

# Quora question pair similarity

## 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. </p>

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 Koh inoor (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 [0]: 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 os
import gc

import re
from nltk.corpus import stopwords
!pip install -q distance
# This package is used for finding longest common subsequence between two strings
# you can write your own dp code for this
import distance
from nltk.stem import PorterStemmer
from bs4 import BeautifulSoup
import warnings
warnings.filterwarnings("ignore")
from nltk.corpus import stopwords
!pip install -q fuzzywuzzy
from fuzzywuzzy import fuzz
from sklearn.manifold import TSNE
# Import the Required Lib packages for WORD-Cloud generation
# https://stackoverflow.com/questions/45625434/how-to-install-wordcloud-in-python3-6
from wordcloud import WordCloud, STOPWORDS
from os import path
from PIL import Image
import nltk
nltk.download('stopwords')

from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
import sys
from tqdm import tqdm
# extract word2vec vectors
# https://github.com/explosion/spaCy/issues/1721
# http://landinghub.visualstudio.com/visual-cpp-build-tools

import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import sqlite3
from sqlalchemy import create_engine # database connection
import csv
import os
warnings.filterwarnings("ignore")
import datetime as dt
import numpy as np
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.linear_model import SGDClassifier
from mlxtend.classifier import StackingClassifier
from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import precision_recall_curve, auc, roc_curve
import spacy
from sklearn.model_selection import StratifiedKFold
from sklearn import preprocessing
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

## 3.1 Reading data and basic stats

```
In [0]: # Load the Drive helper and mount
from google.colab import drive
# This will prompt for authorization.
drive.mount('/content/drive')
```

Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\\_uri=urn%3Aietf%3Awww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response\\_type=code](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3Aietf%3Awww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fdrive.photos.readonly%20https%3A%2F%2Fwww.googleapis.com%2Fauth%2Fpeopleapi.readonly&response_type=code)

Enter your authorization code:  
.....  
Mounted at /content/drive

```
In [0]: #df = pd.read_csv("train.csv")

df=pd.read_csv('/content/drive/My Drive/Colab Notebooks/Quora/train.csv')

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

Number of data points: 404290

```
In [0]: df.head()
```

Out[0]:

	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 [0]: 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:

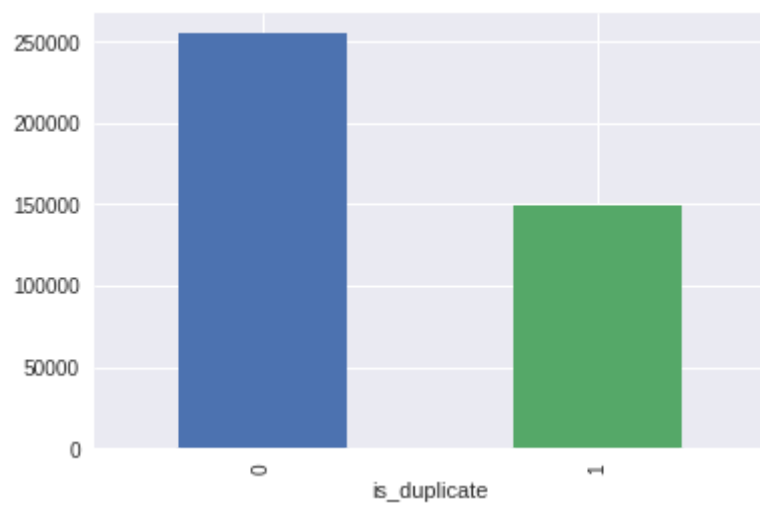
- id: Looks like a simple rowID
- qid{1, 2}: The unique ID of each question in the pair
- question{1, 2}: The actual textual contents of the questions.
- is\_duplicate: The label that we are trying to predict - whether the two questions are duplicates of each other.

## 3.2 Distribution of data points among output classes

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

```
In [0]: df.groupby("is_duplicate")['id'].count().plot.bar()
```

```
Out[0]: <matplotlib.axes._subplots.AxesSubplot at 0x7faa6f5f7828>
```



```
In [0]: print('~> Total number of question pairs for training:\n    {}'.format(len(df)))
```

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

```
In [0]: 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 [0]: 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: {}'.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: {}'.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 [0]: 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 [0]: #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 [0]: plt.figure(figsize=(20, 10))

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

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

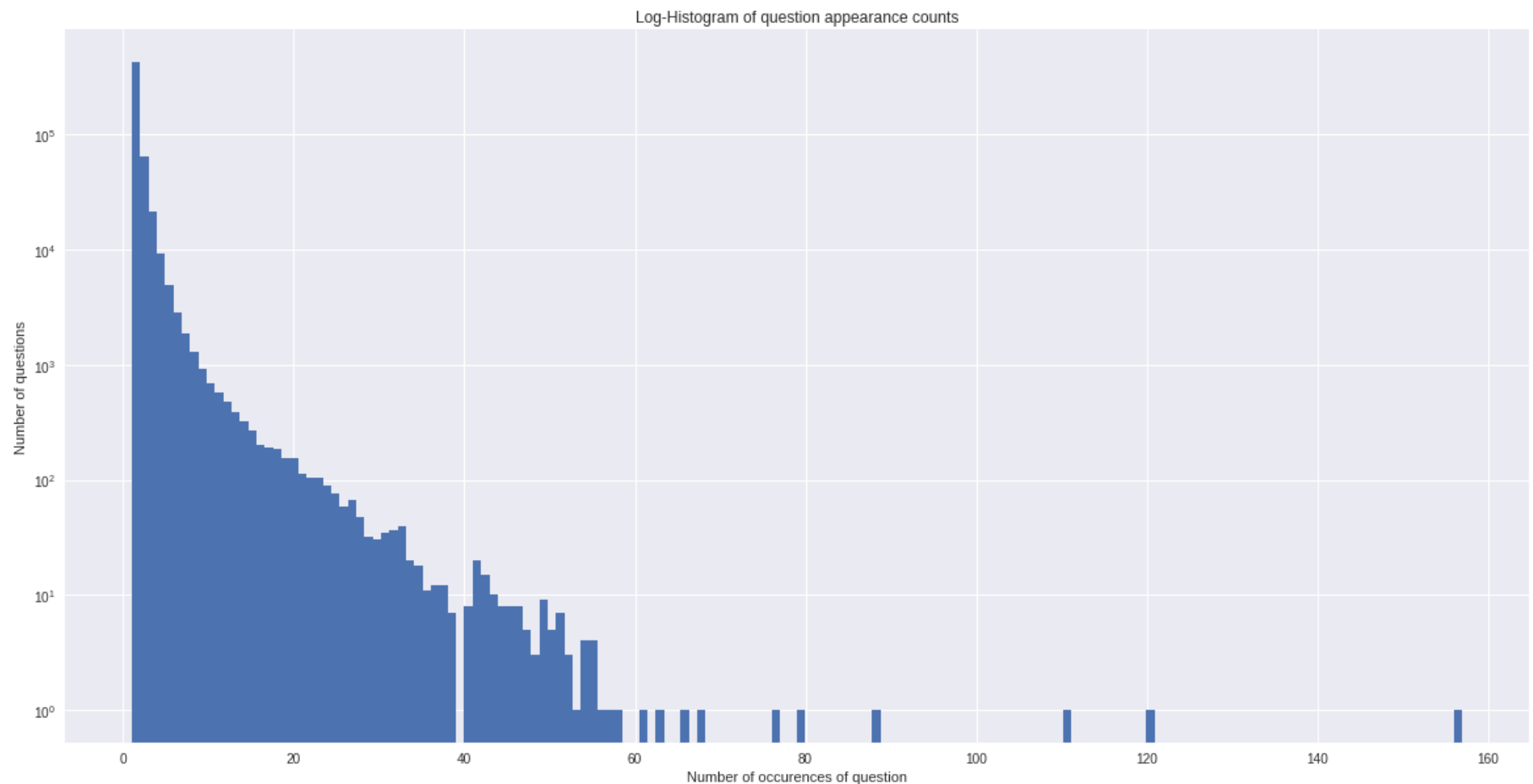
plt.title('Log-Histogram of question appearance counts')

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 [0]: #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	NaN	My Chinese name is Haichao Yu. What English na...	0

- There are two rows with null values in question2

```
In [0]: # 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:(Baisc features)

- **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 [0]: if os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora/df_fe_without_preprocessing_train.csv'):
        df = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora/df_fe_without_preprocessing_train.csv", encoding='latin-1')
        print("from saved files")
    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)

        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))

        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_Common'] = df.apply(normalized_word_Common, axis=1)#aapling it
        df['word_Total'] = df.apply(normalized_word_Total, axis=1)#aapling it
        df['word_share'] = df.apply(normalized_word_share, axis=1)#aapling it
        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)
        print("creating files")

df.head()

```

from saved files

Out[0]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common
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? How can I solve...	Find the remainder when $23^{24}$ $[/math] i...$	0	1	1	50	65	11	9	0.0
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

