



## 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> (<https://www.kaggle.com/c/quora-question-pairs>)

### \_\_ Useful Links \_\_

- Discussions : <https://www.kaggle.com/anokas/data-analysis-xgboost-starter-0-35460-lb/comments> (<https://www.kaggle.com/anokas/data-analysis-xgboost-starter-0-35460-lb/comments>)
- Kaggle Winning Solution and other approaches:  
<https://www.dropbox.com/sh/93968nfnrzh8bp5/AACZdtsApc1QSTQc7X0H3QZ5a?dl=0>  
(<https://www.dropbox.com/sh/93968nfnrzh8bp5/AACZdtsApc1QSTQc7X0H3QZ5a?dl=0>)
- Blog 1 : <https://engineering.quora.com/Semantic-Question-Matching-with-Deep-Learning>  
(<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> (<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>  
(<https://www.kaggle.com/c/quora-question-pairs#evaluation>)

Metric(s):

- log-loss : <https://www.kaggle.com/wiki/LogarithmicLoss>  
(<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 [6]:

```
1 from nltk.corpus import stopwords
2 import distance
3 from nltk.stem import PorterStemmer
4 from bs4 import BeautifulSoup
5 import warnings
6 warnings.filterwarnings("ignore")
7 import numpy as np
8 import pandas as pd
9 import seaborn as sns
10 import matplotlib.pyplot as plt
11 from subprocess import check_output
12 %matplotlib inline
13 import plotly.offline as py
14 py.init_notebook_mode(connected=True)
15 import plotly.graph_objs as go
16 import plotly.tools as tls
17 import os
18 import gc
19 import re
20 from nltk.corpus import stopwords
21 import distance
22 from nltk.stem import PorterStemmer
23 import re
24 # This package is used for finding Longest common subsequence between two strings
25 # you can write your own dp code for this
26 from bs4 import BeautifulSoup
27 from fuzzywuzzy import fuzz
28 from sklearn.manifold import TSNE
29 # Import the Required Lib packages for WORD-Cloud generation
30 # https://stackoverflow.com/questions/45625434/how-to-install-wordcloud-in-python
31 from wordcloud import WordCloud, STOPWORDS
32 from os import path
33 from PIL import Image
34 from sklearn.preprocessing import normalize
35 from sklearn.feature_extraction.text import CountVectorizer
36 from sklearn.feature_extraction.text import TfidfVectorizer
37 from tqdm import tqdm
38 from sklearn.decomposition import TruncatedSVD
39 from sklearn.preprocessing import normalize
40 from sklearn.feature_extraction.text import CountVectorizer
41 from sklearn.manifold import TSNE
42 import seaborn as sns
43 from sklearn.neighbors import KNeighborsClassifier
44 from sklearn.metrics import confusion_matrix
45 from sklearn.metrics.classification import accuracy_score, log_loss
46 from sklearn.feature_extraction.text import TfidfVectorizer
47 from collections import Counter
48 from scipy.sparse import hstack
49 from sklearn.multiclass import OneVsRestClassifier
50 from sklearn.svm import SVC
51 from sklearn.cross_validation import StratifiedKFold
52 from collections import Counter, defaultdict
53 from sklearn.calibration import CalibratedClassifierCV
54 from sklearn.naive_bayes import MultinomialNB
55 from sklearn.naive_bayes import GaussianNB
56 from sklearn.model_selection import train_test_split
```

```

57 from sklearn.model_selection import GridSearchCV
58 import math
59 from sklearn.metrics import normalized_mutual_info_score
60 from sklearn.ensemble import RandomForestClassifier
61 from sklearn.model_selection import cross_val_score
62 from sklearn.linear_model import SGDClassifier
63 from mlxtend.classifier import StackingClassifier
64 from sklearn import model_selection
65 from sklearn.linear_model import LogisticRegression
66 from sklearn.metrics import precision_recall_curve, auc, roc_curve

```

## 3.1 Reading data and basic stats

```

In [7]: 1 df = pd.read_csv("train.csv")
        2
        3 print("Number of data points:",df.shape[0])

```

Number of data points: 404290

```

In [8]: 1 df.head()

```

```

Out[8]:

```

	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 $23^{24}$ is divided by 100	0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

```

In [9]: 1 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         404290 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

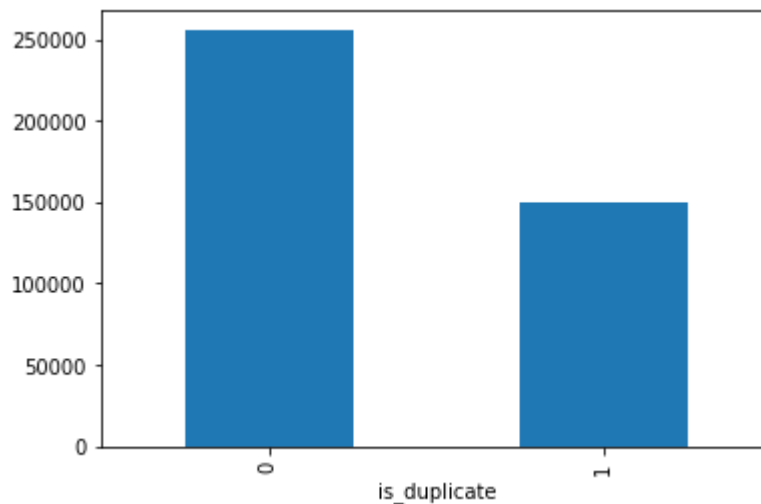
- qid{1, 2}: The unique ID of each question in the pair
- question{1, 2}: The actual textual contents of the questions.
- is\_duplicate: The label that we are trying to predict - whether the two questions are duplicates of each other.

### 3.2.1 Distribution of data points among output classes

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

```
In [10]: 1 df.groupby("is_duplicate")['id'].count().plot.bar()
```

```
Out[10]: <matplotlib.axes._subplots.AxesSubplot at 0x22b00727d30>
```



```
In [11]: 1 print('~> Total number of question pairs for training:\n    {}'.format(len(df)))
```

```
~> Total number of question pairs for training:
404290
```

```
In [12]: 1 print('~> Question pairs are not Similar (is_duplicate = 0):\n    {}'.format(rou
2 print('\n~> Question pairs are Similar (is_duplicate = 1):\n    {}'.format(rou
```

```
~> 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 [13]: 1 qids = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())
2 unique_qs = len(np.unique(qids))
3 qs_morethan_onetime = np.sum(qids.value_counts() > 1)
4 print ('Total number of Unique Questions are: {}'.format(unique_qs))
5 #print len(np.unique(qids))
6
7 print ('Number of unique questions that appear more than one time: {} ({}%)\n'
8
9 print ('Max number of times a single question is repeated: {}'.format(max(qi
10
11 q_vals=qids.value_counts()
12
13 q_vals=q_vals.values

```

Total num of Unique Questions are: 537933

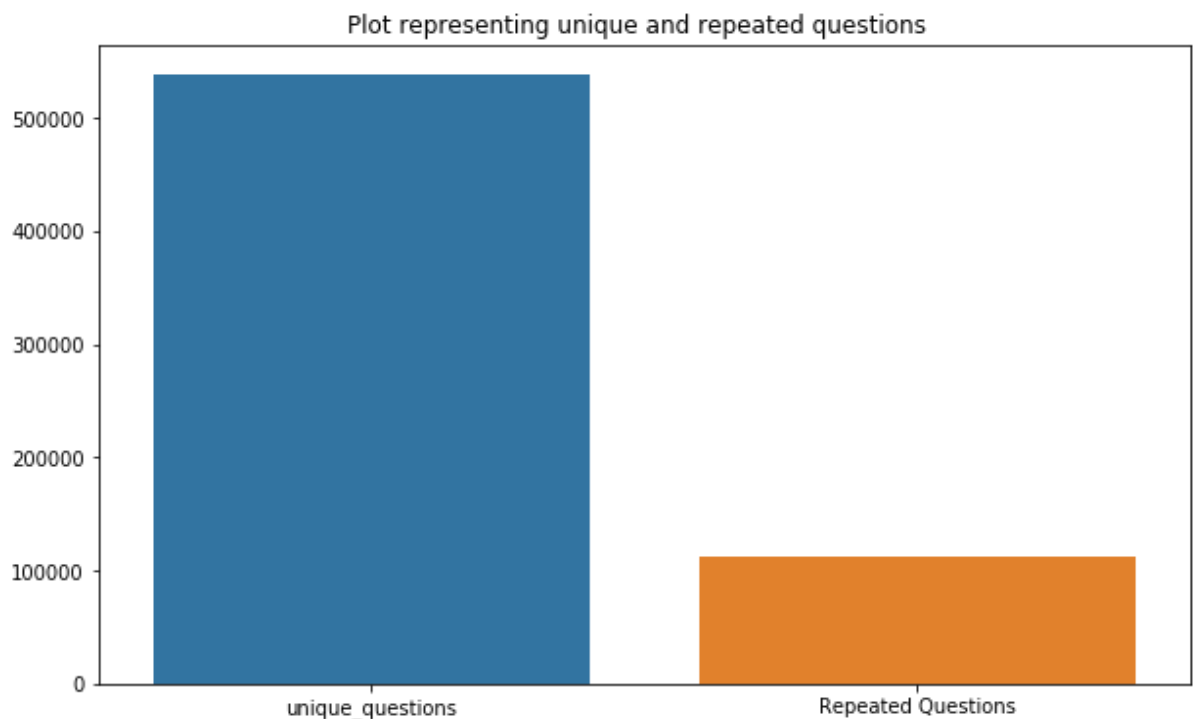
Number of unique questions that appear more than one time: 111780 (20.77953945937505%)

Max number of times a single question is repeated: 157

```

In [14]: 1
2 x = ["unique_questions" , "Repeated Questions"]
3 y = [unique_qs , qs_morethan_onetime]
4
5 plt.figure(figsize=(10, 6))
6 plt.title ("Plot representing unique and repeated questions ")
7 sns.barplot(x,y)
8 plt.show()

```



### 3.2.3 Checking for Duplicates

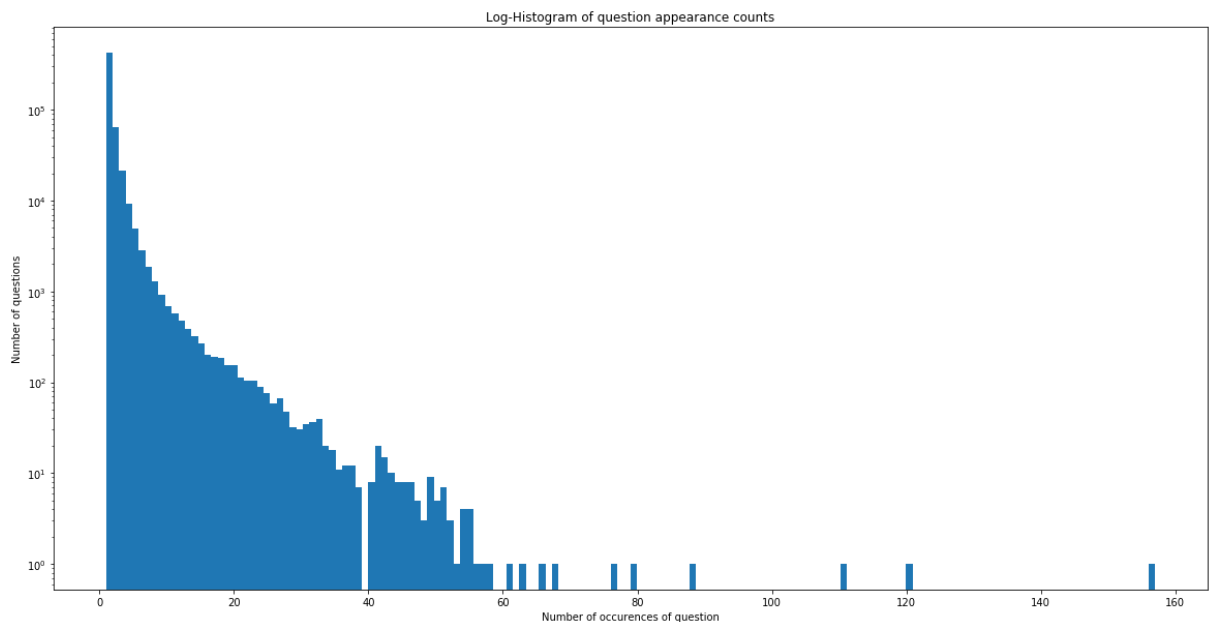
```
In [15]: 1 #checking whether there are any repeated pair of questions
2
3 pair_duplicates = df[['qid1','qid2','is_duplicate']].groupby(['qid1','qid2']).
4
5 print ("Number of duplicate questions",(pair_duplicates).shape[0] - df.shape[0])

Number of duplicate questions 0
```

### 3.2.4 Number of occurrences of each question

```
In [24]: 1 plt.figure(figsize=(20, 10))
2
3 plt.hist(qids.value_counts(), bins=160)
4
5 plt.yscale('log', nonposy='clip')
6
7 plt.title('Log-Histogram of question appearance counts')
8
9 plt.xlabel('Number of occurrences of question')
10
11 plt.ylabel('Number of questions')
12
13 print ('Maximum number of times a single question is repeated: {}'.format(max(qids.value_counts().index)))

