

Quora Question Pairs

1. Business Problem

1.1 Description

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

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

Credits: Kaggle

Problem Statement

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

1.2 Sources/Useful Links

- Source : <https://www.kaggle.com/c/quora-question-pairs> ##### Useful Links
- Discussions : <https://www.kaggle.com/anokas/data-analysis-xgboost-starter-0-35460-lb/comments>
- Kaggle Winning Solution and other approaches:
<https://www.dropbox.com/sh/93968nfnrzh8bp5/AACZdtsApc1QSTQc7X0H3QZ5a?dl=0>
- Blog 1 : <https://engineering.quora.com/Semantic-Question-Matching-with-Deep-Learning>
- Blog 2 : <https://towardsdatascience.com/identifying-duplicate-questions-on-quora-top-12-on-kaggle-4c1cf93f1c30>

1.3 Real world/Business Objectives and Constraints

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

2. Machine Learning Problem

2.1 Data

2.1.1 Data Overview

- Data will be in a file Train.csv
- Train.csv contains 5 columns : qid1, qid2, question1, question2, is_duplicate
- Size of Train.csv - 60MB
- Number of rows in Train.csv = 404,290

2.1.2 Example Data point

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

2.2 Mapping the real world problem to an ML problem

2.2.1 Type of Machine Learning Problem

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

2.2.2 Performance Metric

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

Metric(s):

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

2.3 Train and Test Construction

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

3. Exploratory Data Analysis

In [1]:

```
import warnings
warnings.filterwarnings("ignore")
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
import distance
from nltk.stem import PorterStemmer
from bs4 import BeautifulSoup
import re
from nltk.corpus import stopwords
# 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
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
```

3.1 Reading data and basic stats

In [7]:

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

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

Number of data points: 404290

In [8]:

```
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 $23^{24}/\text{math}$ i...	0
4	4	9	10	Which one dissolve in water quickly sugar, salt...	Which fish would survive in salt water?	0

In [10]:

```
df.info()
```

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

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

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

3.2.1 Distribution of data points among output classes

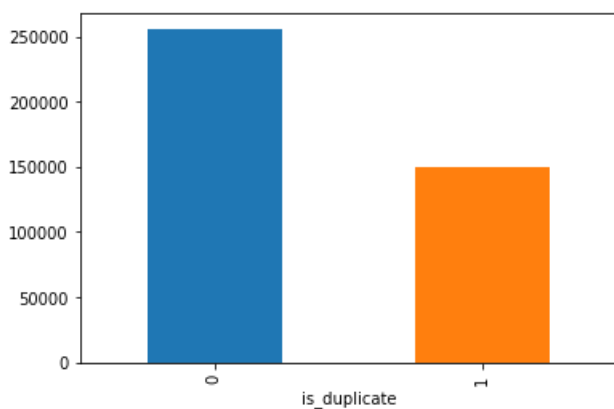
- Number of duplicate(smilar) and non-duplicate(non similar) questions

In [12]:

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

Out[12]:

<matplotlib.axes._subplots.AxesSubplot at 0x7f91bf360128>



In [13]:

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

~> Total number of question pairs for training:
404290

In [14]:

```
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 [15]:

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

Total number of Unique Questions are: 537933

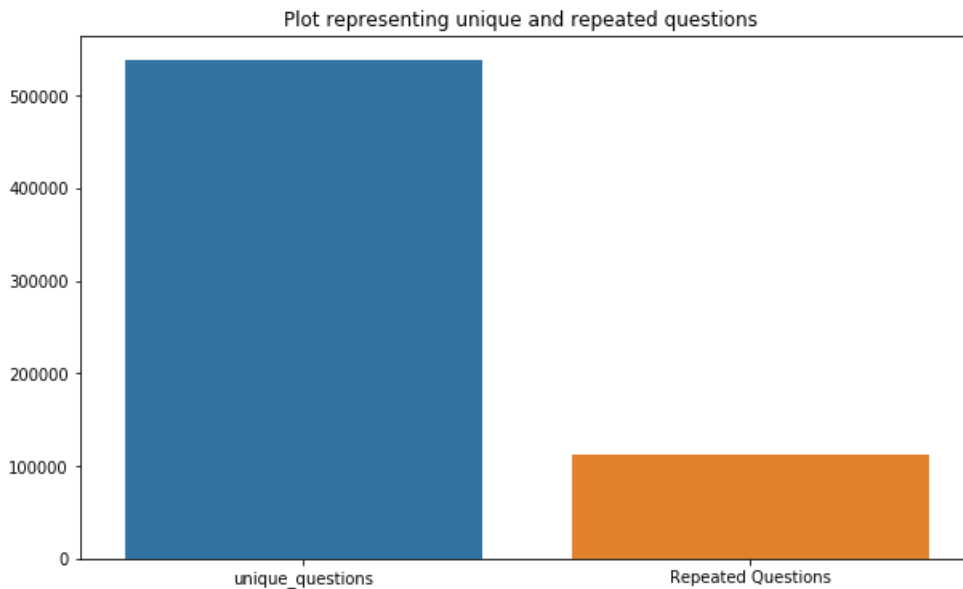
Number of unique questions that appear more than one time: 111780 (20.77953945937505%)

Max number of times a single question is repeated: 157

In [16]:

```
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 [18]:

```
#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 [19]:

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

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

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

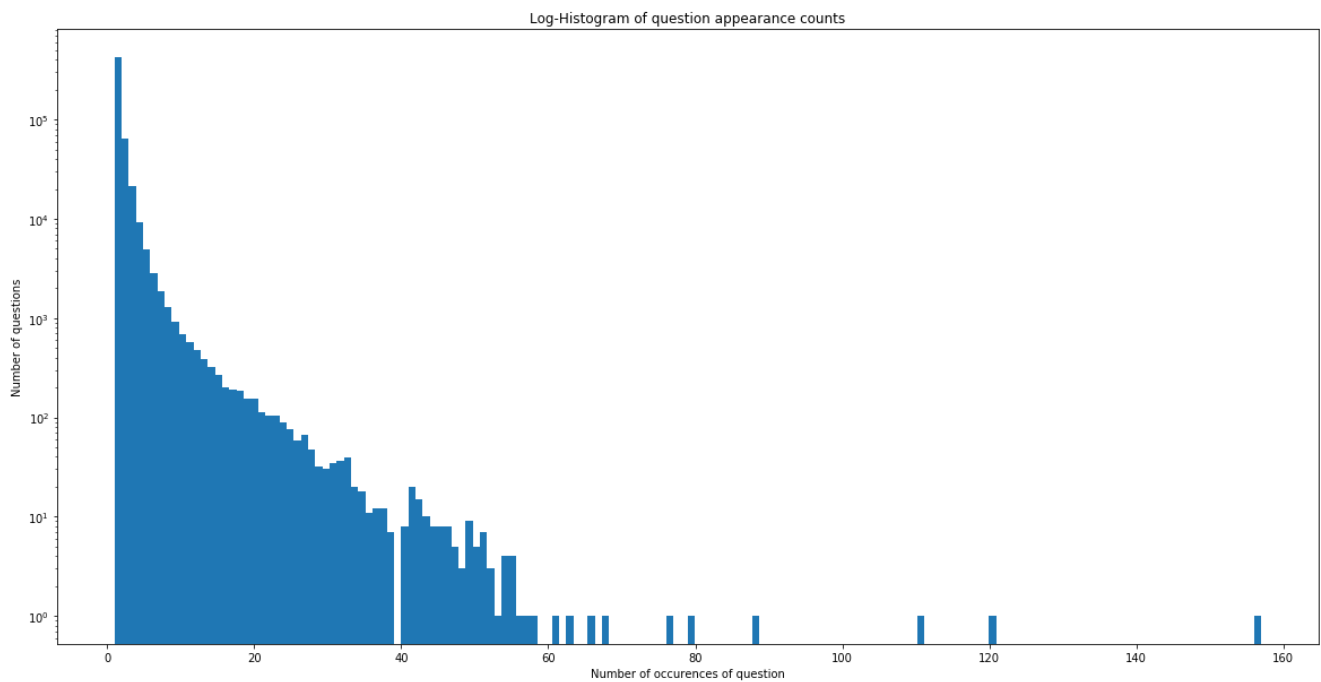
plt.title('Log-Histogram of question appearance counts')

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 [21]:

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

	id	qid1	qid2	question1 \
105780	105780	174363	174364	How can I develop android app?
201841	201841	303951	174364	How can I create an Android app?
363362	363362	493340	493341	NaN

	question2	is_duplicate
105780	NaN	0
201841	NaN	0
363362	My Chinese name is Haichao Yu. What English na...	0

- There are three rows with null values, one in question1 and the other two in question2

In [22]:

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

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

3.3 Basic Feature Extraction (before cleaning)

Let us now construct a few features like:

- **freq_qid1** = Frequency of qid1's
- **freq_qid2** = Frequency of qid2's
- **q1len** = Length of q1
- **q2len** = Length of q2

4	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common
				dissolve in water quickly sugar, salt...	Which fish would survive in salt water?								

3.3.1 Analysis of some of the extracted features

- Here are some questions have only one single words.