### 3.3.1 Analysis of some of the extracted features



- Here are some questions have only one single words.

```
In [0]: 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])

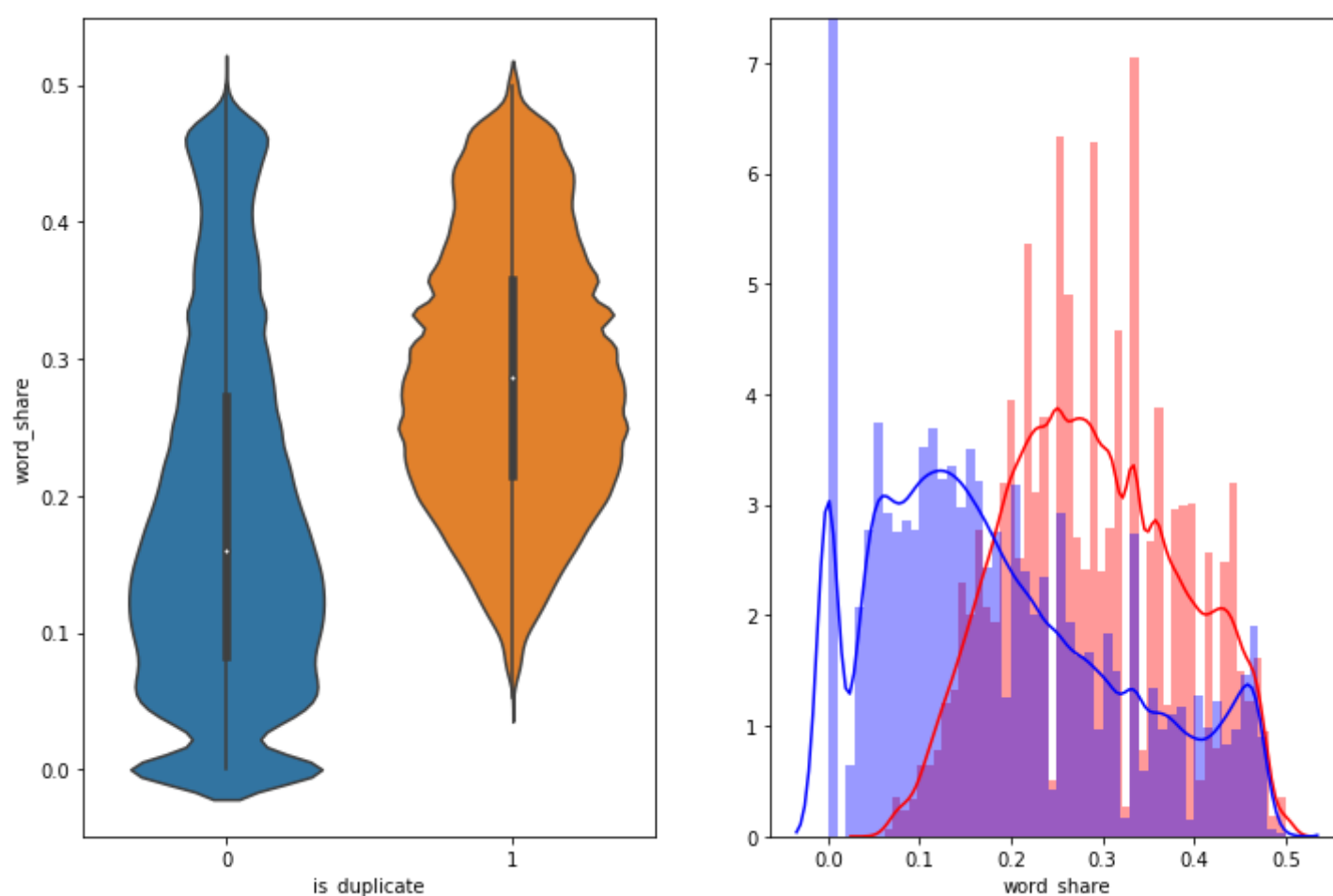
Minimum length of the questions in question1 : 1
Minimum length of the questions in question2 : 1
Number of Questions with minimum length [question1] : 67
Number of Questions with minimum length [question2] : 24
```

### 3.3.1.1 Feature: word\_share

```
In [0]: 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()
```



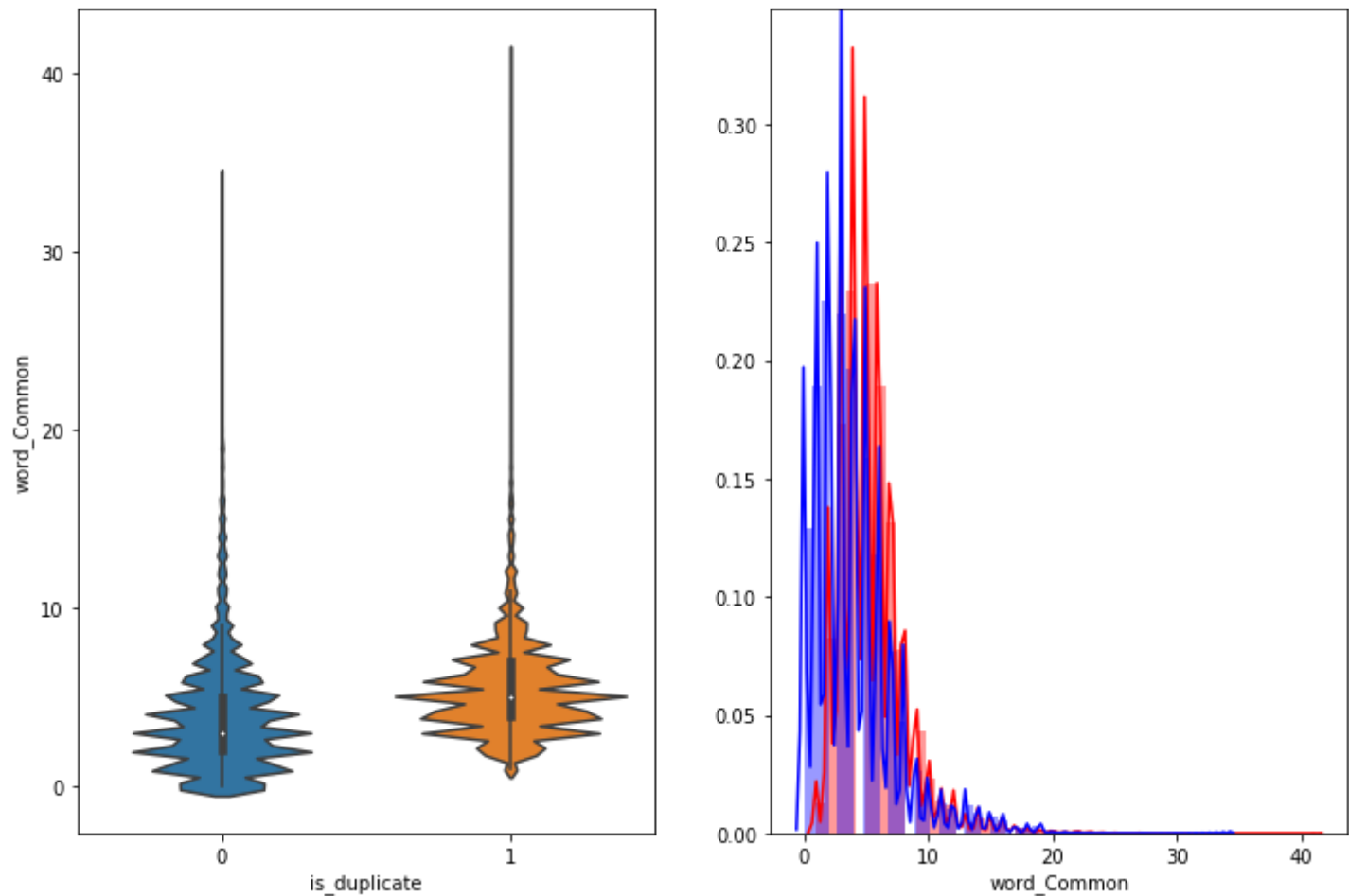
- The distributions for normalized word\_share have some overlap on the far right-hand side, i.e., there are quite a lot of questions with high word similarity
- The average word share and Common no. of words of qid1 and qid2 is more when they are duplicate(Similar)

### 3.3.1.2 Feature: word\_Common

```
In [0]: 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()
```



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

### 3.4 Preprocessing of Text

```
In [0]: #https://stackoverflow.com/questions/12468179/unicodedecodeerror-utf8-codec-cant-decode-byte-0x9c
if os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora/df_fe_without_preprocessing_train.csv'):
    df = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora/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 notebook again!")
```

```
In [0]: df.head(2)
```

Out[0]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	w
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	20
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	20

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

```
In [0]: # To get the results in 4 decemal 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't", "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  
$$\text{cwc\_min} = \text{common\_word\_count} / (\min(\text{len}(\text{q1\_words}), \text{len}(\text{q2\_words})))$$
- **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/> (<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  
$$\text{longest\_substr\_ratio} = \text{len}(\text{longest common substring}) / (\min(\text{len}(\text{q1\_tokens}), \text{len}(\text{q2\_tokens})))$$

```

In [0]: 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 [0]: if os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora/nlp_features_train.csv'):
df = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora/nlp_features_train.csv",encoding='latin-1')
df.fillna('')
print("if")
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)
print("else")
df.head(2)

```

if

Out[0]:

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	first_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	0.999980	0.833319	0.999983	0.999983	...	0.785709	0.0	1.0
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	1.0

2 rows × 21 columns



### 3.5.1 Analysis of extracted features

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

```

In [0]: 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))

#Saving the np array into a text file
np.savetxt('/content/drive/My Drive/Colab Notebooks/Quora/train_p.txt', p, delimiter=' ', fmt='%s')
np.savetxt('/content/drive/My Drive/Colab Notebooks/Quora/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 [0]: # 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 : 16109886  
Total number of words in non duplicate pair questions : 33193130

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

```
In [0]: 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()
```

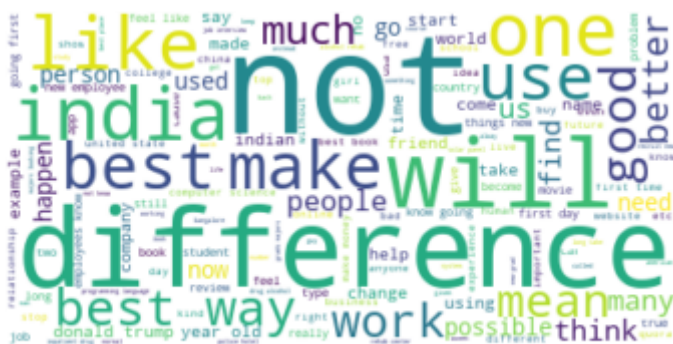
### Word Cloud for Duplicate Question pairs



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

```
In [0]: 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()
```