Maximum number of times a single question is repeated: 157
```



### 3.2.5 Checking for NULL values



```
In [17]: 1 #Checking whether there are any rows with null values
2 nan_rows = df[df.isnull().any(1)]
3 print (nan_rows)
```

	id	qid1	qid2	question1	question2	\
105780	105780	174363	174364	How can I develop android app?	NaN	
201841	201841	303951	174364	How can I create an Android app?	NaN	

	is_duplicate
105780	0
201841	0

- There are two rows with null values in question2

```
In [18]: 1 # Filling the null values with ' '
2 df = df.fillna(' ')
3 nan_rows = df[df.isnull().any(1)]
4 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 [20]: 1 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
2         df = pd.read_csv("df_fe_without_preprocessing_train.csv", encoding='latin-1')
3     else:
4         df['freq_qid1'] = df.groupby('qid1')['qid1'].transform('count')
5         df['freq_qid2'] = df.groupby('qid2')['qid2'].transform('count')
6         df['q1len'] = df['question1'].str.len()
7         df['q2len'] = df['question2'].str.len()
8         df['q1_n_words'] = df['question1'].apply(lambda row: len(row.split(" ")))
9         df['q2_n_words'] = df['question2'].apply(lambda row: len(row.split(" ")))
10
11     def normalized_word_Common(row):
12         w1 = set(map(lambda word: word.lower().strip(), row['question1'].split()))
13         w2 = set(map(lambda word: word.lower().strip(), row['question2'].split()))
14         return 1.0 * len(w1 & w2)
15     df['word_Common'] = df.apply(normalized_word_Common, axis=1)
16
17     def normalized_word_Total(row):
18         w1 = set(map(lambda word: word.lower().strip(), row['question1'].split()))
19         w2 = set(map(lambda word: word.lower().strip(), row['question2'].split()))
20         return 1.0 * (len(w1) + len(w2))
21     df['word_Total'] = df.apply(normalized_word_Total, axis=1)
22
23     def normalized_word_share(row):
24         w1 = set(map(lambda word: word.lower().strip(), row['question1'].split()))
25         w2 = set(map(lambda word: word.lower().strip(), row['question2'].split()))
26         return 1.0 * len(w1 & w2) / (len(w1) + len(w2))
27     df['word_share'] = df.apply(normalized_word_share, axis=1)
28
29     df['freq_q1+q2'] = df['freq_qid1'] + df['freq_qid2']
30     df['freq_q1-q2'] = abs(df['freq_qid1'] - df['freq_qid2'])
31
32     df.to_csv("df_fe_without_preprocessing_train.csv", index=False)
33
34     df.head()

```

Out[20]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_v
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	
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	
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	

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_v
	3	3	7	8	Why am I mentally very lonely? How can I solve... Find the remainder when $23^{24}$ is divided by 1000.	0	1	1	50	65	
	4	4	9	10	Which one dissolves in water quickly sugar, salt... Which fish would survive in salt water?	0	3	1	76	39	

### 3.3.1 Analysis of some of the extracted features

- Here are some questions have only one single words.

In [21]:

```

1 print ("Minimum length of the questions in question1 : " , min(df['q1_n_words']
2
3 print ("Minimum length of the questions in question2 : " , min(df['q2_n_words']
4
5 print ("Number of Questions with minimum length [question1] :", df[df['q1_n_w
6 print ("Number of Questions with minimum length [question2] :", df[df['q2_n_w

```

```

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

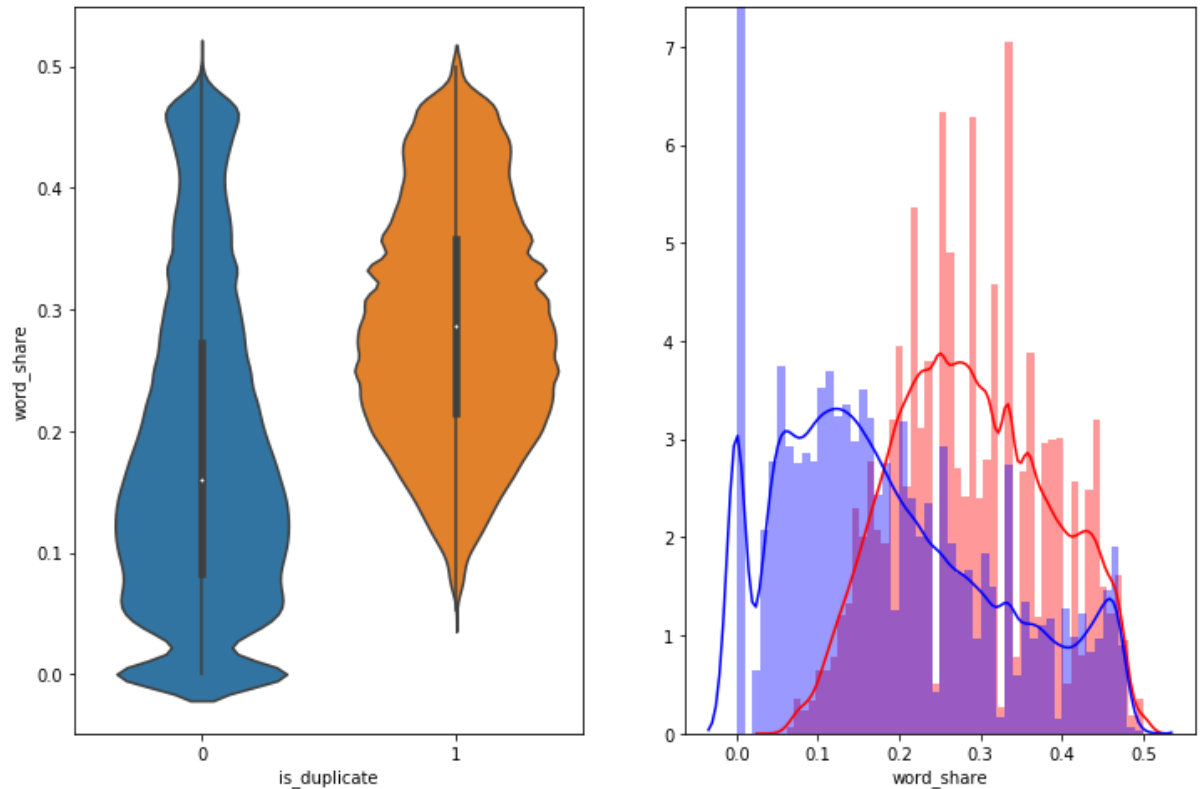
```

#### 3.3.1.1 Feature: word\_share

```

In [25]: 1 plt.figure(figsize=(12, 8))
2
3 plt.subplot(1,2,1)
4 sns.violinplot(x = 'is_duplicate', y = 'word_share', data = df[0:])
5
6 plt.subplot(1,2,2)
7 sns.distplot(df[df['is_duplicate'] == 1.0]['word_share'], label = "1", color = "blue")
8 sns.distplot(df[df['is_duplicate'] == 0.0]['word_share'], label = "0", color = "red")
9 plt.show()

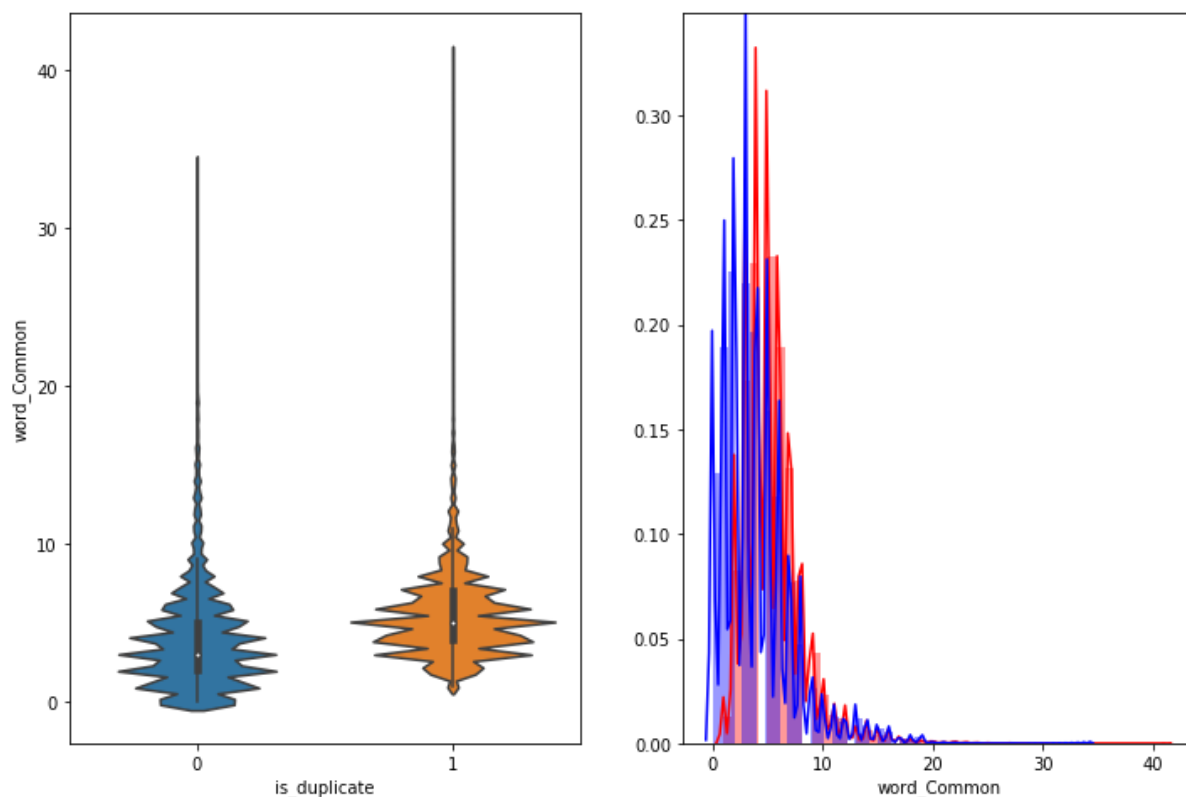
```



- 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 [26]: 1 plt.figure(figsize=(12, 8))
2
3 plt.subplot(1,2,1)
4 sns.violinplot(x = 'is_duplicate', y = 'word_Common', data = df[0:])
5
6 plt.subplot(1,2,2)
7 sns.distplot(df[df['is_duplicate'] == 1.0]['word_Common'][0:], label = "1", c
8 sns.distplot(df[df['is_duplicate'] == 0.0]['word_Common'][0:], label = "0" ,
9 plt.show()
```



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

## Data Pre processing

```
In [7]: 1 #https://stackoverflow.com/questions/12468179/unicodedecodeerror-utf8-codec-co
2 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
3     df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1
4     df = df.fillna('')
5     df.head()
6 else:
7     print("get df_fe_without_preprocessing_train.csv from drive or run the pre
```

In [8]: 1 df.head(2)

Out[8]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_wc
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	
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	

## 3.4 Preprocessing of Text

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

In [9]:

```
1 # To get the results in 4 decemal points
2 SAFE_DIV = 0.0001
3
4 STOP_WORDS = stopwords.words("english")
5
6
7 def preprocess(x):
8     x = str(x).lower()
9     x = x.replace(",000,000", "m").replace(",000", "k").replace("'", "").replace(
10         "won't", "will not").replace("cannot", "can")
11     .replace("n't", " not").replace("what's", "what is")
12     .replace("'ve", " have").replace("i'm", "i am").replace(
13         "he's", "he is").replace("she's", "she is")
14     .replace("%", " percent ").replace("₹", " rupee ").
15     .replace("€", " euro ").replace("'ll", " will")
16     x = re.sub(r"([0-9]+)000000", r"\1m", x)
17     x = re.sub(r"([0-9]+)000", r"\1k", x)
18
19
20     porter = PorterStemmer()
21     pattern = re.compile('\W')
22
23     if type(x) == type(''):
24         x = re.sub(pattern, ' ', x)
25
26
27     if type(x) == type(''):
28         x = porter.stem(x)
29         example1 = BeautifulSoup(x)
30         x = example1.get_text()
31
32
33     return x
34
```

- 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  
$$cwc\_min = common\_word\_count / (\min(\text{len}(q1\_words), \text{len}(q2\_words)))$$

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

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

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

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

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

$$\text{ctc\_max} = \text{common\_token\_count} / (\max(\text{len}(\text{q1\_tokens}), \text{len}(\text{q2\_tokens})))$$
- **last\_word\_eq** : Check if First word of both questions is equal or not  

$$\text{last\_word\_eq} = \text{int}(\text{q1\_tokens}[-1] == \text{q2\_tokens}[-1])$$
- **first\_word\_eq** : Check if First word of both questions is equal or not  

$$\text{first\_word\_eq} = \text{int}(\text{q1\_tokens}[0] == \text{q2\_tokens}[0])$$
- **abs\_len\_diff** : Abs. length difference  

$$\text{abs\_len\_diff} = \text{abs}(\text{len}(\text{q1\_tokens}) - \text{len}(\text{q2\_tokens}))$$
- **mean\_len** : Average Token Length of both Questions  

$$\text{mean\_len} = (\text{len}(\text{q1\_tokens}) + \text{len}(\text{q2\_tokens})) / 2$$
- **fuzz\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>  
<https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)
- **fuzz\_partial\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>  
<https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)
- **token\_sort\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>  
<https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>