In [24]:

```
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
```

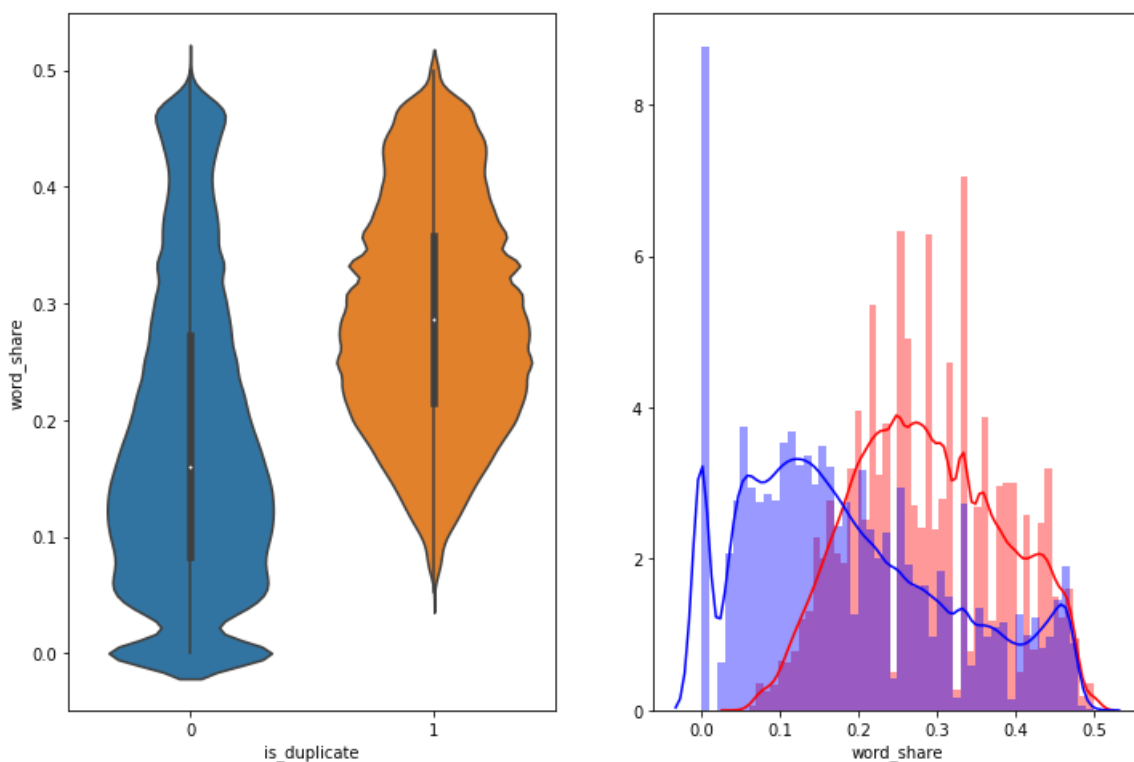
3.3.1.1 Feature: word_share

In [25]:

```
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

with high word similarity

- The average word share and Common no. of words of qid1 and qid2 is more when they are duplicate(Similar)

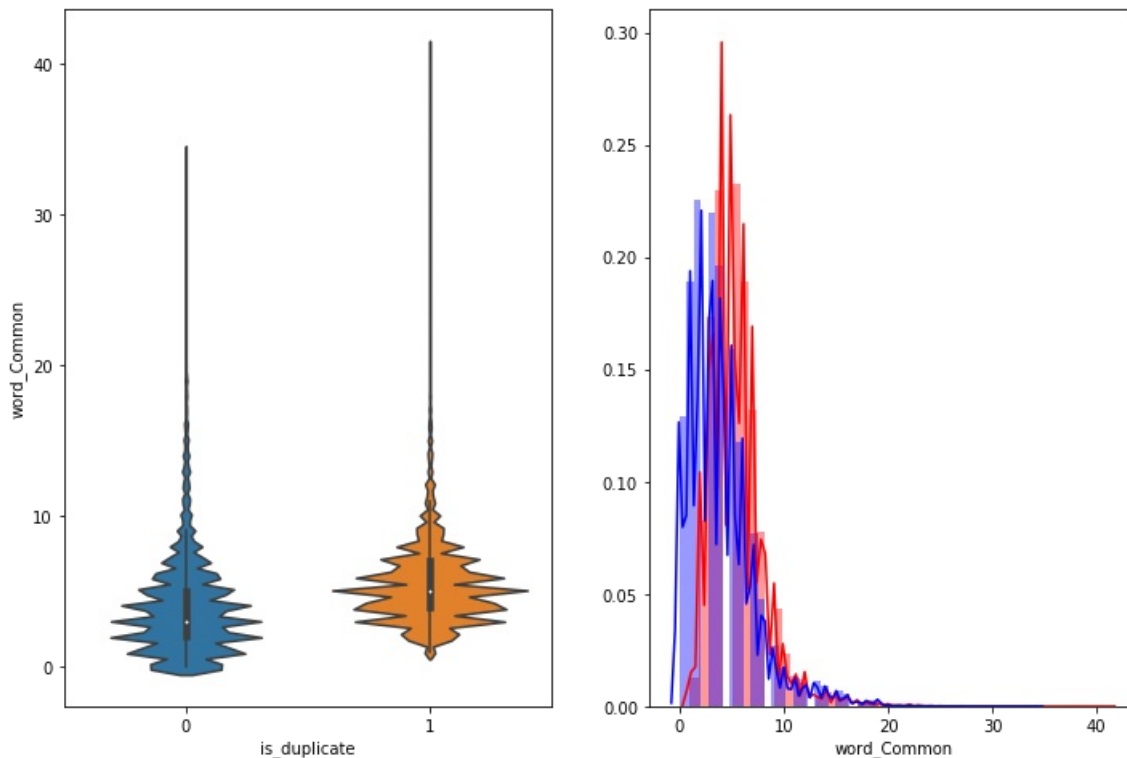
3.3.1.2 Feature: word_Common

In [26]:

```
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

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

In [27]:

```
# 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 "
```

```

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

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

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

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

return x

```

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

3.5 Advanced Feature Extraction (NLP and Fuzzy Features)

Definition:

- **Token**: You get a token by splitting sentence a space
- **Stop_Word** : stop words as per NLTK.
- **Word** : A token that is not a stop_word

Features:

- **cwc_min** : Ratio of common_word_count to min length of word count of Q1 and Q2

$$\text{cwc_min} = \text{common_word_count} / (\min(\text{len}(q1_words), \text{len}(q2_words)))$$
- **cwc_max** : Ratio of common_word_count to max length of word count of Q1 and Q2

$$\text{cwc_max} = \text{common_word_count} / (\max(\text{len}(q1_words), \text{len}(q2_words)))$$
- **csc_min** : Ratio of common_stop_count to min length of stop count of Q1 and Q2

$$\text{csc_min} = \text{common_stop_count} / (\min(\text{len}(q1_stops), \text{len}(q2_stops)))$$
- **csc_max** : Ratio of common_stop_count to max length of stop count of Q1 and Q2

$$\text{csc_max} = \text{common_stop_count} / (\max(\text{len}(q1_stops), \text{len}(q2_stops)))$$
- **ctc_min** : Ratio of common_token_count to min length of token count of Q1 and Q2

$$\text{ctc_min} = \text{common_token_count} / (\min(\text{len}(q1_tokens), \text{len}(q2_tokens)))$$
- **ctc_max** : Ratio of common_token_count to max length of token count of Q1 and Q2

$$\text{ctc_max} = \text{common_token_count} / (\max(\text{len}(q1_tokens), \text{len}(q2_tokens)))$$
- **last_word_eq** : Check if First word of both questions is equal or not

$$\text{last_word_eq} = \text{int}(q1_tokens[-1] == q2_tokens[-1])$$
- **first_word_eq** : Check if First word of both questions is equal or not

$$\text{first_word_eq} = \text{int}(q1_tokens[0] == q2_tokens[0])$$
- **abs_len_diff** : Abs. length difference

$$\text{abs_len_diff} = \text{abs}(\text{len}(q1_tokens) - \text{len}(q2_tokens))$$
- **mean_len** : Average Token Length of both Questions

- **mean_len** : Average Token Length of both Questions

`mean_len = (len(q1_tokens) + len(q2_tokens))/2`

- **fuzz_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **fuzz_partial_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **token_sort_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **token_set_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> <http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>
- **longest_substr_ratio** : Ratio of length longest common substring to min length of token count of Q1 and Q2
`longest_substr_ratio = len(longest common substring) / (min(len(q1_tokens), len(q2_tokens)))`

In [28]:

```
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:
```

```

-- len(strs), ...
    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 [29]:

```

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

```

Extracting features for train:
token features...
fuzzy features..

Out[29]:

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	

id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	first_word
1	1	3	4	story of kohinoor dia...	happen if the indian government sto...	0	0.799984	0.399996	0.749981	0.599988	...	0.466664	0.0

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 [31]:

```
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('train_p.txt', p, delimiter=' ', fmt='%s')
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 [32]:

```
# 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 : 33193067

Word Clouds generated from duplicate pair question's text

In [33]:

```
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()
```

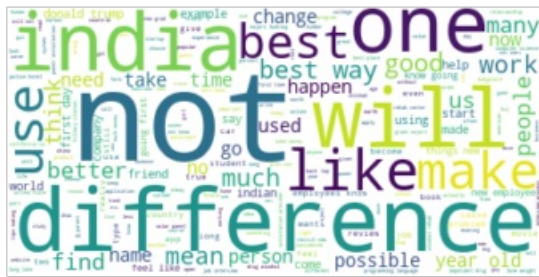
Word Cloud for Duplicate Question pairs



In [34]:

```
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()
```

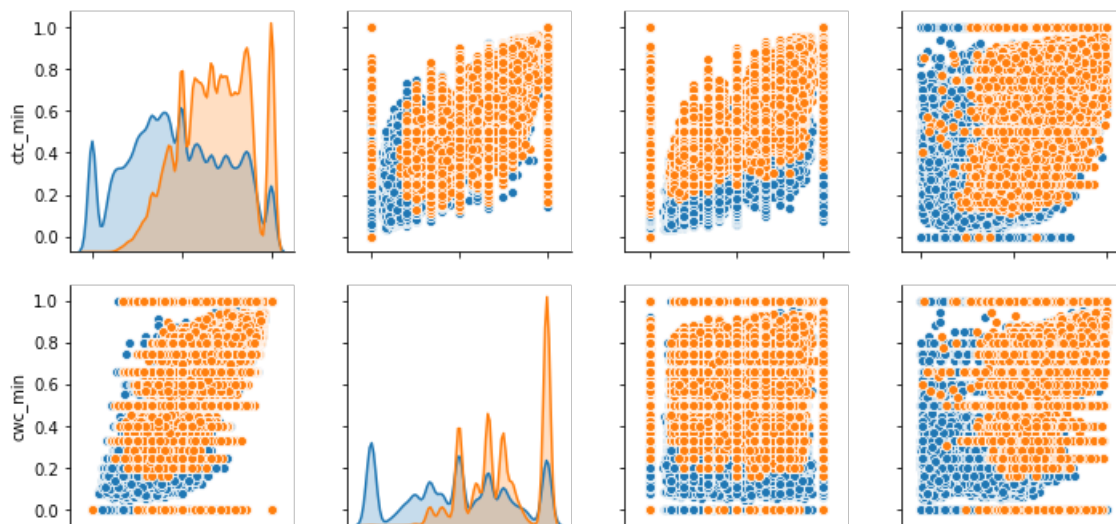
Word Cloud for non-Duplicate Question pairs:

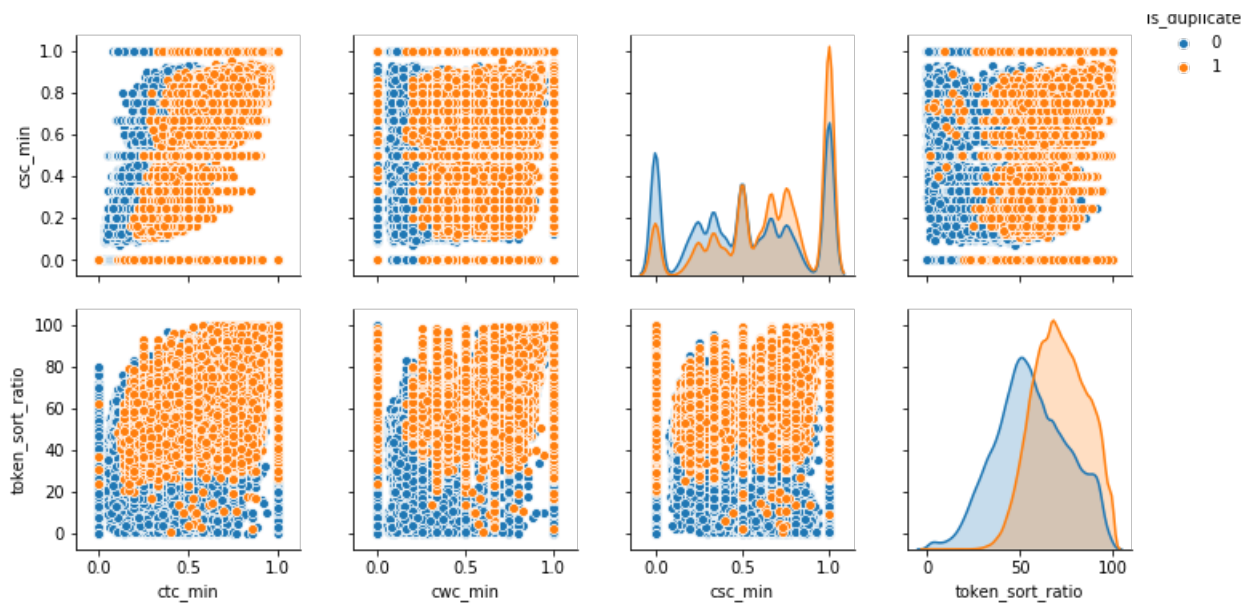


3.5.1.2 Pair plot of features ['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio']

In [35]:

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



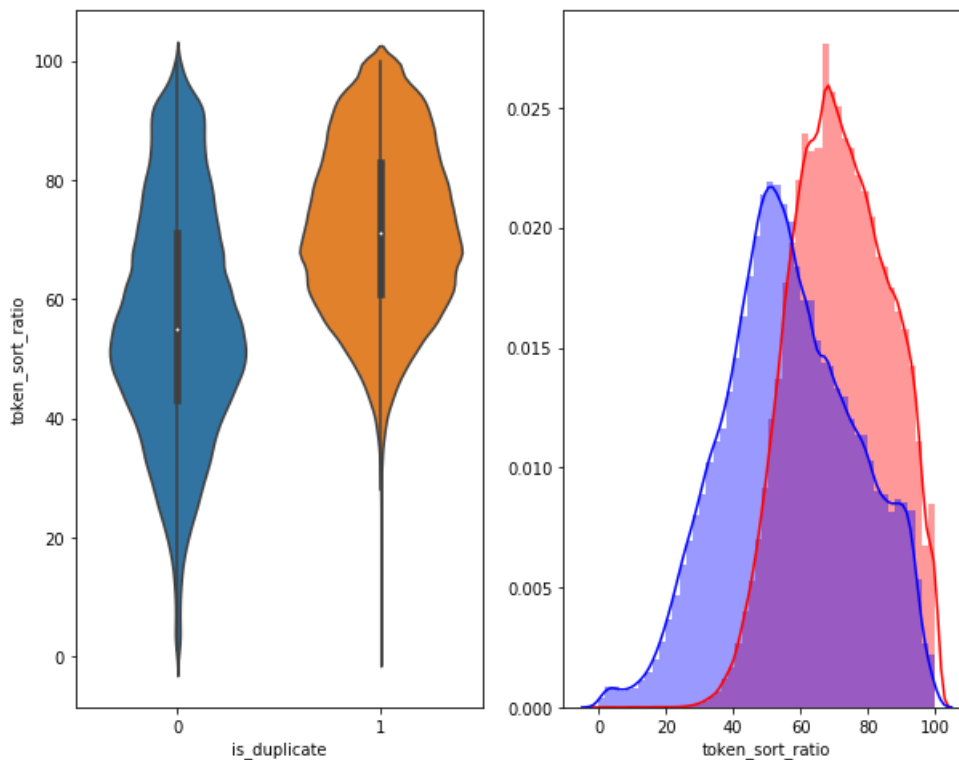


In [36]:

```
# 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 [37]:

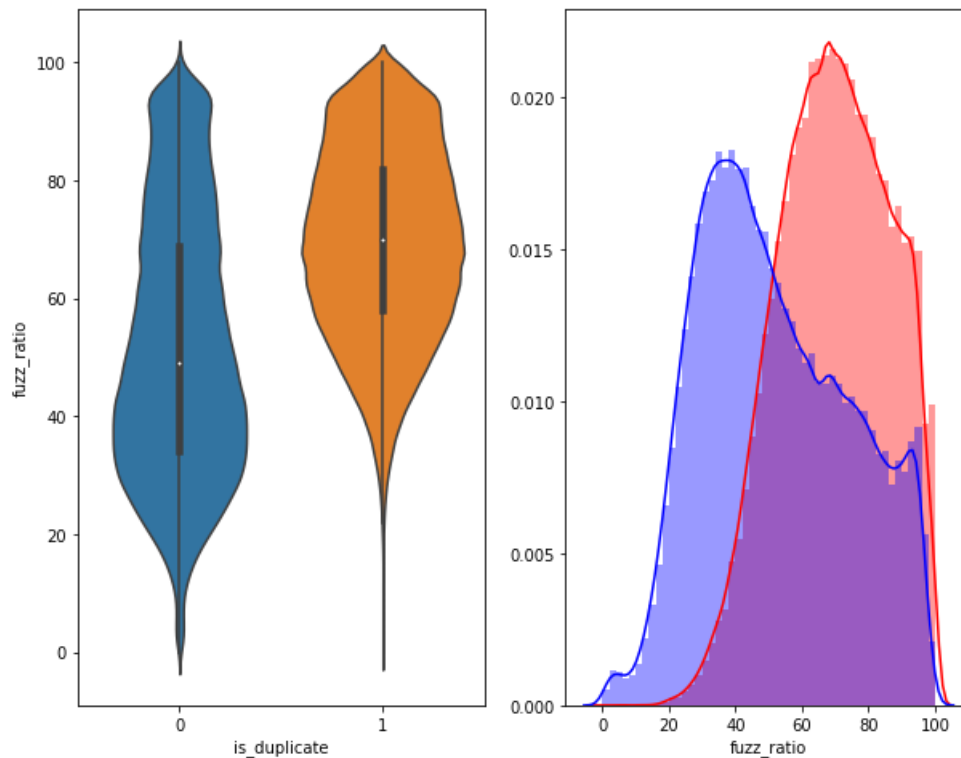
```
plt.figure(figsize=(10, 8))

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 [38]:

```
# Using TSNE for Dimentionality reduction for 15 Features(Generated after cleaning the data) to 2
and 3 dimensions

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 [39]:

```
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 0.028s...
[t-SNE] Computed neighbors for 5000 samples in 0.356s...
[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: 0.130446
[t-SNE] Computed conditional probabilities in 0.489s
[t-SNE] Iteration 50: error = 81.2897949, gradient norm = 0.0455700 (50 iterations in 5.088s)
[t-SNE] Iteration 100: error = 70.6164398, gradient norm = 0.0095177 (50 iterations in 3.951s)
[t-SNE] Iteration 150: error = 68.9172134, gradient norm = 0.0056736 (50 iterations in 3.835s)
```

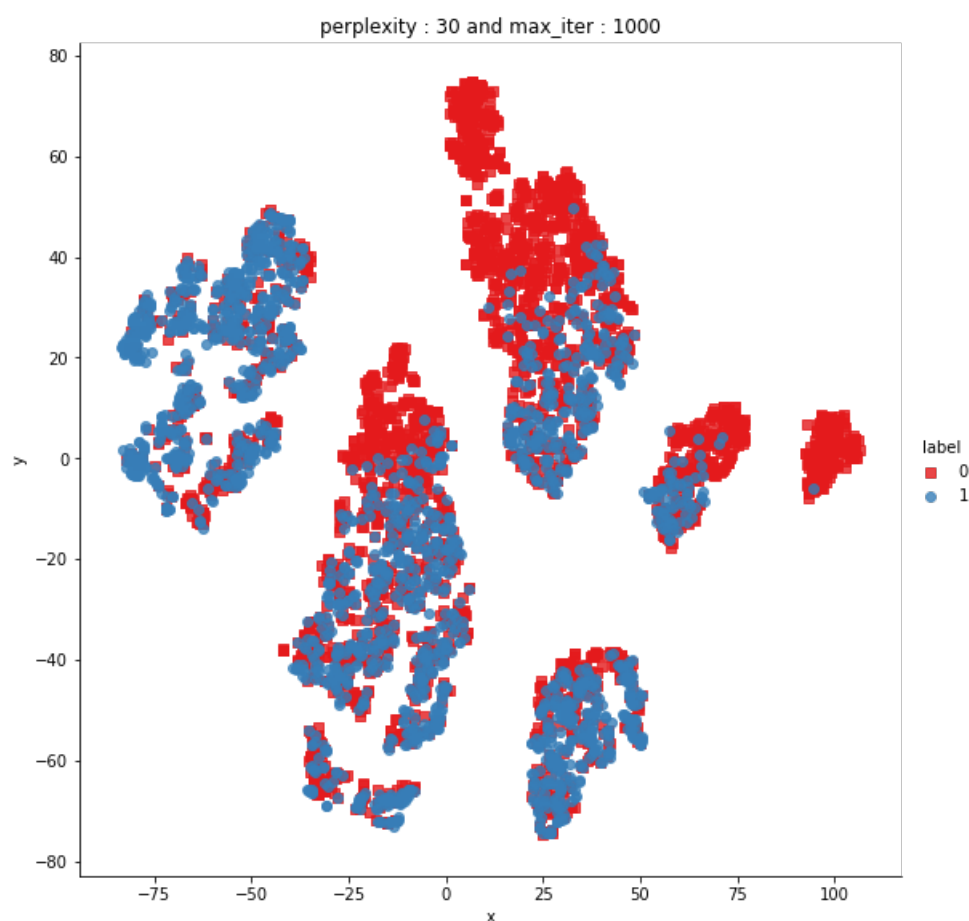


```
[t-SNE] Iteration 200: error = 68.1004639, gradient norm = 0.0049672 (50 iterations in 3.928s)
[t-SNE] Iteration 250: error = 67.5914536, gradient norm = 0.0039700 (50 iterations in 4.025s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.591454
[t-SNE] Iteration 300: error = 1.7926962, gradient norm = 0.0011878 (50 iterations in 4.226s)
[t-SNE] Iteration 350: error = 1.3936826, gradient norm = 0.0004807 (50 iterations in 4.184s)
[t-SNE] Iteration 400: error = 1.2281071, gradient norm = 0.0002778 (50 iterations in 4.203s)
[t-SNE] Iteration 450: error = 1.1385784, gradient norm = 0.0001864 (50 iterations in 4.284s)
[t-SNE] Iteration 500: error = 1.0835493, gradient norm = 0.0001437 (50 iterations in 4.208s)
[t-SNE] Iteration 550: error = 1.0471643, gradient norm = 0.0001152 (50 iterations in 4.199s)
[t-SNE] Iteration 600: error = 1.0231258, gradient norm = 0.0001007 (50 iterations in 4.225s)
[t-SNE] Iteration 650: error = 1.0069925, gradient norm = 0.0000892 (50 iterations in 4.262s)
[t-SNE] Iteration 700: error = 0.9953420, gradient norm = 0.0000804 (50 iterations in 4.267s)
[t-SNE] Iteration 750: error = 0.9866475, gradient norm = 0.0000728 (50 iterations in 4.279s)
[t-SNE] Iteration 800: error = 0.9796536, gradient norm = 0.0000658 (50 iterations in 4.266s)
[t-SNE] Iteration 850: error = 0.9737327, gradient norm = 0.0000618 (50 iterations in 4.250s)
[t-SNE] Iteration 900: error = 0.9688665, gradient norm = 0.0000594 (50 iterations in 4.259s)
[t-SNE] Iteration 950: error = 0.9644679, gradient norm = 0.0000589 (50 iterations in 4.251s)
[t-SNE] Iteration 1000: error = 0.9610358, gradient norm = 0.0000559 (50 iterations in 4.236s)
[t-SNE] Error after 1000 iterations: 0.961036
```

In [40]:

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

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



In [46]:

```
#Trying with different parameters
tsne2d_1 = TSNE(
    n_components=2,
    init='random', # pca
    random_state=101,
    method='barnes_hut',
    n_iter=2000,
```

```

verbose=2,
angle=0.5
).fit_transform(X)