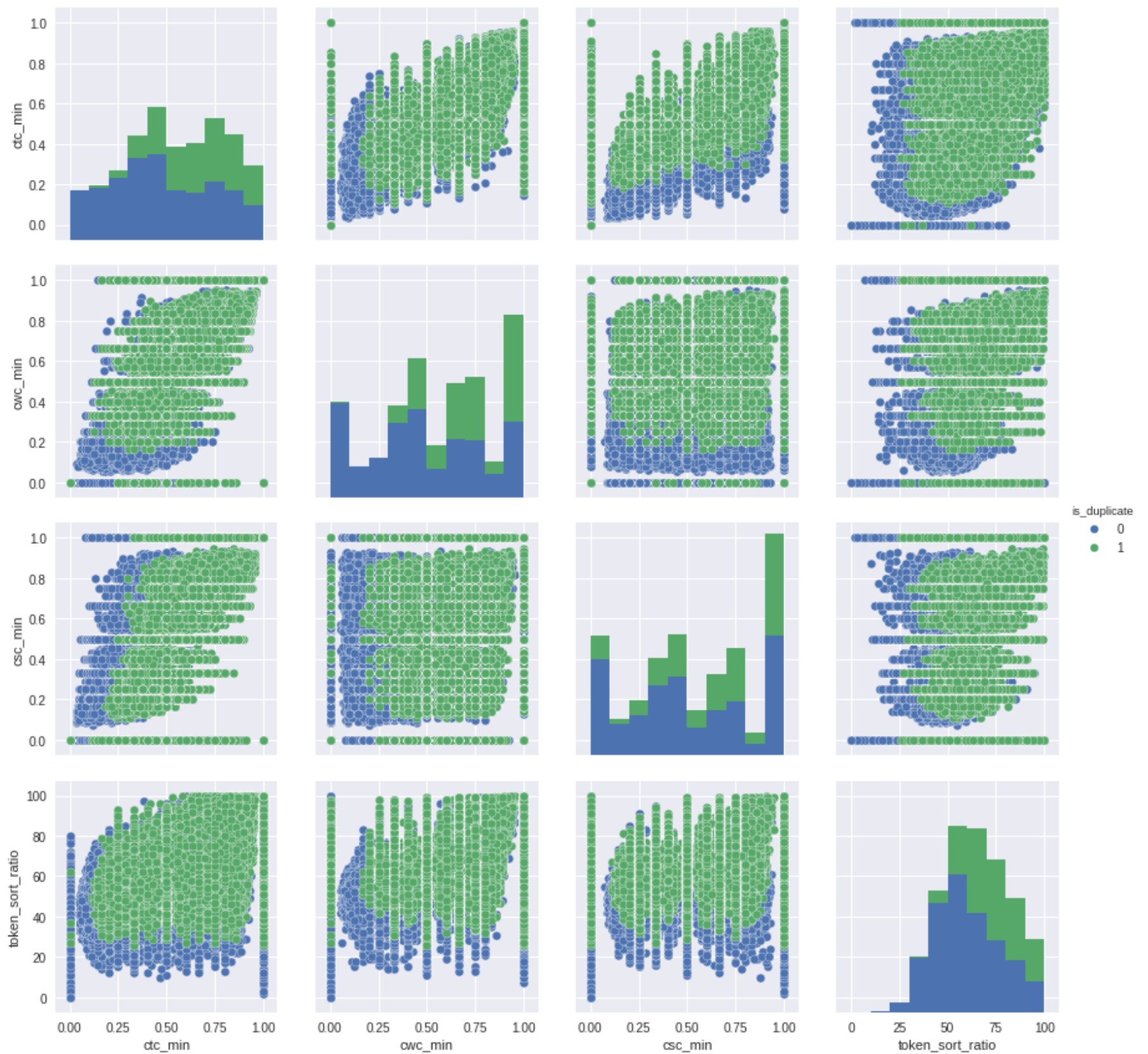
Word Cloud for non-Duplicate Question pairs:



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



```
In [0]: #pair plots
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'],size=3)
plt.show()
```

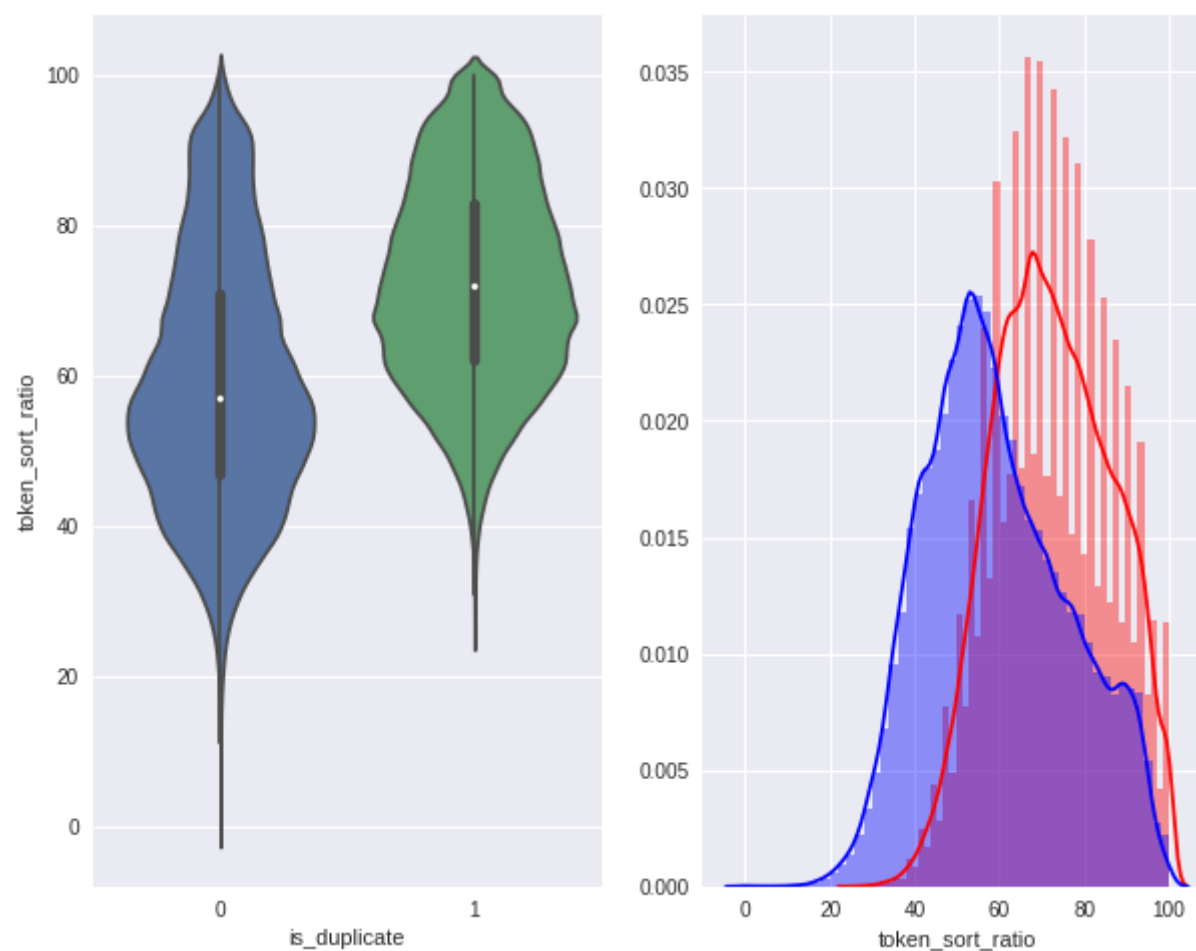




```
In [0]: # Distribution of the token_sort_ratio
plt.figure(figsize=(10, 8))

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

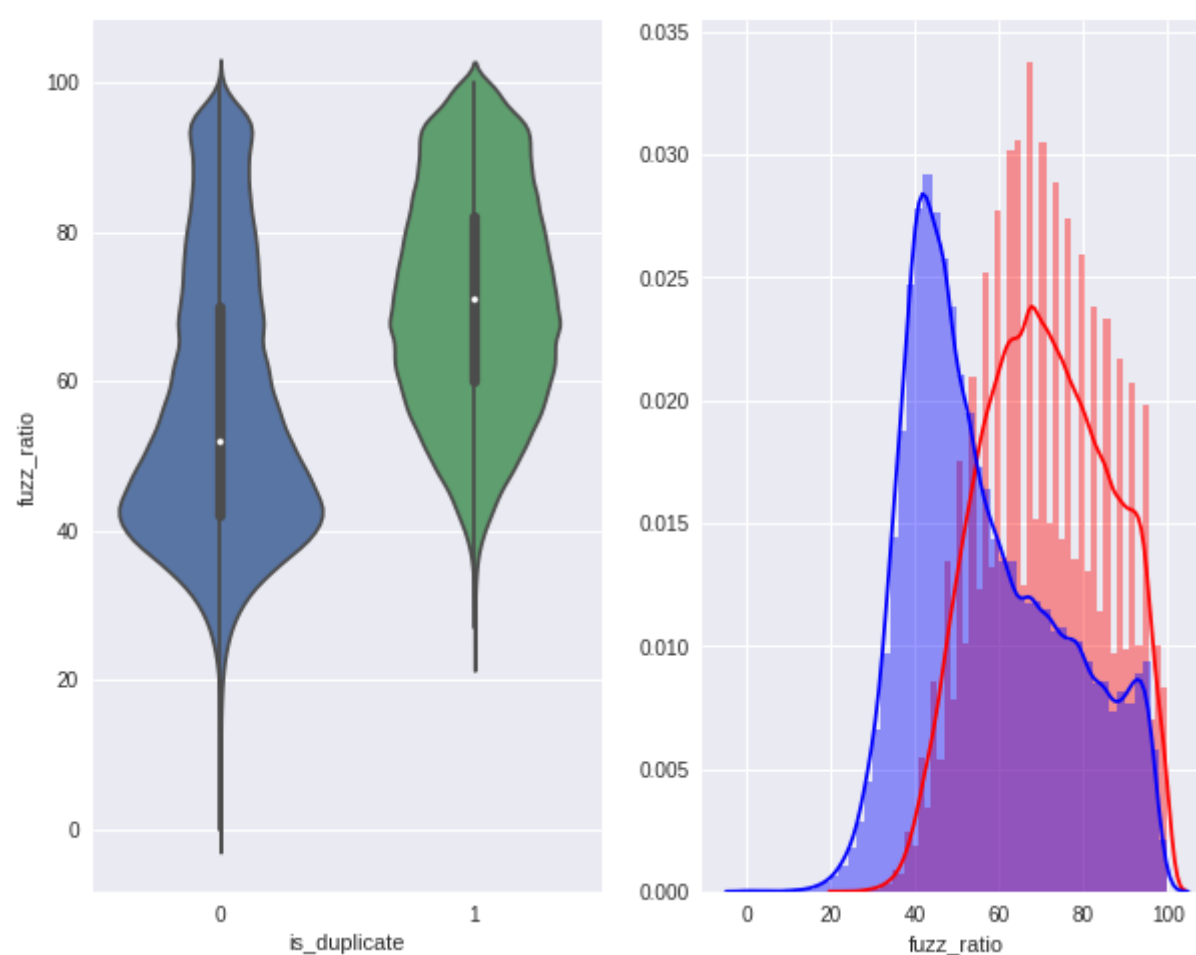
plt.subplot(1,2,2)
sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = "1", color = 'red')
sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = "0" , color = 'blue' )
plt.show()
```



```
In [0]: 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 [0]: # Using TSNE for Dimentionality reduction for 15 Features(Generated after cleaning the data) to 2 dimention

from sklearn.preprocessing import MinMaxScaler

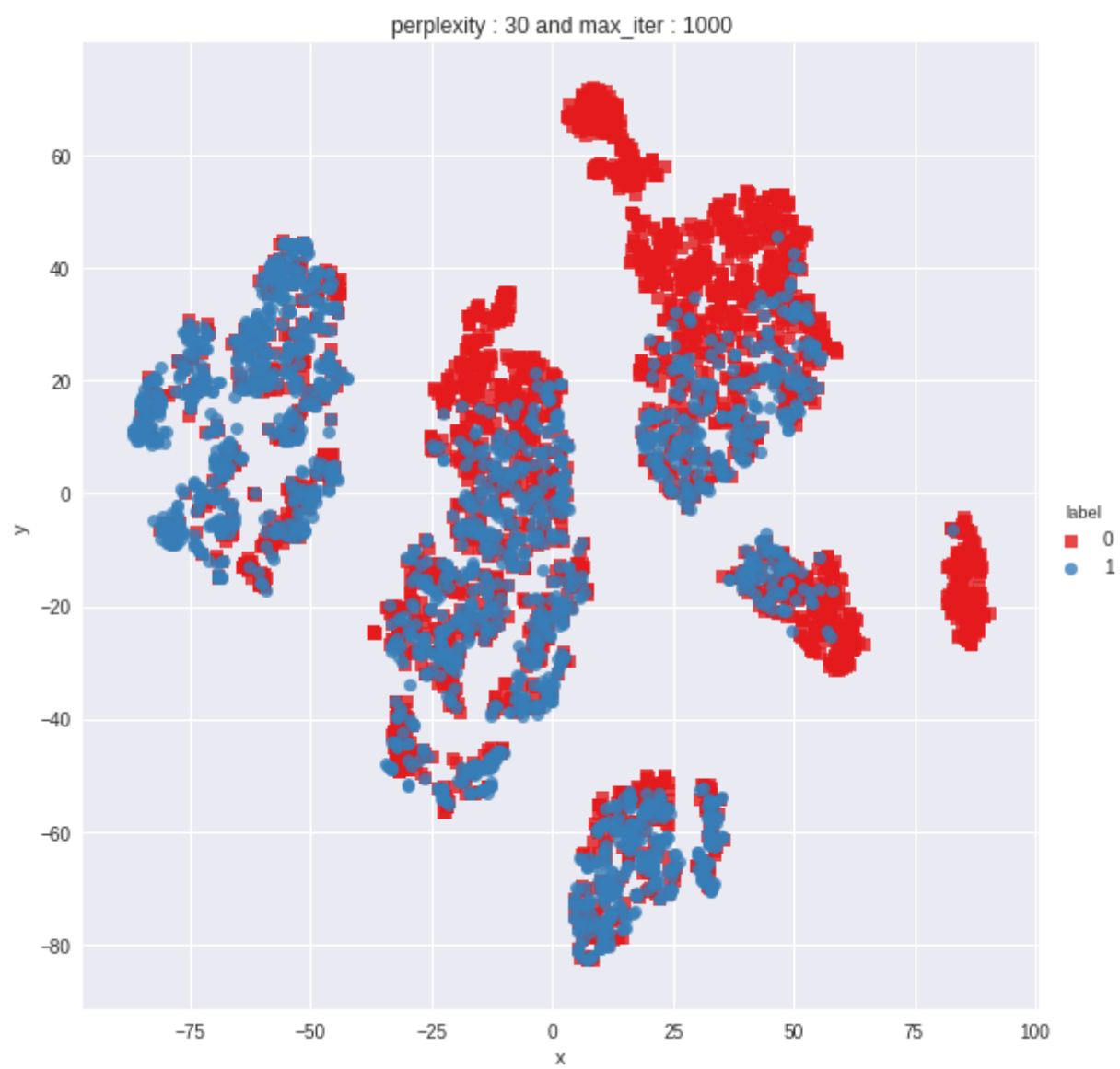
dfp_subsampled = df[0:5000]
X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max',
'last_word_eq', 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_ratio', 'token_sort_ratio', 'fuzz_ratio',
'fuzz_partial_ratio', 'longest_substr_ratio']])
y = dfp_subsampled['is_duplicate'].values
```

```
In [0]: 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 5000 samples in 0.016s...
[t-SNE] Computed neighbors for 5000 samples in 0.369s...
[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 0.316s
[t-SNE] Iteration 50: error = 80.9162369, gradient norm = 0.0427600 (50 iterations in 2.716s)
[t-SNE] Iteration 100: error = 70.3915100, gradient norm = 0.0108003 (50 iterations in 2.060s)
[t-SNE] Iteration 150: error = 68.6126938, gradient norm = 0.0054721 (50 iterations in 2.071s)
[t-SNE] Iteration 200: error = 67.7680206, gradient norm = 0.0042246 (50 iterations in 2.156s)
[t-SNE] Iteration 250: error = 67.2733459, gradient norm = 0.0037275 (50 iterations in 2.125s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.273346
[t-SNE] Iteration 300: error = 1.7734827, gradient norm = 0.0011933 (50 iterations in 2.101s)
[t-SNE] Iteration 350: error = 1.3717980, gradient norm = 0.0004826 (50 iterations in 2.056s)
[t-SNE] Iteration 400: error = 1.2037998, gradient norm = 0.0002772 (50 iterations in 2.113s)
[t-SNE] Iteration 450: error = 1.1133003, gradient norm = 0.0001877 (50 iterations in 2.102s)
[t-SNE] Iteration 500: error = 1.0579894, gradient norm = 0.0001429 (50 iterations in 2.094s)
[t-SNE] Iteration 550: error = 1.0220573, gradient norm = 0.0001178 (50 iterations in 2.102s)
[t-SNE] Iteration 600: error = 0.9990303, gradient norm = 0.0001036 (50 iterations in 2.107s)
[t-SNE] Iteration 650: error = 0.9836842, gradient norm = 0.0000951 (50 iterations in 2.115s)
[t-SNE] Iteration 700: error = 0.9732341, gradient norm = 0.0000860 (50 iterations in 2.135s)
[t-SNE] Iteration 750: error = 0.9649901, gradient norm = 0.0000789 (50 iterations in 2.153s)
[t-SNE] Iteration 800: error = 0.9582695, gradient norm = 0.0000745 (50 iterations in 2.142s)
[t-SNE] Iteration 850: error = 0.9525222, gradient norm = 0.0000732 (50 iterations in 2.140s)
[t-SNE] Iteration 900: error = 0.9479918, gradient norm = 0.0000689 (50 iterations in 2.148s)
[t-SNE] Iteration 950: error = 0.9442031, gradient norm = 0.0000651 (50 iterations in 2.154s)
[t-SNE] Iteration 1000: error = 0.9408465, gradient norm = 0.0000590 (50 iterations in 2.168s)
[t-SNE] KL divergence after 1000 iterations: 0.940847
```

```
In [0]: %%time
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()
```



