

1. Try out models (Logistic regression, Linear-SVM) with simple TF-IDF vectors instead of TD_IDF weighted word2Vec.
2. Perform hyperparameter tuning of XgBoost models using RandomsearchCV with vectorizer as TF-IDF W2V to reduce the log-loss.

In [240]:

```
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 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

import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
warnings.filterwarnings("ignore")
import sys
import os
```

```
import pandas as pd
import numpy as np
from tqdm import tqdm

# extract word2vec vectors
# https://github.com/explosion/spaCy/issues/1721
# http://landinghub.visualstudio.com/visual-cpp-build-tools
import spacy
```

In [241]:

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

from sklearn.model_selection import cross_val_score
from sklearn.linear_model import SGDClassifier
# from mlxtend.classifier import StackingClassifier

from sklearn import model_selection
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import precision_recall_curve, auc, roc_curve
```

In [242]:

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

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

Number of data points: 404290

In [243]:

```
df.head()
```

Out[243]:

id	qid1	qid2	question1	question2	is_duplicate
----	------	------	-----------	-----------	--------------

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

In [244]:

```
df.info()
```

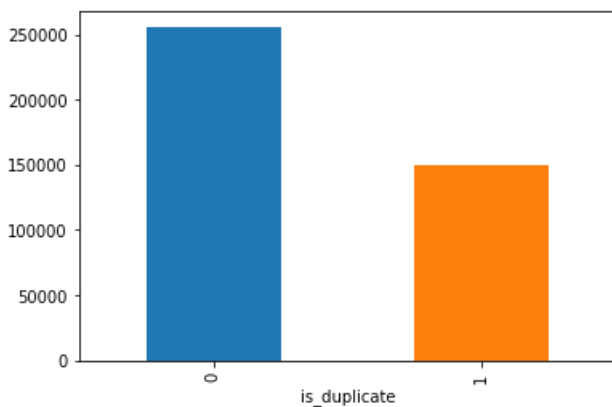
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 404290 entries, 0 to 404289
Data columns (total 6 columns):
id                404290 non-null int64
qid1              404290 non-null int64
qid2              404290 non-null int64
question1         404289 non-null object
question2         404288 non-null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

In [245]:

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

Out[245]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x18ee1d7ab00>
```



In [246]:

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

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

In [247]:

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

```
~> Question pairs that are not Similar (is_duplicate = 0):
63.08%
```

```
~> Question pairs that are Similar (is_duplicate = 1):  
36.92%
```

Number of unique questions

In [248]:

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

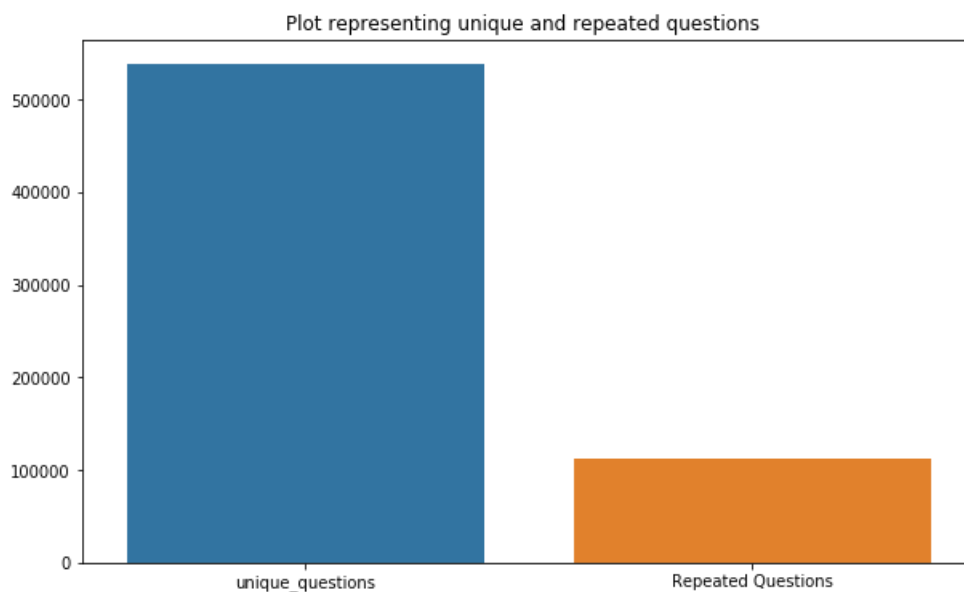
Total number of Unique Questions are: 537933

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

Max number of times a single question is repeated: 157

In [249]:

```
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()
```



Checking for Duplicates

In [250]:

```
#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

Number of occurrences of each question

In [251]:

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

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

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

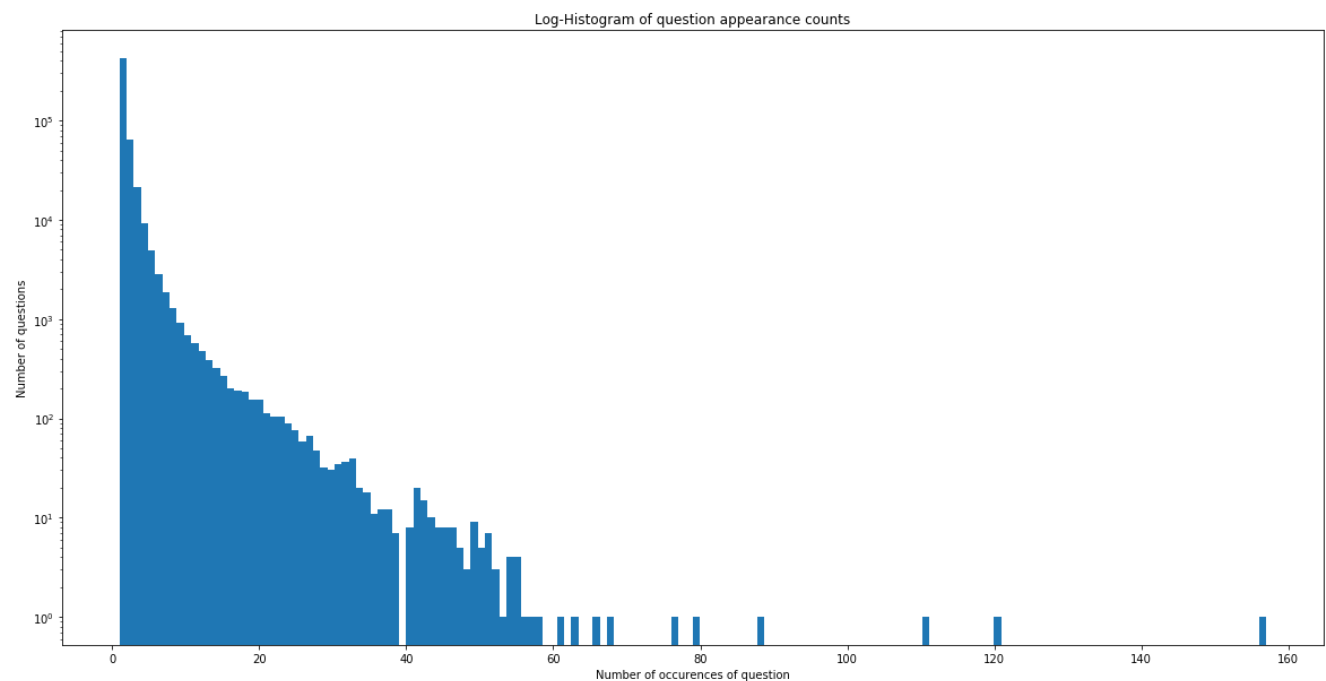
plt.title('Log-Histogram of question appearance counts')

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



Checking for NULL values

In [252]:

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

	id	qid1	qid2	question1 \	question2	is_duplicate
105780	105780	174363	174364	How can I develop android app?	NaN	0
201841	201841	303951	174364	How can I create an Android app?	NaN	0
363362	363362	493340	493341	NaN		

In [253]:

```
# 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: []

Basic Feature Extraction (before cleaning)

Let us now construct a few features like:

- **freq_qid1** = Frequency of qid1's
- **freq_qid2** = Frequency of qid2's
- **q1len** = Length of q1
- **q2len** = Length of q2
- **q1_n_words** = Number of words in Question 1
- **q2_n_words** = Number of words in Question 2
- **word_Common** = (Number of common unique words in Question 1 and Question 2)
- **word_Total** = (Total num of words in Question 1 + Total num of words in Question 2)
- **word_share** = (word_common)/(word_Total)
- **freq_q1+freq_q2** = sum total of frequency of qid1 and qid2
- **freq_q1-freq_q2** = absolute difference of frequency of qid1 and qid2

In [254]:

```

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

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

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

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

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

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

df.head()

```

Out [254] :

[illegible]

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word
	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0	1	1	73	59	14	10	4.0
3	3	7	8	Why am I mentally very lonely? How can I solve...	Find the remainder when 23^{24} $[/math] i...$	0	1	1	50	65	11	9	0.0
4	4	9	10	Which one dissolve in water quikly sugar, salt...	Which fish would survive in salt water?	0	3	1	76	39	13	7	2.0

Analysis of some of the extracted features

- Here are some questions have only one single words.

In [255]:

```
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
```

Feature: word_share

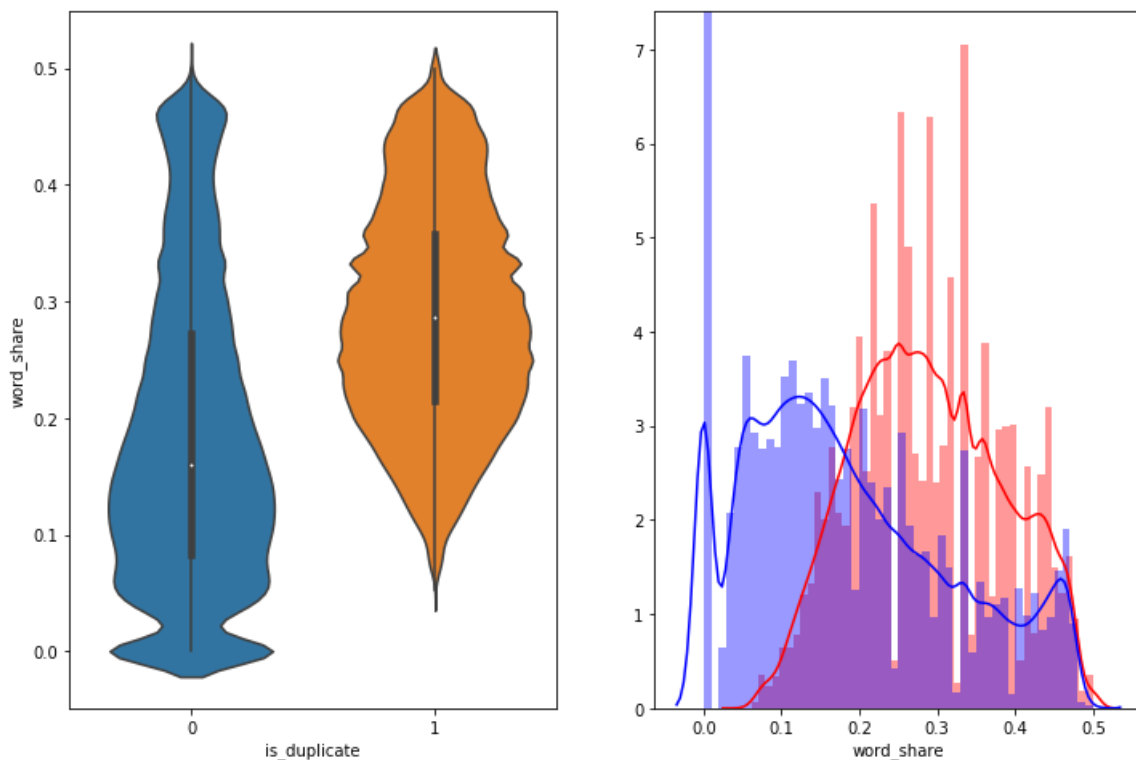
In [256]:

```
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')
```

```
sns.distplot(df[df['is_duplicate'] == 0.0]['word_share'][0:], label = '0', color = 'blue',
plt.show())
```



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

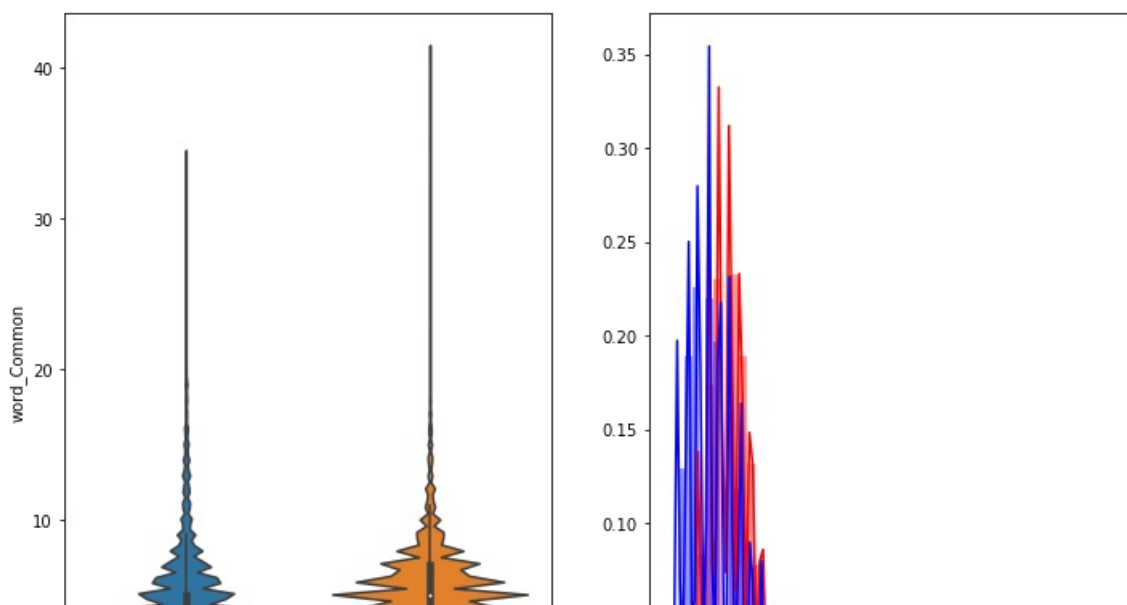
Feature: word_Common

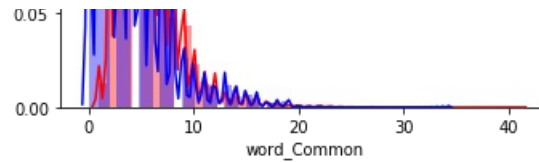
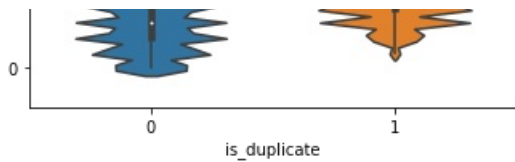
In [257]:

```
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()
```





EDA: Advanced Feature Extraction.

In [258]:

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

In [259]:

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

Out[259]:

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

Preprocessing of Text

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

In [260]:

```
# 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")
```

```

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

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

$$\text{mean_len} = (\text{len}(q1_tokens) + \text{len}(q2_tokens)) / 2$$

- **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

$$\text{longest_substr_ratio} = \text{len}(\text{longest common substring}) / (\min(\text{len}(q1_tokens), \text{len}(q2_tokens)))$$

In [261]:

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

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

    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
        return token_features
    # Get the non-stopwords in Questions
    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])

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

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

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

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

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

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

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

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

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

# get the Longest Common sub string

def get_longest_substr_ratio(a, b):
    strs = list(distance.lcs substrings(a, b))
    if len(strs) == 0:
        return 0
    else:
```

```

        return len(strs[0]) / (min(len(a), len(b)) + 1)

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

    print("token features...")

    # Merging Features with dataset

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

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

    #Computing Fuzzy Features and Merging with Dataset

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

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

```

In [262]:

```

df = extract_features(df)
df.to_csv("nlp_features_train.csv", index=False)

```

token features...
fuzzy features..

In [263]:

```

# 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)

```

Analysis of extracted features

Plotting Word clouds

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

In [264]:

```
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', encoding='utf-8')
```

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

In [265]:

```
# 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))
```

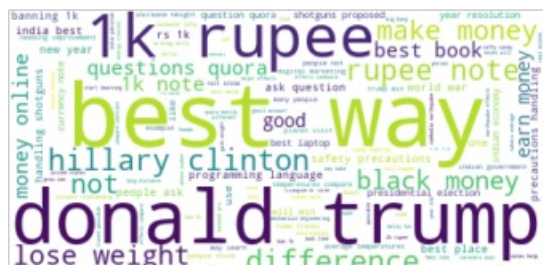
```
Total number of words in duplicate pair questions : 16102776
Total number of words in non duplicate pair questions : 33189630
```

Word Clouds generated from duplicate pair question's text

In [266]:

```
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

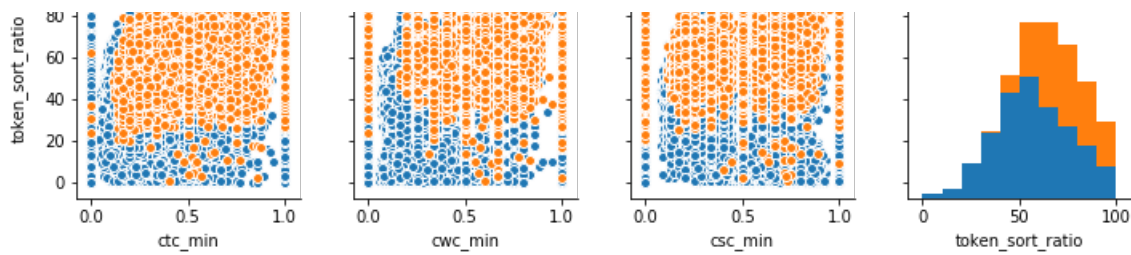


WORK <

```
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()
```



```
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()
```

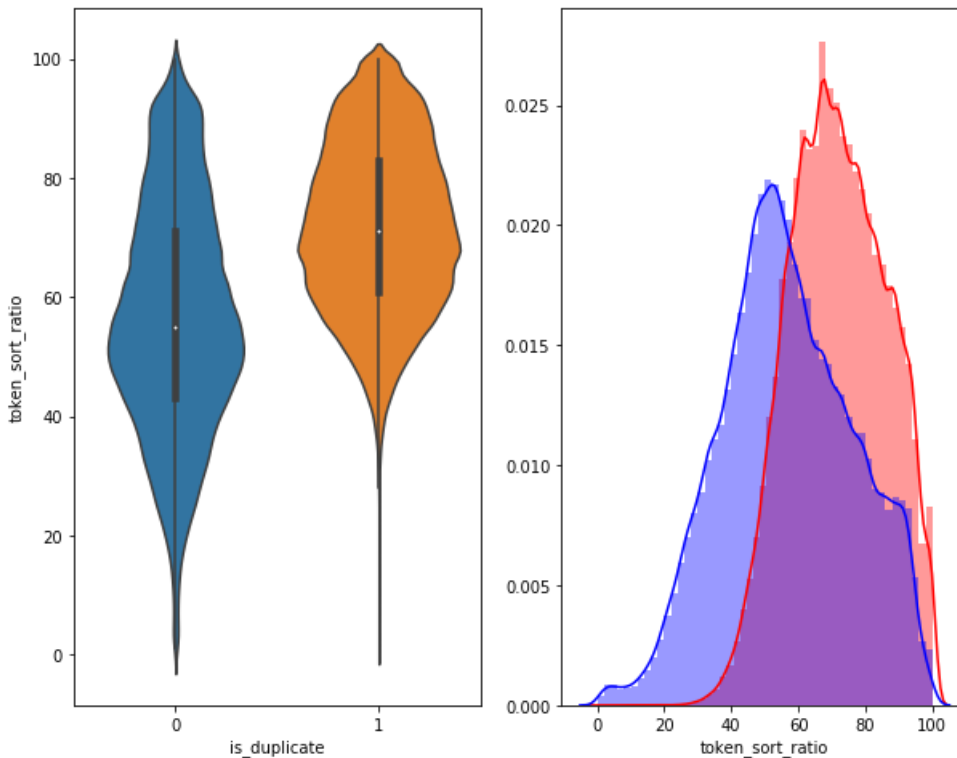


In [269]:

```
# 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()
```

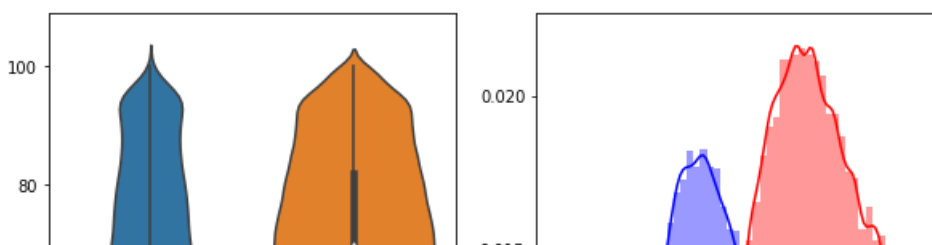


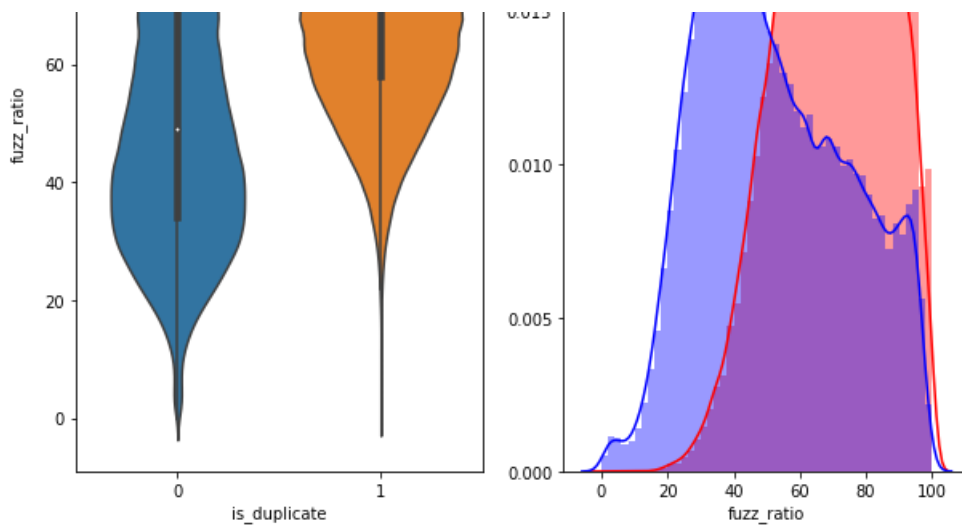
In [270]:

```
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()
```





Visualization

In [271]:

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

from sklearn.preprocessing import MinMaxScaler

dfp_subsampled = df[0:5000]
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 [272]:

```
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 0.014s...
[t-SNE] Computed neighbors for 5000 samples in 0.323s...
[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: 0.130355
[t-SNE] Computed conditional probabilities in 0.177s
[t-SNE] Iteration 50: error = 81.1602478, gradient norm = 0.0414487 (50 iterations in 5.846s)
[t-SNE] Iteration 100: error = 70.6033401, gradient norm = 0.0109017 (50 iterations in 4.538s)
[t-SNE] Iteration 150: error = 68.8870087, gradient norm = 0.0058326 (50 iterations in 4.383s)
[t-SNE] Iteration 200: error = 68.0940781, gradient norm = 0.0041826 (50 iterations in 4.523s)
[t-SNE] Iteration 250: error = 67.5874481, gradient norm = 0.0035499 (50 iterations in 4.630s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.587448
[t-SNE] Iteration 300: error = 1.7909305, gradient norm = 0.0011860 (50 iterations in 4.830s)
[t-SNE] Iteration 350: error = 1.3916727, gradient norm = 0.0004844 (50 iterations in 4.731s)
[t-SNE] Iteration 400: error = 1.2250323, gradient norm = 0.0002794 (50 iterations in 4.702s)
[t-SNE] Iteration 450: error = 1.1362009, gradient norm = 0.0001890 (50 iterations in 4.716s)
[t-SNE] Iteration 500: error = 1.0817227, gradient norm = 0.0001436 (50 iterations in 4.663s)
[t-SNE] Iteration 550: error = 1.0464581, gradient norm = 0.0001171 (50 iterations in 4.758s)
[t-SNE] Iteration 600: error = 1.0233538, gradient norm = 0.0001003 (50 iterations in 4.650s)
[t-SNE] Iteration 650: error = 1.0075148, gradient norm = 0.0000895 (50 iterations in 4.662s)
[t-SNE] Iteration 700: error = 0.9962875, gradient norm = 0.0000833 (50 iterations in 4.670s)
[t-SNE] Iteration 750: error = 0.9878469, gradient norm = 0.0000768 (50 iterations in 4.736s)
```

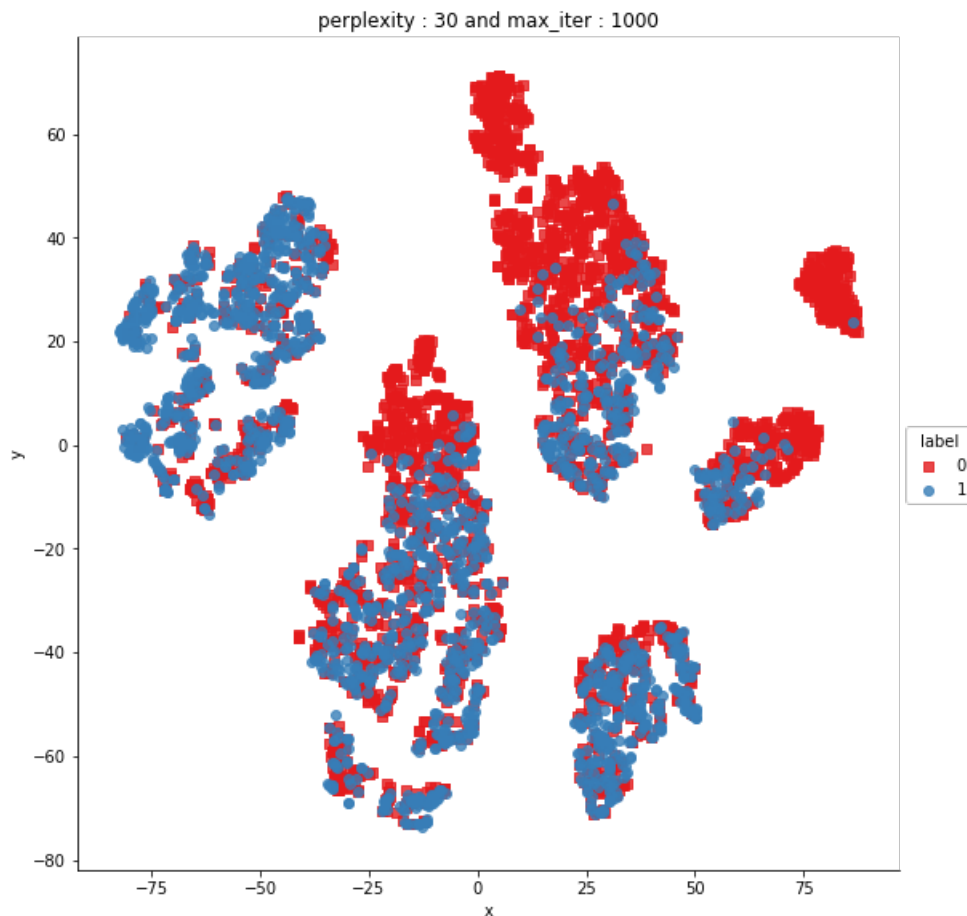


```
[t-SNE] Iteration 800: error = 0.9816954, gradient norm = 0.0000747 (50 iterations in 4.802s)
[t-SNE] Iteration 850: error = 0.9771429, gradient norm = 0.0000694 (50 iterations in 4.734s)
[t-SNE] Iteration 900: error = 0.9733443, gradient norm = 0.0000708 (50 iterations in 4.745s)
[t-SNE] Iteration 950: error = 0.9699866, gradient norm = 0.0000611 (50 iterations in 4.756s)
[t-SNE] Iteration 1000: error = 0.9667592, gradient norm = 0.0000576 (50 iterations in 4.899s)
[t-SNE] Error after 1000 iterations: 0.966759
```

In [273]:

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

# draw the plot in appropriate place in the grid
sns.lmplot(data=df_tsne, 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 [274]:

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

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.009s...
[t-SNE] Computed neighbors for 5000 samples in 0.329s...
[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: 0.130355
```

```

[t-SNE] Computed conditional probabilities in 0.169s
[t-SNE] Iteration 50: error = 80.5703201, gradient norm = 0.0313752 (50 iterations in 11.294s)
[t-SNE] Iteration 100: error = 69.3760147, gradient norm = 0.0032205 (50 iterations in 5.741s)
[t-SNE] Iteration 150: error = 67.9652710, gradient norm = 0.0016286 (50 iterations in 5.319s)
[t-SNE] Iteration 200: error = 67.4106369, gradient norm = 0.0011808 (50 iterations in 5.351s)
[t-SNE] Iteration 250: error = 67.0974884, gradient norm = 0.0008838 (50 iterations in 5.326s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.097488
[t-SNE] Iteration 300: error = 1.5195843, gradient norm = 0.0007023 (50 iterations in 6.660s)
[t-SNE] Iteration 350: error = 1.1824349, gradient norm = 0.0002111 (50 iterations in 8.097s)
[t-SNE] Iteration 400: error = 1.0394241, gradient norm = 0.0000970 (50 iterations in 8.178s)
[t-SNE] Iteration 450: error = 0.9668541, gradient norm = 0.0000614 (50 iterations in 8.179s)
[t-SNE] Iteration 500: error = 0.9275563, gradient norm = 0.0000503 (50 iterations in 8.048s)
[t-SNE] Iteration 550: error = 0.9070587, gradient norm = 0.0000445 (50 iterations in 8.018s)
[t-SNE] Iteration 600: error = 0.8952436, gradient norm = 0.0000391 (50 iterations in 7.959s)
[t-SNE] Iteration 650: error = 0.8866698, gradient norm = 0.0000381 (50 iterations in 7.970s)
[t-SNE] Iteration 700: error = 0.8803310, gradient norm = 0.0000383 (50 iterations in 7.857s)
[t-SNE] Iteration 750: error = 0.8761272, gradient norm = 0.0000379 (50 iterations in 7.730s)
[t-SNE] Iteration 800: error = 0.8730342, gradient norm = 0.0000364 (50 iterations in 7.823s)
[t-SNE] Iteration 850: error = 0.8696711, gradient norm = 0.0000338 (50 iterations in 7.918s)
[t-SNE] Iteration 900: error = 0.8673273, gradient norm = 0.0000354 (50 iterations in 7.912s)
[t-SNE] Iteration 950: error = 0.8650988, gradient norm = 0.0000233 (50 iterations in 7.880s)
[t-SNE] Iteration 1000: error = 0.8614780, gradient norm = 0.0000244 (50 iterations in 7.823s)
[t-SNE] Error after 1000 iterations: 0.861478

```

In [275]:

```

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

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

```

Featurizing text data with tfidf weighted word-vectors

Random train test split(70:30)

In [338]:

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

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

In [339]:

```
y_true = df['is_duplicate']
```

In [340]:

```
X_train, X_test, y_train, y_test = train_test_split(df, y_true, stratify=y_true, test_size=0.3)
```

In [341]:

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

```
Number of data points in train data : (283003, 32)
Number of data points in test data : (121287, 32)
```

In [342]:

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

```
----- Distribution of output variable in train data -----
Class 0:  0.6308025003268517 Class 1:  0.36919749967314835
----- Distribution of output variable in test data -----
Class 0:  0.6308013224830361 Class 1:  0.3691986775169639
```

[illegible]

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

```
# en vectors web lg, which includes over 1 million unique vectors.
```

```
nlp = spacy.load('en_core_web_sm')
```

100% | ██████████ 283003/283003

```
vecs2 = []
```

```
for qu2 in tqdm(list(X_train['question2'])):
```

[illegible]

```
# en vectors web la, which includes over 1 million unique vectors.
```

```
# en_vectors_web_sm, which includes over 1 million unique vectors.
nlp = spacy.load('en_core_web_sm')
# import en_core_web_sm
# nlp = en_core_web_sm.load()

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

X_test['q1_feats_m'] = list(vecs1)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 121287/121287  
[30:00<00:00, 67.36it/s]
```

In [347]:

```
vecs2 = []
for qu2 in tqdm(list(X_test['question2'])):
    doc2 = nlp(qu2)
    mean_vec2 = np.zeros([len(doc1), len(doc2[0].vector)])
    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)
X_test['q2_feats_m'] = list(vecs2)
```

```
100%|██████████████████████████████████████████████████████████████████████████| 121287/121287  
[29:40<00:00, 68.11it/s]
```

In [348]:

```
df_q1_train = pd.DataFrame(X_train.q1_feats_m.values.tolist(), index= X_train.index)
df_q2_train = pd.DataFrame(X_train.q2_feats_m.values.tolist(), index= X_train.index)
```

In [349]:

```
df_q1_test = pd.DataFrame(X_test.q1_feats_m.values.tolist(), index= X_test.index)
df_q2_test = pd.DataFrame(X_test.q2_feats_m.values.tolist(), index= X_test.index)
```

In [350]:

```
df_q1_train['id']=X_train['id']
df_q2_train['id']=X_train['id']
```

In [351]:

```
df_q1_test['id']=X_test['id']
df_q2_test['id']=X_test['id']
```

```
df_q2_test[10]-A_test[10]
```

In [352]:

```
df_train = X_train.drop(['q1_feats_m', 'q2_feats_m'], axis=1)
```

In [353]:

```
df_test = X_test.drop(['q1_feats_m', 'q2_feats_m'], axis=1)
```

In [354]:

```
df_train_ques = df_q1_train.merge(df_q2_train, on='id', how='left')
result_train = df_train.merge(df_train_ques, on='id', how='left')
result_train.head(2)
```

Out[354]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	...	86_y	87_y
0	35850	65451	65452	hypothetical scenarios if obama could run for...	would president obama win if he could run a th...	1	2	1	74	55	...	-19.050019	-5.335857
1	34327	62912	62913	what r d work is being done in capgemini for t...	what are the government initiatives provided t...	0	1	1	67	110	...	1.326836	-23.947595

2 rows × 224 columns

In [355]:

```
df_test_ques = df_q1_test.merge(df_q2_test, on='id', how='left')
result_test = df_test.merge(df_test_ques, on='id', how='left')
result_test.head(2)
```

Out[355]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	...	86_y	87_y
0	190454	289535	289536	what are some great business ideas for 2017	what are some great online business startup id...	0	1	1	44	50	...	-7.488704	-14.754112
1	67929	117479	117480	how do i know whether i am in depression or not	how should i know whether i am in depression o...	1	1	1	48	52	...	22.655746	25.330314

2 rows × 224 columns

In [356]:

```
result_train.to_csv("final_features_train_W2V.csv")
```

In [357]:

```
result_test.to_csv("final_features_test_W2V.csv")
```

Featurizing text data with tfidf

Random train test split(70:30)

In [300]:

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

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

In [301]:

```
df.head()
```

Out[301]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	...	ctc_max	last_word_eq	first_w
0	0	1	2	what is the step by step guide to invest in sh...	what is the step by step guide to invest in sh...	0	1	1	66	57	...	0.785709	0.0	1.0
1	1	3	4	what is the story of kohinoor koh i noor dia...	what would happen if the indian government sto...	0	4	1	51	88	...	0.466664	0.0	1.0
2	2	5	6	how can i increase the speed of my internet co...	how can internet speed be increased by hacking...	0	1	1	73	59	...	0.285712	0.0	1.0
3	3	7	8	why am i mentally very lonely how can i solve...	find the remainder when math 23 24 math i...	0	1	1	50	65	...	0.000000	0.0	0.0
4	4	9	10	which one dissolve in water quickly sugar salt...	which fish would survive in salt water	0	3	1	76	39	...	0.307690	0.0	1.0

5 rows × 32 columns

In [302]:

```
y_true = df['is_duplicate']
```

In [303]:

```
X_train, X_test, y_train, y_test = train_test_split(df, y_true, stratify=y_true, test_size=0.3)
```

In [304]:

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

Number of data points in train data : (283003, 32)

Number of data points in test data : (121287, 32)

In [305]:

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

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

Class 0: 0.6308025003268517 Class 1: 0.36919749967314835

----- Distribution of output variable in test data -----

Class 0: 0.6308013224830361 Class 1: 0.3691986775169639

In [306]:

```
tfidf_vectorizer1 = TfidfVectorizer(lowercase=False,max_features= 20000)
tfidf_train_qes1 = tfidf_vectorizer1.fit_transform(X_train['question1'])
tfidf_test_qes1 = tfidf_vectorizer1.transform(X_test['question1'])
```

In [307]:

```
print(tfidf_train_qes1.shape)
print(tfidf_test_qes1.shape)
```

(283003, 20000)

(121287, 20000)

In [308]:

```
tfidf_vectorizer2 = TfidfVectorizer(lowercase=False,max_features= 20000)
tfidf_train_qes2 = tfidf_vectorizer2.fit_transform(X_train['question2'])
tfidf_test_qes2 = tfidf_vectorizer2.transform(X_test['question2'])
```

In [309]:

```
print(tfidf_train_qes2.shape)
print(tfidf_test_qes2.shape)
```

(283003, 20000)

(121287, 20000)

In [310]:

```
tfidf_train_qes = hstack((tfidf_train_qes1,tfidf_train_qes2))
```

```
tfidf_test_qes = hstack((tfidf_test_qes1,tfidf_test_qes2))
```


In [311]:

```
print(tfidf_train_ques.shape)
print(tfidf_test_ques.shape)
```

```
(283003, 40000)
(121287, 40000)
```

In [312]:

```
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)
```

In [313]:

```
print(train_df.shape)
print(test_df.shape)
```

```
(283003, 28)
(121287, 28)
```

In [314]:

```
import scipy
#converting feature data into sparse matrix
train_sparse = scipy.sparse.csr_matrix(train_df)
test_sparse = scipy.sparse.csr_matrix(test_df)
```

In [315]:

```
# combining features and tfidf

tfidf_train_data = hstack((train_sparse,tfidf_train_ques))

tfidf_test_data = hstack((test_sparse,tfidf_test_ques))
```

In [316]:

```
print(tfidf_train_data.shape)
print(tfidf_test_data.shape)
```

```
(283003, 40028)
(121287, 40028)
```

In [317]:

```
print(y_train.shape)
print(y_test.shape)
```

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

In [318]:

```
X_train = tfidf_train_data
X_test = tfidf_test_data
```

Machine Learning Models

In [319]:

```
# 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 = confusion_matrix(y_test, predicted_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 corresponds 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 corresponds to columns and axis=1 corresponds to rows in two
dimensional array
# C.sum(axix =0) = [[4, 6]]
# (C/C.sum(axis=0)) = [[1/4, 2/6],
#                       [3/4, 4/6]]
plt.figure(figsize=(20,4))

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

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

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

plt.show()

```

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

In [320]:

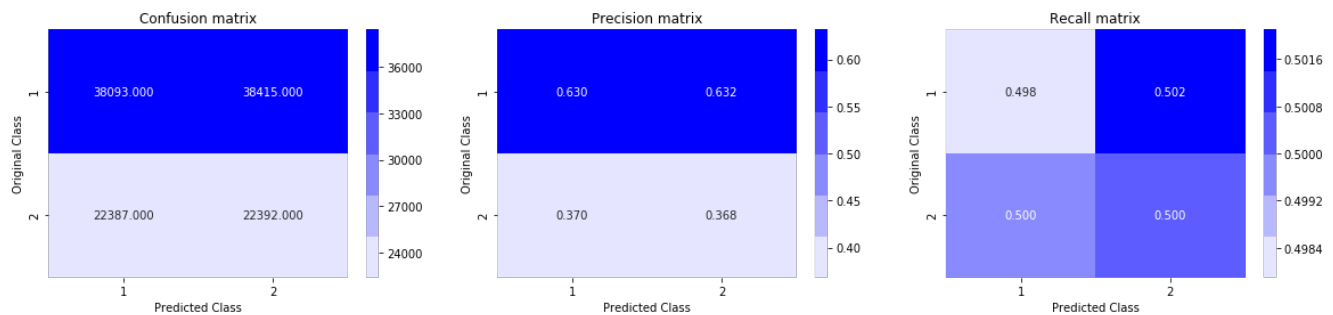
```

# 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 0.8870122114957221



Logistic Regression with hyperparameter tuning

In [321]:

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

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

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

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

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

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

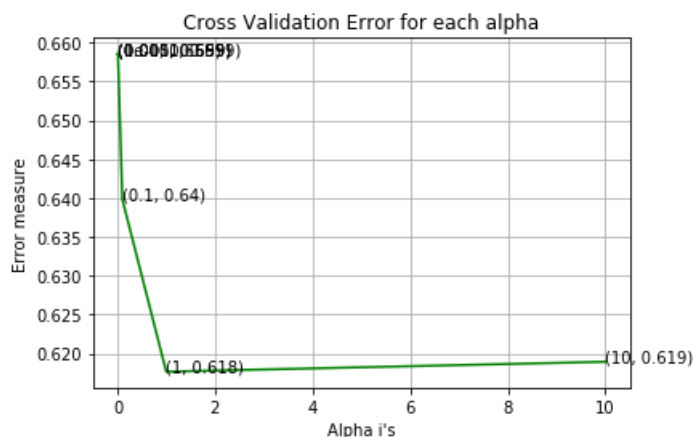
best_alpha = np.argmin(log_error_array)
clf = SGDClassifier(alpha=alpha[best_alpha], penalty='l2', loss='log', random_state=42)
clf.fit(X_train, y_train)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train, y_train)

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

```

For values of alpha = 1e-05 The log loss is: 0.6585278256347588
For values of alpha = 0.0001 The log loss is: 0.6585278256347588
For values of alpha = 0.001 The log loss is: 0.6585278256347588
For values of alpha = 0.01 The log loss is: 0.6585278256347588
For values of alpha = 0.1 The log loss is: 0.6398284021170975
For values of alpha = 1 The log loss is: 0.6176707615389561
For values of alpha = 10 The log loss is: 0.6189418092653223

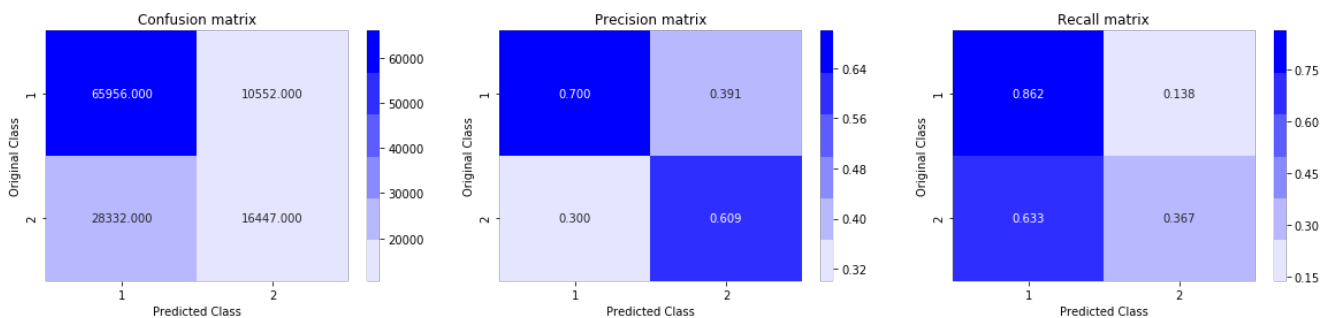
```



```

For values of best alpha = 1 The train log loss is: 0.616651777983376
For values of best alpha = 1 The test log loss is: 0.6176707615389561
Total number of data points : 121287

```



Linear SVM with hyperparameter tuning

In [322]:

```

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

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

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

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

log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
    clf.fit(X train, y train)

```

```

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

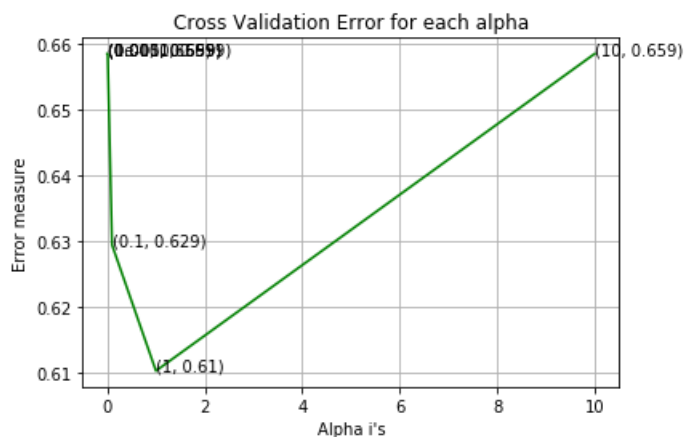
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

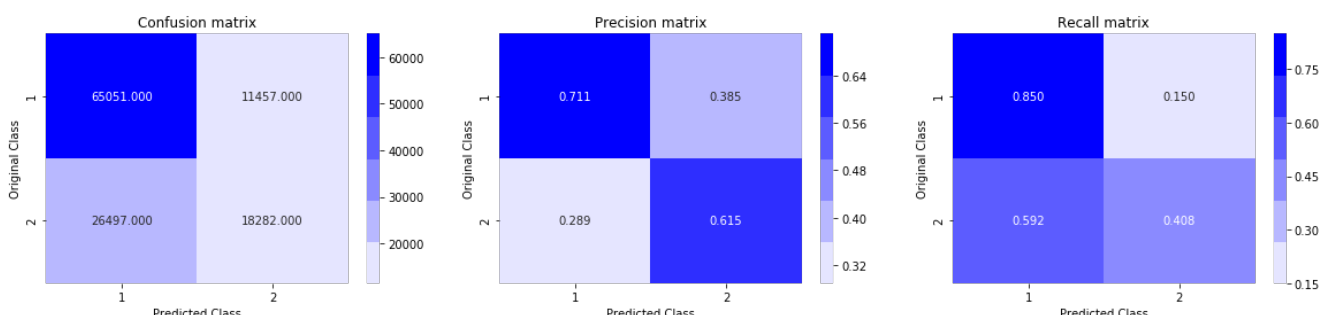
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss 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: 0.6585278256347588
 For values of alpha = 0.0001 The log loss is: 0.6585278256347588
 For values of alpha = 0.001 The log loss is: 0.6585278256347588
 For values of alpha = 0.01 The log loss is: 0.6585278256347588
 For values of alpha = 0.1 The log loss is: 0.6294340108585906
 For values of alpha = 1 The log loss is: 0.6102914391090667
 For values of alpha = 10 The log loss is: 0.6585278256347588



For values of best alpha = 1 The train log loss is: 0.6091544281531446
 For values of best alpha = 1 The test log loss is: 0.6102914391090667
 Total number of data points : 121287



XGBoost

In [369]:

```
result_train = pd.read_csv("final_features_train_W2V.csv")
result_test = pd.read_csv("final_features_test_W2V.csv")
```

In [370]:

```
X_train = result_train[0:15000]
X_test = result_test[0:5000]
```

In [371]:

```
y_train = X_train["is_duplicate"]
y_test = X_test["is_duplicate"]
```

In [372]:

```
X_train = X_train.drop(['id', 'question1', 'question2', 'is_duplicate'], axis=1, inplace=False)
X_test = X_test.drop(['id', 'question1', 'question2', 'is_duplicate'], axis=1, inplace=False)
```

In [375]:

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

```
(15000, 221)
(5000, 221)
```

In [376]:

```
from sklearn.model_selection import RandomizedSearchCV
from xgboost import XGBClassifier
from sklearn.metrics import log_loss

params = {'n_estimators' : [5, 10, 100, 500],
          'max_depth' : [2, 6, 8, 10]}

x_model = RandomizedSearchCV(estimator = XGBClassifier(objective = 'binary:logistic', eval_metric =
'logloss', eta = 0.02),
                             param_distributions = params)

# fit train sets
x_model.fit(X_train, y_train)

# Prediction
predict_y = x_model.predict_proba(X_test)
```

In [377]:

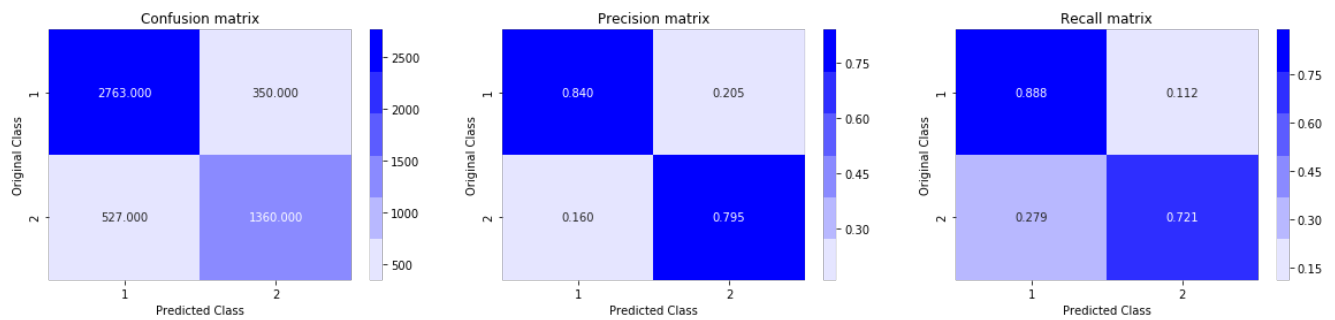
```
print(x_model.best_params_)
```

```
{'n_estimators': 500, 'max_depth': 10}
```

In [378]:

```
print("The test log loss is: ", log_loss(y_test, predict_y, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
plot_confusion_matrix(y_test, predicted_y)
```

```
The test log loss is: 0.4272872944498505
```



In []:

```
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Model", "Vectorizer", "Hyperparameter", "Train Log loss", "Test Log loss"]

x.add_row(["Logistic Regression", "TFIDF", 1, 0.616, 0.617])
x.add_row(["Linear SVM", "TFIDF", 1, 0.609, 0.61])

print(x)
```

In [379]:

```
from prettytable import PrettyTable

x = PrettyTable()

x.field_names = ["Model", "Vectorizer", "n_estimators", "Max Depth", 'Test Log loss']

x.add_row(["XGBoost", "W2V", 500, 10, 0.427])

print(x)
```

```
+-----+-----+-----+-----+-----+
| Model | Vectorizer | n_estimators | Max Depth | Test Log loss |
+-----+-----+-----+-----+-----+
| XGBoost | W2V | 500 | 10 | 0.427 |
+-----+-----+-----+-----+-----+
```