```

```

[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.008s...
[t-SNE] Computed neighbors for 5000 samples in 0.363s...
[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: 0.130446
[t-SNE] Computed conditional probabilities in 0.340s
[t-SNE] Iteration 50: error = 81.2897949, gradient norm = 0.0455700 (50 iterations in 5.166s)
[t-SNE] Iteration 100: error = 70.6164398, gradient norm = 0.0095177 (50 iterations in 3.983s)
[t-SNE] Iteration 150: error = 68.9172134, gradient norm = 0.0056736 (50 iterations in 3.887s)
[t-SNE] Iteration 200: error = 68.1004639, gradient norm = 0.0049672 (50 iterations in 4.047s)
[t-SNE] Iteration 250: error = 67.5914536, gradient norm = 0.0039700 (50 iterations in 4.120s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.591454
[t-SNE] Iteration 300: error = 1.7926962, gradient norm = 0.0011878 (50 iterations in 4.332s)
[t-SNE] Iteration 350: error = 1.3936826, gradient norm = 0.0004807 (50 iterations in 4.278s)
[t-SNE] Iteration 400: error = 1.2281071, gradient norm = 0.0002778 (50 iterations in 4.235s)
[t-SNE] Iteration 450: error = 1.1385784, gradient norm = 0.0001864 (50 iterations in 4.255s)
[t-SNE] Iteration 500: error = 1.0835493, gradient norm = 0.0001437 (50 iterations in 4.276s)
[t-SNE] Iteration 550: error = 1.0471643, gradient norm = 0.0001152 (50 iterations in 4.290s)
[t-SNE] Iteration 600: error = 1.0231258, gradient norm = 0.0001007 (50 iterations in 4.328s)
[t-SNE] Iteration 650: error = 1.0069925, gradient norm = 0.0000892 (50 iterations in 4.400s)
[t-SNE] Iteration 700: error = 0.9953420, gradient norm = 0.0000804 (50 iterations in 4.353s)
[t-SNE] Iteration 750: error = 0.9866475, gradient norm = 0.0000728 (50 iterations in 4.350s)
[t-SNE] Iteration 800: error = 0.9796536, gradient norm = 0.0000658 (50 iterations in 4.336s)
[t-SNE] Iteration 850: error = 0.9737327, gradient norm = 0.0000618 (50 iterations in 4.329s)
[t-SNE] Iteration 900: error = 0.9688665, gradient norm = 0.0000594 (50 iterations in 4.328s)
[t-SNE] Iteration 950: error = 0.9644679, gradient norm = 0.0000589 (50 iterations in 4.321s)
[t-SNE] Iteration 1000: error = 0.9610358, gradient norm = 0.0000559 (50 iterations in 4.349s)
[t-SNE] Iteration 1050: error = 0.9580597, gradient norm = 0.0000548 (50 iterations in 4.334s)
[t-SNE] Iteration 1100: error = 0.9553435, gradient norm = 0.0000529 (50 iterations in 4.339s)
[t-SNE] Iteration 1150: error = 0.9530139, gradient norm = 0.0000485 (50 iterations in 4.322s)
[t-SNE] Iteration 1200: error = 0.9508933, gradient norm = 0.0000530 (50 iterations in 4.315s)
[t-SNE] Iteration 1250: error = 0.9491470, gradient norm = 0.0000486 (50 iterations in 4.326s)
[t-SNE] Iteration 1300: error = 0.9475195, gradient norm = 0.0000466 (50 iterations in 4.357s)
[t-SNE] Iteration 1350: error = 0.9459193, gradient norm = 0.0000445 (50 iterations in 4.337s)
[t-SNE] Iteration 1400: error = 0.9443618, gradient norm = 0.0000420 (50 iterations in 4.323s)
[t-SNE] Iteration 1450: error = 0.9428130, gradient norm = 0.0000435 (50 iterations in 4.326s)
[t-SNE] Iteration 1500: error = 0.9413453, gradient norm = 0.0000414 (50 iterations in 4.379s)
[t-SNE] Iteration 1550: error = 0.9399194, gradient norm = 0.0000399 (50 iterations in 4.339s)
[t-SNE] Iteration 1600: error = 0.9384849, gradient norm = 0.0000405 (50 iterations in 4.314s)
[t-SNE] Iteration 1650: error = 0.9370704, gradient norm = 0.0000394 (50 iterations in 4.310s)
[t-SNE] Iteration 1700: error = 0.9357622, gradient norm = 0.0000386 (50 iterations in 4.361s)
[t-SNE] Iteration 1750: error = 0.9346933, gradient norm = 0.0000400 (50 iterations in 4.337s)
[t-SNE] Iteration 1800: error = 0.9337534, gradient norm = 0.0000398 (50 iterations in 4.430s)
[t-SNE] Iteration 1850: error = 0.9328931, gradient norm = 0.0000358 (50 iterations in 4.554s)
[t-SNE] Iteration 1900: error = 0.9319998, gradient norm = 0.0000359 (50 iterations in 4.521s)
[t-SNE] Iteration 1950: error = 0.9309730, gradient norm = 0.0000353 (50 iterations in 4.344s)
[t-SNE] Iteration 2000: error = 0.9301456, gradient norm = 0.0000403 (50 iterations in 4.318s)
[t-SNE] Error after 2000 iterations: 0.930146

```

In [47]:

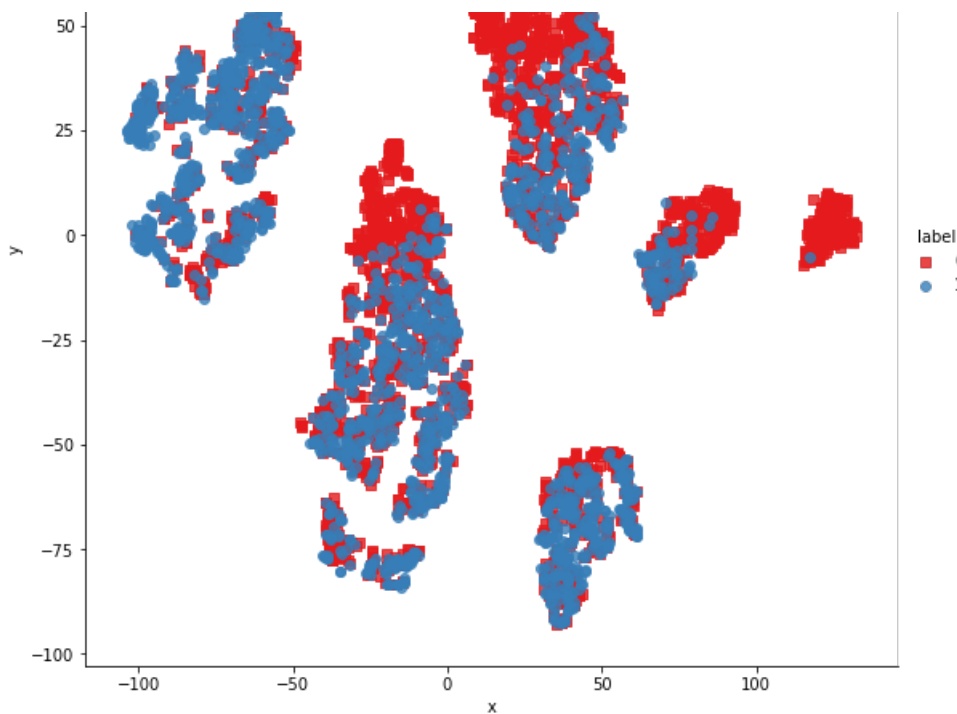
```

df = pd.DataFrame({'x':tsne2d_1[:,0], 'y':tsne2d_1[:,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, 2000))
plt.show()

```





Observation:

There is not much change in the plot after increasing the iterations from 1000 to 2000

In [48]:

```
#
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 0.008s...
[t-SNE] Computed neighbors for 5000 samples in 0.354s...
[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: 0.130446
[t-SNE] Computed conditional probabilities in 0.333s
[t-SNE] Iteration 50: error = 80.5298615, gradient norm = 0.0306586 (50 iterations in 18.491s)
[t-SNE] Iteration 100: error = 69.3777008, gradient norm = 0.0037944 (50 iterations in 10.089s)
[t-SNE] Iteration 150: error = 67.9726028, gradient norm = 0.0017517 (50 iterations in 9.563s)
[t-SNE] Iteration 200: error = 67.4098892, gradient norm = 0.0013384 (50 iterations in 9.555s)
[t-SNE] Iteration 250: error = 67.0977859, gradient norm = 0.0009594 (50 iterations in 9.701s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.097786
[t-SNE] Iteration 300: error = 1.5276405, gradient norm = 0.0007237 (50 iterations in 12.093s)
[t-SNE] Iteration 350: error = 1.1820400, gradient norm = 0.0002119 (50 iterations in 14.811s)
[t-SNE] Iteration 400: error = 1.0407882, gradient norm = 0.0001023 (50 iterations in 14.548s)
[t-SNE] Iteration 450: error = 0.9688321, gradient norm = 0.0000652 (50 iterations in 14.196s)
[t-SNE] Iteration 500: error = 0.9303923, gradient norm = 0.0000554 (50 iterations in 13.928s)
[t-SNE] Iteration 550: error = 0.9110239, gradient norm = 0.0000524 (50 iterations in 13.890s)
[t-SNE] Iteration 600: error = 0.9016075, gradient norm = 0.0000421 (50 iterations in 14.007s)
[t-SNE] Iteration 650: error = 0.8924681, gradient norm = 0.0000360 (50 iterations in 14.145s)
[t-SNE] Iteration 700: error = 0.8837291, gradient norm = 0.0000353 (50 iterations in 14.222s)
[t-SNE] Iteration 750: error = 0.8771634, gradient norm = 0.0000316 (50 iterations in 14.145s)
[t-SNE] Iteration 800: error = 0.8718039, gradient norm = 0.0000295 (50 iterations in 14.066s)
[t-SNE] Iteration 850: error = 0.8669323, gradient norm = 0.0000276 (50 iterations in 14.086s)
[t-SNE] Iteration 900: error = 0.8628623, gradient norm = 0.0000262 (50 iterations in 14.134s)
[t-SNE] Iteration 950: error = 0.8591092, gradient norm = 0.0000241 (50 iterations in 14.114s)
```

```
[t-SNE] Iteration 1000: error = 0.8553245, gradient norm = 0.0000220 (50 iterations in 14.137s)
[t-SNE] Error after 1000 iterations: 0.855325
```

In [49]:

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

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

3.6 Featurizing text data with tfidf weighted word-vectors

In [50]:

```
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
```

In [52]:

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

# encode questions to unicode
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

In [53]:

```
df.head()
```

Out[53]:

	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 quickly sugar, salt...	Which fish would survive in salt water?	0

In [54]:

```
questions = list(df['question1']) + list(df['question2'])

tfidf = TfidfVectorizer(lowercase=False, )
tfidf.fit_transform(questions)

# 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>
- It is trained on Wikipedia and therefore, it is stronger in terms of word semantics.

In [66]:

```
import spacy
from tqdm import tqdm
# 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 qul in tqdm(list(df['question1'])):
    doc1 = nlp(qul)
    # 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 df score
```

```

        # 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%|██████████| 404290/404290 [1:06:55<00:00, 100.69it/s]

In [67]:

```

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%|██████████| 404290/404290 [1:05:40<00:00, 102.59it/s]

In [68]:

```

dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
dfppro = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')

```

In [69]:

```

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 [70]:

```

# dataframe of nlp features
df1.head()

```

Out[70]:

	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
0	0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	2.0	13.0
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0	12.5
2	2	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	4.0	12.0
3	3	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	2.0	12.0
4	4	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	6.0	10.0

In [71]:

```

# data before preprocessing
df2.head()

```

Out[71]:

	id	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	id	freq_q1	freq_q2	q1len	q2len	q1_n_words	q2_n_words	word_Common	word_Total	word_share	freq_q1+q2	freq_q1-q2
0	0	1	1	66	57	14	12	10.0	23.0	0.434783	2	0
1	1	4	1	51	88	8	13	4.0	20.0	0.200000	5	3
2	2	1	1	73	59	14	10	4.0	24.0	0.166667	2	0
3	3	1	1	50	65	11	9	0.0	19.0	0.000000	2	0
4	4	3	1	76	39	13	7	2.0	20.0	0.100000	4	2

In [72]:

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

Out[72]:

	0	1	2	3	4	5	6	7	8	9	...	374
0	121.929932	100.083909	72.497893	115.641794	48.370882	34.619062	172.057790	-92.502619	113.223317	50.562433	...	12.397644
1	-78.070943	54.843787	82.738472	98.191860	51.234856	55.013505	-39.140728	-82.692345	45.161482	-9.556285	...	21.987077
2	-5.355019	73.671803	14.376363	104.130241	1.433539	35.229114	148.519384	-97.124593	41.972190	50.948728	...	3.027700
3	5.778362	-34.712030	48.999616	59.699200	40.661257	41.658736	-36.808593	24.170664	0.235597	29.407294	...	13.100006
4	51.138199	38.587305	123.639489	53.333049	47.062726	37.356195	298.722745	106.421118	106.248904	65.880715	...	13.906532

5 rows × 384 columns

In [73]:

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

Out[73]:

	0	1	2	3	4	5	6	7	8	9	...	374
0	125.983307	95.636488	42.114704	95.449977	37.386311	39.400091	148.116074	-87.851475	110.371982	62.272814	...	16.165594
1	106.871910	80.290334	79.066299	59.302089	42.175328	117.616658	144.364245	127.131503	22.962527	25.397566	...	-4.901128
2	7.072871	15.513378	1.846906	85.937576	33.808816	94.702334	122.256851	114.009526	53.922283	60.131814	...	8.359966
3	39.421531	44.136989	24.010929	85.265870	-0.339022	-9.323147	-60.499652	-37.044766	49.407848	23.350152	...	3.311411
4	31.950102	62.854113	1.778172	36.218766	45.130874	66.674876	106.342344	-22.901006	59.835942	62.663971	...	-2.403870

5 rows × 384 columns

In [74]:

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

Number of features in nlp dataframe : 17
Number of features in preprocessed dataframe : 12
Number of features in question1 w2v dataframe : 384
Number of features in question2 w2v dataframe : 384
Number of features in final dataframe : 797

In [75]:

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

4. Machine Learning Models

4.1 Reading data from file and storing into sql table

In [2]:

```
from sqlalchemy import create_engine # database connection
import datetime as dt
#Creating db file from csv
if not os.path.isfile('train.db'):
    disk_engine = create_engine('sqlite:///train.db')
    start = dt.datetime.now()
    chunksize = 180000
    j = 0
    index_start = 1
    for df in pd.read_csv('final_features.csv', names=['Unnamed: 0','id','is_duplicate','cwc_min','cwc_max','csc_min','csc_max','ctc_min','ctc_max','last_word_eq','first_word_eq','abs_len_diff','mean_len','token_set_ratio','token_sort_ratio','fuzz_ratio','fuzz_partial_ratio','longest_substr_ratio','freq_qid1','freq_qid2','qlen','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','384_y','385_y','386_y','387_y','388_y','389_y','390_y','391_y','392_y','393_y','394_y','395_y','396_y','397_y','398_y','399_y','400_y','401_y','402_y','403_y','404_y','405_y','406_y','407_y','408_y','409_y','410_y','411_y','412_y','413_y','414_y','415_y','416_y','417_y','418_y','419_y','420_y','421_y','422_y','423_y','424_y','425_y','426_y','427_y','428_y','429_y','430_y','431_y','432_y','433_y','434_y','435_y','436_y','437_y','438_y','439_y','440_y','441_y','442_y','443_y','444_y','445_y','446_y','447_y','448_y','449_y','450_y','451_y','452_y','453_y','454_y','455_y','456_y','457_y','458_y','459_y','460_y','461_y','462_y','463_y','464_y','465_y','466_y','467_y','468_y','469_y','470_y','471_y','472_y','473_y','474_y','475_y','476_y','477_y','478_y','479_y','480_y','481_y','482_y','483_y','484_y','485_y','486_y','487_y','488_y','489_y','490_y','491_y','492_y','493_y','494_y','495_y','496_y','497_y','498_y','499_y','500_y','501_y','502_y','503_y','504_y','505_y','506_y','507_y','508_y','509_y','510_y','511_y','512_y','513_y','514_y','515_y','516_y','517_y','518_y','519_y','520_y','521_y','522_y','523_y','524_y','525_y','526_y','527_y','528_y','529_y','530_y','531_y','532_y','533_y','534_y','535_y','536_y','537_y','538_y','539_y','540_y','541_y','542_y','543_y','544_y','545_y','546_y','547_y','548_y','549_y','550_y','551_y','552_y','553_y','554_y','555_y','556_y','557_y','558_y','559_y','560_y','561_y','562_y','563_y','564_y','565_y','566_y','567_y','568_y','569_y','570_y','571_y','572_y','573_y','574_y','575_y','576_y','577_y','578_y','579_y','580_y','581_y','582_y','583_y','584_y','585_y','586_y','587_y','588_y','589_y','590_y','591_y','592_y','593_y','594_y','595_y','596_y','597_y','598_y','599_y','600_y','601_y','602_y','603_y','604_y','605_y','606_y','607_y','608_y','609_y','610_y','611_y','612_y','613_y','614_y','615_y','616_y','617_y','618_y','619_y','620_y','621_y','622_y','623_y','624_y','625_y','626_y','627_y','628_y','629_y','630_y','631_y','632_y','633_y','634_y','635_y','636_y','637_y','638_y','639_y','640_y','641_y','642_y','643_y','644_y','645_y','646_y','647_y','648_y','649_y','650_y','651_y','652_y','653_y','654_y','655_y','656_y','657_y','658_y','659_y','660_y','661_y','662_y','663_y','664_y','665_y','666_y','667_y','668_y','669_y','670_y','671_y','672_y','673_y','674_y','675_y','676_y','677_y','678_y','679_y','680_y','681_y','682_y','683_y','684_y','685_y','686_y','687_y','688_y','689_y','690_y','691_y','692_y','693_y','694_y','695_y','696_y','697_y','698_y','699_y','700_y','701_y','702_y','703_y','704_y','705_y','706_y','707_y','708_y','709_y','710_y','711_y','712_y','713_y','714_y','715_y','716_y','717_y','718_y','719_y','720_y','721_y','722_y','723_y','724_y','725_y','726_y','727_y','728_y','729_y','730_y','731_y','732_y','733_y','734_y','735_y','736_y','737_y','738_y','739_y','740_y','741_y','742_y','743_y','744_y','745_y','746_y','747_y','748_y','749_y','750_y','751_y','752_y','753_y','754_y','755_y','756_y','757_y','758_y','759_y','760_y','761_y','762_y','763_y','764_y','765_y','766_y','767_y','768_y','769_y','770_y','771_y','772_y','773_y','774_y','775_y','776_y','777_y','778_y','779_y','780_y','781_y','782_y','783_y','784_y','785_y','786_y','787_y','788_y','789_y','790_y','791_y','792_y','793_y','794_y','795_y','796_y','797_y','798_y','799_y','800_y','801_y','802_y','803_y','804_y','805_y','806_y','807_y','808_y','809_y','810_y','811_y','812_y','813_y','814_y','815_y','816_y','817_y','818_y','819_y','820_y','821_y','822_y','823_y','824_y','825_y','826_y','827_y','828_y','829_y','830_y','831_y','832_y','833_y','834_y','835_y','836_y','837_y','838_y','839_y','840_y','841_y','842_y','843_y','844_y','845_y','846_y','847_y','848_y','849_y','850_y','851_y','852_y','853_y','854_y','855_y','856_y','857_y','858_y','859_y','860_y','861_y','862_y','863_y','864_y','865_y','866_y','867_y','868_y','869_y','870_y','871_y','872_y','873_y','874_y','875_y','876_y','877_y','878_y','879_y','880_y','881_y','882_y','883_y','884_y','885_y','886_y','887_y','888_y','889_y','890_y','891_y','892_y','893_y','894_y','895_y','896_y','897_y','898_y','899_y','900_y','901_y','902_y','903_y','904_y','905_y','906_y','907_y','908_y','909_y','910_y','911_y','912_y','913_y','914_y','915_y','916_y','917_y','918_y','919_y','920_y','921_y','922_y','923_y','924_y','925_y','926_y','927_y','928_y','929_y','930_y','931_y','932_y','933_y','934_y','935_y','936_y','937_y','938_y','939_y','940_y','941_y','942_y','943_y','944_y','945_y','946_y','947_y','948_y','949_y','950_y','951_y','952_y','953_y','954_y','955_y','956_y','957_y','958_y','959_y','960_y','961_y','962_y','963_y','964_y','965_y','966_y','967_y','968_y','969_y','970_y','971_y','972_y','973_y','974_y','975_y','976_y','977_y','978_y','979_y','980_y','981_y','982_y','983_y','984_y','985_y','986_y','987_y','988_y','989_y','990_y','991_y','992_y','993_y','994_y','995_y','996_y','997_y','998_y','999_y','1000_y']
    df.to_sql('train', disk_engine, if_exists='append', chunksize=chunksize)
```



```

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

```

In [3]:

```

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

    return None

def 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 [4]:

```

import sqlite3
read_db = 'train.db'
conn_r = create_connection(read_db)
checkTableExists(conn_r)
conn_r.close()

```

Tables in the database:
data

In [5]:

```

if os.path.isfile(read_db):
    conn_r = create_connection(read_db)
    if conn_r is not None:
        data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 100001;", conn_r)
        conn_r.commit()
        conn_r.close()

```

In [6]:

```
# 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 [7]:

```
data.head()
```

Out[7]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq
1	0.285710204139941	0.153844970423304	0.499987500312492	0.133332444445037	0.363633057881292	0.105262880887156	0.0
2	0.999985714489793	0.874989062636717	0.999980000399992	0.999980000399992	0.99999166673611	0.923069822539827	1.0
3	0.499987500312492	0.333327777870369	0.33332222259258	0.14285510206997	0.428565306209911	0.230767455634957	0.0
4	0.0	0.0	0.428565306209911	0.29999700003	0.230767455634957	0.107142474491163	0.0
5	0.66664444518516	0.499987500312492	0.599988000239995	0.499991666805553	0.555549382784636	0.454541322351615	1.0

5 rows × 794 columns



4.2 Converting strings to numerics

In [8]:

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

```
cwc_min
cwc_max
csc_min
csc_max
ctc_min
ctc_max
last_word_eq
first_word_eq
abs_len_diff
mean_len
token_set_ratio
token_sort_ratio
fuzz_ratio
fuzz_partial_ratio
longest_substr_ratio
freq_qid1
freq_qid2
q1len
q2len
q1_n_words
q2_n_words
word_Common
word_Total
word_share
freq_q1+q2
freq_q1-q2
0_x
1_x
2_x
3_x
4_x
_
```

5_x
6_x
7_x
8_x
9_x
10_x
11_x
12_x
13_x
14_x
15_x
16_x
17_x
18_x
19_x
20_x
21_x
22_x
23_x
24_x
25_x
26_x
27_x
28_x
29_x
30_x
31_x
32_x
33_x
34_x
35_x
36_x
37_x
38_x
39_x
40_x
41_x
42_x
43_x
44_x
45_x
46_x
47_x
48_x
49_x
50_x
51_x
52_x
53_x
54_x
55_x
56_x
57_x
58_x
59_x
60_x
61_x
62_x
63_x
64_x
65_x
66_x
67_x
68_x
69_x
70_x
71_x
72_x
73_x
74_x
75_x
76_x
77_x
78_x
79_x
80_x
81_x

82_x
83_x
84_x
85_x
86_x
87_x
88_x
89_x
90_x
91_x
92_x
93_x
94_x
95_x
96_x
97_x
98_x
99_x
100_x
101_x
102_x
103_x
104_x
105_x
106_x
107_x
108_x
109_x
110_x
111_x
112_x
113_x
114_x
115_x
116_x
117_x
118_x
119_x
120_x
121_x
122_x
123_x
124_x
125_x
126_x
127_x
128_x
129_x
130_x
131_x
132_x
133_x
134_x
135_x
136_x
137_x
138_x
139_x
140_x
141_x
142_x
143_x
144_x
145_x
146_x
147_x
148_x
149_x
150_x
151_x
152_x
153_x
154_x
155_x
156_x
157_x
158_x

159_x
160_x
161_x
162_x
163_x
164_x
165_x
166_x
167_x
168_x
169_x
170_x
171_x
172_x
173_x
174_x
175_x
176_x
177_x
178_x
179_x
180_x
181_x
182_x
183_x
184_x
185_x
186_x
187_x
188_x
189_x
190_x
191_x
192_x
193_x
194_x
195_x
196_x
197_x
198_x
199_x
200_x
201_x
202_x
203_x
204_x
205_x
206_x
207_x
208_x
209_x
210_x
211_x
212_x
213_x
214_x
215_x
216_x
217_x
218_x
219_x
220_x
221_x
222_x
223_x
224_x
225_x
226_x
227_x
228_x
229_x
230_x
231_x
232_x
233_x
234_x
235_x

236_x
237_x
238_x
239_x
240_x
241_x
242_x
243_x
244_x
245_x
246_x
247_x
248_x
249_x
250_x
251_x
252_x
253_x
254_x
255_x
256_x
257_x
258_x
259_x
260_x
261_x
262_x
263_x
264_x
265_x
266_x
267_x
268_x
269_x
270_x
271_x
272_x
273_x
274_x
275_x
276_x
277_x
278_x
279_x
280_x
281_x
282_x
283_x
284_x
285_x
286_x
287_x
288_x
289_x
290_x
291_x
292_x
293_x
294_x
295_x
296_x
297_x
298_x
299_x
300_x
301_x
302_x
303_x
304_x
305_x
306_x
307_x
308_x
309_x
310_x
311_x
312_x

313_x
314_x
315_x
316_x
317_x
318_x
319_x
320_x
321_x
322_x
323_x
324_x
325_x
326_x
327_x
328_x
329_x
330_x
331_x
332_x
333_x
334_x
335_x
336_x
337_x
338_x
339_x
340_x
341_x
342_x
343_x
344_x
345_x
346_x
347_x
348_x
349_x
350_x
351_x
352_x
353_x
354_x
355_x
356_x
357_x
358_x
359_x
360_x
361_x
362_x
363_x
364_x
365_x
366_x
367_x
368_x
369_x
370_x
371_x
372_x
373_x
374_x
375_x
376_x
377_x
378_x
379_x
380_x
381_x
382_x
383_x
0_y
1_y
2_y
3_y
4_y
5_y

6_y
7_y
8_y
9_y
10_y
11_y
12_y
13_y
14_y
15_y
16_y
17_y
18_y
19_y
20_y
21_y
22_y
23_y
24_y
25_y
26_y
27_y
28_y
29_y
30_y
31_y
32_y
33_y
34_y
35_y
36_y
37_y
38_y
39_y
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 v

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 v

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

In [9]:

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

4.3 Random train test split(70:30)

In [11]:

```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(data, y_true, stratify=y_true, test_size=0.3)
```

In [12]:

```
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 : (70000, 794)

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

In [14]:

```
from collections import Counter, defaultdict
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 train 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:  0.6285714285714286 Class 1:  0.37142857142857144
----- Distribution of output variable in train data -----
Class 0:  0.37143333333333334 Class 1:  0.37143333333333334
```

In [15]:

```
# 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")
```

```

sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=labels)
plt.xlabel('Predicted Class')
plt.ylabel('Original Class')
plt.title("Confusion matrix")

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

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

plt.show()

```

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

In [18]:

```

from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.metrics import confusion_matrix

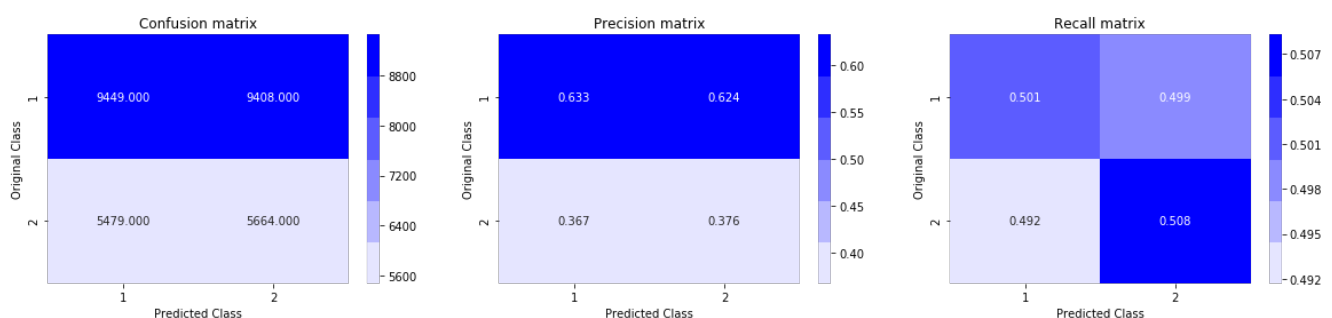
# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their sum
# ref: https://stackoverflow.com/a/18662466/4084039
# we create a output array that has exactly same size as the CV data

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))

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

```

Log loss on Test Data using Random Model 0.8812910323955873



4.4 Logistic Regression with hyperparameter tuning

In [21]:

```

from sklearn.linear_model import SGDClassifier
from sklearn.calibration import CalibratedClassifierCV

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=1000)

```

```

SGDClassifier(alpha=1e-05, 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.cl
asses_, eps=1e-15))

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

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

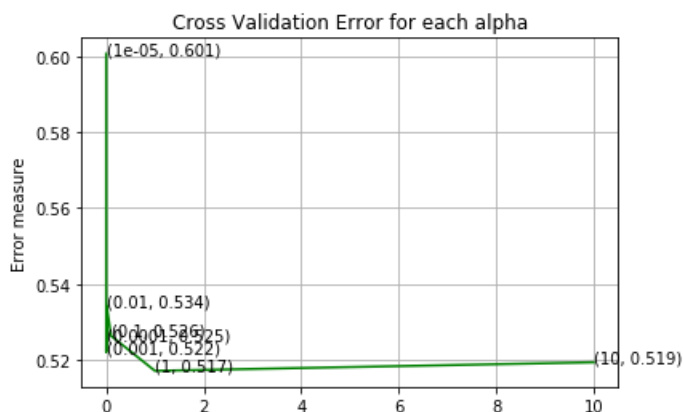
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss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

```

```

For values of alpha = 1e-05 The log loss is: 0.6006655170844521
For values of alpha = 0.0001 The log loss is: 0.5251853551818262
For values of alpha = 0.001 The log loss is: 0.522073497678028
For values of alpha = 0.01 The log loss is: 0.5342984841541478
For values of alpha = 0.1 The log loss is: 0.526364477366014
For values of alpha = 1 The log loss is: 0.5171410985697541
For values of alpha = 10 The log loss is: 0.5193927268268487

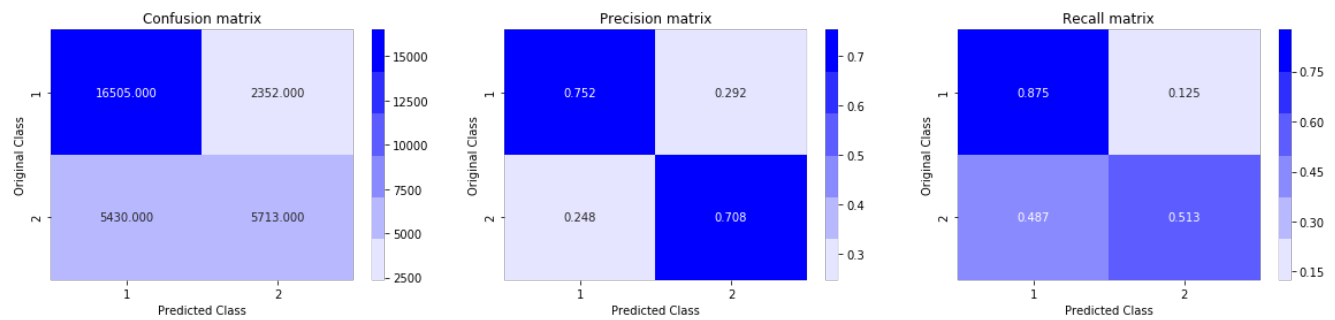
```



For values of best alpha = 1 The train log loss is: 0.5089970061612473

For values of best alpha = 1 The test log loss is: 0.5171410985697541

Total number of data points : 30000



4.5 Linear SVM with hyperparameter tuning

In [22]:

```
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.