CPU times: user 890 ms, sys: 339 ms, total: 1.23 s  
Wall time: 810 ms

```
In [0]: from sklearn.manifold import TSNE
```

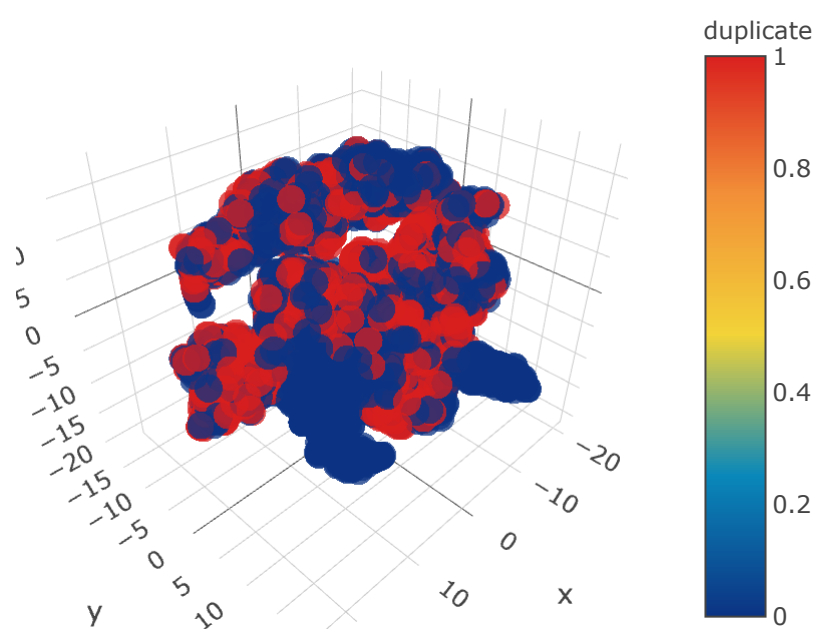
```
tsne3d = TSNE(  
    n_components=3,  
    init='random', # pca  
    random_state=101,  
    method='barnes_hut',  
    n_iter=1000,  
    verbose=2,  
    angle=0.5  
)  
.fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...  
[t-SNE] Indexed 5000 samples in 0.014s...  
[t-SNE] Computed neighbors for 5000 samples in 0.372s...  
[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 0.304s  
[t-SNE] Iteration 50: error = 80.3552017, gradient norm = 0.0329941 (50 iterations in 12.894s)  
[t-SNE] Iteration 100: error = 69.1127167, gradient norm = 0.0036756 (50 iterations in 6.643s)  
[t-SNE] Iteration 150: error = 67.6178818, gradient norm = 0.0017629 (50 iterations in 5.956s)  
[t-SNE] Iteration 200: error = 67.0571747, gradient norm = 0.0011826 (50 iterations in 5.919s)  
[t-SNE] Iteration 250: error = 66.7298050, gradient norm = 0.0008528 (50 iterations in 5.845s)  
[t-SNE] KL divergence after 250 iterations with early exaggeration: 66.729805  
[t-SNE] Iteration 300: error = 1.4963876, gradient norm = 0.0006857 (50 iterations in 7.947s)  
[t-SNE] Iteration 350: error = 1.1549060, gradient norm = 0.0001911 (50 iterations in 10.831s)  
[t-SNE] Iteration 400: error = 1.0083323, gradient norm = 0.0000968 (50 iterations in 10.974s)  
[t-SNE] Iteration 450: error = 0.9356370, gradient norm = 0.0000660 (50 iterations in 10.569s)  
[t-SNE] Iteration 500: error = 0.8982100, gradient norm = 0.0000521 (50 iterations in 10.385s)  
[t-SNE] Iteration 550: error = 0.8778670, gradient norm = 0.0000595 (50 iterations in 10.281s)  
[t-SNE] Iteration 600: error = 0.8642665, gradient norm = 0.0000579 (50 iterations in 10.392s)  
[t-SNE] Iteration 650: error = 0.8558875, gradient norm = 0.0000362 (50 iterations in 10.360s)  
[t-SNE] Iteration 700: error = 0.8492573, gradient norm = 0.0000305 (50 iterations in 10.394s)  
[t-SNE] Iteration 750: error = 0.8432317, gradient norm = 0.0000276 (50 iterations in 10.364s)  
[t-SNE] Iteration 800: error = 0.8378869, gradient norm = 0.0000279 (50 iterations in 10.337s)  
[t-SNE] Iteration 850: error = 0.8331724, gradient norm = 0.0000261 (50 iterations in 10.167s)  
[t-SNE] Iteration 900: error = 0.8291837, gradient norm = 0.0000259 (50 iterations in 10.166s)  
[t-SNE] Iteration 950: error = 0.8255505, gradient norm = 0.0000333 (50 iterations in 9.988s)  
[t-SNE] Iteration 1000: error = 0.8224180, gradient norm = 0.0000235 (50 iterations in 10.077s)  
[t-SNE] KL divergence after 1000 iterations: 0.822418
```

```
In [0]: trace1 = 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=[trace1]
layout=dict(height=500, width=500, title='3d embedding with engineered features')
fig=dict(data=data, layout=layout)
py.iplot(fig, filename='3DBubble')
```

