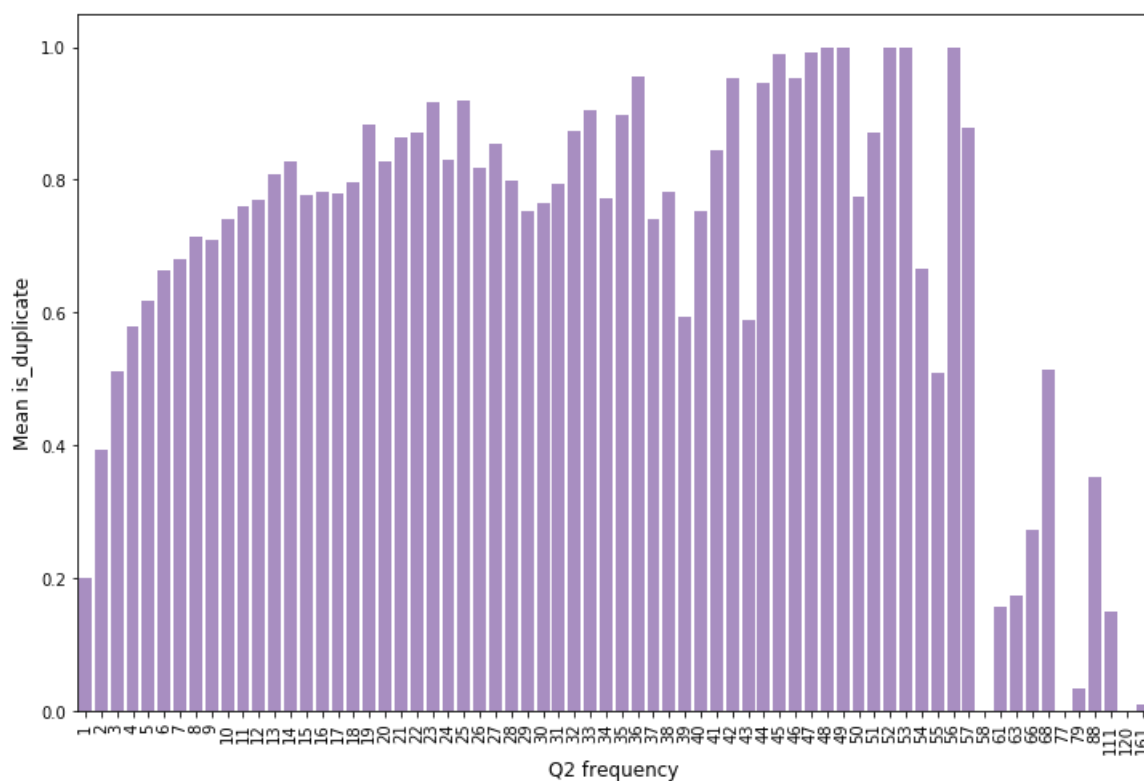
```
cnt_srs = df['q2_freq'].value_counts()

plt.figure(figsize=(12,8))
sns.barplot(cnt_srs.index, cnt_srs.values, alpha=0.8)
plt.xlabel('Q2 frequency', fontsize=12)
plt.ylabel('Number of Occurrences', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



In [66]:

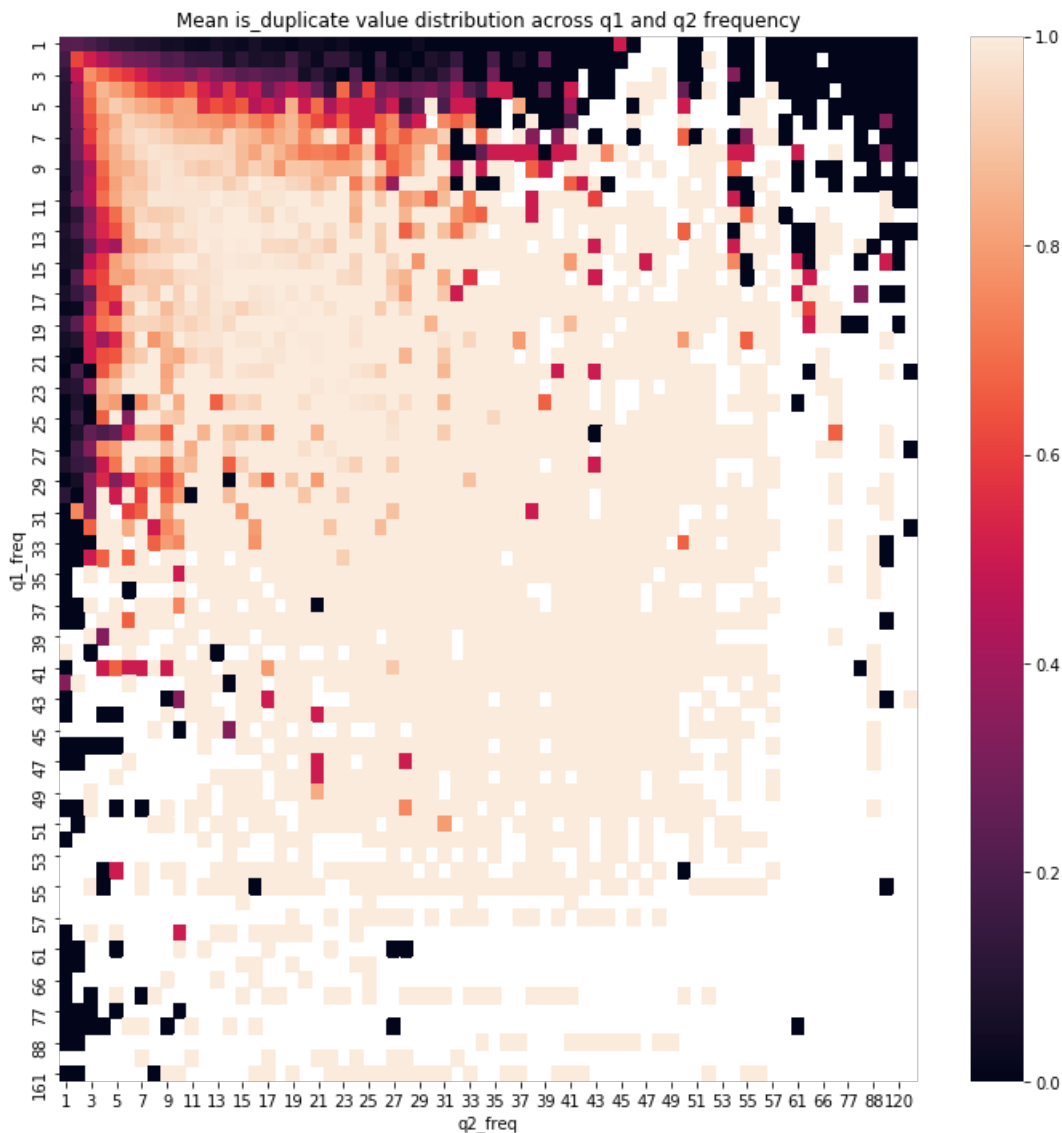
```
plt.figure(figsize=(12,8))
grouped_df = df.groupby('q2_freq')['is_duplicate'].aggregate(np.mean).reset_index()
sns.barplot(grouped_df["q2_freq"].values, grouped_df["is_duplicate"].values, alpha=0.8, color=color
[4])
plt.ylabel('Mean is_duplicate', fontsize=12)
plt.xlabel('Q2 frequency', fontsize=12)
plt.xticks(rotation='vertical')
plt.show()
```



**Plotting heat map between q1 freq and q2 freq to check the target variable distribution**

In [67]:

```
pvt_df = df.pivot_table(index="q1_freq", columns="q2_freq", values="is_duplicate")
plt.figure(figsize=(12,12))
sns.heatmap(pvt_df)
plt.title("Mean is_duplicate value distribution across q1 and q2 frequency")
plt.show()
```

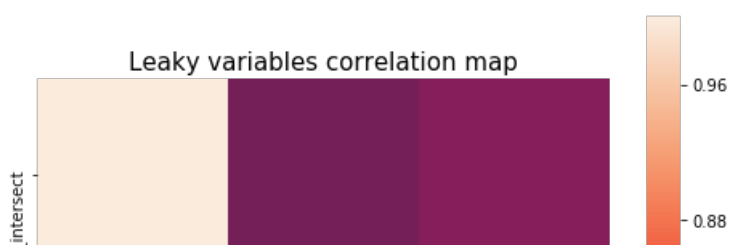


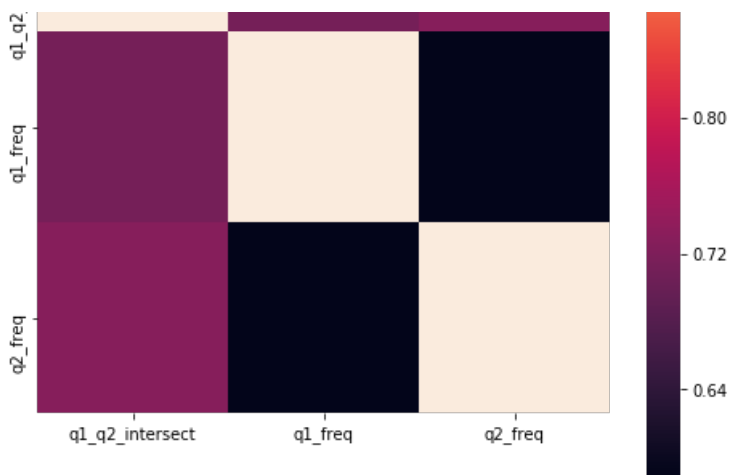
Let us also check the correlation between the three fields

In [68]:

```
cols_to_use = ['q1_q2_intersect', 'q1_freq', 'q2_freq']
temp_df = df[cols_to_use]
corrmat = temp_df.corr(method='spearman')
f, ax = plt.subplots(figsize=(8, 8))

# Draw the heatmap using seaborn
sns.heatmap(corrmat, vmax=1., square=True)
plt.title("Leaky variables correlation map", fontsize=15)
plt.show()
```





### 3.3.1 Analysis of some of the extracted features

- Here are some questions have only one single words.

In [15]:

```
print ("Minimum length of the questions in question1 : " , min(df['q1_n_words']))
print ("Minimum length of the questions in question2 : " , min(df['q2_n_words']))

print ("Number of Questions with minimum length [question1] :", df[df['q1_n_words']== 1].shape[0])
print ("Number of Questions with minimum length [question2] :", df[df['q2_n_words']== 1].shape[0])
```

```
Minimum length of the questions in question1 : 1
Minimum length of the questions in question2 : 1
Number of Questions with minimum length [question1] : 67
Number of Questions with minimum length [question2] : 24
```

#### 3.3.1.1 Feature: word\_share

In [16]:

```
plt.figure(figsize=(12, 8))

plt.subplot(1,2,1)
sns.violinplot(x = 'is_duplicate', y = 'word_share', data = df[0:])

plt.subplot(1,2,2)
sns.distplot(df[df['is_duplicate'] == 1.0]['word_share'][0:] , label = "1", color = 'red')
sns.distplot(df[df['is_duplicate'] == 0.0]['word_share'][0:] , label = "0" , color = 'blue' )
plt.show()
```

C:\Users\Shashank\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning:

Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

C:\Users\Shashank\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning:

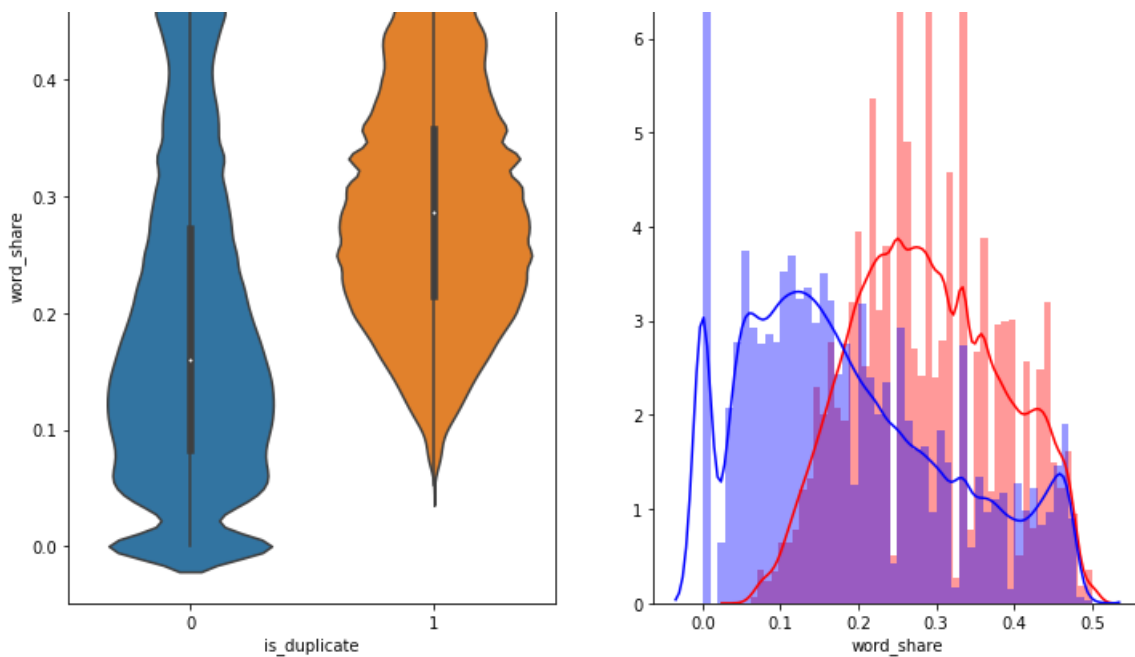
The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.

C:\Users\Shashank\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning:

The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.







- The distributions for normalized word\_share have some overlap on the far right-hand side, i.e., there are quite a lot of questions with high word similarity
- The average word share and Common no. of words of qid1 and qid2 is more when they are duplicate(Similar)

### 3.3.1.2 Feature: word\_Common

In [17]:

```
plt.figure(figsize=(12, 8))

plt.subplot(1,2,1)
sns.violinplot(x = 'is_duplicate', y = 'word_Common', data = df[0:])

plt.subplot(1,2,2)
sns.distplot(df[df['is_duplicate'] == 1.0]['word_Common'][0:], label = "1", color = 'red')
sns.distplot(df[df['is_duplicate'] == 0.0]['word_Common'][0:], label = "0", color = 'blue' )
plt.show()
```

C:\Users\Shashank\Anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning:

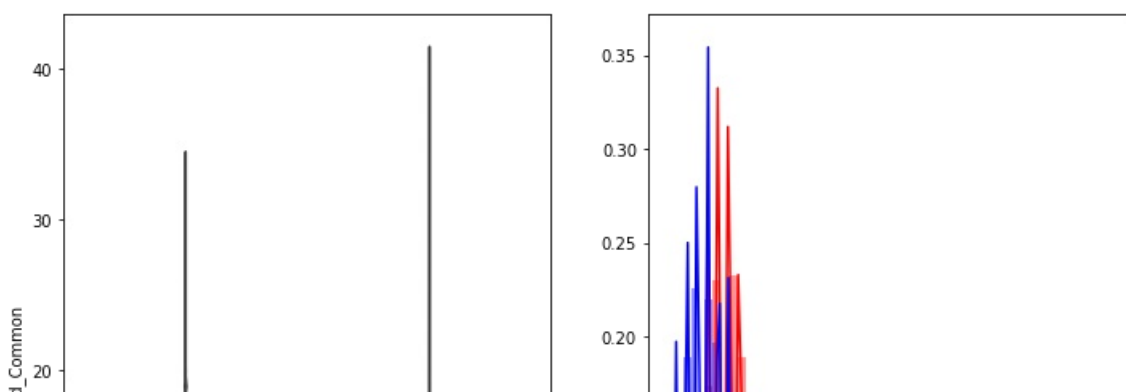
Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.

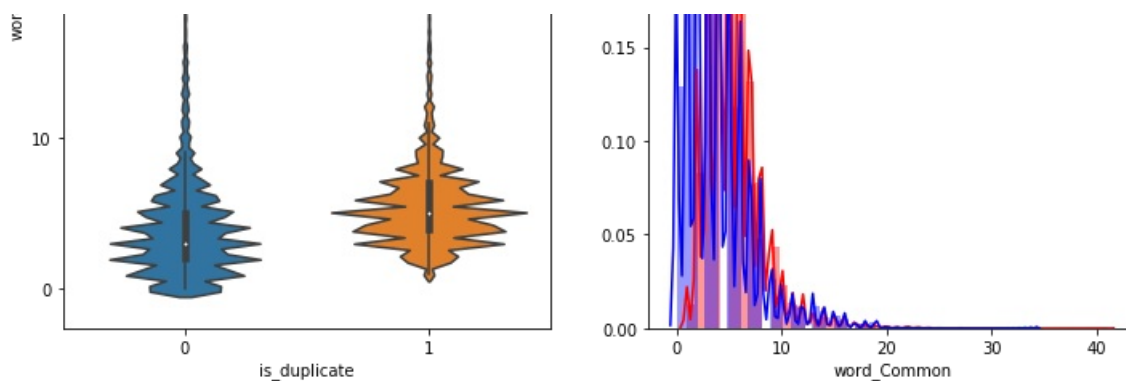
C:\Users\Shashank\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning:

The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.

C:\Users\Shashank\Anaconda3\lib\site-packages\matplotlib\axes\\_axes.py:6462: UserWarning:

The 'normed' kwarg is deprecated, and has been replaced by the 'density' kwarg.





The distributions of the word\_Common feature in similar and non-similar questions are highly overlapping

In [18]:

```
#https://stackoverflow.com/questions/12468179/unicodedecodeerror-utf8-codec-cant-decode-byte-0x9c
if os.path.isfile('df_fe_without_preprocessing_train.csv'):
    df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
    df = df.fillna('')
    df.head()
else:
    print("get df_fe_without_preprocessing_train.csv from drive or run the previous notebook")
```

In [19]:

```
df.head(2)
```

Out[19]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_C
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0

## 3.4 Preprocessing of Text

### - Preprocessing:

- Removing html tags
- Removing Punctuations
- Performing stemming
- Removing Stopwords
- Expanding contractions etc.

In [20]:

```
# m = re.compile('http://www|https://www')
```

```

# To get the results in 4 decimal points
SAFE_DIV = 0.0001

STOP_WORDS = stopwords.words("english")

def preprocess(x):
    x = str(x).lower()
    x = x.replace(",000,000", "m").replace(",000", "k").replace("'", "").replace("/", "").replace("won't", "will not").replace("cannot", "can not").replace("can", "can not")\
        .replace("n't", " not").replace("what's", "what is").replace("it's", "it is")\
        .replace("'ve", " have").replace("i'm", "i am").replace("'re", " are")\
        .replace("he's", "he is").replace("she's", "she is").replace("'s", " own")\
        .replace("%", " percent ").replace("₹", " rupee ").replace("$", " dollar ") \
        .replace("€", " euro ").replace("'ll", " will")
    x = re.sub(r"([0-9]+)000000", r"\1m", x)
    x = re.sub(r"([0-9]+)000", r"\1k", x)

    porter = PorterStemmer()
    pattern = re.compile('\W')

    if type(x) == type(''):
        x = re.sub(pattern, ' ', x)

    if type(x) == type(''):
        x = porter.stem(x)
        example1 = BeautifulSoup(x)
        x = example1.get_text()

    return x

```

Function to Compute and get the features : With 2 parameters of Question 1 and Question 2

## 3.5 Advanced Feature Extraction (NLP and Fuzzy Features)

Definition:

- **Token**: You get a token by splitting sentence a space
- **Stop\_Word** : stop words as per NLTK.
- **Word** : A token that is not a stop\_word

Features:

- **cwc\_min** : Ratio of common\_word\_count to min length of word count of Q1 and Q2  

$$\text{cwc\_min} = \text{common\_word\_count} / (\min(\text{len}(q1\_words), \text{len}(q2\_words)))$$
- **cwc\_max** : Ratio of common\_word\_count to max length of word count of Q1 and Q2  

$$\text{cwc\_max} = \text{common\_word\_count} / (\max(\text{len}(q1\_words), \text{len}(q2\_words)))$$
- **csc\_min** : Ratio of common\_stop\_count to min length of stop count of Q1 and Q2  

$$\text{csc\_min} = \text{common\_stop\_count} / (\min(\text{len}(q1\_stops), \text{len}(q2\_stops)))$$
- **csc\_max** : Ratio of common\_stop\_count to max length of stop count of Q1 and Q2  

$$\text{csc\_max} = \text{common\_stop\_count} / (\max(\text{len}(q1\_stops), \text{len}(q2\_stops)))$$
- **ctc\_min** : Ratio of common\_token\_count to min length of token count of Q1 and Q2  

$$\text{ctc\_min} = \text{common\_token\_count} / (\min(\text{len}(q1\_tokens), \text{len}(q2\_tokens)))$$
- **ctc\_max** : Ratio of common\_token\_count to max length of token count of Q1 and Q2  

$$\text{ctc\_max} = \text{common\_token\_count} / (\max(\text{len}(q1\_tokens), \text{len}(q2\_tokens)))$$

- **last\_word\_eq** : Check if First word of both questions is equal or not  
`last_word_eq = int(q1_tokens[-1] == q2_tokens[-1])`
- **first\_word\_eq** : Check if First word of both questions is equal or not  
`first_word_eq = int(q1_tokens[0] == q2_tokens[0])`
- **abs\_len\_diff** : Abs. length difference  
`abs_len_diff = abs(len(q1_tokens) - len(q2_tokens))`
- **mean\_len** : Average Token Length of both Questions  
`mean_len = (len(q1_tokens) + len(q2_tokens))/2`
- **fuzz\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **fuzz\_partial\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **token\_sort\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **token\_set\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **longest\_substr\_ratio** : Ratio of length longest common substring to min length of token count of Q1 and Q2  
`longest_substr_ratio = len(longest common substring) / (min(len(q1_tokens), len(q2_tokens)))`

In [21]:

```
def get_token_features(q1, q2):
    token_features = [0.0]*10

    # Converting the Sentence into Tokens:
    q1_tokens = q1.split()
    q2_tokens = q2.split()

    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
        return token_features

    # Get the non-stopwords in Questions
    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])

    #Get the stopwords in Questions
    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])

    # Get the common non-stopwords from Question pair
    common_word_count = len(q1_words.intersection(q2_words))

    # Get the common stopwords from Question pair
    common_stop_count = len(q1_stops.intersection(q2_stops))

    # Get the common Tokens from Question pair
    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))

    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)) + SAFE_DIV)
    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)) + SAFE_DIV)
    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)) + SAFE_DIV)
    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)) + SAFE_DIV)
    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)

    # Last word of both question is same or not
    token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])

    # First word of both question is same or not
    token_features[7] = int(q1_tokens[0] == q2_tokens[0])
```

```

token_features[8] = abs(len(q1_tokens) - len(q2_tokens))

#Average Token Length of both Questions
token_features[9] = (len(q1_tokens) + len(q2_tokens))/2
return token_features

# get the Longest Common sub string

def get_longest_substr_ratio(a, b):
    strs = list(distance.lcs substrings(a, b))
    if len(strs) == 0:
        return 0
    else:
        return len(strs[0]) / (min(len(a), len(b)) + 1)

def extract_features(df):
    # preprocessing each question
    df["question1"] = df["question1"].fillna("").apply(preprocess)
    df["question2"] = df["question2"].fillna("").apply(preprocess)

    print("token features...")

    # Merging Features with dataset

    token_features = df.apply(lambda x: get_token_features(x["question1"], x["question2"]), axis=1)

    df["cwc_min"] = list(map(lambda x: x[0], token_features))
    df["cwc_max"] = list(map(lambda x: x[1], token_features))
    df["csc_min"] = list(map(lambda x: x[2], token_features))
    df["csc_max"] = list(map(lambda x: x[3], token_features))
    df["ctc_min"] = list(map(lambda x: x[4], token_features))
    df["ctc_max"] = list(map(lambda x: x[5], token_features))
    df["last_word_eq"] = list(map(lambda x: x[6], token_features))
    df["first_word_eq"] = list(map(lambda x: x[7], token_features))
    df["abs_len_diff"] = list(map(lambda x: x[8], token_features))
    df["mean_len"] = list(map(lambda x: x[9], token_features))

    #Computing Fuzzy Features and Merging with Dataset

    # do read this blog: http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/
    # https://stackoverflow.com/questions/31806695/when-to-use-which-fuzz-function-to-compare-2-strings
    # https://github.com/seatgeek/fuzzywuzzy
    print("fuzzy features..")

    df["token_set_ratio"] = df.apply(lambda x: fuzz.token_set_ratio(x["question1"], x["question2"]), axis=1)
    # The token sort approach involves tokenizing the string in question, sorting the tokens alphabetically, and
    # then joining them back into a string We then compare the transformed strings with a simple ratio().
    df["token_sort_ratio"] = df.apply(lambda x: fuzz.token_sort_ratio(x["question1"], x["question2"]), axis=1)
    df["fuzz_ratio"] = df.apply(lambda x: fuzz.QRatio(x["question1"], x["question2"]), axis=1)
    df["fuzz_partial_ratio"] = df.apply(lambda x: fuzz.partial_ratio(x["question1"], x["question2"]), axis=1)
    df["longest_substr_ratio"] = df.apply(lambda x: get_longest_substr_ratio(x["question1"], x["question2"]), axis=1)
    return df

```

In [22]:

```

import warnings
warnings.filterwarnings("ignore")
if os.path.isfile('nlp_features_train.csv'):
    df = pd.read_csv("nlp_features_train.csv", encoding='latin-1')
    df.fillna('')
else:
    print("Extracting features for train:")
    df = pd.read_csv("train.csv")
    df = extract_features(df)
    df.to_csv("nlp_features_train.csv", index=False)
df.head(2)

```

Out[22]:

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	f
0	0	1	2	what is the step by step guide to invest in sh...	what is the step by step guide to invest in sh...	0	0.999980	0.833319	0.999983	0.999983	...	0.785709	0.0	1
1	1	3	4	what is the story of kohinoor koh i noor dia...	what would happen if the indian government sto...	0	0.799984	0.399996	0.749981	0.599988	...	0.466664	0.0	1

2 rows × 21 columns



## 3.5.1 Analysis of extracted features

### 3.5.1.1 Plotting Word clouds

- Creating Word Cloud of Duplicates and Non-Duplicates Question pairs
- We can observe the most frequent occurring words

In [23]:

```
df_duplicate = df[df['is_duplicate'] == 1]
dfp_nonduplicate = df[df['is_duplicate'] == 0]

# Converting 2d array of q1 and q2 and flatten the array: like {{1,2},{3,4}} to {1,2,3,4}
p = np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
n = np.dstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).flatten()

print ("Number of data points in class 1 (duplicate pairs) :",len(p))
print ("Number of data points in class 0 (non duplicate pairs) :",len(n))
```

Number of data points in class 1 (duplicate pairs) : 298526  
Number of data points in class 0 (non duplicate pairs) : 510054

In [24]:

```
# reading the text files and removing the Stop Words:
d = path.dirname('.')

textp_w = open(path.join(d, 'train_p.txt')).read()
textn_w = open(path.join(d, 'train_n.txt')).read()
stopwords = set(STOPWORDS)
stopwords.add("said")
stopwords.add("br")
stopwords.add(" ")
stopwords.remove("not")

stopwords.remove("no")
#stopwords.remove("good")
#stopwords.remove("love")
stopwords.remove("like")
#stopwords.remove("best")
#stopwords.remove("!")
print ("Total number of words in duplicate pair questions :",len(textp_w))
print ("Total number of words in non duplicate pair questions :",len(textn_w))
```

Total number of words in duplicate pair questions : 891367  
Total number of words in non duplicate pair questions : 33193130

### Word Clouds generated from duplicate pair question's text

In [25]:

```
wc = WordCloud(background_color="white", max_words=len(textp_w), stopwords=stopwords)
wc.generate(textp_w)
print ("Word Cloud for Duplicate Question pairs")
plt.imshow(wc, interpolation='bilinear')
plt.axis("off")
plt.show()
```

Word Cloud for Duplicate Question pairs

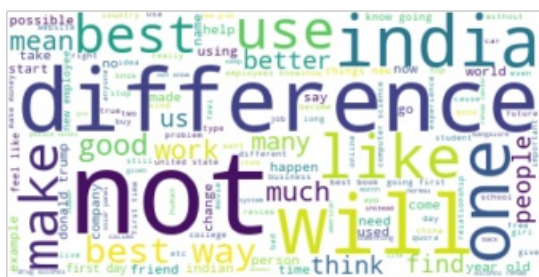


### Word Clouds generated from non duplicate pair question's text

In [26]:

```
wc = WordCloud(background_color="white", max_words=len(textn_w), stopwords=stopwords)
# generate word cloud
wc.generate(textn_w)
print ("Word Cloud for non-Duplicate Question pairs:")
plt.imshow(wc, interpolation='bilinear')
plt.axis("off")
plt.show()
```

Word Cloud for non-Duplicate Question pairs:

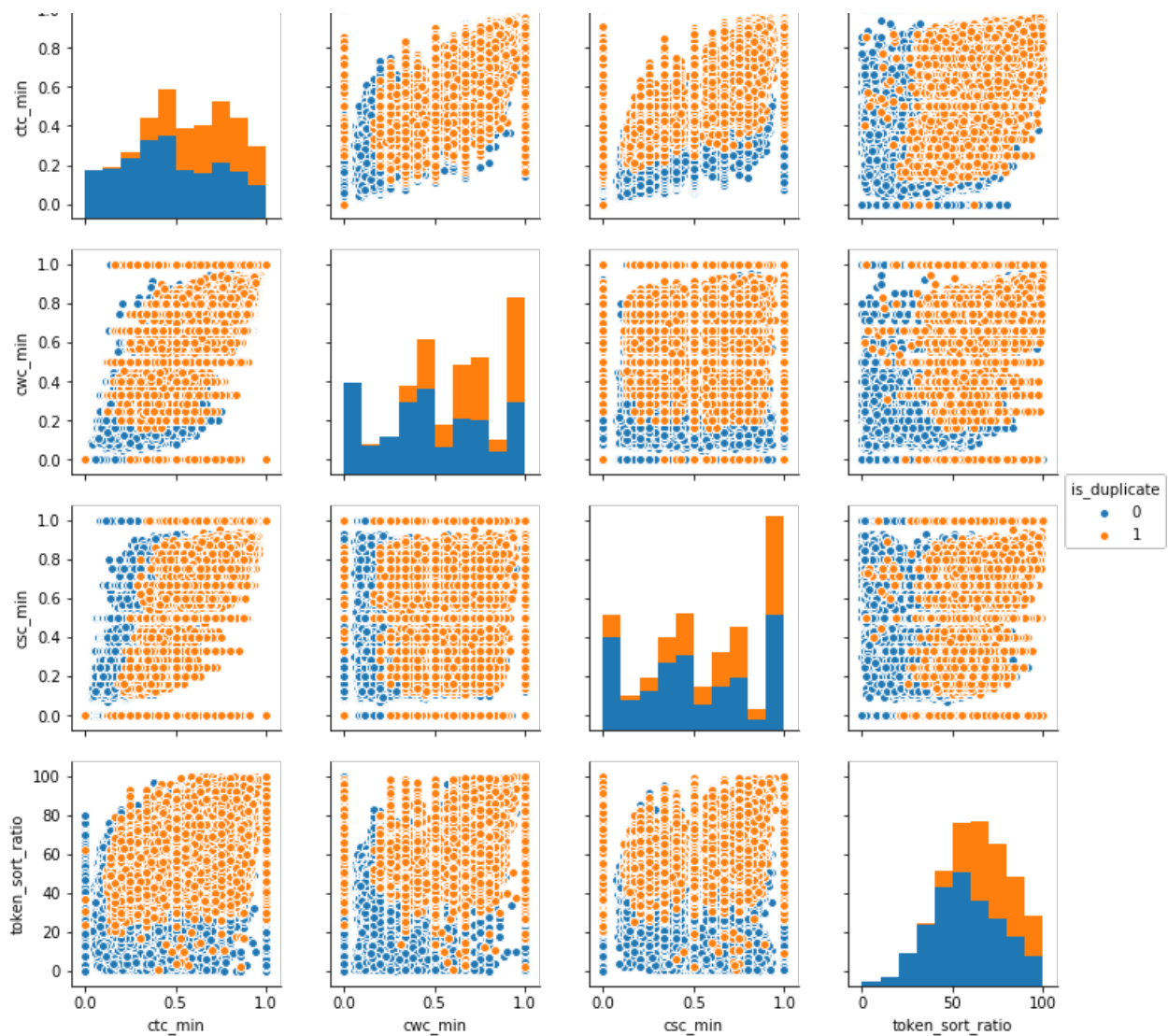


### 3.5.1.2 Pair plot of features ['ctc\_min', 'cwc\_min', 'csc\_min', 'token\_sort\_ratio']

In [27]:

```
n = df.shape[0]
sns.pairplot(df[['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio', 'is_duplicate']][0:n], hue='is_duplicate', vars=['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio'])
plt.show()
```



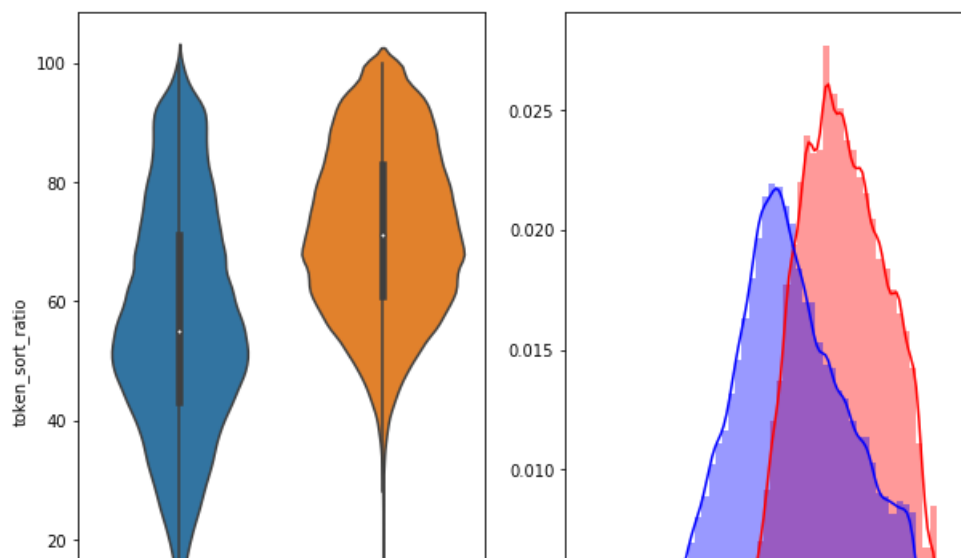


In [28]:

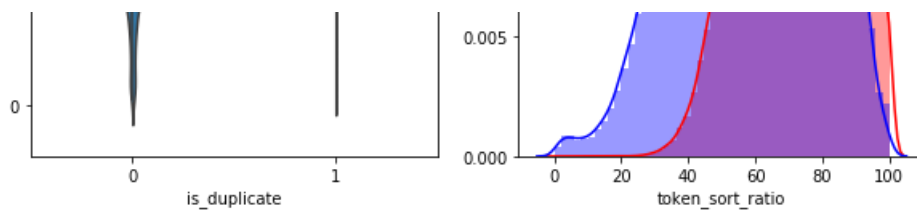
```
# Distribution of the token_sort_ratio
plt.figure(figsize=(10, 8))

plt.subplot(1,2,1)
sns.violinplot(x = 'is_duplicate', y = 'token_sort_ratio', data = df[0:] , )

plt.subplot(1,2,2)
sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = "1", color = 'red')
sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = "0" , color = 'blue' )
plt.show()
```





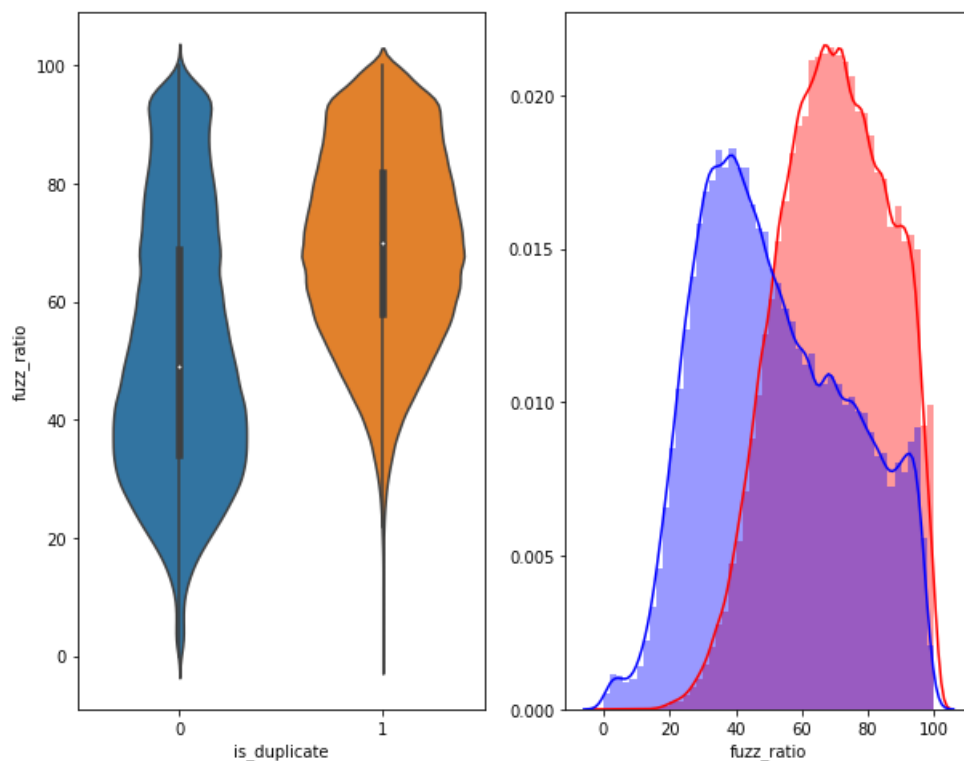


In [29]:

```
plt.figure(figsize=(10, 8))

plt.subplot(1,2,1)
sns.violinplot(x = 'is_duplicate', y = 'fuzz_ratio', data = df[0:] , )

plt.subplot(1,2,2)
sns.distplot(df[df['is_duplicate'] == 1.0]['fuzz_ratio'][0:] , label = "1", color = 'red')
sns.distplot(df[df['is_duplicate'] == 0.0]['fuzz_ratio'][0:] , label = "0" , color = 'blue' )
plt.show()
```



### 3.5.2 Visualization

In [30]:

```
# Using TSNE for Dimentionality reduction for 15 Features(Generated after cleaning the data) to 3
dimention

from sklearn.preprocessing import MinMaxScaler

dfp_subsampled = df[0:3000]
X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_min', 'csc_max',
'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_
ratio', 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio', 'longest_substr_ratio']])
y = dfp_subsampled['is_duplicate'].values
```

In [31]:

```
tsne2d = TSNE(
    n_components=2,
    init='random', # pca
```

```

random_state=101,
method='barnes_hut',
n_iter=1000,
verbose=2,
angle=0.5
).fit_transform(X)

```

```

[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 3000 samples in 0.023s...
[t-SNE] Computed neighbors for 3000 samples in 0.292s...
[t-SNE] Computed conditional probabilities for sample 1000 / 3000
[t-SNE] Computed conditional probabilities for sample 2000 / 3000
[t-SNE] Computed conditional probabilities for sample 3000 / 3000
[t-SNE] Mean sigma: 0.168100
[t-SNE] Computed conditional probabilities in 0.261s
[t-SNE] Iteration 50: error = 74.2438736, gradient norm = 0.0637462 (50 iterations in 2.442s)
[t-SNE] Iteration 100: error = 66.0404434, gradient norm = 0.0192010 (50 iterations in 1.760s)
[t-SNE] Iteration 150: error = 64.6748962, gradient norm = 0.0149074 (50 iterations in 1.649s)
[t-SNE] Iteration 200: error = 64.0936508, gradient norm = 0.0091195 (50 iterations in 1.545s)
[t-SNE] Iteration 250: error = 63.7366714, gradient norm = 0.0074293 (50 iterations in 1.541s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 63.736671
[t-SNE] Iteration 300: error = 1.3798265, gradient norm = 0.0011589 (50 iterations in 1.745s)
[t-SNE] Iteration 350: error = 1.0882488, gradient norm = 0.0004171 (50 iterations in 2.028s)
[t-SNE] Iteration 400: error = 0.9810162, gradient norm = 0.0002293 (50 iterations in 2.003s)
[t-SNE] Iteration 450: error = 0.9283460, gradient norm = 0.0001928 (50 iterations in 2.085s)
[t-SNE] Iteration 500: error = 0.9007203, gradient norm = 0.0001361 (50 iterations in 1.680s)
[t-SNE] Iteration 550: error = 0.8847462, gradient norm = 0.0001154 (50 iterations in 1.641s)
[t-SNE] Iteration 600: error = 0.8747767, gradient norm = 0.0001057 (50 iterations in 1.997s)
[t-SNE] Iteration 650: error = 0.8675045, gradient norm = 0.0000966 (50 iterations in 1.741s)
[t-SNE] Iteration 700: error = 0.8616202, gradient norm = 0.0000919 (50 iterations in 1.597s)
[t-SNE] Iteration 750: error = 0.8572435, gradient norm = 0.0000899 (50 iterations in 1.550s)
[t-SNE] Iteration 800: error = 0.8535899, gradient norm = 0.0000813 (50 iterations in 1.564s)
[t-SNE] Iteration 850: error = 0.8506539, gradient norm = 0.0000810 (50 iterations in 1.698s)
[t-SNE] Iteration 900: error = 0.8481503, gradient norm = 0.0000786 (50 iterations in 1.606s)
[t-SNE] Iteration 950: error = 0.8462221, gradient norm = 0.0000754 (50 iterations in 1.645s)
[t-SNE] Iteration 1000: error = 0.8442046, gradient norm = 0.0000767 (50 iterations in 1.577s)
[t-SNE] KL divergence after 1000 iterations: 0.844205

```

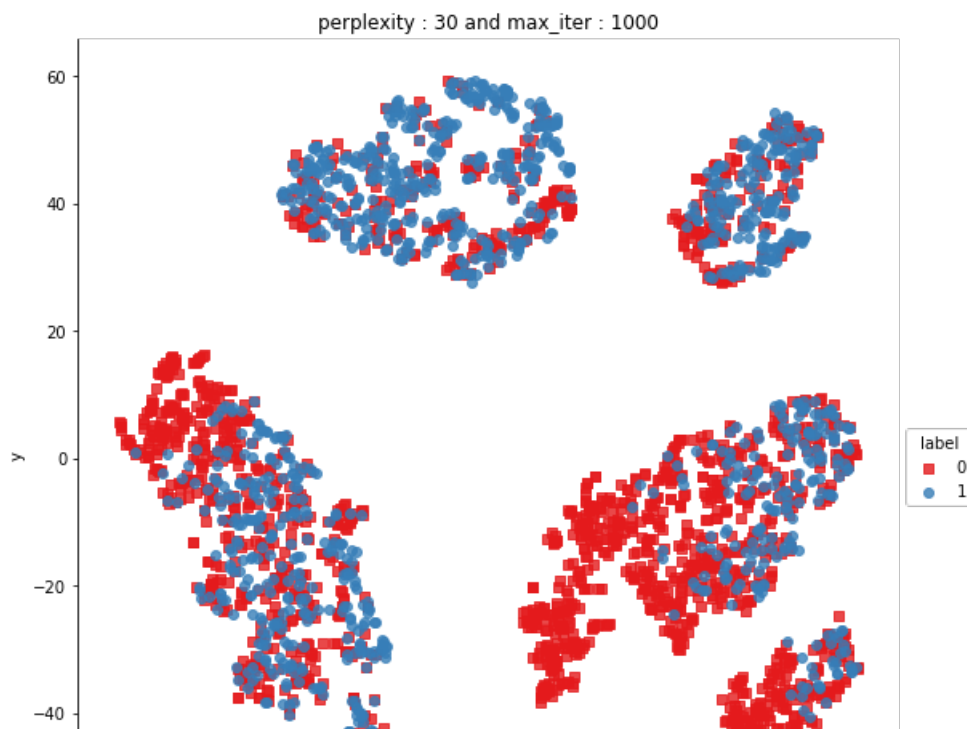
In [32]:

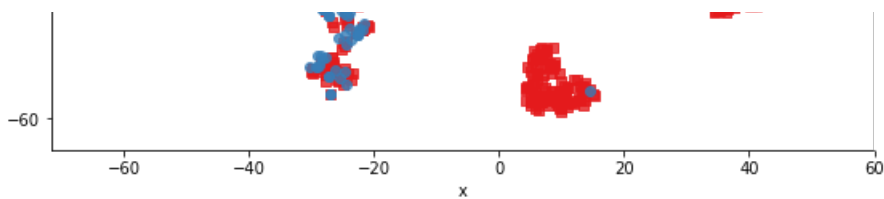
```

df = pd.DataFrame({'x':tsne2d[:,0], 'y':tsne2d[:,1], 'label':y})

# draw the plot in appropriate place in the grid
sns.lmplot(data=df, x='x', y='y', hue='label', fit_reg=False, size=8,palette="Set1",markers=['s','o'])
plt.title("perplexity : {} and max_iter : {}".format(30, 1000))
plt.show()

```





In [33]:

```
from sklearn.manifold import TSNE
tsne3d = TSNE(
    n_components=3,
    init='random', # pca
    random_state=101,
    method='barnes_hut',
    n_iter=1000,
    verbose=2,
    angle=0.5
).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 3000 samples in 0.011s...
[t-SNE] Computed neighbors for 3000 samples in 0.318s...
[t-SNE] Computed conditional probabilities for sample 1000 / 3000
[t-SNE] Computed conditional probabilities for sample 2000 / 3000
[t-SNE] Computed conditional probabilities for sample 3000 / 3000
[t-SNE] Mean sigma: 0.168100
[t-SNE] Computed conditional probabilities in 0.299s
[t-SNE] Iteration 50: error = 73.4699097, gradient norm = 0.0385700 (50 iterations in 9.763s)
[t-SNE] Iteration 100: error = 65.1441422, gradient norm = 0.0062584 (50 iterations in 4.012s)
[t-SNE] Iteration 150: error = 64.2269669, gradient norm = 0.0040364 (50 iterations in 3.570s)
[t-SNE] Iteration 200: error = 63.7893677, gradient norm = 0.0036378 (50 iterations in 3.457s)
[t-SNE] Iteration 250: error = 63.5331802, gradient norm = 0.0024612 (50 iterations in 3.348s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 63.533180
[t-SNE] Iteration 300: error = 1.1797704, gradient norm = 0.0006941 (50 iterations in 4.685s)
[t-SNE] Iteration 350: error = 0.9385557, gradient norm = 0.0001770 (50 iterations in 5.615s)
[t-SNE] Iteration 400: error = 0.8433141, gradient norm = 0.0001146 (50 iterations in 5.987s)
[t-SNE] Iteration 450: error = 0.8030387, gradient norm = 0.0000924 (50 iterations in 5.653s)
[t-SNE] Iteration 500: error = 0.7852126, gradient norm = 0.0000608 (50 iterations in 5.170s)
[t-SNE] Iteration 550: error = 0.7742532, gradient norm = 0.0000479 (50 iterations in 5.870s)
[t-SNE] Iteration 600: error = 0.7652450, gradient norm = 0.0000435 (50 iterations in 5.155s)
[t-SNE] Iteration 650: error = 0.7574351, gradient norm = 0.0000433 (50 iterations in 5.701s)
[t-SNE] Iteration 700: error = 0.7521394, gradient norm = 0.0000360 (50 iterations in 5.456s)
[t-SNE] Iteration 750: error = 0.7471731, gradient norm = 0.0000336 (50 iterations in 5.507s)
[t-SNE] Iteration 800: error = 0.7434089, gradient norm = 0.0000336 (50 iterations in 5.580s)
[t-SNE] Iteration 850: error = 0.7409297, gradient norm = 0.0000338 (50 iterations in 5.314s)
[t-SNE] Iteration 900: error = 0.7386305, gradient norm = 0.0000289 (50 iterations in 5.259s)
[t-SNE] Iteration 950: error = 0.7365760, gradient norm = 0.0000274 (50 iterations in 5.401s)
[t-SNE] Iteration 1000: error = 0.7344688, gradient norm = 0.0000251 (50 iterations in 5.353s)
[t-SNE] KL divergence after 1000 iterations: 0.734469
```

In [34]:

```
tracel = go.Scatter3d(
    x=tsne3d[:,0],
    y=tsne3d[:,1],
    z=tsne3d[:,2],
    mode='markers',
    marker=dict(
        sizemode='diameter',
        color = y,
        colorscale = 'Portland',
        colorbar = dict(title = 'duplicate'),
        line=dict(color='rgb(255, 255, 255)'),
        opacity=0.75
    )
)

data=[tracel]
layout=dict(height=800, width=800, title='3d embedding with engineered features')
fig=dict(data=data, layout=layout)
py.ipplot(fig, filename='3DBubble')
```

In [35]:

```
all_ques_df = pd.DataFrame(pd.concat([df['question1'], df['question2']]))
all_ques_df.columns = ["questions"]

all_ques_df["num of words"] = all_ques_df["questions"].apply(lambda x : len(str(x).split()))
```

```

KeyError                                Traceback (most recent call last)
~\Anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_loc(self, key, method, tolerance)
    3062         try:
-> 3063             return self._engine.get_loc(key)
    3064         except KeyError:

pandas\_libs\index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas\_libs\index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas\_libs\hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

pandas\_libs\hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

KeyError: 'question1'

```

During handling of the above exception, another exception occurred:

**KeyError** Traceback (most recent call last)

```

<ipython-input-35-754d55b51f98> in <module>()
----> 1 all_ques_df = pd.DataFrame(pd.concat([df['question1'], df['question2']]))
      2 all_ques_df.columns = ["questions"]
      3
      4 all_ques_df["num_of_words"] = all_ques_df["questions"].apply(lambda x : len(str(x).split()))
)

~\Anaconda3\lib\site-packages\pandas\core\frame.py in __getitem__(self, key)
    2683         return self._getitem_multilevel(key)
    2684     else:
-> 2685         return self._getitem_column(key)
    2686
    2687     def _getitem_column(self, key):

~\Anaconda3\lib\site-packages\pandas\core\frame.py in _getitem_column(self, key)
    2690         # get column
    2691         if self.columns.is_unique:
-> 2692             return self._get_item_cache(key)
    2693
    2694         # duplicate columns & possible reduce dimensionality

~\Anaconda3\lib\site-packages\pandas\core\generic.py in _get_item_cache(self, item)
    2484         res = cache.get(item)
    2485         if res is None:
-> 2486             values = self._data.get(item)
    2487             res = self._box_item_values(item, values)
    2488             cache[item] = res

~\Anaconda3\lib\site-packages\pandas\core\internals.py in get(self, item, fastpath)
    4113
    4114         if not isna(item):
-> 4115             loc = self.items.get_loc(item)
    4116         else:
    4117             indexer = np.arange(len(self.items))[isna(self.items)]

~\Anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_loc(self, key, method, tolerance)
    3063         return self._engine.get_loc(key)
    3064     except KeyError:
-> 3065         return self._engine.get_loc(self._maybe_cast_indexer(key))
    3066
    3067         indexer = self.get_indexer([key], method=method, tolerance=tolerance)

pandas\_libs\index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas\_libs\index.pyx in pandas._libs.index.IndexEngine.get_loc()

pandas\_libs\hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

pandas\_libs\hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item()

KeyError: 'question1'

```

In [ ]:

```

all_ques_df = pd.DataFrame(pd.concat([train_df['question1'], train_df['question2']]))
all_ques_df.columns = ["questions"]

all_ques_df["num_of_words"] = all_ques_df["questions"].apply(lambda x : len(str(x).split()))

```

## 3.6 Featurizing text data with tfidf weighted word-vectors

In [27]:

```

# avoid decoding problems
df = pd.read_csv("train.csv")

# encode questions to unicode
# https://stackoverflow.com/a/6812069
# ----- python 2 -----
# df['question1'] = df['question1'].apply(lambda x: unicode(str(x), "utf-8"))
# df['question2'] = df['question2'].apply(lambda x: unicode(str(x), "utf-8"))
#
# ----- python 3 -----

```

```
# ----- python 3 -----
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

In [28]:

```
df.head()
```

Out[28]:

	id	qid1	qid2	question1	question2	is_duplicate
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0
3	3	7	8	Why am I mentally very lonely? How can I solve...	Find the remainder when $23^{24}$ i...	0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

In [29]:

```
'''
# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039
def loadGloveModel(gloveFile):
    print ("Loading Glove Model")
    f = open(gloveFile,'r', encoding="utf8")
    model = {}
    for line in tqdm(f):
        splitLine = line.split()
        word = splitLine[0]
        embedding = np.array([float(val) for val in splitLine[1:]])
        model[word] = embedding
    print ("Done.",len(model)," words loaded!")
    return model
model = loadGloveModel('glove.42B.300d.txt')

# =====
Output:

Loading Glove Model
1917495it [06:32, 4879.69it/s]
Done. 1917495 words loaded!

# =====

words = []
for i in preproced_texts:
    words.extend(i.split(' '))

for i in preproced_titles:
    words.extend(i.split(' '))
print("all the words in the coupus", len(words))
words = set(words)
print("the unique words in the coupus", len(words))

inter_words = set(model.keys()).intersection(words)
print("The number of words that are present in both glove vectors and our coupus", \
      len(inter_words), "(",np.round(len(inter_words)/len(words)*100,3),"%")

words_courpus = {}
words_glove = set(model.keys())
for i in words:
    if i in words_glove:
        words_courpus[i] = model[i]
print("word 2 vec length", len(words_courpus))

# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
```

ve and load variables in python:

```
import pickle
with open('glove_vectors', 'wb') as f:
    pickle.dump(words_courpus, f)
```

```
'''
```

Out[29]:

```
'\n# Reading glove vectors in python: https://stackoverflow.com/a/38230349/4084039\ndef
loadGloveModel(gloveFile):\n    print ("Loading Glove Model")\n    f = open(gloveFile,\nencoding="utf8")\n    model = {}\n    for line in tqdm(f):\n        splitLine = line.split()\nword = splitLine[0]\n        embedding = np.array([float(val) for val in splitLine[1:]])\n    model[word] = embedding\n    print ("Done.",len(model)," words loaded!")\n    return model\nmodel =
loadGloveModel('glove.42B.300d.txt')\n\n# =====\n\nOutput:\n\nLoading G
love Model\n1917495it [06:32, 4879.69it/s]\nDone. 1917495 words loaded!\n\n#
=====
\n\nwords = []\nfor i in preprocod_texts:\n    words.extend(i.split('\n\n'))\n\nfor i in preprocod_titles:\n    words.extend(i.split('\n\n'))\n\nprint("all the words in the
coupus", len(words))\n\nwords = set(words)\n\nprint("the unique words in the coupus",
len(words))\n\n\ninter_words = set(model.keys()).intersection(words)\n\nprint("The number of words tha
t are present in both glove vectors and our coupus", len(inter_words),
"\n",np.round(len(inter_words)/len(words)*100,3),"%")\n\n\nwords_courpus = {}\n\nwords_glove =
set(model.keys())\n\nfor i in words:\n    if i in words_glove:\n        words_courpus[i] = model[i]\n\nr
print("word 2 vec length", len(words_courpus))\n\n\n\n# stronging variables into pickle files python
: http://www.jessicayung.com/how-to-use-pickle-to-save-and-load-variables-in-python/\n\n\nimport pic
kle\n\nwith open('glove_vectors', 'wb') as f:\n    pickle.dump(words_courpus, f)\n\n\n'
```

In [38]:

```
# stronging variables into pickle files python: http://www.jessicayung.com/how-to-use-pickle-to-sa
ve-and-load-variables-in-python/
# make sure you have the glove_vectors file
import pickle
with open('glove_vectors', 'rb') as f:
    model = pickle.load(f)
    glove_words = set(model.keys())
```

In [39]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
# merge texts
questions = list(df['question1']) + list(df['question2'])

tfidf = TfidfVectorizer(lowercase=False,)
tfidf.fit_transform(questions)
# we are converting a dictionary with word as a key, and the idf as a value
dictionary = dict(zip(tfidf.get_feature_names(), list(tfidf.idf_)))
tfidf_words = set(tfidf.get_feature_names())
```

In [40]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_qul = []; # the avg-w2v for each sentence/review is stored in this list
for qul in tqdm(list(df['question1'])): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in qul.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(qul.count(word)/len(qul.split())) # getting the tfidf value
for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_qul.append(vector)
print(len(tfidf_w2v_vectors_qul))
```

```
print(len(tfidf_w2v_vectors_qu1[0]))
df['q1_feats_m'] = list(tfidf_w2v_vectors_qu1)
```

100%|



404290  
300

In [41]:

```
# average Word2Vec
# compute average word2vec for each review.
tfidf_w2v_vectors_qu2 = []; # the avg-w2v for each sentence/review is stored in this list
for qu2 in tqdm(list(df['question2'])): # for each review/sentence
    vector = np.zeros(300) # as word vectors are of zero length
    tf_idf_weight = 0; # num of words with a valid vector in the sentence/review
    for word in qu2.split(): # for each word in a review/sentence
        if (word in glove_words) and (word in tfidf_words):
            vec = model[word] # getting the vector for each word
            # here we are multiplying idf value(dictionary[word]) and the tf
            value((sentence.count(word)/len(sentence.split())))
            tf_idf = dictionary[word]*(qu2.count(word)/len(qu2.split())) # getting the tfidf value
            for each word
            vector += (vec * tf_idf) # calculating tfidf weighted w2v
            tf_idf_weight += tf_idf
    if tf_idf_weight != 0:
        vector /= tf_idf_weight
    tfidf_w2v_vectors_qu2.append(vector)
print(len(tfidf_w2v_vectors_qu2))
print(len(tfidf_w2v_vectors_qu2[0]))
df['q2_feats_m'] = list(tfidf_w2v_vectors_qu2)
```

100%|



404290  
300

In [44]:

```
df['question1'].values
```

Out[44]:

```
array(['What is the step by step guide to invest in share market in india?',
      'What is the story of Kohinoor (Koh-i-Noor) Diamond?',
      'How can I increase the speed of my internet connection while using a VPN?',
      ..., 'What is one coin?',
      'What is the approx annual cost of living while studying in UIC Chicago, for an Indian
      student?',
      'What is like to have sex with cousin?'], dtype=object)
```

In [47]:

```
#prepro_features_train.csv (Simple Preprocessing Feartures)
#nlp_features_train.csv (NLP Features)
if os.path.isfile('nlp_features_train.csv'):
    dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
else:
    print("download nlp_features_train.csv from drive or run previous notebook")

if os.path.isfile('df_fe_without_preprocessing_train.csv'):
    dfppro = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
else:
    print("download df_fe_without_preprocessing_train.csv from drive or run previous notebook")
```



In [48]:

```
df1 = dfnlp.drop(['qid1','qid2','question1','question2'],axis=1)
df2 = dfppro.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
df3 = df.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
df3_q1 = pd.DataFrame(df3.q1_feats_m.values.tolist(), index= df3.index)
df3_q2 = pd.DataFrame(df3.q2_feats_m.values.tolist(), index= df3.index)
```

In [49]:

```
# dataframe of nlp features
df1.head()
```

Out[49]:

	id	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff
0	0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	2.0
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0
2	2	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	4.0
3	3	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	2.0
4	4	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	6.0

In [50]:

```
# data before preprocessing
df2.head()
```

Out[50]:

	id	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	word_Total	word_share	freq_q1+q2
0	0	1	1	66	57	14	12	10.0	23.0	0.434783	2
1	1	4	1	51	88	8	13	4.0	20.0	0.200000	5
2	2	1	1	73	59	14	10	4.0	24.0	0.166667	2
3	3	1	1	50	65	11	9	0.0	19.0	0.000000	2
4	4	3	1	76	39	13	7	2.0	20.0	0.100000	4

In [51]:

```
# Questions 1 tfidf weighted word2vec
df3_q1.head()
```

Out[51]:

0	(0, 65736)\t0.07960042526980049\n (0, 31890...
1	(0, 65736)\t0.07960042526980049\n (0, 31890...
2	(0, 65736)\t0.07960042526980049\n (0, 31890...
3	(0, 65736)\t0.07960042526980049\n (0, 31890...
4	(0, 65736)\t0.07960042526980049\n (0, 31890...

In [52]:

```
# Questions 2 tfidf weighted word2vec
df3_q2.head()
```

Out[52]:

۷۵۷ [ ۷۷ ] •

	<b>0</b>
<b>0</b>	(0, 60682)\t0.0829957010020072\n (0, 29680)...
<b>1</b>	(0, 60682)\t0.0829957010020072\n (0, 29680)...
<b>2</b>	(0, 60682)\t0.0829957010020072\n (0, 29680)...
<b>3</b>	(0, 60682)\t0.0829957010020072\n (0, 29680)...
<b>4</b>	(0, 60682)\t0.0829957010020072\n (0, 29680)...

In [53]:

```
print("Number of features in nlp dataframe :", df1.shape[1])
print("Number of features in preprocessed dataframe :", df2.shape[1])
print("Number of features in question1 w2v dataframe :", df3_q1.shape[1])
print("Number of features in question2 w2v dataframe :", df3_q2.shape[1])
print("Number of features in final dataframe :", df1.shape[1]+df2.shape[1]+df3_q1.shape[1]+df3_q2.
shape[1])
```

```
Number of features in nlp dataframe : 17
Number of features in preprocessed dataframe : 12
Number of features in question1 w2v dataframe : 1
Number of features in question2 w2v dataframe : 1
Number of features in final dataframe : 31
```

In [54]:

```
# storing the final features to csv file
if not os.path.isfile('final_features.csv'):
    df3_q1['id']=df1['id']
    df3_q2['id']=df1['id']
    df1 = df1.merge(df2, on='id',how='left')
    df2 = df3_q1.merge(df3_q2, on='id',how='left')
    result = df1.merge(df2, on='id',how='left')
    result.to_csv('final_features.csv')
```

## 4. Machine Learning Models

#### 4.1 Reading data from file and storing into sql table

In [55]:

```
#Creating db file from csv
if not os.path.isfile('train.db'):
    disk_engine = create_engine('sqlite:///train.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('final_features.csv', names=['Unnamed: 0', 'id', 'is_duplicate', 'cwc_min', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio', 'token_sort_ratio', 'fuzz_ratio', 'fuzz_partial_ratio', 'longest_substr_ratio', 'freq_qid1', 'freq_qid2', 'q1len', 'q2len', 'q1_n_words', 'q2_n_words', 'word_Common', 'word_Total', 'word_share', 'freq_q1+q2', 'freq_q1-q2', '0_x', '1_x', '2_x', '3_x', '4_x', '5_x', '6_x', '7_x', '8_x', '9_x', '10_x', '11_x', '12_x', '13_x', '14_x', '15_x', '16_x', '17_x', '18_x', '19_x', '20_x', '21_x', '22_x', '23_x', '24_x', '25_x', '26_x', '27_x', '28_x', '29_x', '30_x', '31_x', '32_x', '33_x', '34_x', '35_x', '36_x', '37_x', '38_x', '39_x', '40_x', '41_x', '42_x', '43_x', '44_x', '45_x', '46_x', '47_x', '48_x', '49_x', '50_x', '51_x', '52_x', '53_x', '54_x', '55_x', '56_x', '57_x', '58_x', '59_x', '60_x', '61_x', '62_x', '63_x', '64_x', '65_x', '66_x', '67_x', '68_x', '69_x', '70_x', '71_x', '72_x', '73_x', '74_x', '75_x', '76_x', '77_x', '78_x', '79_x', '80_x', '81_x', '82_x', '83_x', '84_x', '85_x', '86_x', '87_x', '88_x', '89_x', '90_x', '91_x', '92_x', '93_x', '94_x', '95_x', '96_x', '97_x', '98_x', '99_x', '100_x', '101_x', '102_x', '103_x', '104_x', '105_x', '106_x', '107_x', '108_x', '109_x', '110_x', '111_x', '112_x', '113_x', '114_x', '115_x', '116_x', '117_x', '118_x', '119_x', '120_x', '121_x', '122_x', '123_x', '124_x', '125_x', '126_x', '127_x', '128_x', '129_x', '130_x', '131_x', '132_x', '133_x', '134_x', '135_x', '136_x', '137_x', '138_x', '139_x', '140_x', '141_x', '142_x', '143_x', '144_x', '145_x', '146_x', '147_x', '148_x', '149_x', '150_x', '151_x', '152_x', '153_x', '154_x', '155_x', '156_x', '157_x', '158_x', '159_x', '160_x', '161_x', '162_x', '163_x', '164_x', '165_x', '166_x', '167_x', '168_x', '169_x', '170_x', '171_x', '172_x', '173_x', '174_x', '175_x', '176_x', '177_x', '178_x', '179_x', '180_x', '181_x', '182_x', '183_x', '184_x', '185_x', '186_x', '187_x', '188_x', '189_x', '190_x', '191_x', '192_x', '193_x', '194_x', '195_x', '196_x', '197_x', '198_x', '199_x', '200_x', '201_x', '202_x', '203_x', '204_x', '205_x', '206_x', '207_x', '208_x', '209_x', '210_x', '211_x', '212_x', '213_x', '214_x', '215_x', '216_x', '217_x', '218_x', '219_x', '220_x', '221_x', '222_x', '223_x', '224_x', '225_x', '226_x', '227_x', '228_x', '229_x', '230_x', '231_x', '232_x', '233_x', '234_x', '235_x', '236_x', '237_x', '238_x', '239_x', '240_x', '241_x', '242_x', '243_x', '244_x', '245_x', '246_x', '247_x', '248_x', '249_x', '250_x', '251_x', '252_x', '253_x', '254_x', '255_x', '256_x', '257_x', '258_x', '259_x', '260_x', '261_x', '262_x', '263_x', '264_x', '265_x', '266_x', '267_x', '268_x', '269_x', '270_x', '271_x', '272_x', '273_x', '274_x', '275_x', '276_x', '277_x', '278_x', '279_x', '280_x', '281_x', '282_x', '283_x', '284_x', '285_x', '286_x', '287_x', '288_x', '289_x', '290_x', '291_x', '292_x', '293_x', '294_x', '295_x', '296_x', '297_x', '298_x', '299_x', '300_x', '301_x', '302_x', '303_x', '304_x', '305_x', '306_x', '307_x', '308_x', '309_x', '310_x', '311_x', '312_x', '313_x', '314_x', '315_x', '316_x', '317_x', '318_x', '319_x', '320_x', '321_x', '322_x', '323_x', '324_x', '325_x', '326_x', '327_x', '328_x', '329_x', '330_x', '331_x', '332_x', '333_x', '334_x', '335_x', '336_x', '337_x', '338_x', '339_x', '340_x', '341_x', '342_x', '343_x', '344_x', '345_x', '346_x', '347_x', '348_x', '349_x', '350_x', '351_x', '352_x', '353_x', '354_x', '355_x', '356_x', '357_x', '358_x', '359_x', '360_x', '361_x', '362_x', '363_x', '364_x', '365_x', '366_x', '367_x', '368_x', '369_x', '370_x', '371_x', '372_x', '373_x', '374_x', '375_x', '376_x', '377_x', '378_x', '379_x', '380_x', '381_x', '382_x', '383_x', '384_x', '385_x', '386_x', '387_x', '388_x', '389_x', '390_x', '391_x', '392_x', '393_x', '394_x', '395_x', '396_x', '397_x', '398_x', '399_x', '400_x', '401_x', '402_x', '403_x', '404_x', '405_x', '406_x', '407_x', '408_x', '409_x', '410_x', '411_x', '412_x', '413_x', '414_x', '415_x', '416_x', '417_x', '418_x', '419_x', '420_x', '421_x', '422_x', '423_x', '424_x', '425_x', '426_x', '427_x', '428_x', '429_x', '430_x', '431_x', '432_x', '433_x', '434_x', '435_x', '436_x', '437_x', '438_x', '439_x', '440_x', '441_x', '442_x', '443_x', '444_x', '445_x', '446_x', '447_x', '448_x', '449_x', '450_x', '451_x', '452_x', '453_x', '454_x', '455_x', '456_x', '457_x', '458_x', '459_x', '460_x', '461_x', '462_x', '463_x', '464_x', '465_x', '466_x', '467_x', '468_x', '469_x', '470_x', '471_x', '472_x', '473_x', '474_x', '475_x', '476_x', '477_x', '478_x', '479_x', '480_x', '481_x', '482_x', '483_x', '484_x', '485_x', '486_x', '487_x', '488_x', '489_x', '490_x', '491_x', '492_x', '493_x', '494_x', '495_x', '496_x', '497_x', '498_x', '499_x', '500_x', '501_x', '502_x', '503_x', '504_x', '505_x', '506_x', '507_x', '508_x', '509_x', '510_x', '511_x', '512_x', '513_x', '514_x', '515_x', '516_x', '517_x', '518_x', '519_x', '520_x', '521_x', '522_x', '523_x', '524_x', '525_x', '526_x', '527_x', '528_x', '529_x', '530_x', '531_x', '532_x', '533_x', '534_x', '535_x', '536_x', '537_x', '538_x', '539_x', '540_x', '541_x', '542_x', '543_x', '544_x', '545_x', '546_x', '547_x', '548_x', '549_x', '550_x', '551_x', '552_x', '553_x', '554_x', '555_x', '556_x', '557_x', '558_x', '559_x', '560_x', '561_x', '562_x', '563_x', '564_x', '565_x', '566_x', '567_x', '568_x', '569_x', '570_x', '571_x', '572_x', '573_x', '574_x', '575_x', '576_x', '577_x', '578_x', '579_x', '580_x', '581_x', '582_x', '583_x', '584_x', '585_x', '586_x', '587_x', '588_x', '589_x', '590_x', '591_x', '592_x', '593_x', '594_x', '595_x', '596_x', '597_x', '598_x', '599_x', '600_x', '601_x', '602_x', '603_x', '604_x', '605_x', '606_x', '607_x', '608_x', '609_x', '610_x', '611_x', '612_x', '613_x', '614_x', '615_x', '616_x', '617_x', '618_x', '619_x', '620_x', '621_x', '622_x', '623_x', '624_x', '625_x', '626_x', '627_x', '628_x', '629_x', '630_x', '631_x', '632_x', '633_x', '634_x', '635_x', '636_x', '637_x', '638_x', '639_x', '640_x', '641_x', '642_x', '643_x', '644_x', '645_x', '646_x', '647_x', '648_x', '649_x', '650_x', '651_x', '652_x', '653_x', '654_x', '655_x', '656_x', '657_x', '658_x', '659_x', '660_x', '661_x', '662_x', '
```

```

4_x', '125_x', '126_x', '127_x', '128_x', '129_x', '130_x', '131_x', '132_x', '133_x', '134_x', '135_x', '136_x',
x', '137_x', '138_x', '139_x', '140_x', '141_x', '142_x', '143_x', '144_x', '145_x', '146_x', '147_x', '148_x',
, '149_x', '150_x', '151_x', '152_x', '153_x', '154_x', '155_x', '156_x', '157_x', '158_x', '159_x', '160_x',
161_x', '162_x', '163_x', '164_x', '165_x', '166_x', '167_x', '168_x', '169_x', '170_x', '171_x', '172_x', '173_x',
3_x', '174_x', '175_x', '176_x', '177_x', '178_x', '179_x', '180_x', '181_x', '182_x', '183_x', '184_x', '185_x',
x', '186_x', '187_x', '188_x', '189_x', '190_x', '191_x', '192_x', '193_x', '194_x', '195_x', '196_x', '197_x',
x', '198_x', '199_x', '200_x', '201_x', '202_x', '203_x', '204_x', '205_x', '206_x', '207_x', '208_x', '209_x',
210_x', '211_x', '212_x', '213_x', '214_x', '215_x', '216_x', '217_x', '218_x', '219_x', '220_x', '221_x', '222_x',
2_x', '223_x', '224_x', '225_x', '226_x', '227_x', '228_x', '229_x', '230_x', '231_x', '232_x', '233_x', '234_x',
x', '235_x', '236_x', '237_x', '238_x', '239_x', '240_x', '241_x', '242_x', '243_x', '244_x', '245_x', '246_x',
, '247_x', '248_x', '249_x', '250_x', '251_x', '252_x', '253_x', '254_x', '255_x', '256_x', '257_x', '258_x',
259_x', '260_x', '261_x', '262_x', '263_x', '264_x', '265_x', '266_x', '267_x', '268_x', '269_x', '270_x', '271_x',
1_x', '272_x', '273_x', '274_x', '275_x', '276_x', '277_x', '278_x', '279_x', '280_x', '281_x', '282_x', '283_x',
x', '284_x', '285_x', '286_x', '287_x', '288_x', '289_x', '290_x', '291_x', '292_x', '293_x', '294_x', '295_x',
, '296_x', '297_x', '298_x', '299_x', '300_x', '301_x', '302_x', '303_x', '304_x', '305_x', '306_x', '307_x', '308_x',
308_x', '309_x', '310_x', '311_x', '312_x', '313_x', '314_x', '315_x', '316_x', '317_x', '318_x', '319_x', '320_x',
0_x', '321_x', '322_x', '323_x', '324_x', '325_x', '326_x', '327_x', '328_x', '329_x', '330_x', '331_x', '332_x',
x', '333_x', '334_x', '335_x', '336_x', '337_x', '338_x', '339_x', '340_x', '341_x', '342_x', '343_x', '344_x',
, '345_x', '346_x', '347_x', '348_x', '349_x', '350_x', '351_x', '352_x', '353_x', '354_x', '355_x', '356_x',
357_x', '358_x', '359_x', '360_x', '361_x', '362_x', '363_x', '364_x', '365_x', '366_x', '367_x', '368_x', '369_x',
9_x', '370_x', '371_x', '372_x', '373_x', '374_x', '375_x', '376_x', '377_x', '378_x', '379_x', '380_x', '381_x',
x', '382_x', '383_x', '0_y', '1_y', '2_y', '3_y', '4_y', '5_y', '6_y', '7_y', '8_y', '9_y', '10_y', '11_y', '12_y',
, '13_y', '14_y', '15_y', '16_y', '17_y', '18_y', '19_y', '20_y', '21_y', '22_y', '23_y', '24_y', '25_y', '26_y',
27_y', '28_y', '29_y', '30_y', '31_y', '32_y', '33_y', '34_y', '35_y', '36_y', '37_y', '38_y', '39_y', '40_y',
41_y', '42_y', '43_y', '44_y', '45_y', '46_y', '47_y', '48_y', '49_y', '50_y', '51_y', '52_y', '53_y', '54_y', '55_y',
5_y', '56_y', '57_y', '58_y', '59_y', '60_y', '61_y', '62_y', '63_y', '64_y', '65_y', '66_y', '67_y', '68_y', '69_y',
y', '70_y', '71_y', '72_y', '73_y', '74_y', '75_y', '76_y', '77_y', '78_y', '79_y', '80_y', '81_y', '82_y', '83_y',
y', '84_y', '85_y', '86_y', '87_y', '88_y', '89_y', '90_y', '91_y', '92_y', '93_y', '94_y', '95_y', '96_y', '97_y',
, '98_y', '99_y', '100_y', '101_y', '102_y', '103_y', '104_y', '105_y', '106_y', '107_y', '108_y', '109_y', '110_y',
0_y', '111_y', '112_y', '113_y', '114_y', '115_y', '116_y', '117_y', '118_y', '119_y', '120_y', '121_y', '122_y',
y', '123_y', '124_y', '125_y', '126_y', '127_y', '128_y', '129_y', '130_y', '131_y', '132_y', '133_y', '134_y',
, '135_y', '136_y', '137_y', '138_y', '139_y', '140_y', '141_y', '142_y', '143_y', '144_y', '145_y', '146_y',
147_y', '148_y', '149_y', '150_y', '151_y', '152_y', '153_y', '154_y', '155_y', '156_y', '157_y', '158_y', '159_y',
9_y', '160_y', '161_y', '162_y', '163_y', '164_y', '165_y', '166_y', '167_y', '168_y', '169_y', '170_y', '171_y',
y', '172_y', '173_y', '174_y', '175_y', '176_y', '177_y', '178_y', '179_y', '180_y', '181_y', '182_y', '183_y',
, '184_y', '185_y', '186_y', '187_y', '188_y', '189_y', '190_y', '191_y', '192_y', '193_y', '194_y', '195_y',
196_y', '197_y', '198_y', '199_y', '200_y', '201_y', '202_y', '203_y', '204_y', '205_y', '206_y', '207_y', '208_y',
8_y', '209_y', '210_y', '211_y', '212_y', '213_y', '214_y', '215_y', '216_y', '217_y', '218_y', '219_y', '220_y',
y', '221_y', '222_y', '223_y', '224_y', '225_y', '226_y', '227_y', '228_y', '229_y', '230_y', '231_y', '232_y',
, '233_y', '234_y', '235_y', '236_y', '237_y', '238_y', '239_y', '240_y', '241_y', '242_y', '243_y', '244_y',
245_y', '246_y', '247_y', '248_y', '249_y', '250_y', '251_y', '252_y', '253_y', '254_y', '255_y', '256_y', '257_y',
7_y', '258_y', '259_y', '260_y', '261_y', '262_y', '263_y', '264_y', '265_y', '266_y', '267_y', '268_y', '269_y',
y', '270_y', '271_y', '272_y', '273_y', '274_y', '275_y', '276_y', '277_y', '278_y', '279_y', '280_y', '281_y',
, '282_y', '283_y', '284_y', '285_y', '286_y', '287_y', '288_y', '289_y', '290_y', '291_y', '292_y', '293_y',
294_y', '295_y', '296_y', '297_y', '298_y', '299_y', '300_y', '301_y', '302_y', '303_y', '304_y', '305_y', '306_y',
6_y', '307_y', '308_y', '309_y', '310_y', '311_y', '312_y', '313_y', '314_y', '315_y', '316_y', '317_y', '318_y',
y', '319_y', '320_y', '321_y', '322_y', '323_y', '324_y', '325_y', '326_y', '327_y', '328_y', '329_y', '330_y',
, '331_y', '332_y', '333_y', '334_y', '335_y', '336_y', '337_y', '338_y', '339_y', '340_y', '341_y', '342_y',
343_y', '344_y', '345_y', '346_y', '347_y', '348_y', '349_y', '350_y', '351_y', '352_y', '353_y', '354_y', '355_y',
5_y', '356_y', '357_y', '358_y', '359_y', '360_y', '361_y', '362_y', '363_y', '364_y', '365_y', '366_y', '367_y',
y', '368_y', '369_y', '370_y', '371_y', '372_y', '373_y', '374_y', '375_y', '376_y', '377_y', '378_y', '379_y',
, '380_y', '381_y', '382_y', '383_y'], chunksize=chunksize, iterator=True, encoding='utf-8', ):
    df.index += index_start
    j+=1
    print('{} rows'.format(j*chunksize))
    df.to_sql('data', disk_engine, if_exists='append')
    index_start = df.index[-1] + 1

```

In [56]:

```

#http://www.sqlitetutorial.net/sqlite-python/create-tables/
import sqlite3
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 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))
```

In [57]:

```
read_db = 'train.db'
conn_r = create_connection(read_db)
checkTableExists(conn_r)
conn_r.close()
```

Tables in the database:  
data

In [58]:

```
# try to sample data according to the computing power you have
if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        # for selecting first 1M rows
        # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)

        # for selecting random points
        data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 60000;", conn_r)
        conn_r.commit()

    conn_r.close()
```

In [59]:

```
# remove the first row
data.drop(data.index[0], inplace=True)
y_true = data['is_duplicate']
data.drop(['Unnamed: 0', 'id', 'index', 'is_duplicate'], axis=1, inplace=True)
```

In [60]:

```
data.head()
```

Out[60]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max
1	0.714275510349852	0.714275510349852	0.899991000089999	0.899991000089999	0.636360743814801	0.636360743814801
2	0.499987500312492	0.399992000159997	0.499987500312492	0.285710204139941	0.499993750078124	0.333330555578704
3	0.33332222259258	0.199996000079998	0.66664444518516	0.333327777870369	0.374995312558593	0.333329629670781
4	0.999950002499875	0.499987500312492	0.33332222259258	0.249993750156246	0.599988000239995	0.374995312558593
5	0.66664444518516	0.66664444518516	0.499987500312492	0.499987500312492	0.571420408279882	0.571420408279882

5 rows × 794 columns



## 4.2 Convertina strinas to numerics

In [61]:

```
# after we read from sql table each entry was read it as a string
# we convert all the features into numeric before we apply any model
cols = list(data.columns)
for i in cols:
    data[i] = data[i].apply(pd.to_numeric)
    print(i)
```

```
cwc_min
cwc_max
csc_min
csc_max
ctc_min
ctc_max
last_word_eq
first_word_eq
abs_len_diff
mean_len
token_set_ratio
token_sort_ratio
fuzz_ratio
fuzz_partial_ratio
longest_substr_ratio
freq_qid1
freq_qid2
q1len
q2len
q1_n_words
q2_n_words
word_Common
word_Total
word_share
freq_q1+q2
freq_q1-q2
0_x
1_x
2_x
3_x
4_x
5_x
6_x
7_x
8_x
9_x
10_x
11_x
12_x
13_x
14_x
15_x
16_x
17_x
18_x
19_x
20_x
21_x
22_x
23_x
24_x
25_x
26_x
27_x
28_x
29_x
30_x
31_x
32_x
33_x
34_x
35_x
36_x
37_x
38_x
```

39\_x  
40\_x  
41\_x  
42\_x  
43\_x  
44\_x  
45\_x  
46\_x  
47\_x  
48\_x  
49\_x  
50\_x  
51\_x  
52\_x  
53\_x  
54\_x  
55\_x  
56\_x  
57\_x  
58\_x  
59\_x  
60\_x  
61\_x  
62\_x  
63\_x  
64\_x  
65\_x  
66\_x  
67\_x  
68\_x  
69\_x  
70\_x  
71\_x  
72\_x  
73\_x  
74\_x  
75\_x  
76\_x  
77\_x  
78\_x  
79\_x  
80\_x  
81\_x  
82\_x  
83\_x  
84\_x  
85\_x  
86\_x  
87\_x  
88\_x  
89\_x  
90\_x  
91\_x  
92\_x  
93\_x  
94\_x  
95\_x  
96\_x  
97\_x  
98\_x  
99\_x  
100\_x  
101\_x  
102\_x  
103\_x  
104\_x  
105\_x  
106\_x  
107\_x  
108\_x  
109\_x  
110\_x  
111\_x  
112\_x  
113\_x  
114\_x  
115\_x

116\_x  
117\_x  
118\_x  
119\_x  
120\_x  
121\_x  
122\_x  
123\_x  
124\_x  
125\_x  
126\_x  
127\_x  
128\_x  
129\_x  
130\_x  
131\_x  
132\_x  
133\_x  
134\_x  
135\_x  
136\_x  
137\_x  
138\_x  
139\_x  
140\_x  
141\_x  
142\_x  
143\_x  
144\_x  
145\_x  
146\_x  
147\_x  
148\_x  
149\_x  
150\_x  
151\_x  
152\_x  
153\_x  
154\_x  
155\_x  
156\_x  
157\_x  
158\_x  
159\_x  
160\_x  
161\_x  
162\_x  
163\_x  
164\_x  
165\_x  
166\_x  
167\_x  
168\_x  
169\_x  
170\_x  
171\_x  
172\_x  
173\_x  
174\_x  
175\_x  
176\_x  
177\_x  
178\_x  
179\_x  
180\_x  
181\_x  
182\_x  
183\_x  
184\_x  
185\_x  
186\_x  
187\_x  
188\_x  
189\_x  
190\_x  
191\_x  
192\_x

193\_x  
194\_x  
195\_x  
196\_x  
197\_x  
198\_x  
199\_x  
200\_x  
201\_x  
202\_x  
203\_x  
204\_x  
205\_x  
206\_x  
207\_x  
208\_x  
209\_x  
210\_x  
211\_x  
212\_x  
213\_x  
214\_x  
215\_x  
216\_x  
217\_x  
218\_x  
219\_x  
220\_x  
221\_x  
222\_x  
223\_x  
224\_x  
225\_x  
226\_x  
227\_x  
228\_x  
229\_x  
230\_x  
231\_x  
232\_x  
233\_x  
234\_x  
235\_x  
236\_x  
237\_x  
238\_x  
239\_x  
240\_x  
241\_x  
242\_x  
243\_x  
244\_x  
245\_x  
246\_x  
247\_x  
248\_x  
249\_x  
250\_x  
251\_x  
252\_x  
253\_x  
254\_x  
255\_x  
256\_x  
257\_x  
258\_x  
259\_x  
260\_x  
261\_x  
262\_x  
263\_x  
264\_x  
265\_x  
266\_x  
267\_x  
268\_x  
269\_x



270\_x  
271\_x  
272\_x  
273\_x  
274\_x  
275\_x  
276\_x  
277\_x  
278\_x  
279\_x  
280\_x  
281\_x  
282\_x  
283\_x  
284\_x  
285\_x  
286\_x  
287\_x  
288\_x  
289\_x  
290\_x  
291\_x  
292\_x  
293\_x  
294\_x  
295\_x  
296\_x  
297\_x  
298\_x  
299\_x  
300\_x  
301\_x  
302\_x  
303\_x  
304\_x  
305\_x  
306\_x  
307\_x  
308\_x  
309\_x  
310\_x  
311\_x  
312\_x  
313\_x  
314\_x  
315\_x  
316\_x  
317\_x  
318\_x  
319\_x  
320\_x  
321\_x  
322\_x  
323\_x  
324\_x  
325\_x  
326\_x  
327\_x  
328\_x  
329\_x  
330\_x  
331\_x  
332\_x  
333\_x  
334\_x  
335\_x  
336\_x  
337\_x  
338\_x  
339\_x  
340\_x  
341\_x  
342\_x  
343\_x  
344\_x  
345\_x  
346\_x

347\_x  
348\_x  
349\_x  
350\_x  
351\_x  
352\_x  
353\_x  
354\_x  
355\_x  
356\_x  
357\_x  
358\_x  
359\_x  
360\_x  
361\_x  
362\_x  
363\_x  
364\_x  
365\_x  
366\_x  
367\_x  
368\_x  
369\_x  
370\_x  
371\_x  
372\_x  
373\_x  
374\_x  
375\_x  
376\_x  
377\_x  
378\_x  
379\_x  
380\_x  
381\_x  
382\_x  
383\_x  
0\_y  
1\_y  
2\_y  
3\_y  
4\_y  
5\_y  
6\_y  
7\_y  
8\_y  
9\_y  
10\_y  
11\_y  
12\_y  
13\_y  
14\_y  
15\_y  
16\_y  
17\_y  
18\_y  
19\_y  
20\_y  
21\_y  
22\_y  
23\_y  
24\_y  
25\_y  
26\_y  
27\_y  
28\_y  
29\_y  
30\_y  
31\_y  
32\_y  
33\_y  
34\_y  
35\_y  
36\_y  
37\_y  
38\_y  
39\_y

39\_y  
40\_y  
41\_y  
42\_y  
43\_y  
44\_y  
45\_y  
46\_y  
47\_y  
48\_y  
49\_y  
50\_y  
51\_y  
52\_y  
53\_y  
54\_y  
55\_y  
56\_y  
57\_y  
58\_y  
59\_y  
60\_y  
61\_y  
62\_y  
63\_y  
64\_y  
65\_y  
66\_y  
67\_y  
68\_y  
69\_y  
70\_y  
71\_y  
72\_y  
73\_y  
74\_y  
75\_y  
76\_y  
77\_y  
78\_y  
79\_y  
80\_y  
81\_y  
82\_y  
83\_y  
84\_y  
85\_y  
86\_y  
87\_y  
88\_y  
89\_y  
90\_y  
91\_y  
92\_y  
93\_y  
94\_y  
95\_y  
96\_y  
97\_y  
98\_y  
99\_y  
100\_y  
101\_y  
102\_y  
103\_y  
104\_y  
105\_y  
106\_y  
107\_y  
108\_y  
109\_y  
110\_y  
111\_y  
112\_y  
113\_y  
114\_y  
115\_y  
116\_y

116\_y  
117\_y  
118\_y  
119\_y  
120\_y  
121\_y  
122\_y  
123\_y  
124\_y  
125\_y  
126\_y  
127\_y  
128\_y  
129\_y  
130\_y  
131\_y  
132\_y  
133\_y  
134\_y  
135\_y  
136\_y  
137\_y  
138\_y  
139\_y  
140\_y  
141\_y  
142\_y  
143\_y  
144\_y  
145\_y  
146\_y  
147\_y  
148\_y  
149\_y  
150\_y  
151\_y  
152\_y  
153\_y  
154\_y  
155\_y  
156\_y  
157\_y  
158\_y  
159\_y  
160\_y  
161\_y  
162\_y  
163\_y  
164\_y  
165\_y  
166\_y  
167\_y  
168\_y  
169\_y  
170\_y  
171\_y  
172\_y  
173\_y  
174\_y  
175\_y  
176\_y  
177\_y  
178\_y  
179\_y  
180\_y  
181\_y  
182\_y  
183\_y  
184\_y  
185\_y  
186\_y  
187\_y  
188\_y  
189\_y  
190\_y  
191\_y  
192\_y  
193\_y

193\_y  
194\_y  
195\_y  
196\_y  
197\_y  
198\_y  
199\_y  
200\_y  
201\_y  
202\_y  
203\_y  
204\_y  
205\_y  
206\_y  
207\_y  
208\_y  
209\_y  
210\_y  
211\_y  
212\_y  
213\_y  
214\_y  
215\_y  
216\_y  
217\_y  
218\_y  
219\_y  
220\_y  
221\_y  
222\_y  
223\_y  
224\_y  
225\_y  
226\_y  
227\_y  
228\_y  
229\_y  
230\_y  
231\_y  
232\_y  
233\_y  
234\_y  
235\_y  
236\_y  
237\_y  
238\_y  
239\_y  
240\_y  
241\_y  
242\_y  
243\_y  
244\_y  
245\_y  
246\_y  
247\_y  
248\_y  
249\_y  
250\_y  
251\_y  
252\_y  
253\_y  
254\_y  
255\_y  
256\_y  
257\_y  
258\_y  
259\_y  
260\_y  
261\_y  
262\_y  
263\_y  
264\_y  
265\_y  
266\_y  
267\_y  
268\_y  
269\_y  
270\_y

270\_y  
271\_y  
272\_y  
273\_y  
274\_y  
275\_y  
276\_y  
277\_y  
278\_y  
279\_y  
280\_y  
281\_y  
282\_y  
283\_y  
284\_y  
285\_y  
286\_y  
287\_y  
288\_y  
289\_y  
290\_y  
291\_y  
292\_y  
293\_y  
294\_y  
295\_y  
296\_y  
297\_y  
298\_y  
299\_y  
300\_y  
301\_y  
302\_y  
303\_y  
304\_y  
305\_y  
306\_y  
307\_y  
308\_y  
309\_y  
310\_y  
311\_y  
312\_y  
313\_y  
314\_y  
315\_y  
316\_y  
317\_y  
318\_y  
319\_y  
320\_y  
321\_y  
322\_y  
323\_y  
324\_y  
325\_y  
326\_y  
327\_y  
328\_y  
329\_y  
330\_y  
331\_y  
332\_y  
333\_y  
334\_y  
335\_y  
336\_y  
337\_y  
338\_y  
339\_y  
340\_y  
341\_y  
342\_y  
343\_y  
344\_y  
345\_y  
346\_y  
347\_y

34 / \_y  
348 \_y  
349 \_y  
350 \_y  
351 \_y  
352 \_y  
353 \_y  
354 \_y  
355 \_y  
356 \_y  
357 \_y  
358 \_y  
359 \_y  
360 \_y  
361 \_y  
362 \_y  
363 \_y  
364 \_y  
365 \_y  
366 \_y  
367 \_y  
368 \_y  
369 \_y  
370 \_y  
371 \_y  
372 \_y  
373 \_y  
374 \_y  
375 \_y  
376 \_y  
377 \_y  
378 \_y  
379 \_y  
380 \_y  
381 \_y  
382 \_y  
383 \_y

In [62]:

```
# https://stackoverflow.com/questions/7368789/convert-all-strings-in-a-list-to-int
y_true = list(map(int, y_true.values))
```

4.3 Random train test split( 70:30)

In [30]:

```
nlp_features_train=pd.read_csv('nlp_features_train.csv',encoding='latin-1')
```

In [31]:

```
nlp_features_train.head()
```

Out[31]:

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	fi
0	0	1	2	what is the step by step guide to invest in sh...	what is the step by step guide to invest in sh...	0	0.999980	0.833319	0.999983	0.999983	...	0.785709	0.0	1
1	1	3	4	what is the story of kohinoor ...	what would happen if the indian government	0	0.799984	0.399996	0.749981	0.599988	...	0.466664	0.0	1

	id	qid1	qid2	KOHINOOR question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	f
2	2	5	6	how can i increase the speed of my internet co...	how can internet speed be increased by hacking...	0	0.399992	0.333328	0.399992	0.249997	...	0.285712	0.0	1
3	3	7	8	why am i mentally very lonely how can i solve...	find the remainder when math 23 24 math i...	0	0.000000	0.000000	0.000000	0.000000	...	0.000000	0.0	0
4	4	9	10	which one dissolve in water quikly sugar salt...	which fish would survive in salt water	0	0.399992	0.199998	0.999950	0.666644	...	0.307690	0.0	1

5 rows × 21 columns

◀		▶
---	--	---

In [32]:

```
df_fe_without_preprocessing_train=pd.read_csv('df_fe_without_preprocessing_train.csv',encoding='latin-1')
```

In [33]:

```
df_fe_without_preprocessing_train.head(3)
```

Out[33]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_C
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0
1	1	3	4	What is the story of Kohinoor (Koh-i- Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0	1	1	73	59	14	10	4.0

◀		▶
---	--	---

In [34]:

```
d3=nlp_features_train.drop(['qid1','qid2'],axis=1)
d4=df_fe_without_preprocessing_train.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
data=d3.merge(d4,on='id',how='left')
```



In [35]:

```
y_class = data['is_duplicate']
data.drop(['id','is_duplicate'], axis=1, inplace=True)
```

In [37]:

```
#splitting the data into train and test dataset
from sklearn.model_selection import train_test_split
X_train,X_test, y_train, y_test = train_test_split(data, y_class, stratify=y_class, test_size=0.3)
```

In [43]:

```
print(X_train.shape)
print(X_test.shape)
```

```
(283003, 26)
(121287, 26)
```

In [40]:

```
#transforming question1 and question2 of train dataset
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf1 = TfidfVectorizer()
train_q1 = tfidf1.fit_transform(X_train['question1'].values.astype('U'))
test_q1 =tfidf1.transform(X_test['question1'].values.astype('U'))
```

In [41]:

```
#transforming question1 and question2 of test dataset
tfidf2=TfidfVectorizer()
train_q2 = tfidf2.fit_transform(X_train['question2'].values.astype('U'))
test_q2 =tfidf2.transform(X_test['question2'].values.astype('U'))
```

In [42]:

```
#Dropping question-1 and question-2 and Replacing with tfidf values
X_train.drop(['question1','question2'], axis=1, inplace=True)
X_test.drop(['question1','question2'], axis=1, inplace=True)
```

In [45]:

```
#Combining Question-1 and Question-2
from scipy.sparse import coo_matrix, hstack
train_tfidf = hstack((train_q1,train_q2))
test_tfidf = hstack((test_q1,test_q2))
```

In [46]:

```
#combining all basic,advanced and tfidf features
X_train = hstack((X_train, train_tfidf)).tocsr()
X_test = hstack((X_test, test_tfidf)).tocsr()
```

In [47]:

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

```
Number of data points in train data : (283003, 112084)
Number of data points in test data : (121287, 112084)
```

In [68]:

```
# This function plots the confusion matrices given y_i, y_i_hat.
from sklearn.metrics import confusion_matrix
def plot_confusion_matrix(test_y, predict_y):
```

```

C = confusion_matrix(test_y, predict_y)
# C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

A = ((C.T)/(C.sum(axis=1))).T)
#divid each element of the confusion matrix with the sum of elements in that column

# C = [[1, 2],
#       [3, 4]]
# C.T = [[1, 3],
#         [2, 4]]
# C.sum(axis=1) axis=0 corresonds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axix =1) = [[3, 7]]
# ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7],
#                             [2/3, 4/7]]

# ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3],
#                               [3/7, 4/7]]
# sum of row elements = 1

B = (C/C.sum(axis=0))
#divid each element of the confusion matrix with the sum of elements in that row
# C = [[1, 2],
#       [3, 4]]
# C.sum(axis=0) axis=0 corresonds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axix =0) = [[4, 6]]
# (C/C.sum(axis=0)) = [[1/4, 2/6],
#                       [3/4, 4/6]]
plt.figure(figsize=(20,4))

labels = [1,2]
# representing A in heatmap format
cmap=sns.light_palette("blue")
plt.subplot(1, 3, 1)
sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.title("Confusion matrix")

plt.subplot(1, 3, 2)
sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.title("Precision matrix")

plt.subplot(1, 3, 3)
# representing B in heatmap format
sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.title("Recall matrix")

plt.show()

```

## 4.4 Building a random model (Finding worst-case log-loss)

In [69]:

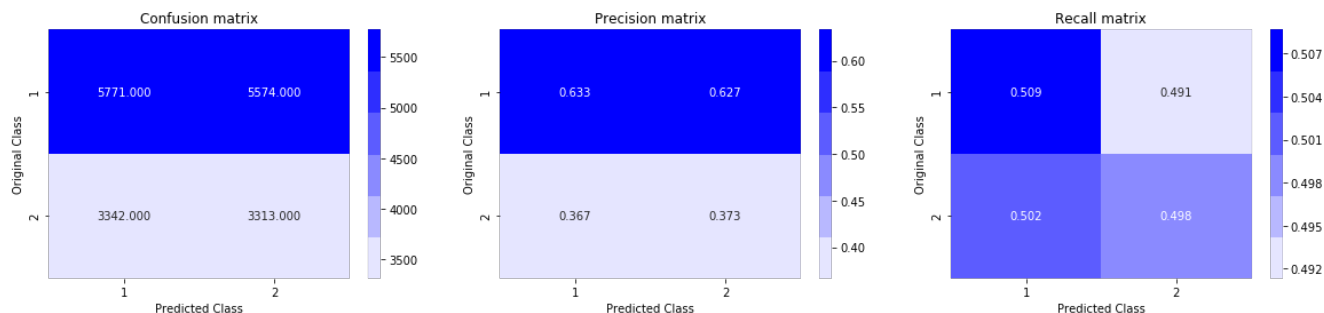
```

# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their sum
# ref: https://stackoverflow.com/a/18662466/4084039
# we create a output array that has exactly same size as the CV data
from sklearn.metrics import log_loss
predicted_y = np.zeros((test_len,2))
for i in range(test_len):
    rand_probs = np.random.rand(1,2)
    predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-15))

predicted_y =np.argmax(predicted_y, axis=1)
plot_confusion_matrix(y_test, predicted_y)

```

Log loss on Test Data using Random Model 0.8837470522382065



## 4.4 Logistic Regression with hyperparameter tuning

In [72]:

```
alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----

from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV
log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

fig, ax = plt.subplots()
ax.plot(alpha, log_error_array,c='g')
for i, txt in enumerate(np.round(log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train, y_train)

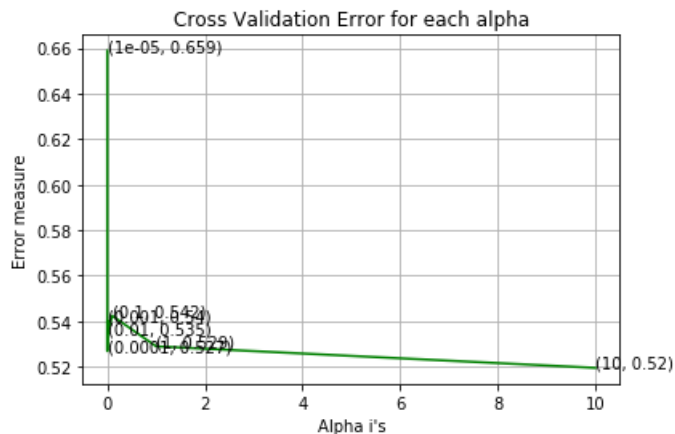
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
```

```

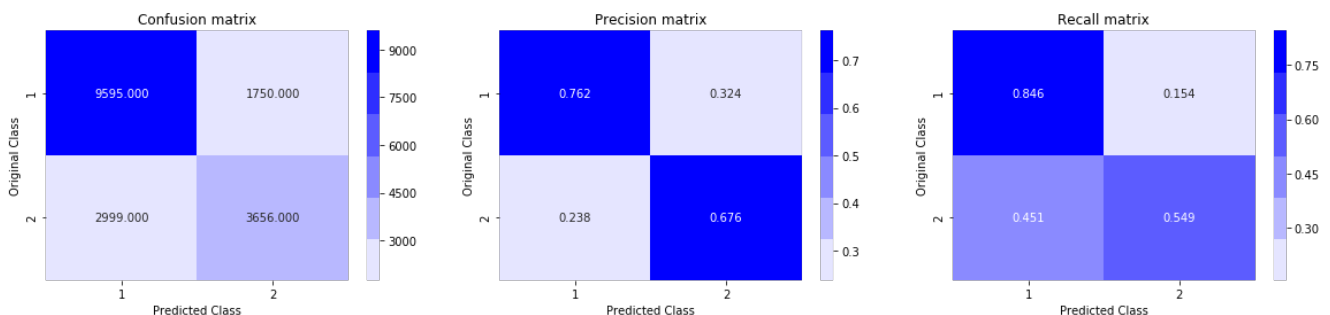
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

```

For values of alpha = 1e-05 The log loss is: 0.658807677966553  
 For values of alpha = 0.0001 The log loss is: 0.5268692295041751  
 For values of alpha = 0.001 The log loss is: 0.5402688688457457  
 For values of alpha = 0.01 The log loss is: 0.5347628254014818  
 For values of alpha = 0.1 The log loss is: 0.5424147491669886  
 For values of alpha = 1 The log loss is: 0.5289762417303445  
 For values of alpha = 10 The log loss is: 0.5195220613606596



For values of best alpha = 10 The train log loss is: 0.5148048328037605  
 For values of best alpha = 10 The test log loss is: 0.5195220613606596  
 Total number of data points : 18000



## 4.5 Linear SVM with hyperparameter tuning

In [74]:

```

alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-
learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_i
ter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:

```

```

# VIDEO LINK:
#-----

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
    clf.fit(X_train, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

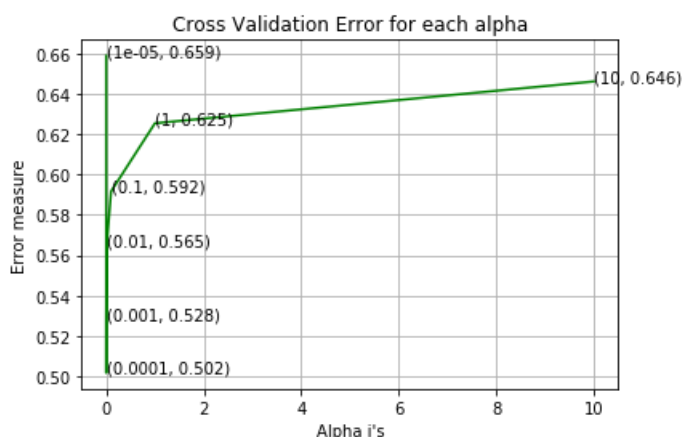
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', random_state=42)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train, y_train)

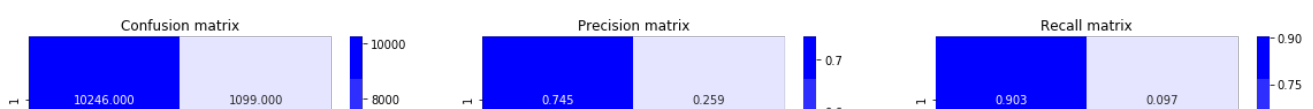
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

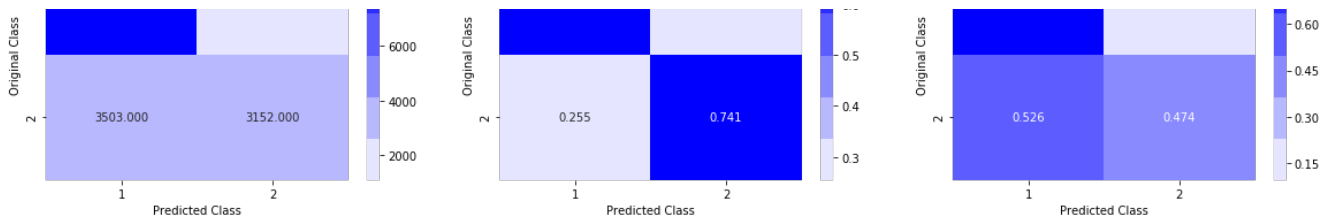
```

For values of alpha = 1e-05 The log loss is: 0.658807677966553  
 For values of alpha = 0.0001 The log loss is: 0.501758546876382  
 For values of alpha = 0.001 The log loss is: 0.5282452062377582  
 For values of alpha = 0.01 The log loss is: 0.5646235090120997  
 For values of alpha = 0.1 The log loss is: 0.5916062464007461  
 For values of alpha = 1 The log loss is: 0.6253664924902231  
 For values of alpha = 10 The log loss is: 0.64605795435085



For values of best alpha = 0.0001 The train log loss is: 0.4961221923159768  
 For values of best alpha = 0.0001 The test log loss is: 0.501758546876382  
 Total number of data points : 18000





## 4.6 XGBoost

In [81]:

```
import xgboost as xgb
params = {}
params['objective'] = 'binary:logistic'
params['eval_metric'] = 'logloss'
params['eta'] = 0.02
params['max_depth'] = 4

d_train = xgb.DMatrix(new_x_train, label=new_y_train)
d_test = xgb.DMatrix(new_x_test, label=new_y_test)

watchlist = [(d_train, 'train'), (d_test, 'valid')]

bst = xgb.train(params, d_train, 400, watchlist, early_stopping_rounds=20, verbose_eval=10)

xgdmatrix = xgb.DMatrix(new_x_train, new_y_train)
predict_y = bst.predict(d_test)
print("The test log loss is:", log_loss(new_y_test, predict_y, labels=clf.classes_, eps=1e-15))
```

```
[00:18:49] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[0] train-logloss:0.6848 valid-logloss:0.685058
Multiple eval metrics have been passed: 'valid-logloss' will be used for early stopping.
```

Will train until valid-logloss hasn't improved in 20 rounds.

```
[00:18:50] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:18:50] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:51] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:18:51] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:51] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:18:52] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:52] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:53] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:53] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:53] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[10] train-logloss:0.61446 valid-logloss:0.617895
[00:18:54] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:54] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:55] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:55] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:56] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:56] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:18:56] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
```

[illegible]

[illegible]



[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

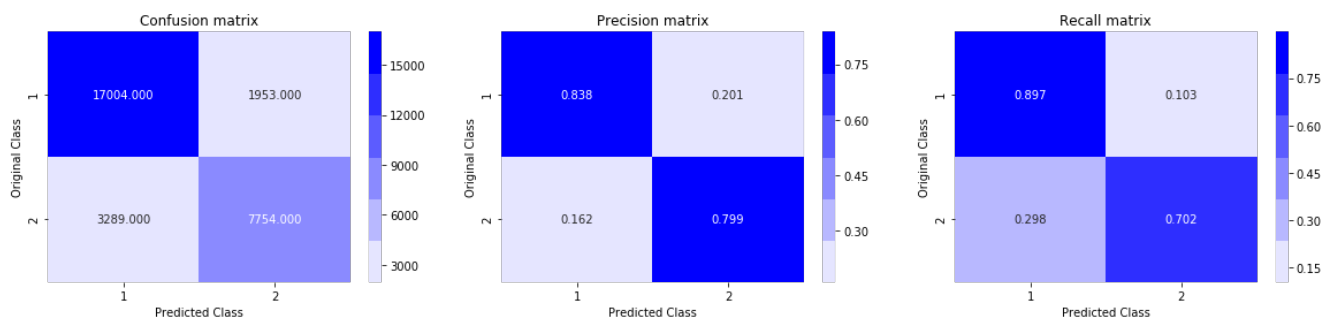


```
[00:21:22] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:21:22] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:23] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 22 extra nodes, 0 pruned nodes, max_depth=4
[00:21:23] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:21:23] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 28 extra nodes, 0 pruned nodes, max_depth=4
[00:21:24] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 22 extra nodes, 0 pruned nodes, max_depth=4
[00:21:24] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[390] train-logloss:0.301228 valid-logloss:0.381765
[00:21:25] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:25] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:25] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:26] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 18 extra nodes, 0 pruned nodes, max_depth=4
[00:21:26] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 22 extra nodes, 0 pruned nodes, max_depth=4
[00:21:26] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:27] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 26 extra nodes, 0 pruned nodes, max_depth=4
[00:21:27] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[00:21:27] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
[399] train-logloss:0.298798 valid-logloss:0.381676
The test log loss is: 0.381676321641348
```

In [179]:

```
predicted_y = np.array(predict_y>0.5, dtype=int)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(new_y_test, predicted_y)
```

Total number of data points : 30000



## 5. Assignments

1. Try out models (Logistic regression, Linear-SVM) with simple TF-IDF vectors instead of TD\_IDF weighted word2Vec.
2. Hyperparameter tune XgBoost using RandomSearch to reduce the log-loss.

## Logistic Regression with TFIDF and hyperparameter tuning

In [61]:

```
alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
from sklearn.metrics import log_loss
```

```

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0=0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----
from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV
log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(X_train,y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train, y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

fig, ax = plt.subplots()
ax.plot(alpha, log_error_array,c='g')
for i, txt in enumerate(np.round(log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train, y_train)

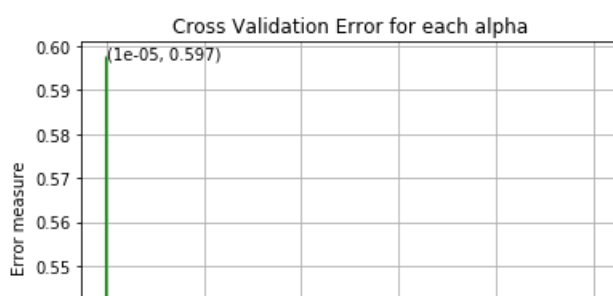
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(new_y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y,axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

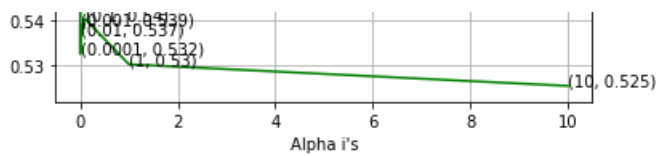
```

```

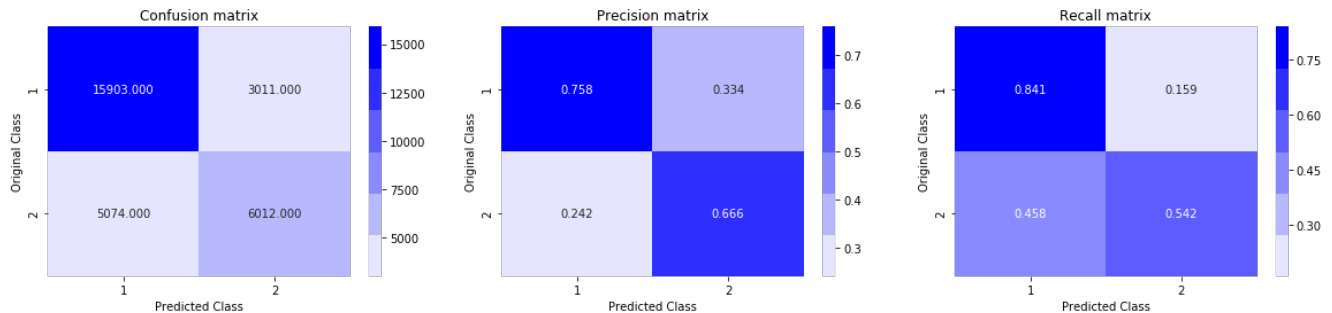
For values of alpha = 1e-05 The log loss is: 0.5973379930195214
For values of alpha = 0.0001 The log loss is: 0.5324688775826292
For values of alpha = 0.001 The log loss is: 0.5391307902511334
For values of alpha = 0.01 The log loss is: 0.5370251185379298
For values of alpha = 0.1 The log loss is: 0.5404253950867678
For values of alpha = 1 The log loss is: 0.5301213032440313
For values of alpha = 10 The log loss is: 0.525223471496805

```





For values of best alpha = 10 The train log loss is: 0.5163528679900335  
 For values of best alpha = 10 The test log loss is: 0.525223471496805  
 Total number of data points : 30000



## Linear SVM with TFIDF and hyperparameter tuning

In [82]:

```
alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.

# read more about SGDClassifier() at http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.SGDClassifier.html
# -----
# default parameters
# SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_intercept=True, max_iter=None, tol=None,
# shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_rate='optimal', eta0
# =0.0, power_t=0.5,
# class_weight=None, warm_start=False, average=False, n_iter=None)

# some of methods
# fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic Gradient Descent.
# predict(X) Predict class labels for samples in X.

#-----
# video link:
#-----

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
    clf.fit(new_x_train, new_y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(new_x_train, new_y_train)
    predict_y = sig_clf.predict_proba(X_test)
    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.cl
asses_, eps=1e-15))

fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', random_state=42)
clf.fit(new_x_train, new_y_train)
```

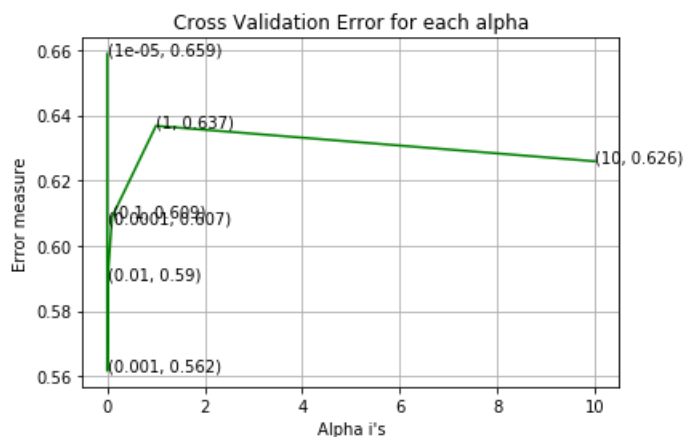
```

sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(new_x_train, new_y_train)

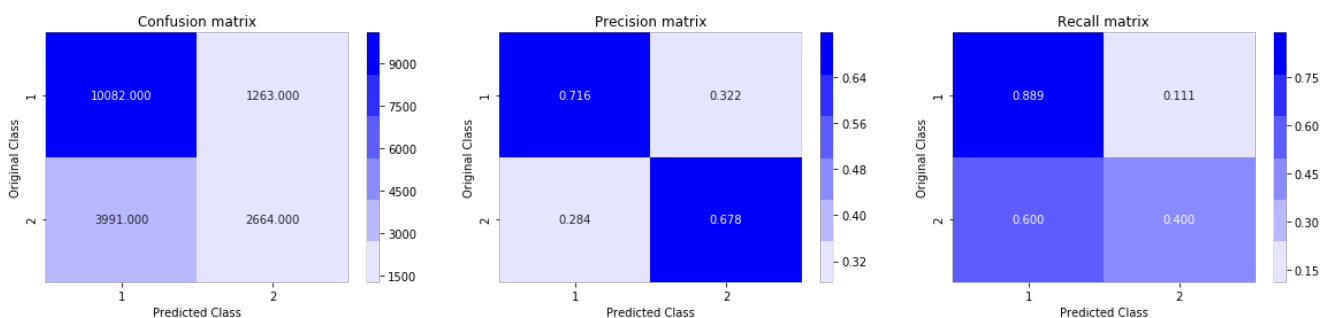
predict_y = sig_clf.predict_proba(new_x_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss
is:", log_loss(new_y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, p
redict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

```

For values of alpha = 1e-05 The log loss is: 0.6588378264676555  
 For values of alpha = 0.0001 The log loss is: 0.6071160913315697  
 For values of alpha = 0.001 The log loss is: 0.5616058796517264  
 For values of alpha = 0.01 The log loss is: 0.5896569821761471  
 For values of alpha = 0.1 The log loss is: 0.6088181085142541  
 For values of alpha = 1 The log loss is: 0.6367239445519107  
 For values of alpha = 10 The log loss is: 0.6258621473423245



For values of best alpha = 0.001 The train log loss is: 0.538992339428178  
 For values of best alpha = 0.001 The test log loss is: 0.5616058796517264  
 Total number of data points : 18000



## Hyperparameter tune for Xgboost

In [75]:

```

# https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier
from sklearn.datasets import make_classification
from xgboost import XGBClassifier

model = XGBClassifier()
param_grid = {
    'n_estimators': [100, 200, 300],
    'max_depth' : [5, 10]
}

```

```

}

clf = GridSearchCV(model, param_grid, scoring='roc_auc', cv = 10)

clf.fit(new_x_train, new_y_train)

train_auc= clf.cv_results_['mean_train_score']
train_auc_std= clf.cv_results_['std_train_score']

best_parameters=clf.best_params_

```

In [76]:

```

for param_name in sorted(best_parameters.keys()):
    print("%s: %r" % (param_name, best_parameters[param_name]))

```

```

max_depth: 10
n_estimators: 300

```

In [78]:

```

from sklearn.metrics import log_loss
import xgboost as xgb
params = {}
params['objective'] = 'binary:logistic'
params['eval_metric'] = 'logloss'
params['max_depth'] = 10
params['n_estimators'] = 300

d_train = xgb.DMatrix(final_x, label=final_y)
d_test = xgb.DMatrix(X_test, label=y_test)

watchlist = [(d_train, 'train'), (d_test, 'valid')]

bst = xgb.train(params, d_train, 400, watchlist, early_stopping_rounds=20, verbose_eval=10)

xgdmatrix = xgb.DMatrix(new_x_train, new_y_train)
predict_y = bst.predict(d_test)
print("The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))

```

```

[13:57:18] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 570 extra nodes, 0 pruned nodes, max_depth=10
[0] train-logloss:0.537612 valid-logloss:0.570369
Multiple eval metrics have been passed: 'valid-logloss' will be used for early stopping.

```

Will train until valid-logloss hasn't improved in 20 rounds.

```

[13:57:22] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 594 extra nodes, 0 pruned nodes, max_depth=10
[13:57:26] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 568 extra nodes, 0 pruned nodes, max_depth=10
[13:57:29] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 604 extra nodes, 0 pruned nodes, max_depth=10
[13:57:33] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 522 extra nodes, 0 pruned nodes, max_depth=10
[13:57:37] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 564 extra nodes, 0 pruned nodes, max_depth=10
[13:57:41] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 640 extra nodes, 0 pruned nodes, max_depth=10
[13:57:45] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 350 extra nodes, 0 pruned nodes, max_depth=10
[13:57:49] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 306 extra nodes, 0 pruned nodes, max_depth=10
[13:57:52] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 294 extra nodes, 0 pruned nodes, max_depth=10
[13:57:56] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 362 extra nodes, 0 pruned nodes, max_depth=10
[10] train-logloss:0.183503 valid-logloss:0.375455
[13:58:00] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 360 extra nodes, 0 pruned nodes, max_depth=10
[13:58:03] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 360 extra nodes, 0 pruned nodes, max_depth=10

```

```

[13:58:03] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 440 extra nodes, 0 pruned nodes, max_depth=10
[13:58:07] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 464 extra nodes, 0 pruned nodes, max_depth=10
[13:58:11] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 272 extra nodes, 0 pruned nodes, max_depth=10
[13:58:15] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 434 extra nodes, 0 pruned nodes, max_depth=10
[13:58:19] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 274 extra nodes, 0 pruned nodes, max_depth=10
[13:58:22] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 402 extra nodes, 0 pruned nodes, max_depth=10
[13:58:27] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 220 extra nodes, 0 pruned nodes, max_depth=10
[13:58:31] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 488 extra nodes, 0 pruned nodes, max_depth=10
[13:58:35] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 470 extra nodes, 0 pruned nodes, max_depth=10
[20] train-logloss:0.096059 valid-logloss:0.375361
[13:58:38] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 338 extra nodes, 0 pruned nodes, max_depth=10
[13:58:42] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 408 extra nodes, 0 pruned nodes, max_depth=10
[13:58:46] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 496 extra nodes, 0 pruned nodes, max_depth=10
[13:58:49] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 250 extra nodes, 0 pruned nodes, max_depth=10
[13:58:52] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 324 extra nodes, 0 pruned nodes, max_depth=10
[13:58:56] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 600 extra nodes, 0 pruned nodes, max_depth=10
[13:59:00] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 616 extra nodes, 0 pruned nodes, max_depth=10
[13:59:04] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 354 extra nodes, 0 pruned nodes, max_depth=10
[13:59:07] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 256 extra nodes, 0 pruned nodes, max_depth=10
[13:59:11] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 486 extra nodes, 0 pruned nodes, max_depth=10
[30] train-logloss:0.049011 valid-logloss:0.382999
[13:59:15] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 530 extra nodes, 0 pruned nodes, max_depth=10
[13:59:18] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 284 extra nodes, 0 pruned nodes, max_depth=10
[13:59:21] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 436 extra nodes, 0 pruned nodes, max_depth=10
[13:59:25] d:\build\xgboost\xgboost-0.80.git\src\tree\updater_prune.cc:74: tree pruning end, 1 roots, 264 extra nodes, 0 pruned nodes, max_depth=10
Stopping. Best iteration:
[14] train-logloss:0.142144 valid-logloss:0.372167

```

The test log loss is: 0.3848602981221957

In [81]:

```

# http://zetcode.com/python/prettytable/
#Xgboost observation
from prettytable import PrettyTable

#If you get a ModuleNotFoundError error , install prettytable using: pip3 install prettytable

x = PrettyTable()
x.field_names = ["Featurization","MODEL","test log loss"]
x.add_row(["TFIDF_W2V","Logistic Regression", 0.51])
x.add_row(["TFIDF_W2V","Linear SVM",0.50])
x.add_row(["TFIDF_W2V","XGBoost",0.40])
x.add_row(["TFIDF","Logistic Regression",0.52])
x.add_row(["TFIDF","Linear SVM",0.49])
x.add_row(["TFIDF","XGBoost(with hyperparameter tuning)",0.38])

print(x)

```

```

+-----+-----+-----+
| Featurization |          MODEL          | test log loss |
+-----+-----+-----+
|  TFIDF_W2V   |  Logistic Regression    |      0.51     |

```

	TFIDF_W2V		Linear SVM		0.5	
	TFIDF_W2V		XGBoost		0.4	
	TFIDF		Logistic Regression		0.52	
	TFIDF		Linear SVM		0.49	
	TFIDF		XGBoost(with hyperparameter tuning)		0.38	
+-----+-----+-----+-----+-----+						

## Conclusion

1.As we know that linear model works well using high dimension data, so we use tfidf featurization for linear models like logistic regression and linear svm.

2.We have used both the featurization (tfidf and tfidf avg ) for xgboost.

3.XGBOOST using tfidf avg seems overfit as high difference between train log loss and test log loss. So we will discard this.

4.We get low log-loss value using xgboost model with tfidf featurization and not getting high difference between train and test log loss as well.

5.So we can use XGBOOST model with tfidf featurization for this problem.