[fuzzy-string-matching-in-python/](http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/) (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)

- **token\_set\_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> (<https://github.com/seatgeek/fuzzywuzzy#usage>) <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<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  
$$\text{longest\_substr\_ratio} = \text{len}(\text{longest common substring}) / (\min(\text{len}(q1\_tokens), \text{len}(q2\_tokens)))$$

In [10]:

```
1 def get_token_features(q1, q2):
2     token_features = [0.0]*10
3
4     # Converting the Sentence into Tokens:
5     q1_tokens = q1.split()
6     q2_tokens = q2.split()
7
8     if len(q1_tokens) == 0 or len(q2_tokens) == 0:
9         return token_features
10    # Get the non-stopwords in Questions
11    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
12    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])
13
14    #Get the stopwords in Questions
15    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
16    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])
17
18    # Get the common non-stopwords from Question pair
19    common_word_count = len(q1_words.intersection(q2_words))
20
21    # Get the common stopwords from Question pair
22    common_stop_count = len(q1_stops.intersection(q2_stops))
23
24    # Get the common Tokens from Question pair
25    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))
26
27
28    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)))
29    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)))
30    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)))
31    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)))
32    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)))
33    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)))
34
35    # Last word of both question is same or not
36    token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])
37
38    # First word of both question is same or not
39    token_features[7] = int(q1_tokens[0] == q2_tokens[0])
40
41    token_features[8] = abs(len(q1_tokens) - len(q2_tokens))
42
43    #Average Token Length of both Questions
44    token_features[9] = (len(q1_tokens) + len(q2_tokens))/2
45    return token_features
46
47    # get the Longest Common sub string
48
49    def get_longest_substr_ratio(a, b):
50        strs = list(distance.lcs substrings(a, b))
51        if len(strs) == 0:
52            return 0
53        else:
54            return len(strs[0]) / (min(len(a), len(b)) + 1)
55
56    def extract_features(df):
```

```

57 # preprocessing each question
58 df["question1"] = df["question1"].fillna("").apply(preprocess)
59 df["question2"] = df["question2"].fillna("").apply(preprocess)
60
61 print("token features...")
62
63 # Merging Features with dataset
64
65 token_features = df.apply(lambda x: get_token_features(x["question1"], x["question2"]), axis=1)
66
67 df["cwc_min"] = list(map(lambda x: x[0], token_features))
68 df["cwc_max"] = list(map(lambda x: x[1], token_features))
69 df["csc_min"] = list(map(lambda x: x[2], token_features))
70 df["csc_max"] = list(map(lambda x: x[3], token_features))
71 df["ctc_min"] = list(map(lambda x: x[4], token_features))
72 df["ctc_max"] = list(map(lambda x: x[5], token_features))
73 df["last_word_eq"] = list(map(lambda x: x[6], token_features))
74 df["first_word_eq"] = list(map(lambda x: x[7], token_features))
75 df["abs_len_diff"] = list(map(lambda x: x[8], token_features))
76 df["mean_len"] = list(map(lambda x: x[9], token_features))
77
78 #Computing Fuzzy Features and Merging with Dataset
79
80 # do read this blog: http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string
81 # https://stackoverflow.com/questions/31806695/when-to-use-which-fuzz-func
82 # https://github.com/seatgeek/fuzzywuzzy
83 print("fuzzy features..")
84
85 df["token_set_ratio"] = df.apply(lambda x: fuzz.token_set_ratio(x["question1"], x["question2"]), axis=1)
86 # The token sort approach involves tokenizing the string in question, sort
87 # then joining them back into a string We then compare the transformed strings
88 df["token_sort_ratio"] = df.apply(lambda x: fuzz.token_sort_ratio(x["question1"], x["question2"]), axis=1)
89 df["fuzz_ratio"] = df.apply(lambda x: fuzz.QRatio(x["question1"], x["question2"]), axis=1)
90 df["fuzz_partial_ratio"] = df.apply(lambda x: fuzz.partial_ratio(x["question1"], x["question2"]), axis=1)
91 df["longest_substr_ratio"] = df.apply(lambda x: fuzz.longest_common_substring_ratio(x["question1"], x["question2"]), axis=1)
92 return df

```

```
In [12]: 1 if os.path.isfile('nlp_features_train.csv'):
2         df = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
3         df.fillna('')
4     else:
5         print("Extracting features for train:")
6         df = pd.read_csv("train.csv")
7         df = extract_features(df)
8         df.to_csv("nlp_features_train.csv", index=False)
9     df.head(2)
```

```
Out[12]:
```

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...
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	...
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	...

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 [17]: 1 df_duplicate = df[df['is_duplicate'] == 1]
2 dfp_nonduplicate = df[df['is_duplicate'] == 0]
3
4 # Converting 2d array of q1 and q2 and flatten the array: Like {{1,2},{3,4}} t
5 p = np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten(
6 n = np.dstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).
7
8 print ("Number of data points in class 1 (duplicate pairs) :",len(p))
9 print ("Number of data points in class 0 (non duplicate pairs) :",len(n))
10
11 #Saving the np array into a text file
12 np.savetxt('train_p.txt', p, delimiter=' ', fmt='%s')
13 np.savetxt('train_n.txt', n, delimiter=' ', fmt='%s')
```

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

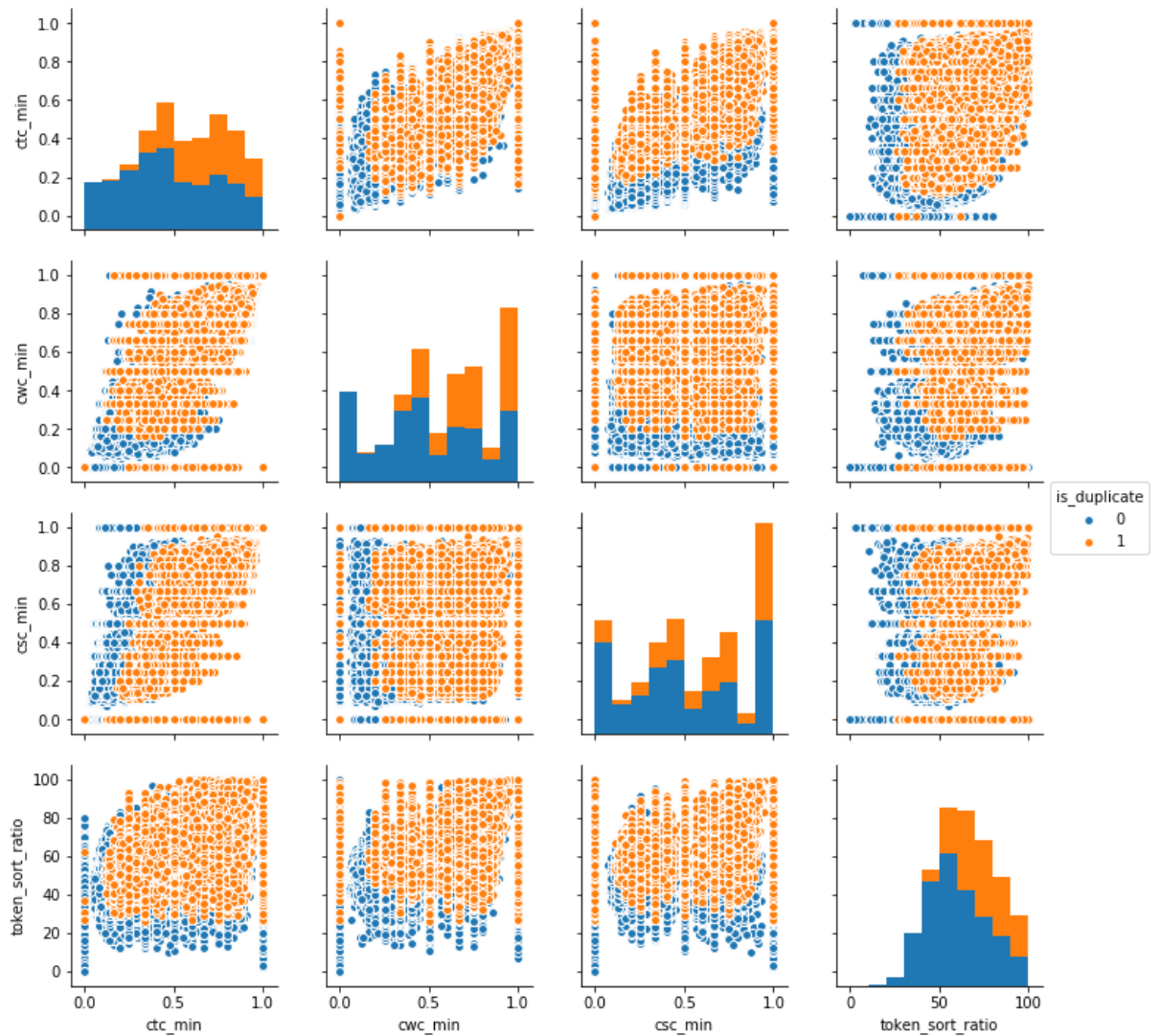
Number of data points in class 0 (non duplicate pairs) : 510054





In [33]:

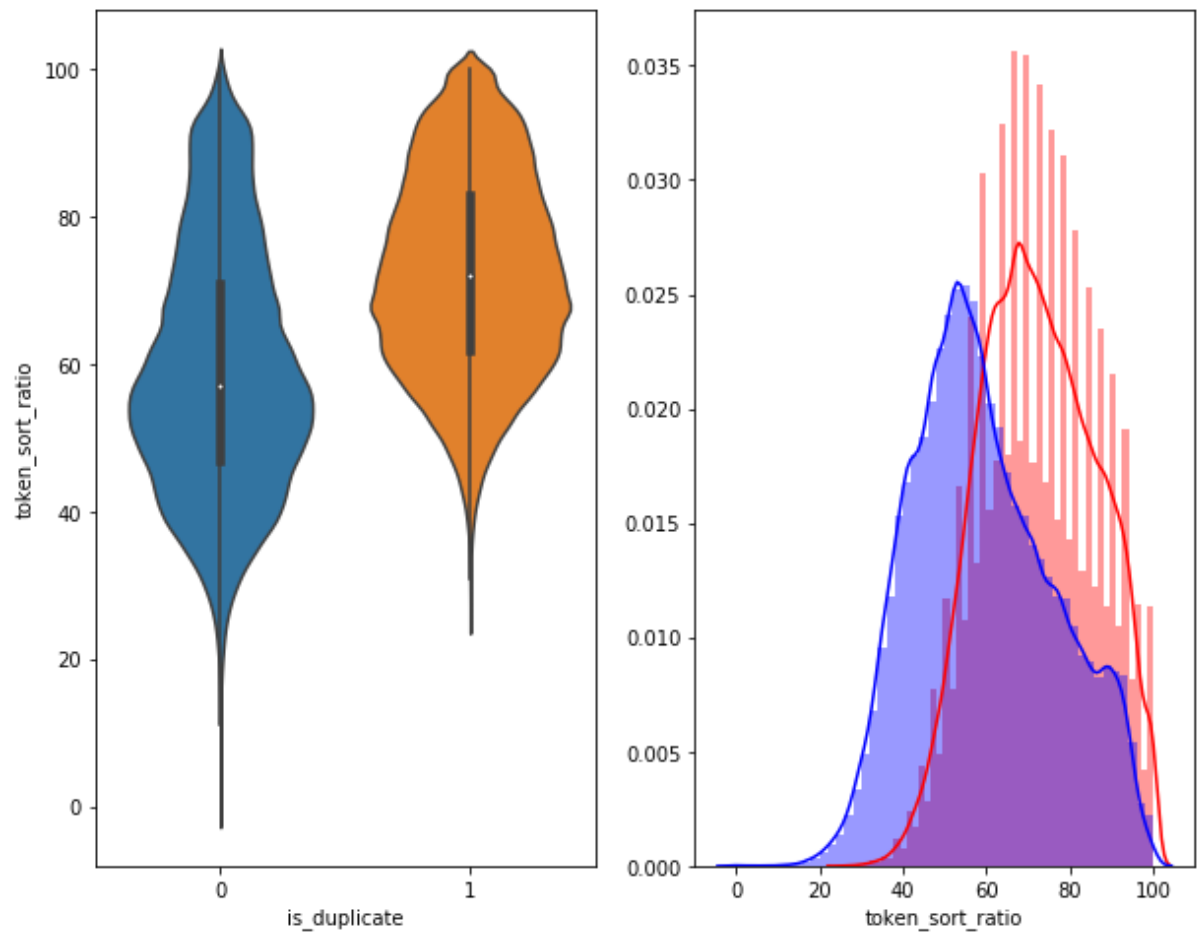
```
1 n = df.shape[0]
2 sns.pairplot(df[['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio', 'is_duplicate']
3 plt.show())
```



```

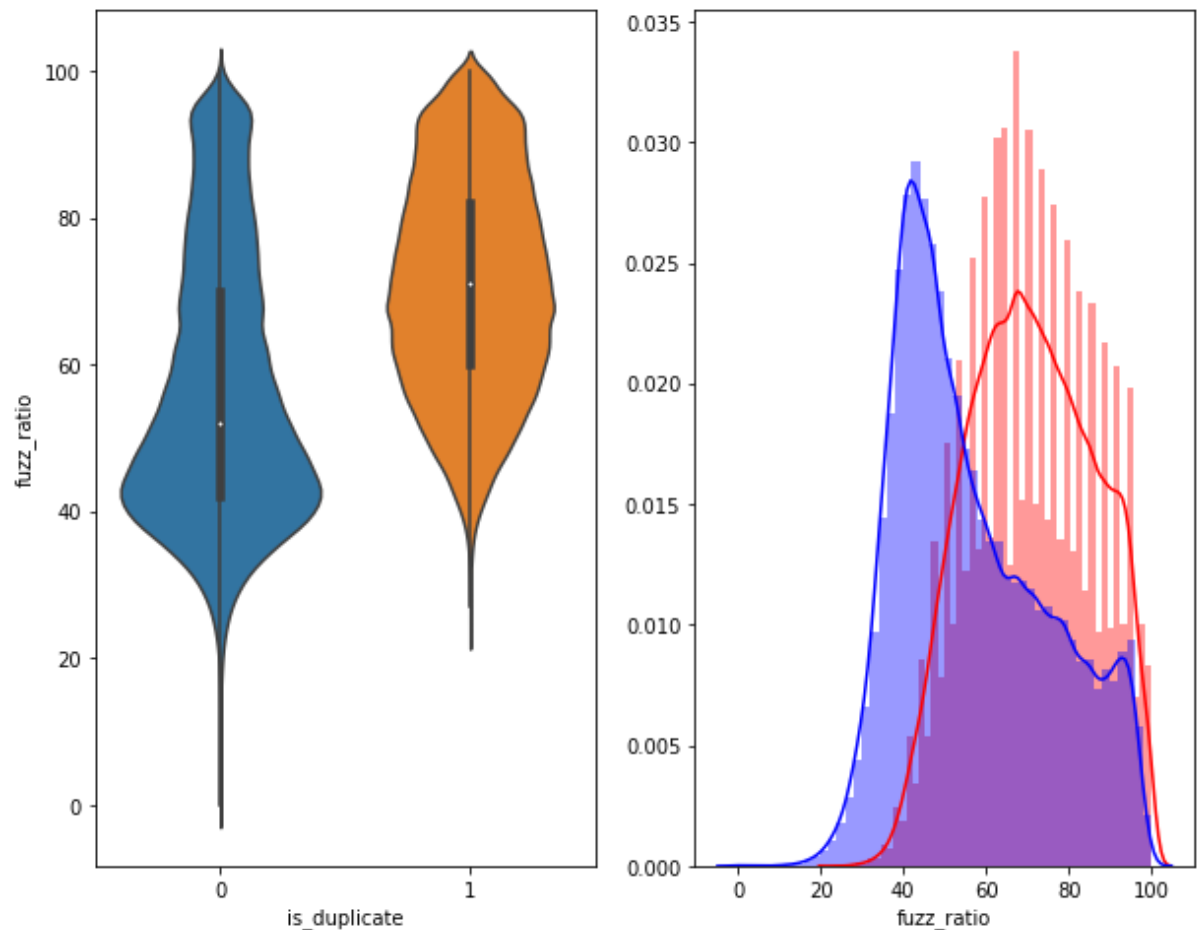
In [32]: 1 # Distribution of the token_sort_ratio
2 plt.figure(figsize=(10, 8))
3
4 plt.subplot(1,2,1)
5 sns.violinplot(x = 'is_duplicate', y = 'token_sort_ratio', data = df[0:] , )
6
7 plt.subplot(1,2,2)
8 sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = '
9 sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = '
10 plt.show()

```





```
In [31]: 1 plt.figure(figsize=(10, 8))
2
3 plt.subplot(1,2,1)
4 sns.violinplot(x = 'is_duplicate', y = 'fuzz_ratio', data = df[0:] , )
5
6 plt.subplot(1,2,2)
7 sns.distplot(df[df['is_duplicate'] == 1.0]['fuzz_ratio'][0:] , label = "1", color = "#1f77b4")
8 sns.distplot(df[df['is_duplicate'] == 0.0]['fuzz_ratio'][0:] , label = "0", color = "#d62728")
9 plt.show()
```



### 3.5.2 Visualization

```
In [34]: 1 # Using TSNE for Dimentionalty reduction for 15 Features(Generated after clea
2
3 from sklearn.preprocessing import MinMaxScaler
4
5 dfp_subsampled = df[0:5000]
6 X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_mi
7 y = dfp_subsampled['is_duplicate'].values
```

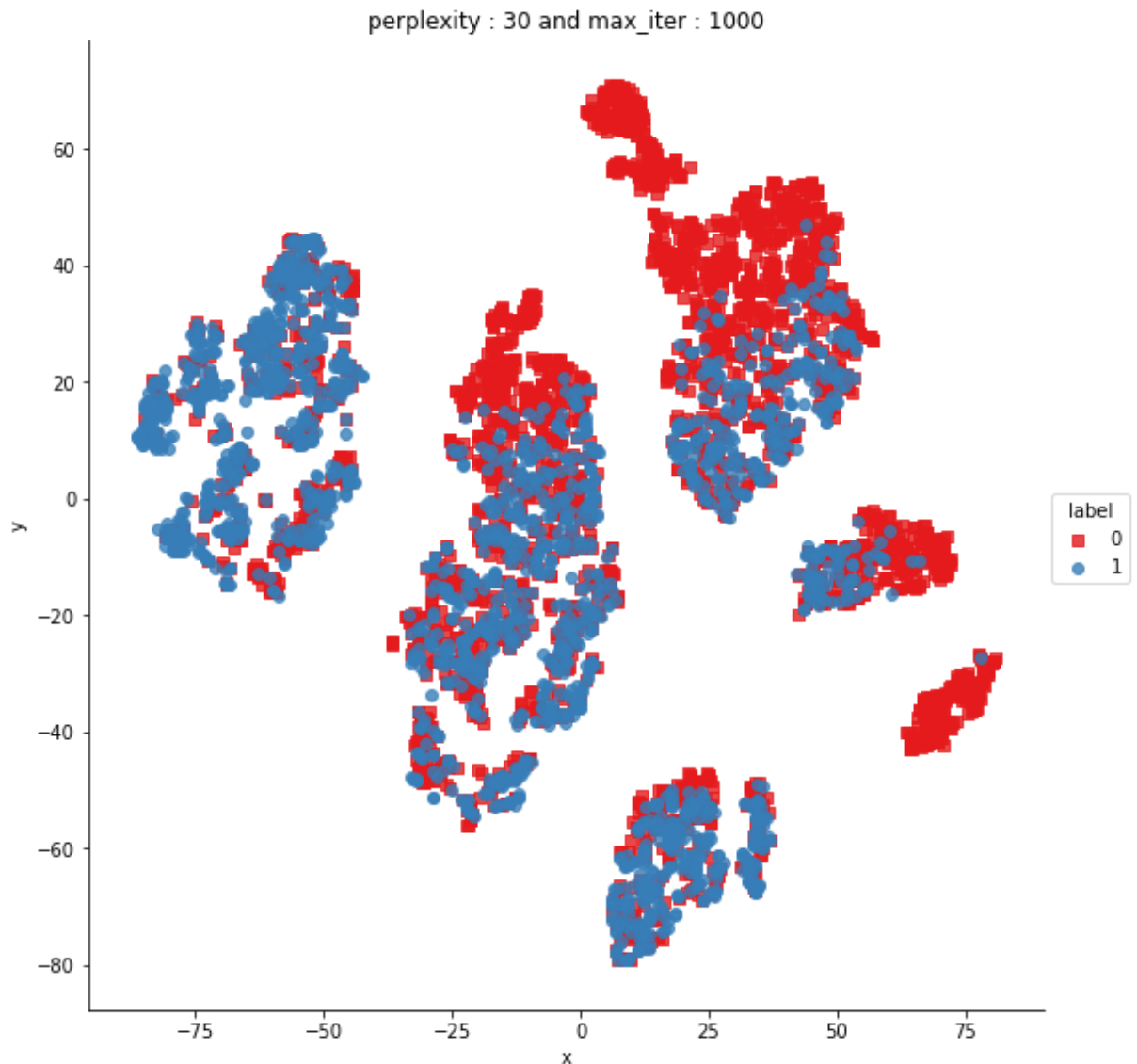
In [42]:

```
1 tsne2d = TSNE(  
2     n_components=2,  
3     init='random', # pca  
4     random_state=101,  
5     method='barnes_hut',  
6     n_iter=1000,  
7     verbose=2,  
8     angle=0.5  
9 ).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...  
[t-SNE] Indexed 5000 samples in 0.011s...  
[t-SNE] Computed neighbors for 5000 samples in 0.912s...  
[t-SNE] Computed conditional probabilities for sample 1000 / 5000  
[t-SNE] Computed conditional probabilities for sample 2000 / 5000  
[t-SNE] Computed conditional probabilities for sample 3000 / 5000  
[t-SNE] Computed conditional probabilities for sample 4000 / 5000  
[t-SNE] Computed conditional probabilities for sample 5000 / 5000  
[t-SNE] Mean sigma: 0.116557  
[t-SNE] Computed conditional probabilities in 0.433s  
[t-SNE] Iteration 50: error = 80.9244080, gradient norm = 0.0428133 (50 iterations in 13.099s)  
[t-SNE] Iteration 100: error = 70.3858795, gradient norm = 0.0100968 (50 iterations in 9.067s)  
[t-SNE] Iteration 150: error = 68.6138382, gradient norm = 0.0058392 (50 iterations in 9.602s)  
[t-SNE] Iteration 200: error = 67.7700119, gradient norm = 0.0036596 (50 iterations in 9.121s)  
[t-SNE] Iteration 250: error = 67.2725067, gradient norm = 0.0034962 (50 iterations in 11.305s)  
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.272507  
[t-SNE] Iteration 300: error = 1.7737305, gradient norm = 0.0011918 (50 iterations in 8.289s)  
[t-SNE] Iteration 350: error = 1.3720417, gradient norm = 0.0004822 (50 iterations in 10.526s)  
[t-SNE] Iteration 400: error = 1.2039998, gradient norm = 0.0002768 (50 iterations in 9.600s)  
[t-SNE] Iteration 450: error = 1.1133438, gradient norm = 0.0001881 (50 iterations in 11.827s)  
[t-SNE] Iteration 500: error = 1.0579143, gradient norm = 0.0001434 (50 iterations in 8.941s)  
[t-SNE] Iteration 550: error = 1.0221983, gradient norm = 0.0001164 (50 iterations in 11.092s)  
[t-SNE] Iteration 600: error = 0.9987167, gradient norm = 0.0001039 (50 iterations in 11.467s)  
[t-SNE] Iteration 650: error = 0.9831534, gradient norm = 0.0000938 (50 iterations in 11.799s)  
[t-SNE] Iteration 700: error = 0.9722011, gradient norm = 0.0000858 (50 iterations in 12.028s)  
[t-SNE] Iteration 750: error = 0.9643636, gradient norm = 0.0000799 (50 iterations in 12.120s)  
[t-SNE] Iteration 800: error = 0.9584482, gradient norm = 0.0000785 (50 iterations in 11.867s)  
[t-SNE] Iteration 850: error = 0.9538348, gradient norm = 0.0000739 (50 iterations in 11.461s)  
[t-SNE] Iteration 900: error = 0.9496906, gradient norm = 0.0000712 (50 iterations in 11.461s)
```

ions in 11.023s)  
[t-SNE] Iteration 950: error = 0.9463405, gradient norm = 0.0000673 (50 iterations in 11.755s)  
[t-SNE] Iteration 1000: error = 0.9432716, gradient norm = 0.0000662 (50 iterations in 11.493s)  
[t-SNE] Error after 1000 iterations: 0.943272

```
In [44]: 1 df = pd.DataFrame({'x':tsne2d[:,0], 'y':tsne2d[:,1] , 'label':y})  
2  
3 # draw the plot in appropriate place in the grid  
4 sns.lmplot(data=df, x='x', y='y', hue='label', fit_reg=False, size=8,palette='  
5 plt.title("perplexity : {} and max_iter : {}".format(30, 1000))  
6 plt.show()
```



In [45]:

```
1 from sklearn.manifold import TSNE
2 tsne3d = TSNE(
3     n_components=3,
4     init='random', # pca
5     random_state=101,
6     method='barnes_hut',
7     n_iter=1000,
8     verbose=2,
9     angle=0.5
10 ).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.010s...
[t-SNE] Computed neighbors for 5000 samples in 0.935s...
[t-SNE] Computed conditional probabilities for sample 1000 / 5000
[t-SNE] Computed conditional probabilities for sample 2000 / 5000
[t-SNE] Computed conditional probabilities for sample 3000 / 5000
[t-SNE] Computed conditional probabilities for sample 4000 / 5000
[t-SNE] Computed conditional probabilities for sample 5000 / 5000
[t-SNE] Mean sigma: 0.116557
[t-SNE] Computed conditional probabilities in 0.363s
[t-SNE] Iteration 50: error = 77.7944183, gradient norm = 0.1014017 (50 iterations in 34.931s)
[t-SNE] Iteration 100: error = 69.2682266, gradient norm = 0.0248657 (50 iterations in 15.147s)
[t-SNE] Iteration 150: error = 67.7877655, gradient norm = 0.0150941 (50 iterations in 13.761s)
[t-SNE] Iteration 200: error = 67.1991119, gradient norm = 0.0126559 (50 iterations in 13.425s)
[t-SNE] Iteration 250: error = 66.8560715, gradient norm = 0.0074975 (50 iterations in 12.904s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 66.856071
[t-SNE] Iteration 300: error = 1.2356015, gradient norm = 0.0007033 (50 iterations in 13.302s)
[t-SNE] Iteration 350: error = 0.9948602, gradient norm = 0.0001997 (50 iterations in 18.898s)
[t-SNE] Iteration 400: error = 0.9168936, gradient norm = 0.0001430 (50 iterations in 13.397s)
[t-SNE] Iteration 450: error = 0.8863022, gradient norm = 0.0000975 (50 iterations in 16.379s)
[t-SNE] Iteration 500: error = 0.8681002, gradient norm = 0.0000854 (50 iterations in 17.791s)
[t-SNE] Iteration 550: error = 0.8564141, gradient norm = 0.0000694 (50 iterations in 17.060s)
[t-SNE] Iteration 600: error = 0.8470711, gradient norm = 0.0000640 (50 iterations in 15.454s)
[t-SNE] Iteration 650: error = 0.8389117, gradient norm = 0.0000561 (50 iterations in 17.562s)
[t-SNE] Iteration 700: error = 0.8325295, gradient norm = 0.0000529 (50 iterations in 13.443s)
[t-SNE] Iteration 750: error = 0.8268463, gradient norm = 0.0000528 (50 iterations in 17.981s)
[t-SNE] Iteration 800: error = 0.8219477, gradient norm = 0.0000477 (50 iterations in 17.448s)
[t-SNE] Iteration 850: error = 0.8180174, gradient norm = 0.0000490 (50 iterations in 18.376s)
```

[t-SNE] Iteration 900: error = 0.8150476, gradient norm = 0.0000456 (50 iterations in 17.778s)  
[t-SNE] Iteration 950: error = 0.8122067, gradient norm = 0.0000472 (50 iterations in 16.983s)  
[t-SNE] Iteration 1000: error = 0.8095787, gradient norm = 0.0000489 (50 iterations in 18.581s)  
[t-SNE] Error after 1000 iterations: 0.809579

In [46]:

```
1 trace1 = go.Scatter3d(  
2     x=tsne3d[:,0],  
3     y=tsne3d[:,1],  
4     z=tsne3d[:,2],  
5     mode='markers',  
6     marker=dict(  
7         sizemode='diameter',  
8         color = y,  
9         colorscale = 'Portland',  
10        colorbar = dict(title = 'duplicate'),  
11        line=dict(color='rgb(255, 255, 255)'),  
12        opacity=0.75  
13    )  
14 )  
15  
16 data=[trace1]  
17 layout=dict(height=800, width=800, title='3d embedding with engineered feature  
18 fig=dict(data=data, layout=layout)  
19 py.iplot(fig, filename='3DBubble')
```

## Computing Word Vectors (Average W2V)

```
In [2]: 1 # avoid decoding problems
2 df = pd.read_csv("train.csv")
3
4 # encode questions to unicode
5 # https://stackoverflow.com/a/6812069
6 # ----- python 2 -----
7 # df['question1'] = df['question1'].apply(lambda x: unicode(str(x), "utf-8"))
8 # df['question2'] = df['question2'].apply(lambda x: unicode(str(x), "utf-8"))
9 # ----- python 3 -----
10 df['question1'] = df['question1'].apply(lambda x: str(x))
11 df['question2'] = df['question2'].apply(lambda x: str(x))
```

```
In [3]: 1 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 $23^{24}$ is divided by 100	0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

```
In [4]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.feature_extraction.text import CountVectorizer
3 # merge texts
4 questions = list(df['question1']) + list(df['question2'])
5
6 tfidf = TfidfVectorizer(lowercase=False, )
```








```
In [10]: 1 # data before preprocessing
         2 df2.head()
```

```
Out[10]:
```

	id	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	word_Total	w
0	0	1	1	66	57	14	12	10.0	23.0	
1	1	4	1	51	88	8	13	4.0	20.0	
2	2	1	1	73	59	14	10	4.0	24.0	
3	3	1	1	50	65	11	9	0.0	19.0	
4	4	3	1	76	39	13	7	2.0	20.0	

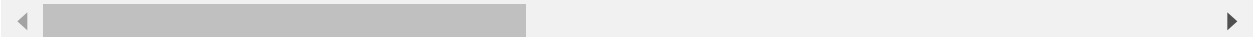


```
In [11]: 1 # Questions 1 tfidf weighted word2vec
         2 df3_q1.head()
```

```
Out[11]:
```

	0	1	2	3	4	5	6	7
0	121.929927	100.083900	72.497894	115.641800	-48.370870	34.619058	-172.057787	-92.502617
1	-78.070939	54.843781	82.738482	98.191872	-51.234859	55.013510	-39.140730	-82.692352
2	-5.355015	73.671810	14.376365	104.130241	1.433537	35.229116	-148.519385	-97.124595
3	5.778359	-34.712038	48.999631	59.699204	40.661263	-41.658731	-36.808594	24.170655
4	51.138220	38.587312	123.639488	53.333041	-47.062739	37.356212	-298.722753	-106.421119

5 rows × 384 columns




```
In [12]: 1 # Questions 2 tfidf weighted word2vec
         2 df3_q2.head()
```

```
Out[12]:
```

	0	1	2	3	4	5	6	7
0	125.983301	95.636485	42.114702	95.449980	-37.386295	39.400078	-148.116070	-87.851475
1	-106.871904	80.290331	79.066297	59.302092	-42.175328	117.616655	-144.364237	-127.131513
2	7.072875	15.513378	1.846914	85.937583	-33.808811	94.702337	-122.256856	-114.009530
3	39.421531	44.136989	-24.010929	85.265863	-0.339022	-9.323137	-60.499651	-37.044763
4	31.950101	62.854106	1.778164	36.218768	-45.130875	66.674880	-106.342341	-22.901008

5 rows × 384 columns



```
In [13]: 1 print("Number of features in nlp dataframe :", df1.shape[1])
2 print("Number of features in preprocessed dataframe :", df2.shape[1])
3 print("Number of features in question1 w2v dataframe :", df3_q1.shape[1])
4 print("Number of features in question2 w2v dataframe :", df3_q2.shape[1])
5 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 : 384
Number of features in question2 w2v dataframe : 384
Number of features in final dataframe : 794
```

```
In [14]: 1 # storing the final features to csv file
2 if not os.path.isfile('final_features.csv'):
3     df3_q1['id']=df1['id']
4     df3_q2['id']=df1['id']
5     df1 = df1.merge(df2, on='id',how='left')
6     df2 = df3_q1.merge(df3_q2, on='id',how='left')
7     result = df1.merge(df2, on='id',how='left')
8     result.to_csv('final_features.csv')
```

## Applying Machine Learning Algorithms

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

```
In [4]: 1 #Creating db file from csv
2 if not os.path.isfile('train.db'):
3     disk_engine = create_engine('sqlite:///train.db')
4     start = dt.datetime.now()
5     chunksize = 180000
6     j = 0
7     index_start = 1
8     for df in pd.read_csv('final_features.csv', names=['Unnamed: 0','id','is_c
9         df.index += index_start
10        j+=1
11        print('{} rows'.format(j*chunksize))
12        df.to_sql('data', disk_engine, if_exists='append')
13        index_start = df.index[-1] + 1
```

```

In [5]: 1 #http://www.sqlitetutorial.net/sqlite-python/create-tables/
2 def create_connection(db_file):
3     """ create a database connection to the SQLite database
4         specified by db_file
5     :param db_file: database file
6     :return: Connection object or None
7     """
8     try:
9         conn = sqlite3.connect(db_file)
10        return conn
11    except Error as e:
12        print(e)
13
14    return None
15
16
17 def checkTableExists(dbcon):
18     cursr = dbcon.cursor()
19     str = "select name from sqlite_master where type='table'"
20     table_names = cursr.execute(str)
21     print("Tables in the databse:")
22     tables = table_names.fetchall()
23     print(tables[0][0])
24     return(len(tables))

```

```

In [6]: 1 read_db = 'train.db'
2 conn_r = create_connection(read_db)
3 checkTableExists(conn_r)
4 conn_r.close()

```

Tables in the databse:  
data

```

In [7]: 1 # try to sample data according to the computing power you have
2 if os.path.isfile(read_db):
3     conn_r = create_connection(read_db)
4     if conn_r is not None:
5         # for selecting first 1M rows
6         # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)
7
8         # for selecting random points
9         data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 100001", conn_r)
10        conn_r.commit()
11        conn_r.close()

```

```

In [8]: 1 # remove the first row
2 data.drop(data.index[0], inplace=True)
3 y_true = data['is_duplicate']
4 data.drop(['Unnamed: 0', 'id', 'index', 'is_duplicate'], axis=1, inplace=True)

```

In [9]:

```
1 data.head()
```

Out[9]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_mi
1	0.199996000079998	0.166663888935184	0.0	0.0	0.1428551020699
2	0.399992000159997	0.399992000159997	0.499987500312492	0.499987500312492	0.44443950622770
3	0.833319444675922	0.714275510349852	0.999983333611106	0.857130612419823	0.68749570315185
4	0.0	0.0	0.599988000239995	0.499991666805553	0.24999791668402
5	0.749981250468738	0.749981250468738	0.499987500312492	0.499987500312492	0.62499218759765

5 rows × 794 columns

## Converting strings to numerics

In [10]:

```
1 # after we read from sql table each entry was read it as a string
2 # we convert all the features into numeric before we apply any model
3 cols = list(data.columns)
4 for i in cols:
5     data[i] = data[i].apply(pd.to_numeric)
6     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
-1 in words
```

In [18]:

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

## Random train test split( 70:30)

In [19]:

```
1 X_train,X_test, y_train, y_test = train_test_split(data, y_true, stratify=y_tr
```

In [20]:

```
1 f
```

Number of data points in train data : (70000, 794)

Number of data points in test data : (30000, 794)

In [25]:

```
1 print("-"*10, "Distribution of output variable in train data", "-"*10)
2 train_distr = Counter(y_train)
3 train_len = len(y_train)
4 print("Class 0: ",int(train_distr[0])/train_len,"Class 1: ", int(train_distr[1])/train_len)
5 print("-"*10, "Distribution of output variable in train data", "-"*10)
6 test_distr = Counter(y_test)
7 test_len = len(y_test)
8 print("Class 0: ",int(test_distr[0])/test_len, "Class 1: ",int(test_distr[1])/test_len)
```

----- Distribution of output variable in train data -----

Class 0: 0.6324857142857143 Class 1: 0.36751428571428574

----- Distribution of output variable in train data -----

Class 0: 0.3675 Class 1: 0.3675

In [34]:

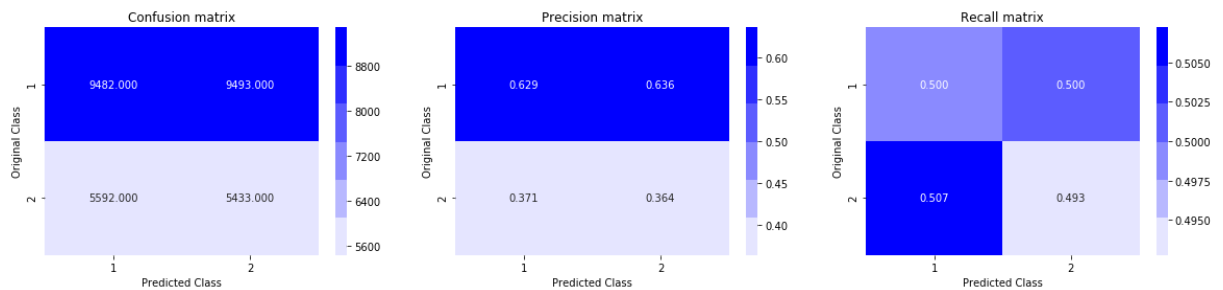
```
1  # This function plots the confusion matrices given y_i, y_i_hat.
2  def plot_confusion_matrix(test_y, predict_y):
3      C = confusion_matrix(test_y, predict_y)
4      # C = 9,9 matrix, each cell (i,j) represents number of points of class i c
5
6      A = (((C.T)/(C.sum(axis=1))).T)
7      #divid each element of the confusion matrix with the sum of elements in th
8
9      # C = [[1, 2],
10     #      [3, 4]]
11     # C.T = [[1, 3],
12     #        [2, 4]]
13     # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to
14     # C.sum(axix =1) = [[3, 7]]
15     # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
16     #                             [2/3, 4/7]]
17
18     # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
19     #                             [3/7, 4/7]]
20     # sum of row elements = 1
21
22     B = (C/C.sum(axis=0))
23     #divid each element of the confusion matrix with the sum of elements in th
24     # C = [[1, 2],
25     #      [3, 4]]
26     # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to
27     # C.sum(axix =0) = [[4, 6]]
28     # (C/C.sum(axis=0)) = [[1/4, 2/6],
29     #                       [3/4, 4/6]]
30     plt.figure(figsize=(20,4))
31
32     labels = [1,2]
33     # representing A in heatmap format
34     cmap=sns.light_palette("blue")
35     plt.subplot(1, 3, 1)
36     sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
37     plt.xlabel('Predicted Class')
38     plt.ylabel('Original Class')
39     plt.title("Confusion matrix")
40
41     plt.subplot(1, 3, 2)
42     sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
43     plt.xlabel('Predicted Class')
44     plt.ylabel('Original Class')
45     plt.title("Precision matrix")
46
47     plt.subplot(1, 3, 3)
48     # representing B in heatmap format
49     sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, ytick
50     plt.xlabel('Predicted Class')
51     plt.ylabel('Original Class')
52     plt.title("Recall matrix")
53
54     plt.show()
```

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

In [35]:

```
1 # we need to generate 9 numbers and the sum of numbers should be 1
2 # one solution is to generate 9 numbers and divide each of the numbers by their
3 # ref: https://stackoverflow.com/a/18662466/4084039
4 # we create a output array that has exactly same size as the CV data
5 predicted_y = np.zeros((test_len,2))
6 for i in range(test_len):
7     rand_probs = np.random.rand(1,2)
8     predicted_y[i] = ((rand_probs/sum(sum(rand_probs))))[0])
9 print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y),
10
11 predicted_y = np.argmax(predicted_y, axis=1)
12 plot_confusion_matrix(y_test, predicted_y)
```

Log loss on Test Data using Random Model 0.887242646958



## Logistic Regression with hyperparameter tuning

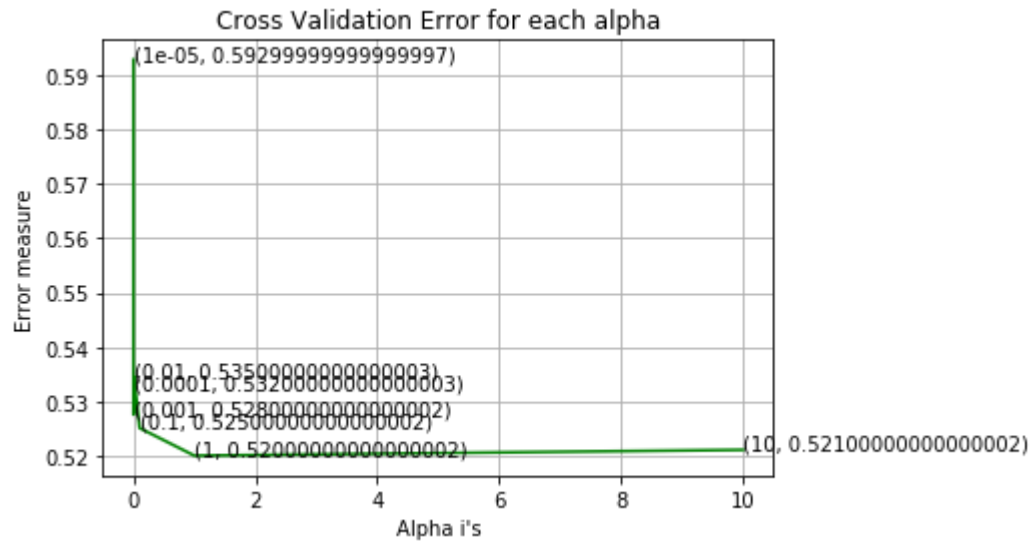


In [50]:

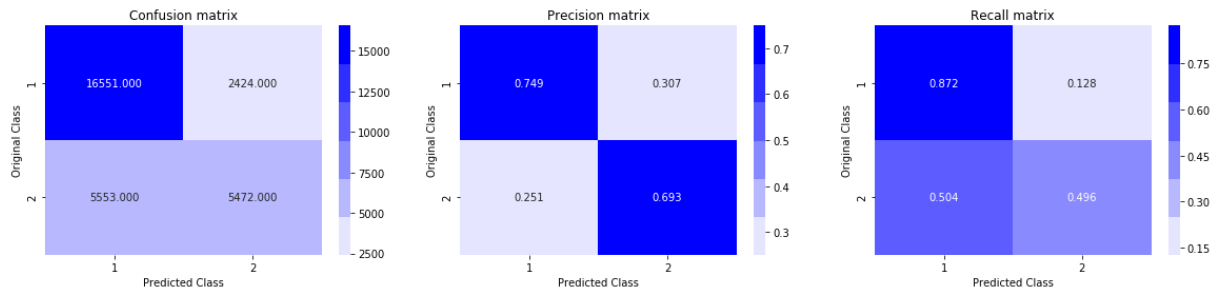
```
1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/ge
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_i
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic G
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, ep
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, pre
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")
35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)
```

For values of alpha = 1e-05 The log loss is: 0.592800211149  
For values of alpha = 0.0001 The log loss is: 0.532351700629

For values of alpha = 0.001 The log loss is: 0.527562275995  
 For values of alpha = 0.01 The log loss is: 0.534535408885  
 For values of alpha = 0.1 The log loss is: 0.525117052926  
 For values of alpha = 1 The log loss is: 0.520035530431  
 For values of alpha = 10 The log loss is: 0.521097925307



For values of best alpha = 1 The train log loss is: 0.513842874233  
 For values of best alpha = 1 The test log loss is: 0.520035530431  
 Total number of data points : 30000



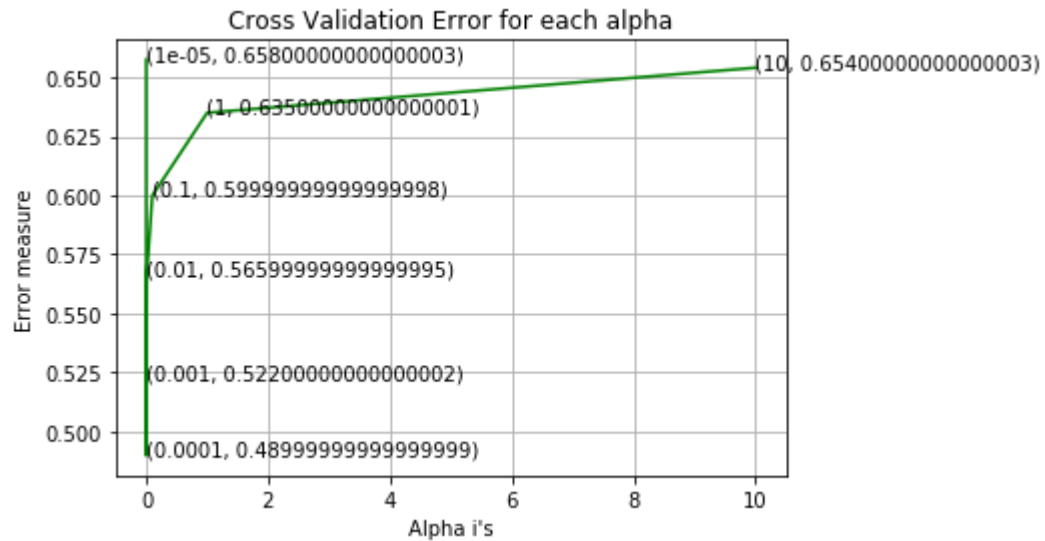
## Linear SVM with hyperparameter tuning

In [51]:

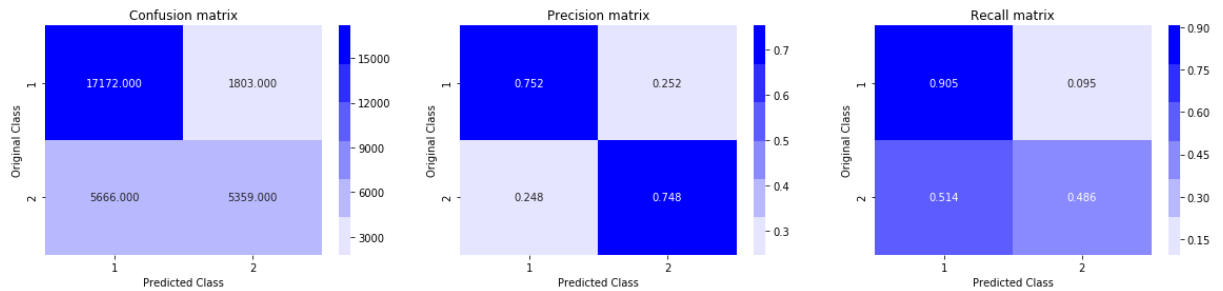
```
1  alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3  # read more about SGDClassifier() at http://scikit-learn.org/stable/modules/ge
4  # -----
5  # default parameters
6  # SGDClassifier(loss='hinge', penalty='l2', alpha=0.0001, l1_ratio=0.15, fit_i
7  # shuffle=True, verbose=0, epsilon=0.1, n_jobs=1, random_state=None, learning_
8  # class_weight=None, warm_start=False, average=False, n_iter=None)
9
10 # some of methods
11 # fit(X, y[, coef_init, intercept_init, ...]) Fit linear model with Stochastic G
12 # predict(X) Predict class labels for samples in X.
13
14 #-----
15 # video link:
16 #-----
17
18
19 log_error_array=[]
20 for i in alpha:
21     clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
22     clf.fit(X_train, y_train)
23     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
24     sig_clf.fit(X_train, y_train)
25     predict_y = sig_clf.predict_proba(X_test)
26     log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, ep
27     print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, pre
28
29 fig, ax = plt.subplots()
30 ax.plot(alpha, log_error_array, c='g')
31 for i, txt in enumerate(np.round(log_error_array, 3)):
32     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
33 plt.grid()
34 plt.title("Cross Validation Error for each alpha")
35 plt.xlabel("Alpha i's")
36 plt.ylabel("Error measure")
37 plt.show()
38
39
40 best_alpha = np.argmin(log_error_array)
41 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', rand
42 clf.fit(X_train, y_train)
43 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
44 sig_clf.fit(X_train, y_train)
45
46 predict_y = sig_clf.predict_proba(X_train)
47 print('For values of best alpha = ', alpha[best_alpha], "The train log loss is
48 predict_y = sig_clf.predict_proba(X_test)
49 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:
50 predicted_y = np.argmax(predict_y, axis=1)
51 print("Total number of data points :", len(predicted_y))
52 plot_confusion_matrix(y_test, predicted_y)
```

For values of alpha = 1e-05 The log loss is: 0.657611721261  
For values of alpha = 0.0001 The log loss is: 0.489669093534

For values of alpha = 0.001 The log loss is: 0.521829068562  
 For values of alpha = 0.01 The log loss is: 0.566295616914  
 For values of alpha = 0.1 The log loss is: 0.599957866217  
 For values of alpha = 1 The log loss is: 0.635059427016  
 For values of alpha = 10 The log loss is: 0.654159467907



For values of best alpha = 0.0001 The train log loss is: 0.478054677285  
 For values of best alpha = 0.0001 The test log loss is: 0.489669093534  
 Total number of data points : 30000



## XGBoost

In [52]:

```
1 import xgboost as xgb
2 params = {}
3 params['objective'] = 'binary:logistic'
4 params['eval_metric'] = 'logloss'
5 params['eta'] = 0.02
6 params['max_depth'] = 4
7
8 d_train = xgb.DMatrix(X_train, label=y_train)
9 d_test = xgb.DMatrix(X_test, label=y_test)
10
11 watchlist = [(d_train, 'train'), (d_test, 'valid')]
12
13 bst = xgb.train(params, d_train, 400, watchlist, early_stopping_rounds=20, ver
14
15 xgdmatrix = xgb.DMatrix(X_train, y_train)
16 predict_y = bst.predict(d_test)
17 print("The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_,
```

```
[0]    train-logloss:0.684819  valid-logloss:0.684845
```

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.

```
[10]    train-logloss:0.61583  valid-logloss:0.616104
[20]    train-logloss:0.564616 valid-logloss:0.565273
[30]    train-logloss:0.525758 valid-logloss:0.52679
[40]    train-logloss:0.496661 valid-logloss:0.498021
[50]    train-logloss:0.473563 valid-logloss:0.475182
[60]    train-logloss:0.455315 valid-logloss:0.457186
[70]    train-logloss:0.440442 valid-logloss:0.442482
[80]    train-logloss:0.428424 valid-logloss:0.430795
[90]    train-logloss:0.418803 valid-logloss:0.421447
[100]   train-logloss:0.41069  valid-logloss:0.413583
[110]   train-logloss:0.403831 valid-logloss:0.40693
[120]   train-logloss:0.398076 valid-logloss:0.401402
[130]   train-logloss:0.393305 valid-logloss:0.396851
[140]   train-logloss:0.38913  valid-logloss:0.392952
[150]   train-logloss:0.385469 valid-logloss:0.389521
[160]   train-logloss:0.382327 valid-logloss:0.386667
[170]   train-logloss:0.379541 valid-logloss:0.384148
[180]   train-logloss:0.377014 valid-logloss:0.381932
[190]   train-logloss:0.374687 valid-logloss:0.379883
[200]   train-logloss:0.372585 valid-logloss:0.378068
[210]   train-logloss:0.370615 valid-logloss:0.376367
[220]   train-logloss:0.368559 valid-logloss:0.374595
[230]   train-logloss:0.366545 valid-logloss:0.372847
[240]   train-logloss:0.364708 valid-logloss:0.371311
[250]   train-logloss:0.363021 valid-logloss:0.369886
[260]   train-logloss:0.36144  valid-logloss:0.368673
[270]   train-logloss:0.359899 valid-logloss:0.367421
[280]   train-logloss:0.358465 valid-logloss:0.366395
[290]   train-logloss:0.357128 valid-logloss:0.365361
[300]   train-logloss:0.355716 valid-logloss:0.364315
[310]   train-logloss:0.354425 valid-logloss:0.363403
[320]   train-logloss:0.353276 valid-logloss:0.362595
[330]   train-logloss:0.352084 valid-logloss:0.361823
```

```

[340] train-logloss:0.351051 valid-logloss:0.361167
[350] train-logloss:0.349867 valid-logloss:0.36043
[360] train-logloss:0.348829 valid-logloss:0.359773
[370] train-logloss:0.347689 valid-logloss:0.359019
[380] train-logloss:0.346607 valid-logloss:0.358311
[390] train-logloss:0.345568 valid-logloss:0.357674
The test log loss is: 0.357054433715

```

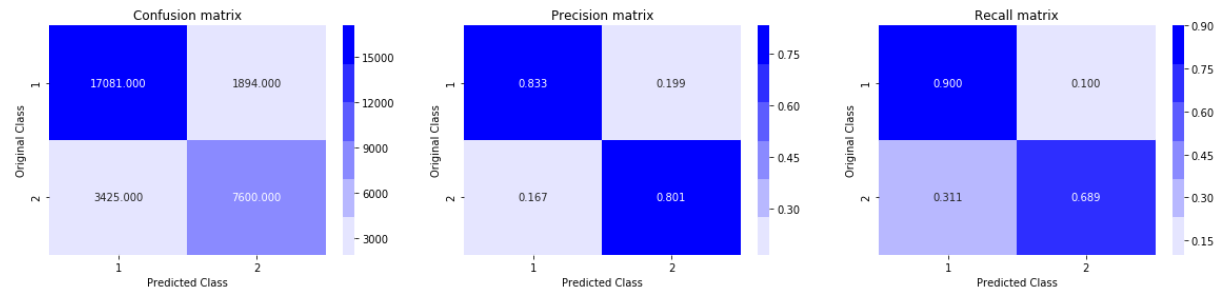
In [60]:

```

1 predicted_y = np.array(predict_y>0.5,dtype=int)
2 print("Total number of data points :", len(predicted_y))
3 plot_confusion_matrix(y_test, predicted_y)

```

Total number of data points : 30000



## Applying Logistic regression and linear SVM on Simple TFIDF features

In [2]:

```

1 # avoid decoding problems
2 df = pd.read_csv("train.csv")
3
4 # encode questions to unicode
5 # https://stackoverflow.com/a/6812069
6 # ----- python 2 -----
7 # df['question1'] = df['question1'].apply(lambda x: unicode(str(x),"utf-8"))
8 # df['question2'] = df['question2'].apply(lambda x: unicode(str(x),"utf-8"))
9 # ----- python 3 -----
10 df['question1'] = df['question1'].apply(lambda x: str(x))
11 df['question2'] = df['question2'].apply(lambda x: str(x))

```

In [3]:

```
1 df.shape
```

Out[3]: (404290, 6)

In [25]:

```
1 final_df = pd.read_csv('final_features.csv')
```

In [26]:

```
1 final_col = list(final_df.columns)
```

In [29]:

```

1 final_fil_df = final_df[final_col[1:29]]
2 final_fil_df.to_csv('final_features_fil.csv', index=False)

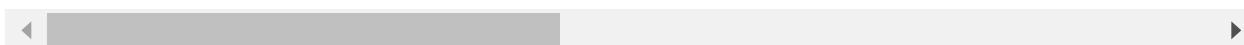
```

```
In [30]: 1 final_fil_df.head()
```

```
Out[30]:
```

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

5 rows × 28 columns



```
In [33]: 1 comb_df = df.merge(final_fil_df, on=["id"])
```

```
In [35]: 1 comb_df.drop(['is_duplicate_x'], axis=1, inplace=True)
```

```
In [37]: 1 comb_df.rename(index=str, columns={'is_duplicate_y':'is_duplicate'}, inplace=True)
```

```
In [40]: 1 comb_df.to_csv('comb_fil_df.csv', index=False)
```

```
In [42]: 1 comb_df.shape
```

```
Out[42]: (404290, 32)
```

\*\* Computing on all the data points is raising "MemoryError". So, I have used 10000 points from data. \*\*

```
In [44]: 1 sam_comb_df = comb_df[:10000]
2
3 sam_comb_df.to_csv('sample_comb_fil_df.csv', index=False)
```

```
In [103]: 1 sam_comb_df = pd.read_csv('sample_comb_fil_df.csv', encoding='cp1252')
```

```
In [45]: 1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.feature_extraction.text import CountVectorizer
3 # merge texts
4 questions = list(sam_comb_df['question1']) + list(sam_comb_df['question2'])
5
6 tfidf = TfidfVectorizer(lowercase=False, )
```

```
In [46]: 1 tfidf_vec = tfidf.fit_transform(questions)
2
```

```
In [49]: 1 # Assigning the Sparse Matrices of Q1 and Q2 Sparse Matrices
        2 q1_vec = tfidf_vec[:sam_comb_df.shape[0]]
        3 q2_vec = tfidf_vec[sam_comb_df.shape[0]:]
```

```
In [55]: 1 q1_vec = q1_vec.toarray()
        2 q2_vec = q2_vec.toarray()
```

```
In [58]: 1 q1_vec.shape
```

```
Out[58]: (10000, 17418)
```

```
In [61]: 1 index_x = ["x_" + str(i) for i in range(0, q1_vec.shape[1])]
        2 index_y = ["y_" + str(i) for i in range(0, q2_vec.shape[1])]
        3 df3_q1 = pd.DataFrame(data=q1_vec, columns=index_x)
        4 df3_q2 = pd.DataFrame(data=q2_vec, columns=index_y)
```

**\*\* Making X and Y datasets. \*\***

```
In [104]: 1 sam_comb_df_features = sam_comb_df.iloc[:,6:]
```

```
In [92]: 1 sam_tfidf_features = pd.concat([df3_q1, df3_q2], axis=1)
```

```
In [93]: 1 sam_tfidf_features.shape
```

```
Out[93]: (10000, 34836)
```

```
In [94]: 1 sam_comb_df_features.shape
```

```
Out[94]: (10000, 26)
```

```
In [105]: 1 sam_tfidf_features_arr = sam_tfidf_features.values
        2 sam_com_df_features_arr = sam_comb_df_features.values
```

```
In [109]: 1 com_data = np.hstack((sam_tfidf_features_arr, sam_com_df_features_arr))
```

```
In [111]: 1 np.save('com_data.npy', com_data)
```

```
In [85]: 1 y = sam_comb_df['is_duplicate'].values
```

```
In [121]: 1 np.save('com_data_y.npy', y)
```

**\*\* Train and test splitting of data. \*\***

```
In [87]: 1 from sklearn.model_selection import train_test_split
```



```
In [112]: 1 X_train,X_test, y_train, y_test = train_test_split(com_data, y, stratify=y, t
```

```
** Logistic Regression. **
```

```

In [118]: 1 alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3 log_error_array=[]
4 for i in alpha:
5     print("For Alpha: {} \n".format(i))
6     clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
7     clf.fit(X_train, y_train)
8     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
9     sig_clf.fit(X_train, y_train)
10    predict_y = sig_clf.predict_proba(X_test)
11    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, e
12    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, pr
13
14 fig, ax = plt.subplots()
15 ax.plot(alpha, log_error_array, c='g')
16 for i, txt in enumerate(np.round(log_error_array, 3)):
17     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
18 plt.grid()
19 plt.title("Cross Validation Error for each alpha")
20 plt.xlabel("Alpha i's")
21 plt.ylabel("Error measure")
22 plt.show()
23
24
25 best_alpha = np.argmin(log_error_array)
26 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random
27 clf.fit(X_train, y_train)
28 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
29 sig_clf.fit(X_train, y_train)
30
31 predict_y = sig_clf.predict_proba(X_train)
32 print('For values of best alpha = ', alpha[best_alpha], "The train log loss i
33 predict_y = sig_clf.predict_proba(X_test)
34 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
35 predicted_y = np.argmax(predict_y, axis=1)
36 print("Total number of data points :", len(predicted_y))
37 plot_confusion_matrix(y_test, predicted_y)

```

For Alpha: 1e-05

For values of alpha = 1e-05 The log loss is: 0.458121867813

For Alpha: 0.0001

For values of alpha = 0.0001 The log loss is: 0.464271808076

For Alpha: 0.001

For values of alpha = 0.001 The log loss is: 0.46098815888

For Alpha: 0.01

For values of alpha = 0.01 The log loss is: 0.461200475162

For Alpha: 0.1

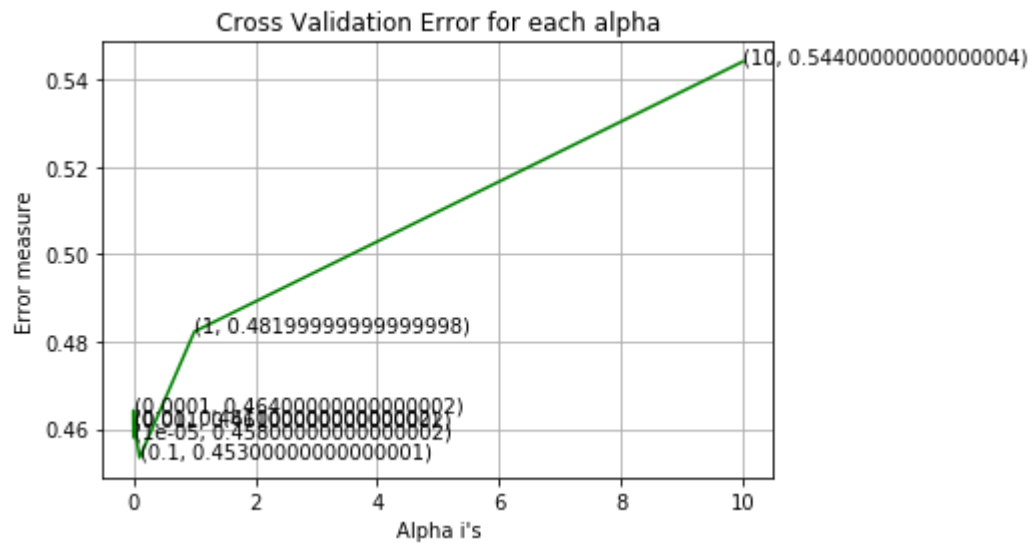
For values of alpha = 0.1 The log loss is: 0.453438911955

For Alpha: 1

For values of alpha = 1 The log loss is: 0.482321302379

For Alpha: 10

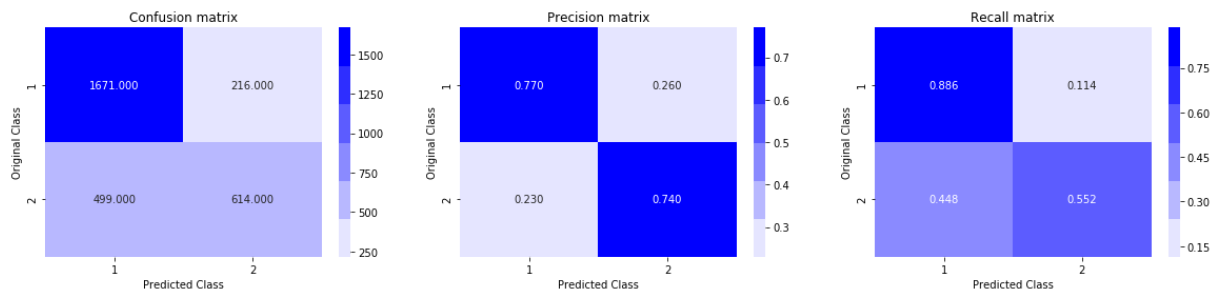
For values of alpha = 10 The log loss is: 0.544186641282



For values of best alpha = 0.1 The train log loss is: 0.469279314863

For values of best alpha = 0.1 The test log loss is: 0.453438911955

Total number of data points : 3000



**\*\* Support Vector Machine.\*\***

```

In [120]: 1 alpha = [10 ** x for x in range(-5, 2)] # hyperparam for SGD classifier.
2
3 log_error_array=[]
4 for i in alpha:
5     print("For Alpha: {} \n".format(i))
6     clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
7     clf.fit(X_train, y_train)
8     sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
9     sig_clf.fit(X_train, y_train)
10    predict_y = sig_clf.predict_proba(X_test)
11    log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, e
12    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, pr
13
14 fig, ax = plt.subplots()
15 ax.plot(alpha, log_error_array, c='g')
16 for i, txt in enumerate(np.round(log_error_array, 3)):
17     ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
18 plt.grid()
19 plt.title("Cross Validation Error for each alpha")
20 plt.xlabel("Alpha i's")
21 plt.ylabel("Error measure")
22 plt.show()
23
24
25 best_alpha = np.argmin(log_error_array)
26 clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l1', loss='hinge', rand
27 clf.fit(X_train, y_train)
28 sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
29 sig_clf.fit(X_train, y_train)
30
31 predict_y = sig_clf.predict_proba(X_train)
32 print('For values of best alpha = ', alpha[best_alpha], "The train log loss i
33 predict_y = sig_clf.predict_proba(X_test)
34 print('For values of best alpha = ', alpha[best_alpha], "The test log loss is
35 predicted_y = np.argmax(predict_y, axis=1)
36 print("Total number of data points :", len(predicted_y))
37 plot_confusion_matrix(y_test, predicted_y)

```

For Alpha: 1e-05

For values of alpha = 1e-05 The log loss is: 0.468595574651

For Alpha: 0.0001

For values of alpha = 0.0001 The log loss is: 0.461435202833

For Alpha: 0.001

For values of alpha = 0.001 The log loss is: 0.478022693837

For Alpha: 0.01

For values of alpha = 0.01 The log loss is: 0.481395005323

For Alpha: 0.1

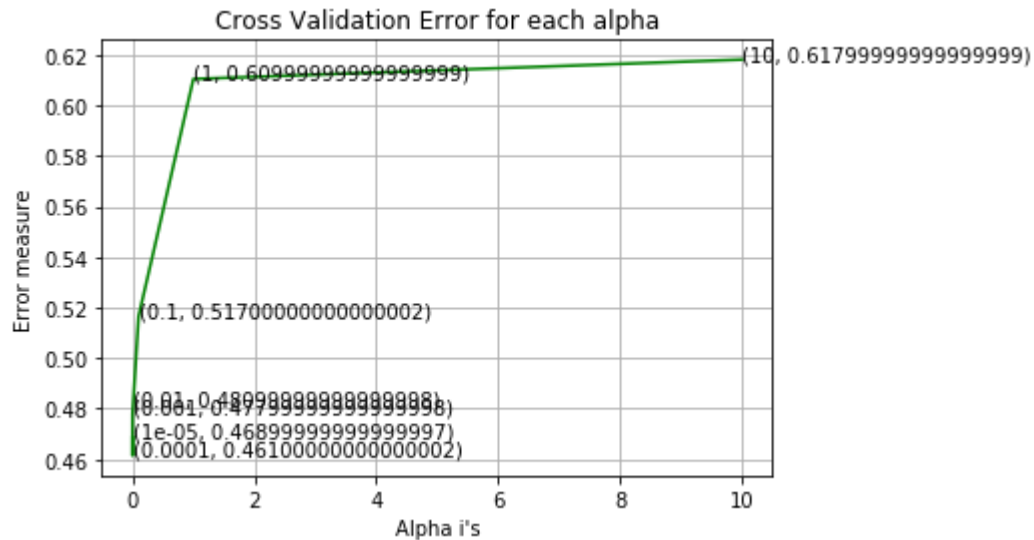
For values of alpha = 0.1 The log loss is: 0.516529462333

For Alpha: 1

For values of alpha = 1 The log loss is: 0.610460609251

For Alpha: 10

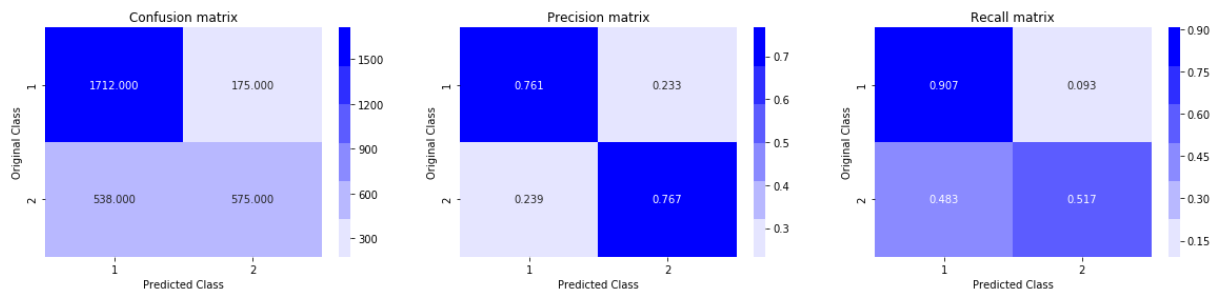
For values of alpha = 10 The log loss is: 0.618137656884



For values of best alpha = 0.0001 The train log loss is: 0.47891025282

For values of best alpha = 0.0001 The test log loss is: 0.461435202833

Total number of data points : 3000



## XGBoost Hyper Parameter Tuning

```
In [63]: 1 def xgboost(x_train, x_test, y_train, y_test):
2         alpha = {
3             'learning_rate' : [0.1, 0.2, 0.3],
4             'n_estimators':[i*10 for i in range(1, 10)],
5             'max_depth':[i for i in range(1,5)],
6         }
7
8         x_model = xgb.XGBClassifier(objective='binary:logistic', eval_metric='logl
9         xgb_model = GridSearchCV(estimator=x_model, param_grid=alpha, cv=10, scor
10        xgb_model.fit(x_train, y_train)
11        print('The Best Parameters are: \n')
12        print(xgb_model.best_params_)
13        predict_y = xgb_model.predict(x_test)
14        print("The test log loss is:",log_loss(y_test, predict_y, eps=1e-15))
15        plot_confusion_matrix(y_test, predict_y)
16        return None
```

**\*\* The model is fitted on the sample of the data using 10000 data points. \*\***

```
In [64]: 1 sam_df = pd.read_csv('sample_comb_fil_df.csv', encoding='cp1252')
```

```
In [65]: 1 data = sam_df.drop(['is_duplicate', 'question1', 'question2', 'id'], axis=1)
2 y = sam_df['is_duplicate'].values
```

```
In [66]: 1 X_train, X_test, y_train, y_test = train_test_split(data, y, test_size=0.33, s
```

```
In [67]: 1 X_train.shape
```

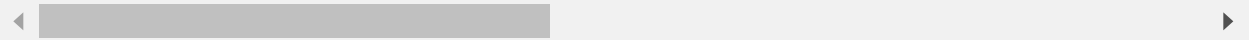
Out[67]: (6700, 28)

```
In [68]: 1 X_train.head()
```

Out[68]:

	qid1	qid2	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	fir
<b>5883</b>	11552	11553	0.000000	0.000000	0.499994	0.444440	0.210525	0.210525	0.0	
<b>2889</b>	5729	5730	0.749981	0.374995	0.666644	0.499988	0.714276	0.416663	0.0	
<b>6069</b>	11901	11902	0.499988	0.333328	0.799984	0.799984	0.599994	0.545450	1.0	
<b>6392</b>	12526	12527	0.999950	0.666644	0.000000	0.000000	0.499988	0.499988	1.0	
<b>66</b>	133	134	0.499988	0.499988	0.999967	0.499992	0.714276	0.499995	0.0	

5 rows × 28 columns



## 4.6 XGBoost

```
► In [69]: 1 xgboost(X_train, X_test, y_train, y_test)
```

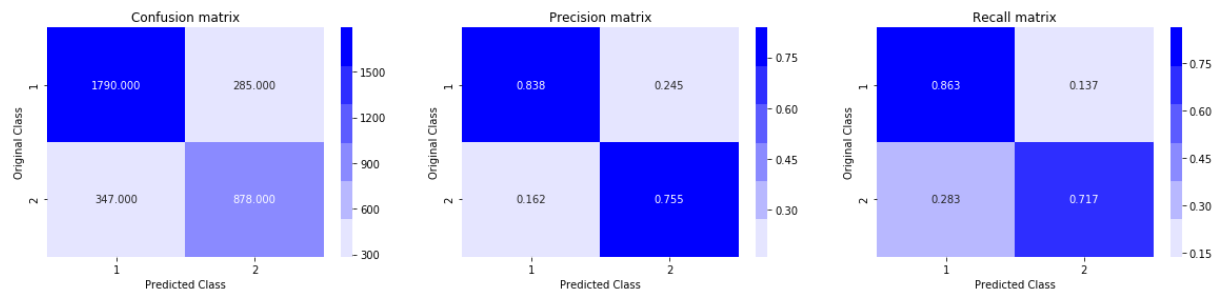
Fitting 10 folds for each of 108 candidates, totalling 1080 fits

```
[Parallel(n_jobs=-1)]: Done    5 tasks      | elapsed:    15.1s
[Parallel(n_jobs=-1)]: Done   10 tasks      | elapsed:    16.0s
[Parallel(n_jobs=-1)]: Done   17 tasks      | elapsed:    17.2s
[Parallel(n_jobs=-1)]: Done   24 tasks      | elapsed:    18.5s
[Parallel(n_jobs=-1)]: Done   33 tasks      | elapsed:    20.6s
[Parallel(n_jobs=-1)]: Done   42 tasks      | elapsed:    22.7s
[Parallel(n_jobs=-1)]: Done   53 tasks      | elapsed:    25.6s
[Parallel(n_jobs=-1)]: Done   64 tasks      | elapsed:    28.6s
[Parallel(n_jobs=-1)]: Done   77 tasks      | elapsed:    33.0s
[Parallel(n_jobs=-1)]: Done   90 tasks      | elapsed:    38.0s
[Parallel(n_jobs=-1)]: Done  105 tasks      | elapsed:    40.8s
[Parallel(n_jobs=-1)]: Done  120 tasks      | elapsed:    44.6s
[Parallel(n_jobs=-1)]: Done  137 tasks      | elapsed:    50.4s
[Parallel(n_jobs=-1)]: Done  154 tasks      | elapsed:    57.7s
[Parallel(n_jobs=-1)]: Done  173 tasks      | elapsed:   1.1min
[Parallel(n_jobs=-1)]: Done  192 tasks      | elapsed:   1.2min
[Parallel(n_jobs=-1)]: Done  213 tasks      | elapsed:   1.4min
[Parallel(n_jobs=-1)]: Done  234 tasks      | elapsed:   1.5min
[Parallel(n_jobs=-1)]: Done  257 tasks      | elapsed:   1.8min
[Parallel(n_jobs=-1)]: Done  280 tasks      | elapsed:   2.0min
[Parallel(n_jobs=-1)]: Done  305 tasks      | elapsed:   2.2min
[Parallel(n_jobs=-1)]: Done  330 tasks      | elapsed:   2.5min
[Parallel(n_jobs=-1)]: Done  357 tasks      | elapsed:   3.0min
[Parallel(n_jobs=-1)]: Done  384 tasks      | elapsed:   3.0min
[Parallel(n_jobs=-1)]: Done  413 tasks      | elapsed:   3.1min
[Parallel(n_jobs=-1)]: Done  442 tasks      | elapsed:   3.3min
[Parallel(n_jobs=-1)]: Done  473 tasks      | elapsed:   3.4min
[Parallel(n_jobs=-1)]: Done  504 tasks      | elapsed:   3.5min
[Parallel(n_jobs=-1)]: Done  537 tasks      | elapsed:   3.8min
[Parallel(n_jobs=-1)]: Done  570 tasks      | elapsed:   4.0min
[Parallel(n_jobs=-1)]: Done  605 tasks      | elapsed:   4.3min
[Parallel(n_jobs=-1)]: Done  640 tasks      | elapsed:   4.7min
[Parallel(n_jobs=-1)]: Done  677 tasks      | elapsed:   5.0min
[Parallel(n_jobs=-1)]: Done  714 tasks      | elapsed:   5.5min
[Parallel(n_jobs=-1)]: Done  753 tasks      | elapsed:   5.6min
[Parallel(n_jobs=-1)]: Done  792 tasks      | elapsed:   5.8min
[Parallel(n_jobs=-1)]: Done  833 tasks      | elapsed:   5.9min
[Parallel(n_jobs=-1)]: Done  874 tasks      | elapsed:   6.2min
[Parallel(n_jobs=-1)]: Done  917 tasks      | elapsed:   6.4min
[Parallel(n_jobs=-1)]: Done  960 tasks      | elapsed:   6.7min
[Parallel(n_jobs=-1)]: Done 1005 tasks      | elapsed:   7.2min
[Parallel(n_jobs=-1)]: Done 1050 tasks      | elapsed:   7.6min
[Parallel(n_jobs=-1)]: Done 1080 out of 1080 | elapsed:   8.0min finished
```

The Best Parameters are:

```
{'learning_rate': 0.2, 'max_depth': 3, 'n_estimators': 90}
```

The test log loss is: 6.61476805056



The value of the log\_loss is more as the model training is done on the sample data.

In [ ]:

1