3d embedding with engineered features



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

```
In [0]: # avoid decoding problems
df = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/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 -----
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

```
In [0]: questions = list(df['question1']) + list(df['question2'])
questions=questions[:50000]#taking 1m data points
```

```
In [0]: df=df.head(50000)#taking 1m data points
print(len(df['is_duplicate'].values))

50000
```

```
In [0]: #we have splited data into 70 :30 ratio
X_train,X_test, y_train, y_test = train_test_split(questions, df['is_duplicate'].values,test_size=0.30,shuffle=False)
```

```
In [0]: 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(X_train)

# dict key:word and value:tf-idf score
word2tfidf = dict(zip(tfidf.get_feature_names(), tfidf.idf_))
```

- After we find TF-IDF scores, we convert each question to a weighted average of word2vec vectors by these scores.
- here we use a pre-trained GLOVE model which comes free with "Spacy". <https://spacy.io/usage/vectors-similarity> (<https://spacy.io/usage/vectors-similarity>)
- It is trained on Wikipedia and therefore, it is stronger in terms of word semantics.

```
In [0]: %%time
# en_vectors_web_lg, which includes over 1 million unique vectors.
nlp = spacy.load('en_core_web_sm')

vecs1 = []
# https://github.com/noamraph/tqdm
# tqdm is used to print the progress bar
for qu1 in tqdm(list(df['question1'])):
    doc1 = nlp(qu1)
    # 384 is the number of dimensions of vectors
    mean_vec1 = np.zeros([len(doc1), 384])
    for word1 in doc1:
        # word2vec
        vec1 = word1.vector
        # fetch df score
        try:
            idf = word2tfidf[str(word1)]
        except:
            idf = 0
        # compute final vec
        mean_vec1 += vec1 * idf
    mean_vec1 = mean_vec1.mean(axis=0)
    vecs1.append(mean_vec1)
df['q1_feats_m'] = list(vecs1)

100%|██████████| 50000/50000 [13:54<00:00, 59.90it/s]

CPU times: user 18min 33s, sys: 8min 22s, total: 26min 56s
Wall time: 13min 57s
```

```
In [0]: vecs2 = []
for qu2 in tqdm(list(df['question2'])):
    doc2 = nlp(qu2)
    mean_vec2 = np.zeros([len(doc2), 384])
    for word2 in doc2:
        # word2vec
        vec2 = word2.vector
        # fetch df score
        try:
            idf = word2tfidf[str(word2)]
        except:
            #print word
            idf = 0
        # compute final vec
        mean_vec2 += vec2 * idf
    mean_vec2 = mean_vec2.mean(axis=0)
    vecs2.append(mean_vec2)
df['q2_feats_m'] = list(vecs2)
```

```
100%|██████████| 50000/50000 [14:03<00:00, 59.24it/s]
```

```
In [0]: #prepro_features_train.csv (Simple Preprocessing Feartures)
#nlp_features_train.csv (NLP Features)
if os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/nlp_features_train.csv'):
    dfnlp = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/nlp_features_train.csv",encoding='latin-1')
else:
    print("download nlp_features_train.csv from drive or run previous notebook")

if os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/df_fe_without_preprocessing_train.csv'):
    dfppro = pd.read_csv("/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/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 [0]: dfnlp=dfnlp.head(50000)
```

```
In [0]: dfppro=dfppro.head(50000)
```

```
In [0]: 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 [0]: # 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')#i just need to take this result ...do not need to store in file..just
perform ur task in one go
    result.to_csv('/content/drive/My Drive/Colab Notebooks/Quora/final_features.csv')
```

```
In [0]: result.columns
```

```
Out[0]: Index(['id', 'is_duplicate', 'cwc_min', 'cwc_max', 'csc_min', 'csc_max',
              'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq',
              ...,
              '374_y', '375_y', '376_y', '377_y', '378_y', '379_y', '380_y', '381_y',
              '382_y', '383_y'],
              dtype='object', length=796)
```

## 4. Machine Learning Models

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

```

In [0]: #Creating db file from csv
if not os.path.isfile('/content/drive/My Drive/Colab Notebooks/Quora/train.db'):
    disk_engine = create_engine('sqlite:///content/drive/My Drive/Colab Notebooks/Quora/train.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('/content/drive/My Drive/Colab Notebooks/Quora/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', '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', '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', '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', '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', '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', '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'], 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
else:
    print("else")

```

```

else
CPU times: user 459 µs, sys: 172 µs, total: 631 µs
Wall time: 1.57 ms

```



```
In [0]: #http://www.sqlitetutorial.net/sqlite-python/create-tables/
def create_connection(db_file):
    """ create a database connection to the SQLite database
        specified by db_file
    :param db_file: database file
    :return: Connection object or None
    """
    try:
        conn = sqlite3.connect(db_file)
        return conn

    except Error as e:
        print(e)

    return None
```

```
def 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 [0]: %%time
read_db = '/content/drive/My Drive/Colab Notebooks/Quora/train.db'
conn_r = create_connection(read_db)
checkTableExists(conn_r)
conn_r.close()
```

Tables in the database:  
data  
CPU times: user 3.31 ms, sys: 2.39 ms, total: 5.7 ms  
Wall time: 434 ms

```
In [0]: # 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 1000;", conn_r)
        conn_r.commit()
        conn_r.close()
```

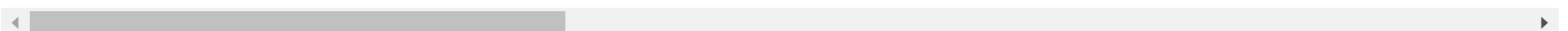
```
In [0]: # remove the first row
data.drop(data.index[0], inplace=True)
y_true = data['is_duplicate']
data.drop(['Unnamed: 0', 'id', 'index', 'is_duplicate'], axis=1, inplace=True)
```

```
In [0]: data.head()
```

```
Out[0]:
```

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_work
1	0.374995312558593	0.29999700003	0.249993750156246	0.166663888935184	0.333330555578704	0.22222098766118	0.0
2	0.499975001249937	0.249993750156246	0.999966667777741	0.599988000239995	0.571420408279882	0.571420408279882	1.0
3	0.249993750156246	0.14285510206997	0.66664444518516	0.66664444518516	0.428565306209911	0.29999700003	0.0
4	0.999966667777741	0.333329629670781	0.0	0.0	0.428565306209911	0.166665740745885	0.0
5	0.833319444675922	0.833319444675922	0.799984000319994	0.799984000319994	0.818174380232907	0.818174380232907	1.0

5 rows × 794 columns



## 4.2 Converting strings to numerics

```
In [0]: %%time
# 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(result.columns)
#j=0
for i in cols:
    result[i] = result[i].apply(pd.to_numeric)
    # j=j+1
    #if j<=5:
    #print(i)
```

```
In [0]: y_true = result['is_duplicate']
```

```
In [0]: # https://stackoverflow.com/questions/7368789/convert-all-strings-in-a-list-to-int
y_true = list(map(int, y_true))
```

```
In [0]: print(len(y_true))

50000
```

### 4.3 Random train test split( 70:30)

```
In [0]: print(result.shape)
#y_true = result['is_duplicate']

(50000, 796)
```

```
In [0]: result.head()
```

Out[0]:

	id	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	...	374_y	375_y
0	0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	...	16.158090	32.718838
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	...	-2.299923	-3.723903
2	2	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	...	12.546578	1.177362
3	3	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	...	3.131514	3.592912
4	4	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	...	-2.500425	11.748801

5 rows × 796 columns

```
In [0]: X_train,X_test, y_train, y_test = train_test_split(result, y_true, stratify=y_true, test_size=0.3)#use result
```

```
In [0]: 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 : (35000, 796)
Number of data points in test data : (15000, 796)
```

```
In [0]: print("-"*10, "Distribution of output variable in train data", "-"*10)
train_distr = Counter(y_train)
train_len = len(y_train)
print("Class 0: ",int(train_distr[0])/train_len,"Class 1: ", int(train_distr[1])/train_len)
```

```
print("-"*10, "Distribution of output variable in test data", "-"*10)
test_distr = Counter(y_test)
test_len = len(y_test)
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:  0.6270285714285714 Class 1:  0.37297142857142856
----- Distribution of output variable in test data -----
Class 0:  0.627 Class 1:  0.373
```

```

In [0]: # This function plots the confusion matrices given y_i, y_i_hat.
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, 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, 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, fmt=".3f", xticklabels=labels, yticklabels=labels)
    plt.xlabel('Predicted Class')
    plt.ylabel('Original Class')
    plt.title("Recall matrix")

    plt.show()

```

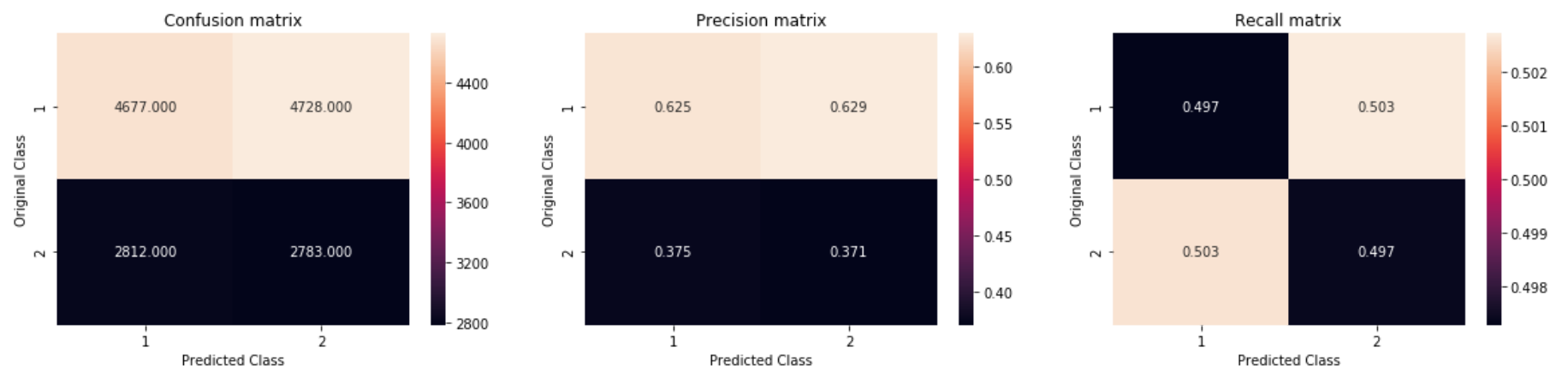
#### 4.4 Building a random model (Finding worst-case log-loss)

```
In [0]: # 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
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))#shows error

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

Log loss on Test Data using Random Model 0.8844677880953877



## 4.5 Logistic Regression with hyperparameter tuning

```

In [0]: 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='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))

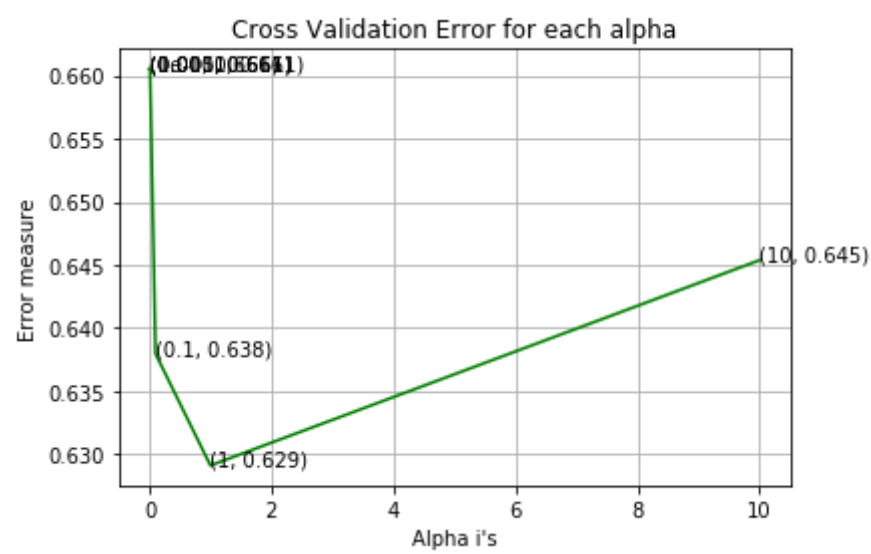
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

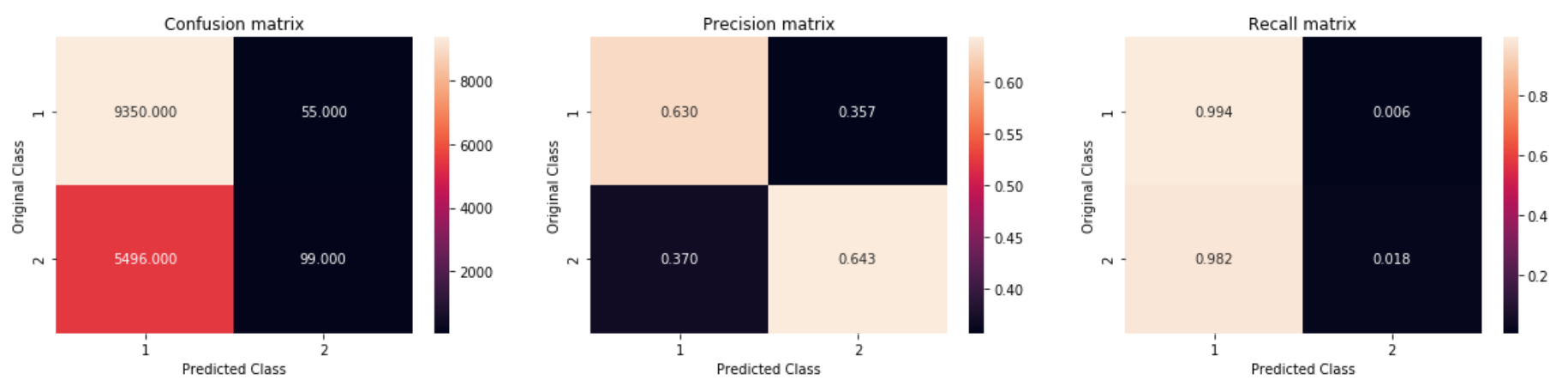
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, 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: 0.6605330475770012  
 For values of alpha = 0.0001 The log loss is: 0.6605330475770012  
 For values of alpha = 0.001 The log loss is: 0.6605330475770012  
 For values of alpha = 0.01 The log loss is: 0.6605330475770012  
 For values of alpha = 0.1 The log loss is: 0.6379625001226136  
 For values of alpha = 1 The log loss is: 0.6291150675941164  
 For values of alpha = 10 The log loss is: 0.6453690663497577



For values of best alpha = 1 The train log loss is: 0.6282913651481622  
 For values of best alpha = 1 The test log loss is: 0.6291150675941164  
 Total number of data points : 15000



## 4.6 Linear SVM with hyperparameter tuning

```

In [0]: 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(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))

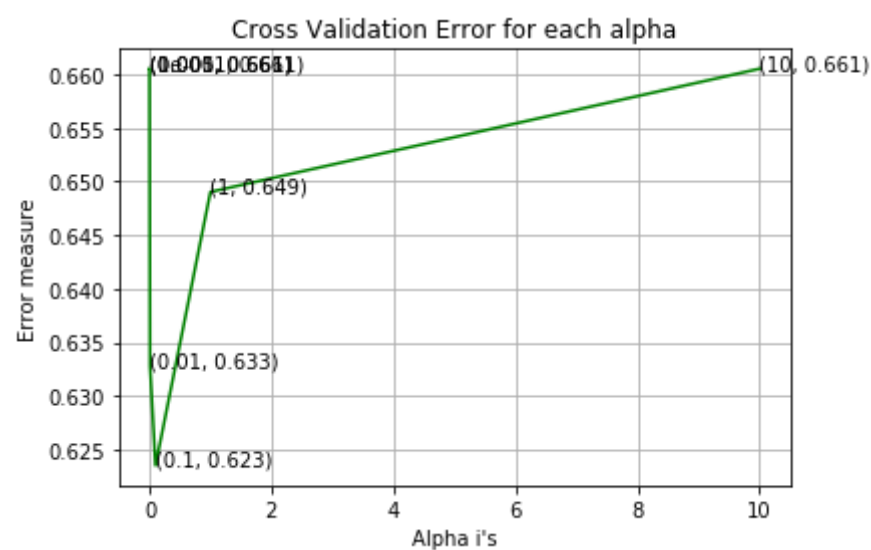
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

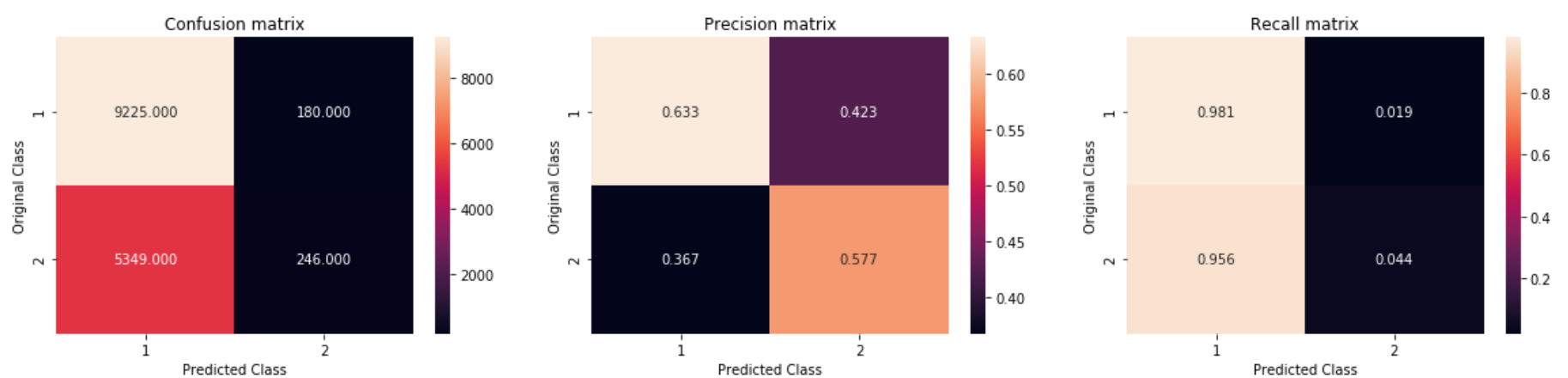
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss 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: 0.6605330475770012  
 For values of alpha = 0.0001 The log loss is: 0.6605330475770012  
 For values of alpha = 0.001 The log loss is: 0.6605330475770012  
 For values of alpha = 0.01 The log loss is: 0.6327511870707146  
 For values of alpha = 0.1 The log loss is: 0.6234992921121445  
 For values of alpha = 1 The log loss is: 0.6490305987749726  
 For values of alpha = 10 The log loss is: 0.6605312690926457



For values of best alpha = 0.1 The train log loss is: 0.6246086457684077  
 For values of best alpha = 0.1 The test log loss is: 0.6234992921121445  
 Total number of data points : 15000



## 5. Additional Task to Perform

### 5.1 Pre Processing

In [0]: `#loading from nlp file  
nlp_features_train=pd.read_csv('/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/nlp_features_train.csv',encoding='latin-1')`

In [0]: `nlp_features_train.head(2)`

Out[0]:

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	first_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	0.999980	0.833319	0.999983	0.999983	...	0.785709	0.0	1.0
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	1.0

2 rows × 21 columns

In [0]: `#loading form df_fe_without_preprocessing_train file  
df_fe_without_preprocessing_train=pd.read_csv('/content/drive/My Drive/Colab Notebooks/Quora_ravi/Quora/df_fe_without_preprocessing_train.csv',encoding='latin-1')`



```
In [0]: df_fe_without_preprocessing_train.head(2)
```

Out[0]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	w
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	2
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

```
In [0]: #removing basics coloumns and merging both the files as data
```

```
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 [0]: #taking 1m points
data=data.head(200000)
```

```
In [0]: # Label for y class
y_class = data['is_duplicate']
data.drop(['id','is_duplicate'], axis=1, inplace=True)
```

```
In [0]: print(data.shape)

(200000, 28)
```

```
In [0]: #spliting the data into train and test dataset in 70:30 ratio

X_train,X_test, y_train, y_test = train_test_split(data, y_class, shuffle=False, test_size=0.3)
```

```
In [0]: #astype is used for string to float
#standardized the data
tfidf1=TfidfVectorizer()
train_q1 = tfidf1.fit_transform(X_train['question1'].values.astype('U'))
test_q1 =tfidf1.transform(X_test['question1'].values.astype('U'))
```

```
In [0]: #transforming question1 and question2 of test dataset
#standardized the data
tfidf2=TfidfVectorizer()
train_q2 = tfidf2.fit_transform(X_train['question2'].values.astype('U'))
test_q2 =tfidf2.transform(X_test['question2'].values.astype('U'))
```

```
In [0]: #Dropping question1 and question2 and Replacing with tfidf values
X_train.drop(['question1','question2'], axis=1, inplace=True)
X_test.drop(['question1','question2'], axis=1, inplace=True)
```

```
In [0]: #Combining Question1 and Question2
from scipy.sparse import coo_matrix, hstack
train_tfidf = hstack((train_q1,train_q2))#hstack is used for combining two dataFrames
test_tfidf = hstack((test_q1,test_q2))
```

```
In [0]: #combining all basic,advanced and tfidf features
#.tocsr() = Convert this matrix to Compressed Sparse Row format

X_train = hstack((X_train, train_tfidf)).tocsr()
X_test = hstack((X_test, test_tfidf)).tocsr()
```

```
In [0]: 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 : (140000, 81696)
Number of data points in test data : (60000, 81696)
```

**5.2 Logistic Regression with TFIDF and hyperparameter tuning**

In [0]: *#performing hyperparametertuning , modeling , plotting error plot and plotting confusion matix*

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

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l2', loss='log', random_state=42)
    clf.fit(X_train, y_train)#values from TFIDF
    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))

print("----"*44)
print("----"*44)

#error plot
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()

print("----"*44)
print("----"*44)

#bulding model
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, 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))

train_predicted = sig_clf.predict_proba(X_train)[:, 1]
test_predicted = sig_clf.predict_proba(X_test)[:, 1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, train_predicted)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, test_predicted)

# test plotting
plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("Alpha")
plt.plot([0, 1], [0, 1], 'k--')
plt.ylabel("AUC")
plt.title("TEST ERROR PLOTS")
plt.legend(loc='best')
plt.show()

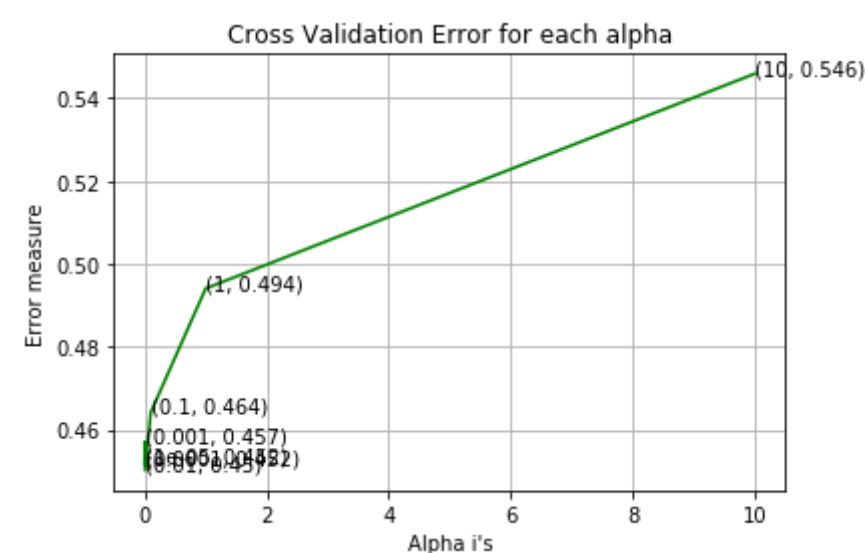
#confusion matrix
plot_confusion_matrix(y_test, predicted_y)
```

For values of alpha = 1e-05 The log loss is: 0.452166522828434  
 For values of alpha = 0.0001 The log loss is: 0.45155161794026755  
 For values of alpha = 0.001 The log loss is: 0.4570964981295094  
 For values of alpha = 0.01 The log loss is: 0.45018599170785584  
 For values of alpha = 0.1 The log loss is: 0.46419309390805763  
 For values of alpha = 1 The log loss is: 0.4940887866823525  
 For values of alpha = 10 The log loss is: 0.5458906040136738

-----

-----

-----

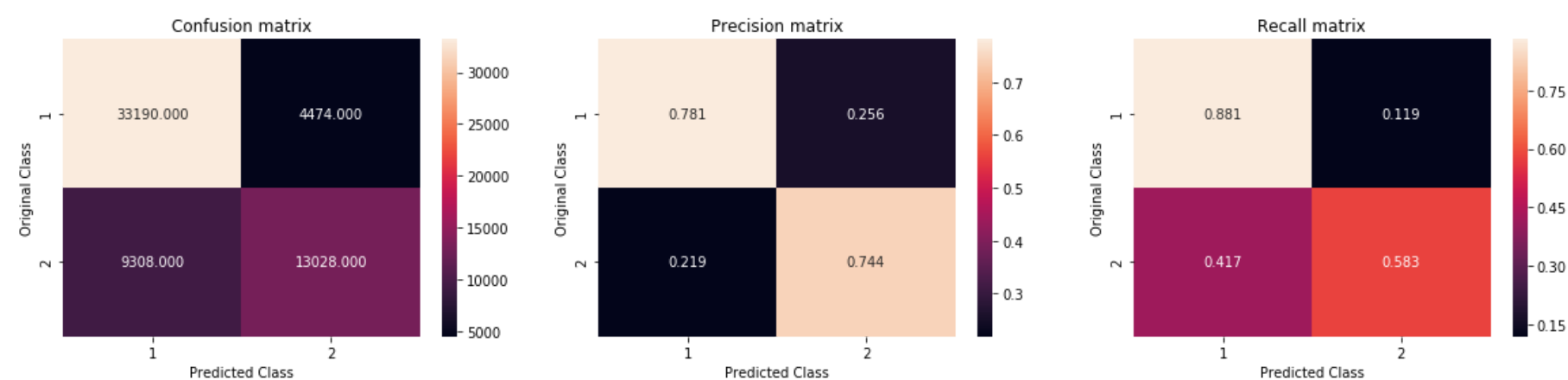
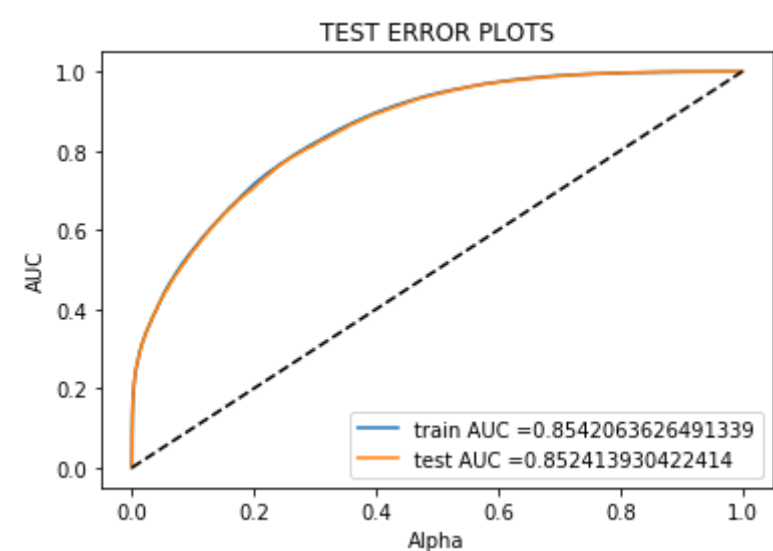


-----

-----

-----

For values of best alpha = 0.01 The train log loss is: 0.44890756801054355  
 For values of best alpha = 0.01 The test log loss is: 0.45018599170785584  
 Total number of data points : 60000



## 5.2 Linear SVM with TFIDF and hyperparameter tuning

In [0]: *#performing hyperparametertuning , modeling , plotting error plot and plotting confusion matix*

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

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))

print("----"*44)
print("----"*44)

#Error Plot
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()

print("----"*44)
print("----"*44)

#Building Model
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))

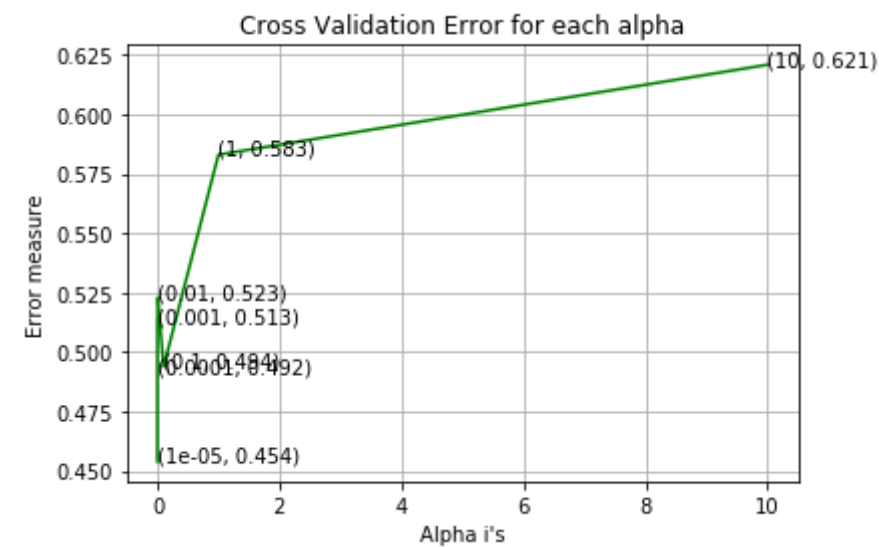
train_predicted = sig_clf.predict_proba(X_train)[:, 1]
test_predicted = sig_clf.predict_proba(X_test)[:, 1]

train_fpr, train_tpr, tr_thresholds = roc_curve(y_train, train_predicted)
test_fpr, test_tpr, te_thresholds = roc_curve(y_test, test_predicted)

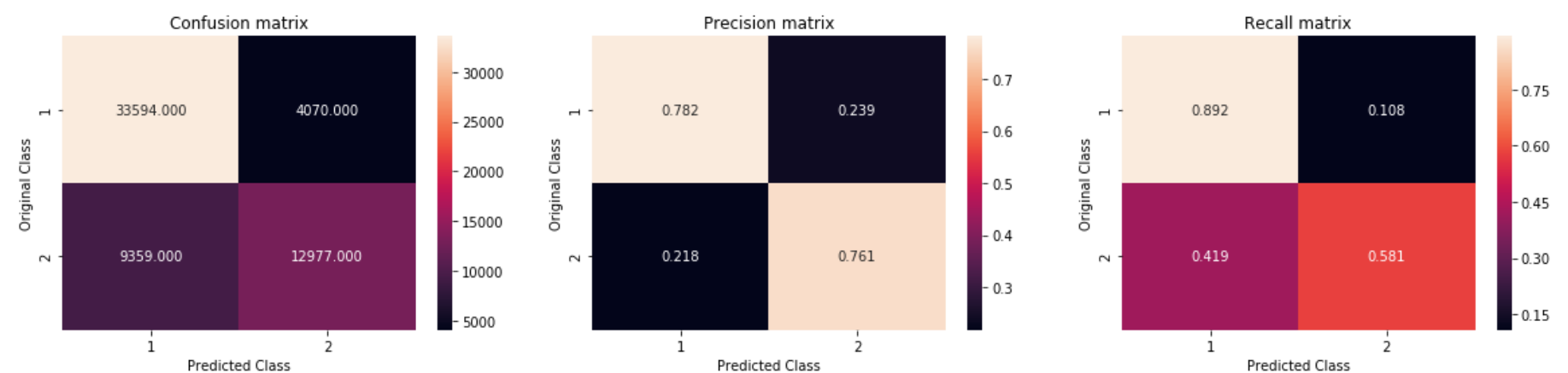
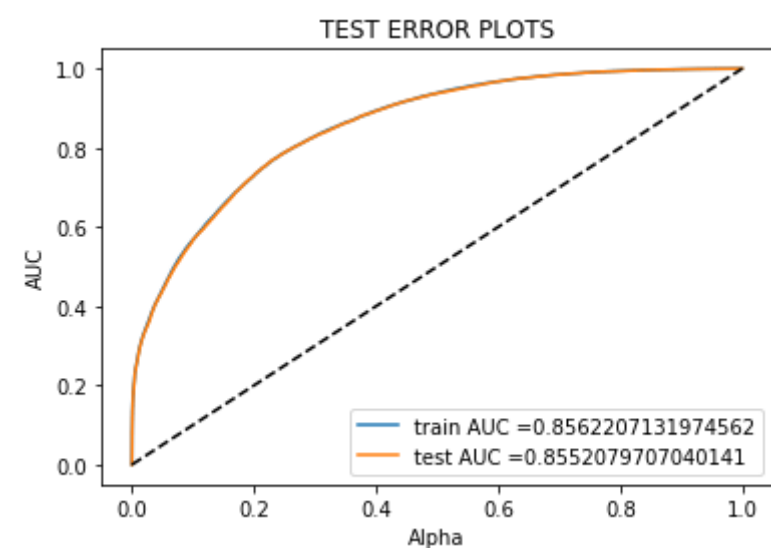
#test plotting
plt.plot(train_fpr, train_tpr, label="train AUC =" +str(auc(train_fpr, train_tpr)))
plt.plot(test_fpr, test_tpr, label="test AUC =" +str(auc(test_fpr, test_tpr)))
plt.legend()
plt.xlabel("Alpha")
plt.plot([0, 1], [0, 1], 'k--')
plt.ylabel("AUC")
plt.title("TEST ERROR PLOTS")
plt.legend(loc='best')
plt.show()

#Confusion Matrix
plot_confusion_matrix(y_test, predicted_y)
```

For values of alpha = 1e-05 The log loss is: 0.4539358164754859  
 For values of alpha = 0.0001 The log loss is: 0.49167466825173906  
 For values of alpha = 0.001 The log loss is: 0.5125347414444757  
 For values of alpha = 0.01 The log loss is: 0.5229077281816606  
 For values of alpha = 0.1 The log loss is: 0.49367715423273384  
 For values of alpha = 1 The log loss is: 0.5830223884422053  
 For values of alpha = 10 The log loss is: 0.6209083109445622



For values of best alpha = 1e-05 The train log loss is: 0.45253939539505317  
 For values of best alpha = 1e-05 The test log loss is: 0.4539358164754859  
 Total number of data points : 60000



## 5.3 Xgboost with TFIDF and hyperparameter tuning Using RandomizedSearchCV

```

In [0]: %%time
#HyperparameterTunning on XGboost

import xgboost as xgb
from sklearn.model_selection import TimeSeriesSplit
from sklearn.model_selection import RandomizedSearchCV
#Finding the best 'max_depth' using Forward Chaining Cross Validation or Time Series CV
#https://stackoverflow.com/questions/40005795/how-does-sp-randint-work
#g = sp_randint(2,300)
xgb_clf= xgb.XGBClassifier()

n_estimators=[1, 5, 10, 15, 20]
max_depth=[1, 5, 10, 50, 100]

param_grid = {'n_estimators':n_estimators, 'max_depth':max_depth} #params we need to try on classifier
tscv = TimeSeriesSplit(n_splits=5) #For time based splitting
gsv = RandomizedSearchCV(xgb_clf,param_grid,cv=tscv,verbose=1,scoring='roc_auc',n_jobs=-1)
gsv.fit(X_train, y_train)
#savetofile(gsv,"Decision Trees on BOW")
print("Best HyperParameter: ",gsv.best_params_)
print("Best Accuracy: %.2f%%"%(gsv.best_score_*100))

```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

```

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done 46 tasks      | elapsed: 21.1min
[Parallel(n_jobs=-1)]: Done 50 out of 50 | elapsed: 26.0min finished

```

```

Best HyperParameter: {'n_estimators': 20, 'max_depth': 50}
Best Accuracy: 90.77%
CPU times: user 3min 20s, sys: 291 ms, total: 3min 21s
Wall time: 29min 17s

```

```

In [0]: #Building XGBoost Model

xgb_clf = xgb.XGBClassifier(max_depth=50,n_estimators=20)#best values we got after RandomizedSearchCV

xgb_clf.fit(X_train,y_train)
y_pred = xgb_clf.predict(X_test)

```

```

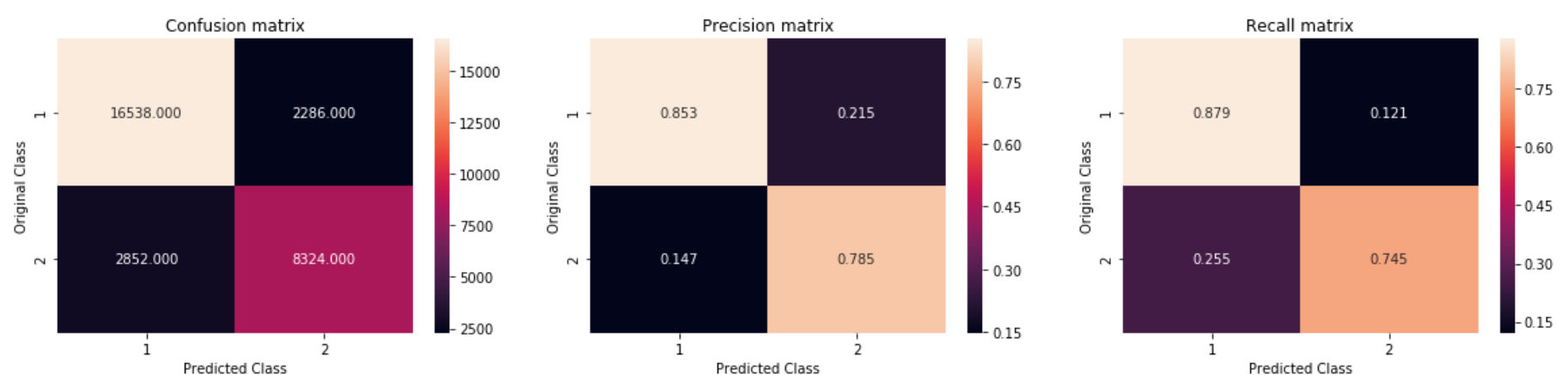
In [0]: #performing hyperparametertuning , modeling , plotting error plot and plotting confusion matix

predict_y = xgb_clf.predict_proba(X_test)
print("The test log loss is:",log_loss(y_test, predict_y))

#Confusion Matrix
plot_confusion_matrix(y_test, predicted_y)

```

The test log loss is: 0.3715554483974973



## 6.Conclusion

S.no	Featurization	Model Algo	alpha	max_depth	n_estimator	log loss
1	Word2Vec	Random Model	-	-	-	.8844
2	Word2Vec	Logistic Regression	1	-	-	.6291
3	Word2Vec	SVM	0.1	-	-	.6234
4	TFIDF	Logistic Regression	0.001	-	-	.4539
5	TFIDF	SVM	0.00001	-	-	.4501
6	TFIDF	XGBoost	-	50	20	.3715

1. As very first we have perform Random model for getting worst results of log loss, and we know our classified models can not exceed this limits, if they does then that model is totally not perfect.
2. we have used Word2Vec featurization on our linear models( Logistic Regression and SVM), they show great results , reducing log loss to .5198.
3. But still there is some chance of improving this models So we go for TFIDF featurization.
4. Again we have applied TFIDF on Logistic Regression (0.45018599170785584 logg loss) and SVM (.4539358164754859 log loss ) models they gives us much satisfactory results then previous model .
5. Now At last we have used Most complex model i.e XGBoost for this optomization and we got fantastic results.
6. XGboost takes much time to perform but gives best results till, Log Loss : .3715554483974973 which is least values till now.
7. So we conclude for this type of 'Question pair similarity' problem we must go for XGBoost.

---xxx---