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))
```

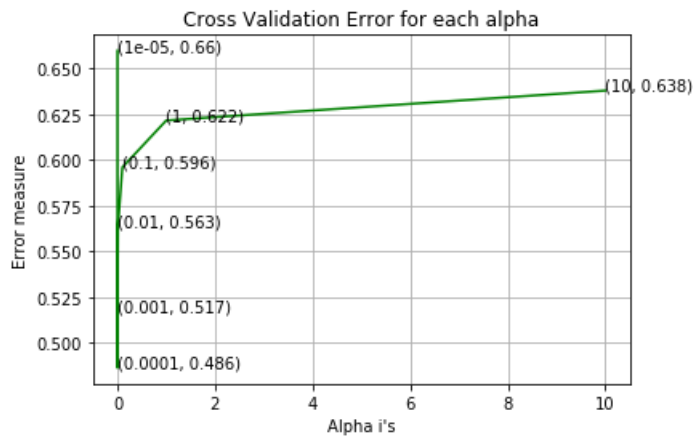


```

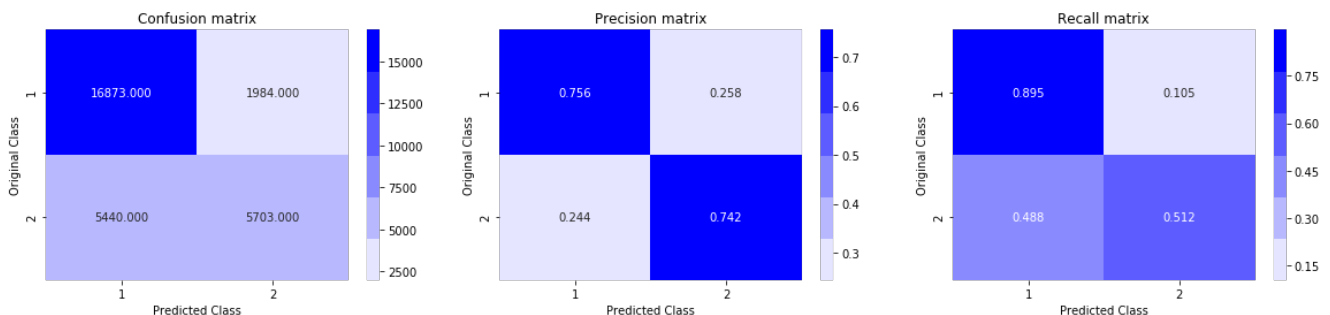
predict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y,axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)

```

For values of alpha = 1e-05 The log loss is: 0.6597141204189816
 For values of alpha = 0.0001 The log loss is: 0.48614706251166956
 For values of alpha = 0.001 The log loss is: 0.5172527985436652
 For values of alpha = 0.01 The log loss is: 0.5633373620425512
 For values of alpha = 0.1 The log loss is: 0.5958782537226457
 For values of alpha = 1 The log loss is: 0.6215492842471588
 For values of alpha = 10 The log loss is: 0.637883887273433



For values of best alpha = 0.0001 The train log loss is: 0.47763640515097006
 For values of best alpha = 0.0001 The test log loss is: 0.48614706251166956
 Total number of data points : 30000



4.6 XGBoost

In [24]:

```

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

d_train = xgb.DMatrix(X_train, label=y_train)
d_test = xgb.DMatrix(X_test, label=y_test)

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

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

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

```

[17:18:50] /workspace/src/tree/updater_prune.cc:74: tree pruning end, 1 roots, 30 extra nodes, 0 pruned nodes, max_depth=4
 [0] train-logloss:0.68495 valid-logloss:0.684973

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.

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

[illegible]

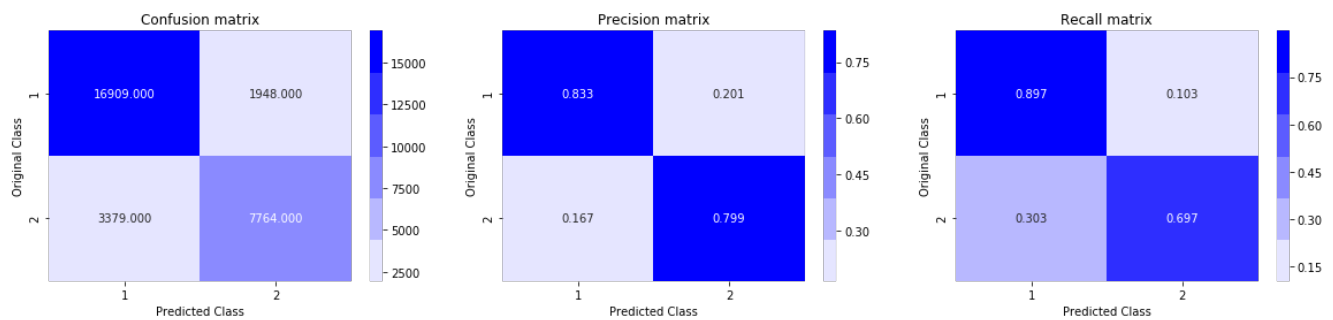
[illegible]

```

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

```

Total number of data points : 30000



5.1 XGBoost with hyperparameter tuning

In [35]:

```

from xgboost import XGBClassifier
from sklearn.model_selection import RandomizedSearchCV

params = {'n_estimators': [10,20,50,100,500], 'max_depth':[1, 5, 10, 15, 20]}

clf = XGBClassifier(objective='binary:logistic', eval_metric='logloss', n_jobs=-1)
model = RandomizedSearchCV(clf, params, cv=3, scoring='neg_log_loss', n_jobs=-1)
model.fit(X_train,y_train)

#Optimal value of depth
optimal_depth = model.best_estimator_.max_depth
print("\nThe optimal value of depth is : ",optimal_depth)

#Optimal value of number of estimators
optimal_estimators = model.best_estimator_.n_estimators
print("\nThe optimal value of depth is : ",optimal_estimators)

print('Log loss is :', model.best_score_)

```

The optimal value of depth is : 5

The optimal value of depth is : 500

Log loss is : -0.3337186777321428

Correction: The optimal value of depth is : 5 and The optimal value of estimators* are : 500

5.1.1 XGBoost with optimal hyperparameters

In [37]:

```

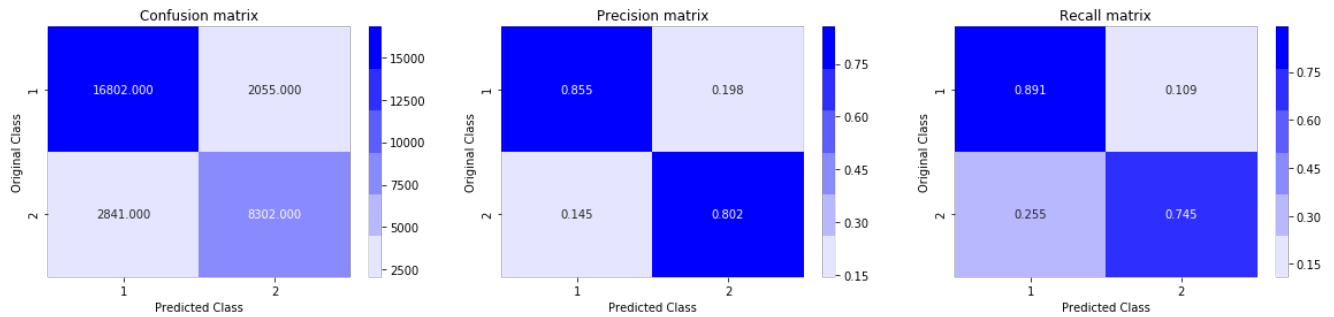
clf = XGBClassifier(max_depth=5,objective='binary:logistic',eval_metric='logloss',
n_estimators=500, n_jobs=-1)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train, y_train)
predict_y = sig_clf.predict_proba(X_train)
print("Log loss for tfidf train data : ",log_loss(y_train, predict_y, eps=1e-15))

predict_y = sig_clf.predict_proba(X_test)
print("Log loss for tfidf test data : ",log_loss(y_test, predict_y, eps=1e-15))
predicted_y = np.argmax(predict_y,axis=1)
plot_confusion_matrix(y_test, predicted_y)

```

Log loss for tfidf train data : 0.1867460755323528

Log loss for tfidf test data: 0.343841158138813



5.2 TFIDF Vectorization

In [52]:

```
df = pd.read_csv("train.csv", usecols= ['id','question1','question2','is_duplicate'])

# encode questions to unicode
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

In [53]:

```
print(df.shape)
df.head()
```

(404290, 4)

Out[53]:

	id	question1	question2	is_duplicate
0	0	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	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0
2	2	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0
3	3	Why am I mentally very lonely? How can I solve...	Find the remainder when $23^{24}/\text{math}$ i...	0
4	4	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0

In [54]:

```
col_names1 = ['id','cwc_min','cwc_max','csc_min','csc_max','ctc_min','ctc_max','last_word_eq','fir
st_word_eq',
              'abs_len_diff','mean_len','token_set_ratio','token_sort_ratio','fuzz_ratio','fuzz_par
tial_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']
```

```
data = pd.read_csv('final_features.csv',usecols=col_names1)
```

In [55]:

```
print("Rows: ",data.shape[0])
print("Dimensions: ",data.shape[1])
data.head()
```

Rows: 404290
Dimensions: 27

Out[55]:

	id	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	...	freq_qid2	q1len	q2len
	id	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	...	freq_qid2	q1len	q2len
0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	2.0	...	1	66	
1	1	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0	...	1	51	
2	2	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	4.0	...	1	73	
3	3	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	2.0	...	1	50	
4	4	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	6.0	...	1	76	

5 rows × 27 columns

In [56]:

```
df = df.merge(data, on='id', how='inner')
print(df.shape)
df.head()
```

(404290, 30)

Out[56]:

	id	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	...	freq_qid2	q1len	q2len
0	0	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	0.916659	0.785709	...	1	66	57
1	1	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	0.699993	0.466664	...	1	51	88
2	2	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	...	1	73	59
3	3	Why am I mentally very lonely? How can I solve...	Find the remainder when 23^{24} $[/math> [/math> i...$	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	...	1	50	65
4	4	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	...	1	76	39

5 rows × 30 columns

In [57]:

```
#Taking a random sample of 100k data points

sample = df.sample(n=100000, random_state = 42)
sample.shape
```

Out[57]:

(100000, 30)

In [58]:

```
#Taking the class lables from the data fram
y_true = sample['is_duplicate']
y_true.shape[0]
```

```
y_true.shape[0]
```

```
Out[58]:
```

```
100000
```

5.2.1 Random train test split(70:30)

```
In [59]:
```

```
from sklearn.model_selection import train_test_split
X_train,X_test, y_train, y_test = train_test_split(sample, y_true, stratify=y_true, test_size=0.3)
```

```
In [60]:
```

```
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 : (70000, 30)
```

```
Number of data points in test data : (30000, 30)
```

```
In [61]:
```

```
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 train data", "-"*10)
test_distr = Counter(y_test)
test_len = len(y_test)
print("Class 0: ",int(test_distr[1])/test_len, "Class 1: ",int(test_distr[1])/test_len)
```

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

```
Class 0:  0.6287857142857143 Class 1:  0.3712142857142857
```

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

```
Class 0:  0.3712 Class 1:  0.3712
```

5.2.2 Featurizing text data with tfidf

```
In [82]:
```

```
vect = TfidfVectorizer(lowercase=False)
train_q1 = vect.fit_transform(X_train['question1'])
test_q1  = vect.transform(X_test['question1'])
train_q2 = vect.fit_transform(X_train['question2'])
test_q2  = vect.transform(X_test['question2'])
```

```
In [83]:
```

```
print(train_q1.shape)
print(test_q1.shape)
print("-"*15)
print(train_q2.shape)
print(test_q2.shape)
```

```
(70000, 37120)
```

```
(30000, 37120)
```

```
*****
```

```
(70000, 34496)
```

```
(30000, 34496)
```

```
In [84]:
```

```
from scipy.sparse import hstack
Train_data = hstack((train_q1,train_q2))
Test_data = hstack((test_q1,test_q2))
```



```
print(Train_data.shape)
print(Test_data.shape)
```

```
(70000, 71616)
(30000, 71616)
```

In [87]:

```
#Extraction 26 features from train and test data frame
train_df = X_train.drop(['id','question1','question2','is_duplicate'], axis=1, inplace=False)
test_df = X_test.drop(['id','question1','question2','is_duplicate'], axis=1, inplace=False)
print("train Shape:",train_df.shape)
print("test shape:",test_df.shape)
```

```
train Shape: (70000, 26)
test shape: (30000, 26)
```

In [89]:

```
X_train_tfidf = hstack((train_df,Train_data))
X_test_tfidf = hstack((test_df,Test_data))

print(X_train_tfidf.shape)
print(X_test_tfidf.shape)
```

```
(70000, 71642)
(30000, 71642)
```

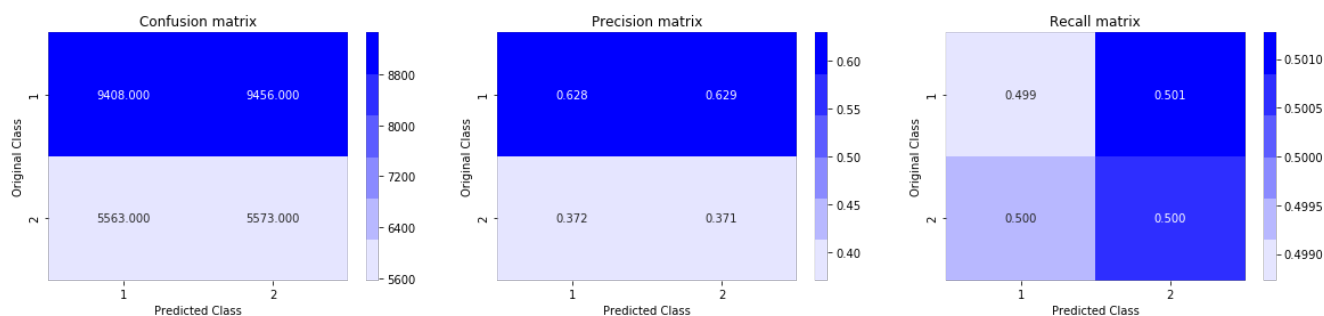
5.3 Building a random model (Finding worst-case log-loss)

In [97]:

```
predicted_y = np.zeros((test_len,2))
for i in range(test_len):
    rand_probs = np.random.rand(1,2)
    predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-15))

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

Log loss on Test Data using Random Model 0.8849848790191922



5.3.1 Logistic Regression with hyperparameter tuning

In [99]:

```
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_tfidf, y_train)
```

```

sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train_tfidf, y_train)
predict_y = sig_clf.predict_proba(X_test_tfidf)
log_error_array.append(log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
print('For values of alpha = ', i, "The log loss is:", log_loss(y_test, predict_y, labels=clf.cl
asses_, eps=1e-15))

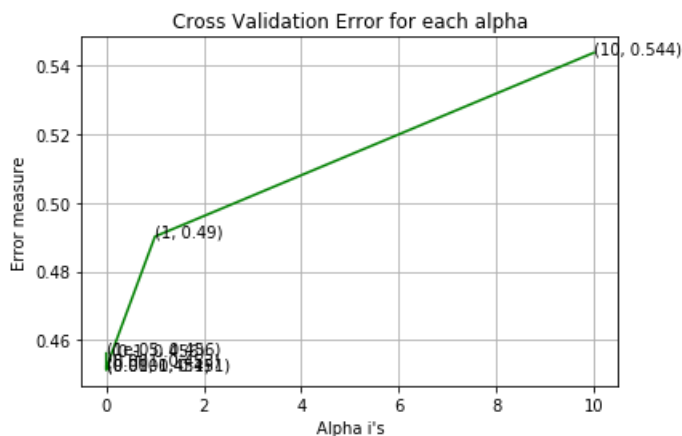
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

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

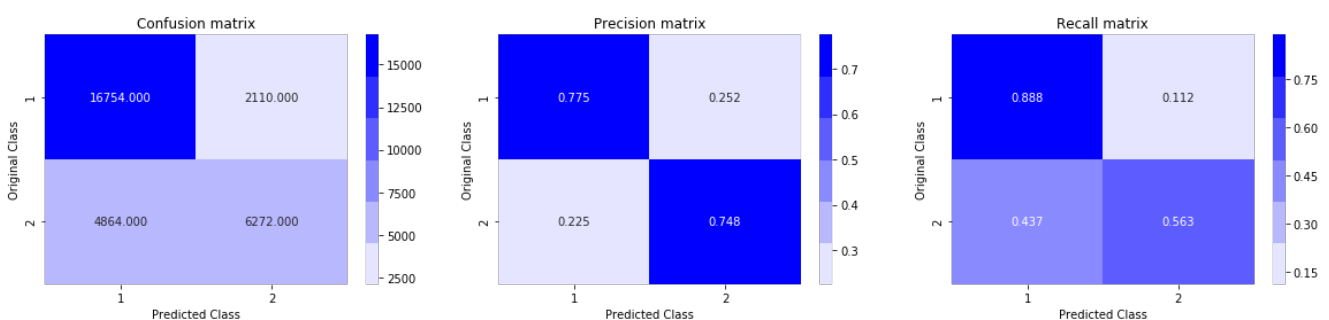
predict_y = sig_clf.predict_proba(X_train_tfidf)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test_tfidf)
print('For values of best alpha = ', alpha[best_alpha], "The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

```

For values of alpha = 1e-05 The log loss is: 0.4561824066110501
 For values of alpha = 0.0001 The log loss is: 0.4513461331823313
 For values of alpha = 0.001 The log loss is: 0.45264151688397863
 For values of alpha = 0.01 The log loss is: 0.4514760278903491
 For values of alpha = 0.1 The log loss is: 0.45577846438226766
 For values of alpha = 1 The log loss is: 0.4902087836358565
 For values of alpha = 10 The log loss is: 0.5437718913436526



For values of best alpha = 0.0001 The train log loss is: 0.45488707473791223
 For values of best alpha = 0.0001 The test log loss is: 0.4513461331823313
 Total number of data points : 30000



5.3.2 Linear SVM with hyperparameter tuning

In [101]:

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

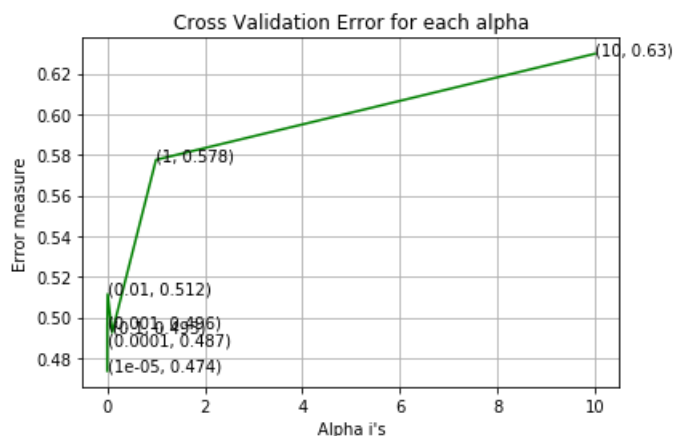
log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
    clf.fit(X_train_tfidf, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train_tfidf, y_train)
    predict_y = sig_clf.predict_proba(X_test_tfidf)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

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

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

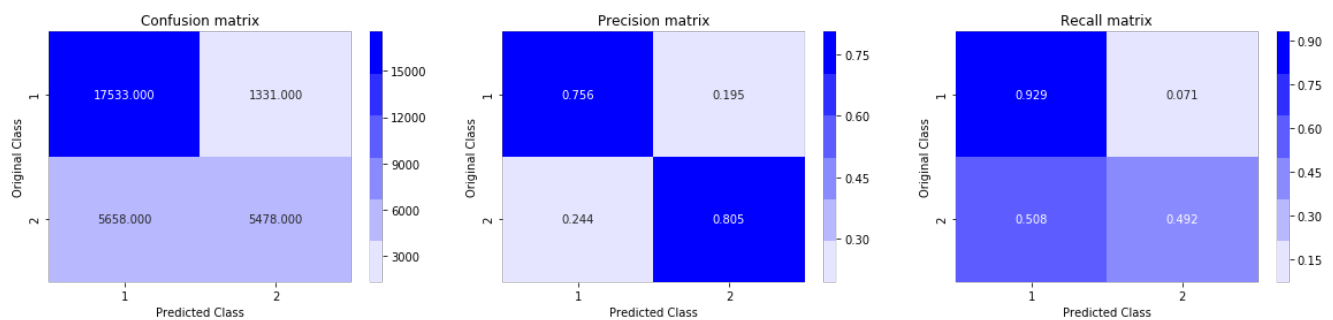
predict_y = sig_clf.predict_proba(X_train_tfidf)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test_tfidf)
print('For values of best alpha = ', alpha[best_alpha], "The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

```
For values of alpha = 1e-05 The log loss is: 0.4737451413724052
For values of alpha = 0.0001 The log loss is: 0.4868747392634236
For values of alpha = 0.001 The log loss is: 0.495560071960926
For values of alpha = 0.01 The log loss is: 0.5115242943592478
For values of alpha = 0.1 The log loss is: 0.49258573926068977
For values of alpha = 1 The log loss is: 0.5775389917454128
For values of alpha = 10 The log loss is: 0.6297178940309919
```



```
For values of best alpha = 1e-05 The train log loss is: 0.47827727803155495
For values of best alpha = 1e-05 The test log loss is: 0.4737451413724052
```

Total number of data points : 30000



6. Conclusion:

1. Initially, after importing the data we've performed the EDA to understand what all the features are important.
2. Before pre-processing the data, we've performed some basic feature engineering and extracted 11 additional features and analysed the features using violin plots, PDF plots to understand the how the features have impact on data.
3. Now, we've performed the pre-processing steps like Removing html tags, Removing Punctuations, Performing stemming, Removing Stopwords, Expanding contractions etc., on the data.
4. In the next step, we've performed advanced feature engineering using NLP and fuzzy features where we've extracted 15 additional features.
5. We've analysed these features using word clouds, pair plots, violin plots and PDF plots to understand the behaviour of data with specific features.
6. Using all these features, we have visualized the data using T-SNE in both 2-D and 3-D.
7. We now have featurized the data with TF-IDF weighted word2vec and TF-IDF and trained different machine learning models like Logistic regression, Linear SVM and XG Boost on top of it using log-loss as our metric.

In [104]:

```
from prettytable import PrettyTable
x = PrettyTable()

x.field_names = ["Model", "Vectorizer", "Train Loss", "Test Loss"]

x.add_row(['Logistic Regression', 'TF-IDF weighted W2V', 0.508, 0.517])
x.add_row(['Linear SVM', 'TF-IDF weighted W2V', 0.477, 0.486])
x.add_row(['XGBoost', 'TF-IDF weighted W2V', 0.344, 0.357])
x.add_row(['XGBoost(hyperparameter tuning)', 'TF-IDF weighted W2V', 0.186, 0.343])
x.add_row(['Logistic Regression', 'TF-IDF', 0.454, 0.451])
x.add_row(['Linear SVM', 'TF-IDF', 0.478, 0.473])

print(x)
```

Model	Vectorizer	Train Loss	Test Loss
Logistic Regression	TF-IDF weighted W2V	0.508	0.517
Linear SVM	TF-IDF weighted W2V	0.477	0.486
XGBoost	TF-IDF weighted W2V	0.344	0.357
XGBoost(hyperparameter tuning)	TF-IDF weighted W2V	0.186	0.343
Logistic Regression	TF-IDF	0.454	0.451
Linear SVM	TF-IDF	0.478	0.473

In []: