

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

- 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
```

```

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

```

```

C:\ProgramData\Anaconda3\lib\site-packages\fuzzywuzzy\fuzz.py:11: UserWarning:
Using slow pure-python SequenceMatcher. Install python-Levenshtein to remove this warning

```

3.1 Reading data and basic stats

In [6]:

```

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

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

```

Number of data points: 404290

In [7]:

```
df.head()
```

Out[7]:

	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 1000	0
4	4	9	10	Which one dissolve in water quickly sugar, salt...	Which fish would survive in salt water?	0

In [8]:

```
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 on:

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

3.2.1 Distribution of data points among output classes

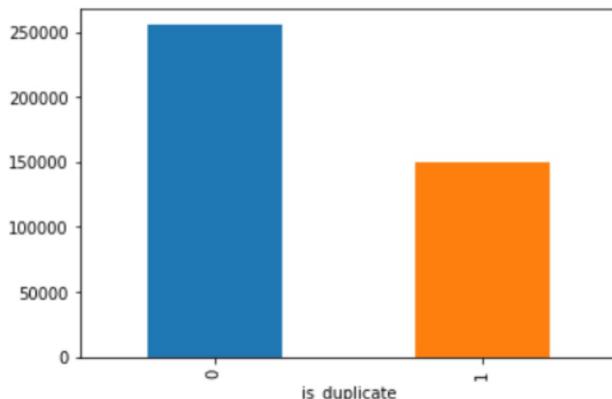
- Number of duplicate(similar) and non-duplicate(non similar) questions

In [9]:

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

Out[9]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x2a82125be10>
```



In [10]:

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

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

In [11]:

```
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 [12]:

```
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: {}')
```

```
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%)

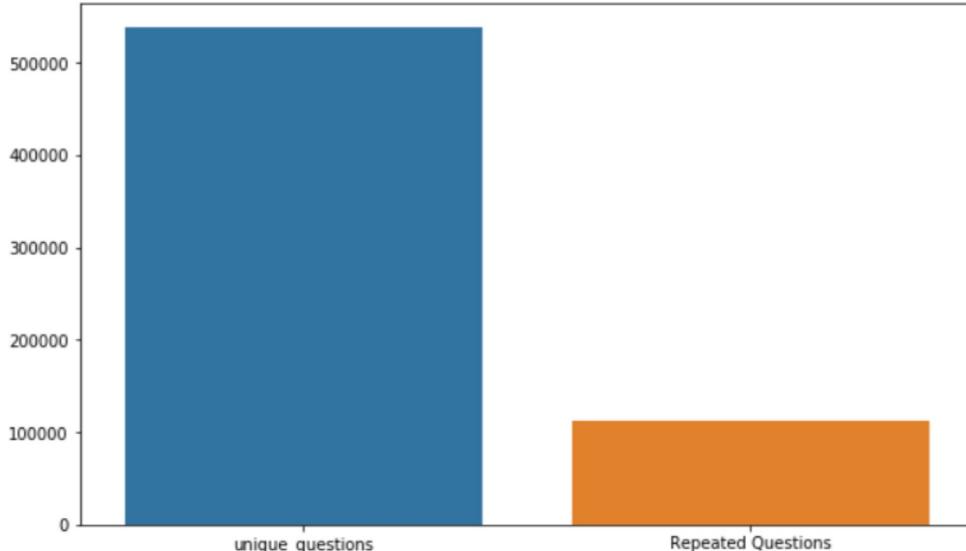
Max number of times a single question is repeated: 157
```

In [13]:

```
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()
```

Plot representing unique and repeated questions



3.2.3 Checking for Duplicates

In [14]:

```
#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 [15]:

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

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

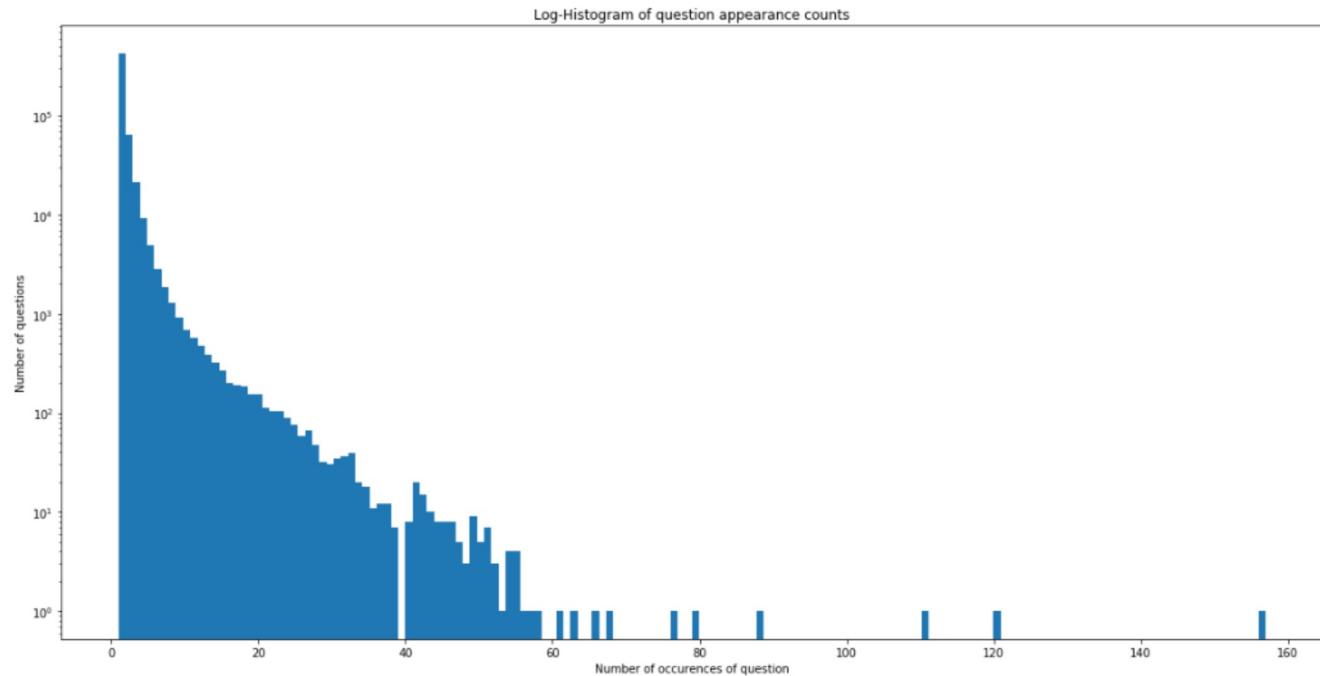
```

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: {} \n'.format(max(qids.value_counts())))

```

Maximum number of times a single question is repeated: 157



3.2.5 Checking for NULL values

In [16]:

```

#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	0

- There are two rows with null values in question2 and one row with null value in question1

In [17]:

```

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

```

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

3.3 Basic Feature Extraction (before cleaning)

Let us now construct a few features like:

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

In [18]:

```
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[18]:

	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	What would happen if the Indian	0	4	1	51	88	8	13	4.0

				Dia question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word
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 $[math]23^{24}[/math] i...$	0	1	1	50	65	11	9	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	13	7	2.0

◀ ▶

3.3.1 Analysis of some of the extracted features

- Here are some questions have only one single words.

In [20]:

```
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

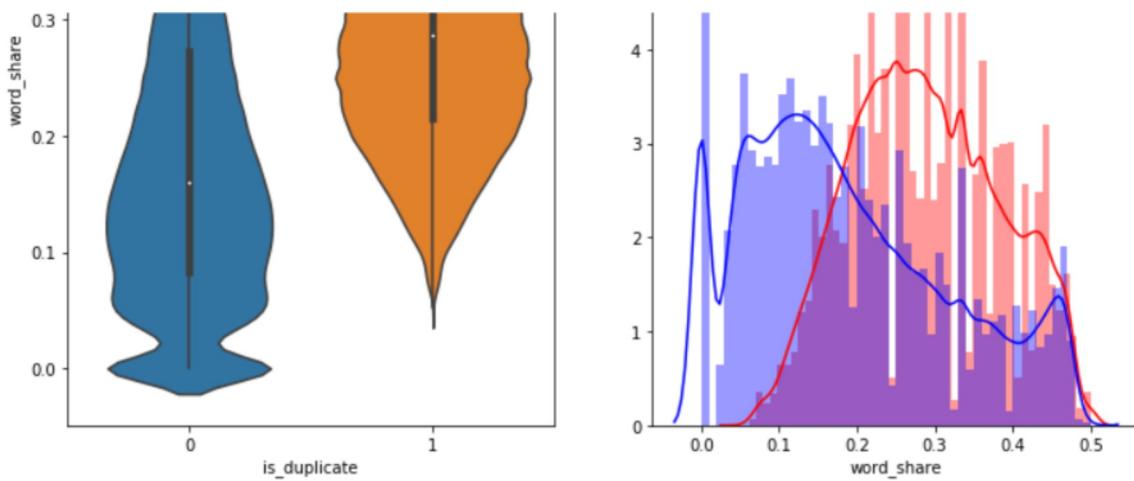
In [22]:

```
import warnings
warnings.filterwarnings("ignore")
plt.figure(figsize=(12, 8))

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

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





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

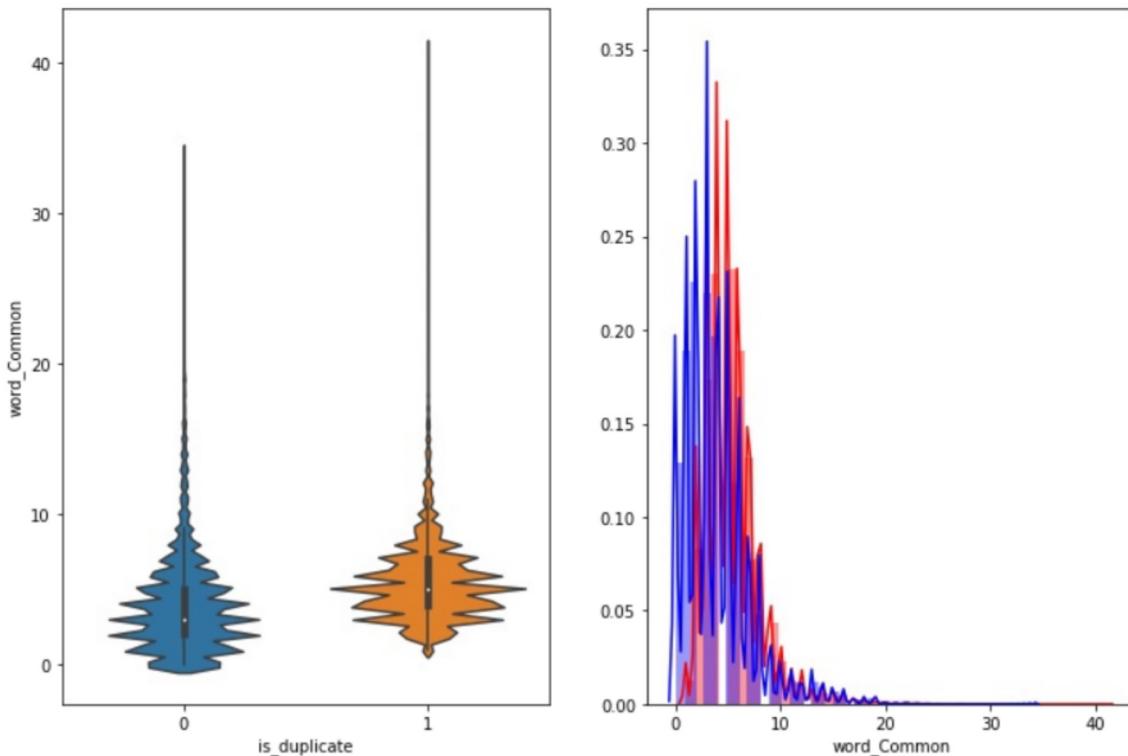
3.3.1.2 Feature: word_Common

In [23]:

```
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

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

In [2] :

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

STOP_WORDS = stopwords.words("english")

def preprocess(x):
    x = str(x).lower()
    x = x.replace(",000,000", "m").replace("", " ").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

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

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

$$csc_min = \text{common_stop_count} / (\min(\text{len}(q1_stops), \text{len}(q2_stops)))$$

- **csc_max** : Ratio of common_stop_count to max length of stop count of Q1 and Q2

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

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

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

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

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

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

$$\text{mean_len} = (\text{len}(\text{q1_tokens}) + \text{len}(\text{q2_tokens}))/2$$
- **fuzz_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage 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

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

In [31]:

```
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[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.lcsubstrings(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 [32]:

```
if os.path.isfile('nlp_features_train.csv'):
```

```

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[32]:

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

2 rows × 21 columns

3.5.1 Analysis of extracted features

3.5.1.1 Plotting Word clouds

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

In [35]:

```

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.vstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
n = np.vstack([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',encoding='utf-8')
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 [38]:

```

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

textp_w = open(path.join(d, 'train_p.txt'),encoding='utf-8').read()
textn_w = open(path.join(d, 'train_n.txt'),encoding='utf-8').read()
stopwords = set(STOPWORDS)
stopwords.add("said")

```

```
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 [39]:

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

Word Cloud for Duplicate Question pairs



Word Clouds generated from non duplicate pair question's text

In [40]:

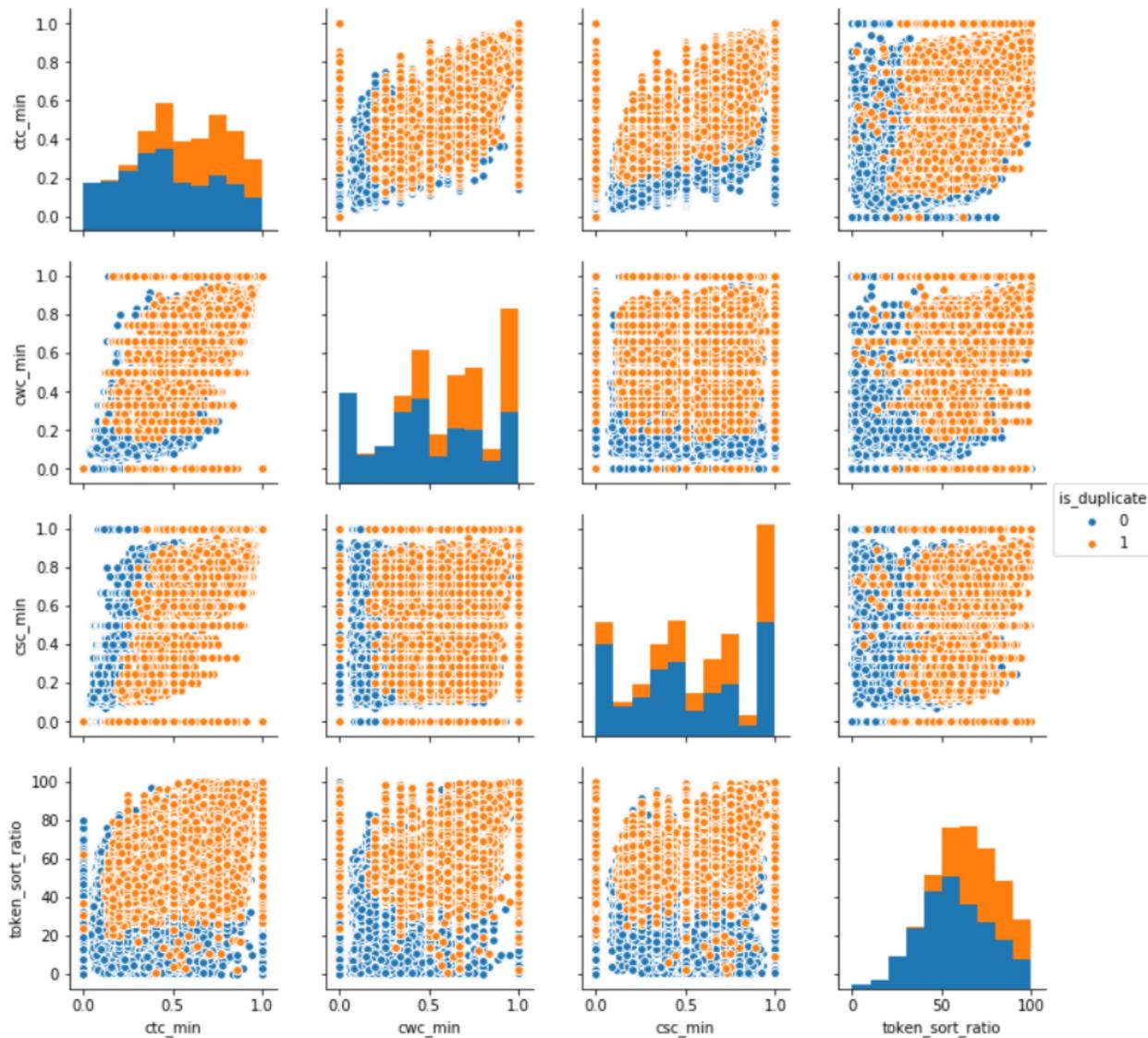
```
wc = WordCloud(background_color="white", max_words=len(textn_w),stopwords=stopwords)
# generate word cloud
wc.generate(textn_w)
print ("Word Cloud for non-Duplicate Question pairs:")
plt.imshow(wc, interpolation='bilinear')
plt.axis("off")
plt.show()
```

Word Cloud for non-Duplicate Question pairs:



In [41]:

```
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()
```

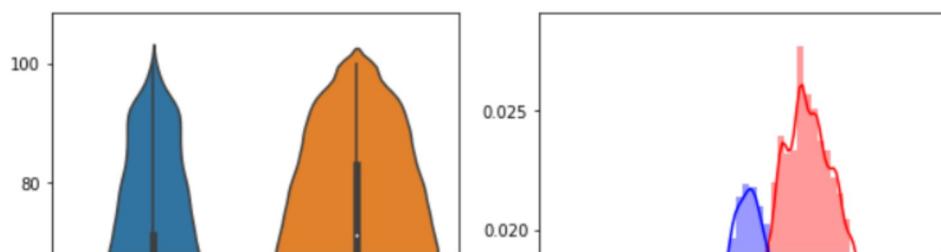


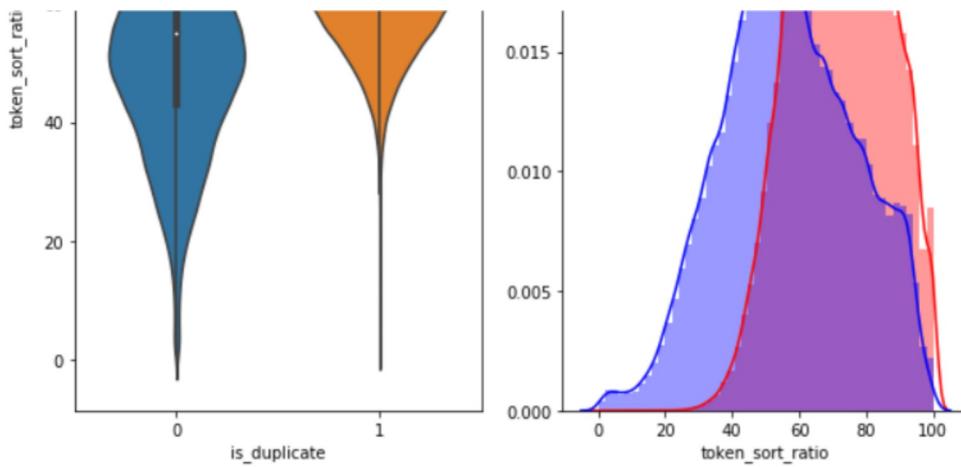
In [42]:

```
# 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()
```



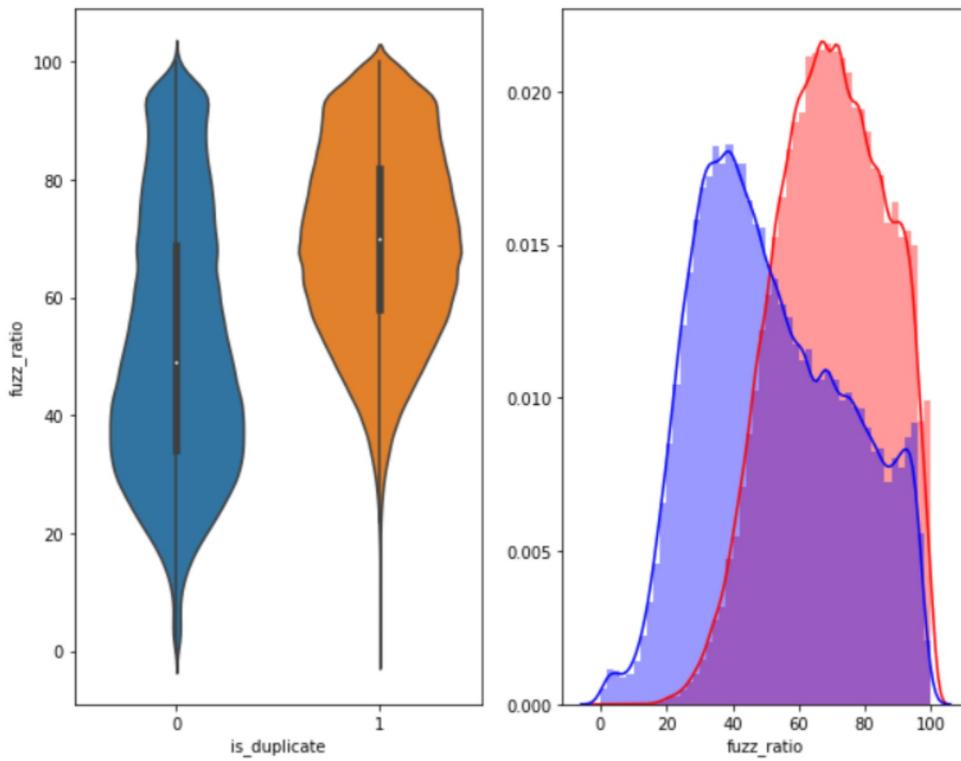


In [43]:

```
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 [45]:

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

from sklearn.preprocessing import MinMaxScaler

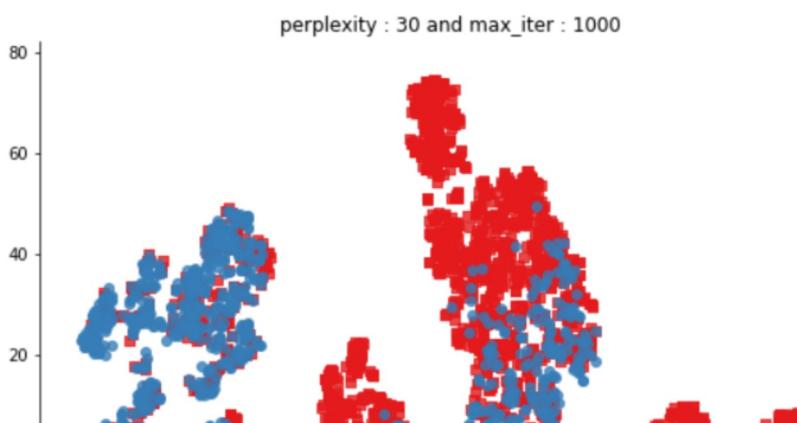
dfp_subsampled = df[0:5000]
X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_min', 'csc_max',
'ctc_min', 'ctc_max', 'last_word_eq', 'first_word_eq', 'abs_len_diff', 'mean_len', 'token_set_size',
'len_ratio', 'len_ratio_sq', 'len_sqrt', 'len_log', 'len_log_sqrt', 'len_log_ratio', 'len_log_ratio_sq']]]
```

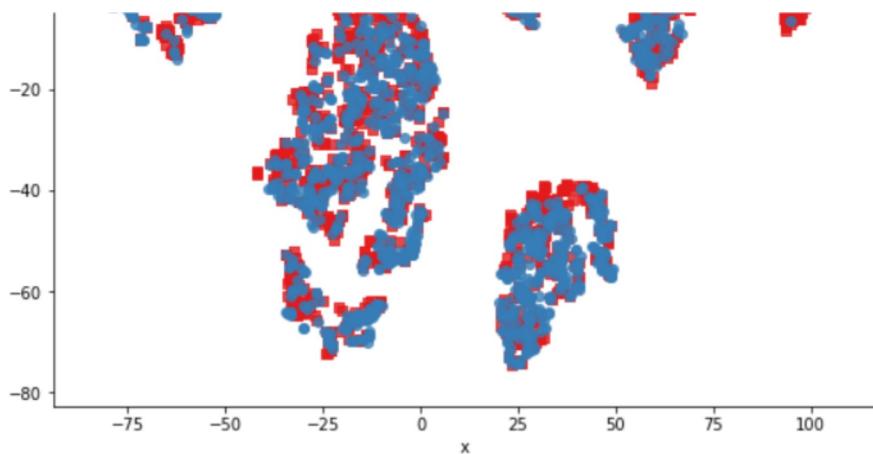
```
In [46]:
```

```
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 0.019s...  
[t-SNE] Computed neighbors for 5000 samples in 0.598s...  
[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 0.406s  
[t-SNE] Iteration 50: error = 81.2911148, gradient norm = 0.0457501 (50 iterations in 4.737s)  
[t-SNE] Iteration 100: error = 70.6044159, gradient norm = 0.0086692 (50 iterations in 3.102s)  
[t-SNE] Iteration 150: error = 68.9124908, gradient norm = 0.0056016 (50 iterations in 2.993s)  
[t-SNE] Iteration 200: error = 68.1010742, gradient norm = 0.0047585 (50 iterations in 3.264s)  
[t-SNE] Iteration 250: error = 67.5907974, gradient norm = 0.0033576 (50 iterations in 3.477s)  
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.590797  
[t-SNE] Iteration 300: error = 1.7929677, gradient norm = 0.0011899 (50 iterations in 3.345s)  
[t-SNE] Iteration 350: error = 1.3937442, gradient norm = 0.0004817 (50 iterations in 3.151s)  
[t-SNE] Iteration 400: error = 1.2280033, gradient norm = 0.0002773 (50 iterations in 3.059s)  
[t-SNE] Iteration 450: error = 1.1383208, gradient norm = 0.0001865 (50 iterations in 3.163s)  
[t-SNE] Iteration 500: error = 1.0834006, gradient norm = 0.0001423 (50 iterations in 3.197s)  
[t-SNE] Iteration 550: error = 1.0474092, gradient norm = 0.0001144 (50 iterations in 3.249s)  
[t-SNE] Iteration 600: error = 1.0231259, gradient norm = 0.0000995 (50 iterations in 3.193s)  
[t-SNE] Iteration 650: error = 1.0066353, gradient norm = 0.0000895 (50 iterations in 3.230s)  
[t-SNE] Iteration 700: error = 0.9954656, gradient norm = 0.0000805 (50 iterations in 3.256s)  
[t-SNE] Iteration 750: error = 0.9871529, gradient norm = 0.0000719 (50 iterations in 3.301s)  
[t-SNE] Iteration 800: error = 0.9801921, gradient norm = 0.0000657 (50 iterations in 3.190s)  
[t-SNE] Iteration 850: error = 0.9743395, gradient norm = 0.0000631 (50 iterations in 3.271s)  
[t-SNE] Iteration 900: error = 0.9693972, gradient norm = 0.0000606 (50 iterations in 3.242s)  
[t-SNE] Iteration 950: error = 0.9654404, gradient norm = 0.0000594 (50 iterations in 3.398s)  
[t-SNE] Iteration 1000: error = 0.9622302, gradient norm = 0.0000565 (50 iterations in 3.328s)  
[t-SNE] KL divergence after 1000 iterations: 0.962230
```

```
In [47]:
```

```
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 [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 0.014s...
[t-SNE] Computed neighbors for 5000 samples in 0.639s...
[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 0.309s
[t-SNE] Iteration 50: error = 80.5316772, gradient norm = 0.0296611 (50 iterations in 16.542s)
[t-SNE] Iteration 100: error = 69.3815765, gradient norm = 0.0033166 (50 iterations in 7.897s)
[t-SNE] Iteration 150: error = 67.9724655, gradient norm = 0.0018542 (50 iterations in 7.240s)
[t-SNE] Iteration 200: error = 67.4176865, gradient norm = 0.0012513 (50 iterations in 7.215s)
[t-SNE] Iteration 250: error = 67.1036377, gradient norm = 0.0009096 (50 iterations in 7.483s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.103638
[t-SNE] Iteration 300: error = 1.5251231, gradient norm = 0.0007399 (50 iterations in 9.248s)
[t-SNE] Iteration 350: error = 1.1820215, gradient norm = 0.0002076 (50 iterations in 12.185s)
[t-SNE] Iteration 400: error = 1.0389463, gradient norm = 0.0000969 (50 iterations in 11.220s)
[t-SNE] Iteration 450: error = 0.9659566, gradient norm = 0.0000635 (50 iterations in 11.109s)
[t-SNE] Iteration 500: error = 0.9267892, gradient norm = 0.0000482 (50 iterations in 10.961s)
[t-SNE] Iteration 550: error = 0.9053178, gradient norm = 0.0000406 (50 iterations in 10.987s)
[t-SNE] Iteration 600: error = 0.8915660, gradient norm = 0.0000349 (50 iterations in 11.089s)
[t-SNE] Iteration 650: error = 0.8804696, gradient norm = 0.0000345 (50 iterations in 10.934s)
[t-SNE] Iteration 700: error = 0.8723292, gradient norm = 0.0000358 (50 iterations in 10.984s)
[t-SNE] Iteration 750: error = 0.8668707, gradient norm = 0.0000314 (50 iterations in 11.328s)
[t-SNE] Iteration 800: error = 0.8626194, gradient norm = 0.0000250 (50 iterations in 10.915s)
[t-SNE] Iteration 850: error = 0.8584315, gradient norm = 0.0000253 (50 iterations in 11.090s)
[t-SNE] Iteration 900: error = 0.8547347, gradient norm = 0.0000261 (50 iterations in 10.555s)
[t-SNE] Iteration 950: error = 0.8517873, gradient norm = 0.0000263 (50 iterations in 10.724s)
[t-SNE] Iteration 1000: error = 0.8493521, gradient norm = 0.0000250 (50 iterations in 10.768s)
[t-SNE] KL divergence after 1000 iterations: 0.849352
```

In [49]:

```
trace1 = go.Scatter3d(
    x=tsne3d[:,0],
    y=tsne3d[:,1],
    z=tsne3d[:,2],
    mode='markers',
    marker=dict(
        sizemode='diameter',
        color = v,
```

```
        line=dict(color='rgb(255, 255, 255)'),
        opacity=0.75
    )
)

data=[trace1]
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 [44]:

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

In [98]:

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

```
def question2_x = df['question2'].apply(lambda x: set(x))
```

In [55]:

```
X_train=df[0:70000]
X_test=df[70000:100000]
```

In [104]:

```
X_train = X_train.fillna('')
X_test= X_test.fillna('')
```

In [105]:

```
nan_rows = X_train[X_train.isnull().any(1)]
print (nan_rows)
```

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

In [106]:

```
nan_rows = X_test[X_test.isnull().any(1)]
print (nan_rows)
```

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

Using only X_train to produce TFIDF scores

In [107]:

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

TFIDF word2vec vectorization of Train data

In [109]:

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

vecs1 = []

for qul in list(X_train['question1']):
    doc1 = nlp(qul)
    # 384 is the number of dimensions of vectors
    mean_vec1 = np.zeros([len(doc1), 96])
    for word1 in doc1:
        # word2vec
        # word2vec
```

```

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_train['q1_feats_m'] = list(vecs1)

```

In [110]:

```

vecs2 = []
for qu2 in list(X_train['question2']):
    doc2 = nlp(qu2)
    mean_vec2 = np.zeros([len(doc2), 96])
    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_train['q2_feats_m'] = list(vecs2)

```

TFIDF word2vec vectorization of Test Data

In [111]:

```

# en_vectors_web_lg, which includes over 1 million unique vectors.

vecs1 = []

for qu1 in list(X_test['question1']):
    doc1 = nlp(qu1)
    # 384 is the number of dimensions of vectors
    mean_vec1 = np.zeros([len(doc1), 96])
    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)

```

In [112]:

```

vecs2 = []
for qu2 in list(X_test['question2']):
    doc2 = nlp(qu2)
    mean_vec2 = np.zeros([len(doc2), 96])
    for word2 in doc2:
        # word2vec
        vec2 = word2.vector
        # fetch df score
        try:
            idf = word2tfidf[str(word2)]
        except:
            #print word

```

```

    mean_vec2 += vec2 * idf
mean_vec2 = mean_vec2.mean(axis=0)
vecs2.append(mean_vec2)
X_test['q2_feats_m'] = list(vecs2)

```

Loading previously created features

In [116]:

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

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

In [117]:

```

df1 = dfnlp.drop(['qid1','qid2','question1','question2'],axis=1)
df2 = dfppro.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)

df3 = X_train.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
df4 = X_test.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)

df3_q1_train = pd.DataFrame(df3.q1_feats_m.values.tolist(), index= df3.index)
df3_q2_train = pd.DataFrame(df3.q2_feats_m.values.tolist(), index= df3.index)

df4_q1_test = pd.DataFrame(df4.q1_feats_m.values.tolist(), index= df4.index)
df4_q2_test = pd.DataFrame(df4.q2_feats_m.values.tolist(), index= df4.index)

```

In [132]:

```
# Splitting df1 i.e. NLP features in test and train
df1_train=df1[0:70000]
df1_test=df1[70000:100000]

# splitting df2 i.e preprocessed features in test and train
df2_train=df2[0:70000]
df2_test=df2[70000:100000]
```

In [120]:

```
# dataframe of nlp features  
df1 train.head(2)
```

Out[120]:

	id	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff
0	0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	2.0
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0

In [121]:

```
# data before preprocessing  
df2 train.head(2)
```

Out[121]:

0	0	1	freq_qid1	1	freq_qid2	66	q1len	57	q2len	14	n_words	12	n_words	10.0	word_Common	23.0	Word_Total	0.434783	2	freq_q1+q2	f
1	1	4		1		51	88	8		13		4.0			20.0		0.200000	5		3	

In [122]:

```
# Questions 1 tfidf weighted word2vec
df3_q1_train.head(2)
```

Out[122]:

	0	1	2	3	4	5	6	7	8	9	...		
0	199.005999	-	-	-	-	2.381461	128.083333	46.694098	-	62.432003	49.942147	...	3
1	113.862048	-99.744038	-	89.539876	-82.767201	-	9.223255	39.254904	80.029392	10.840831	146.492736	...	5

2 rows × 96 columns

In [126]:

```
# Questions 2 tfidf weighted word2vec
df3_q2_train.head(2)
```

Out[126]:

	0	1	2	3	4	5	6	7	8	9	...		
0	141.072898	-	-	-	-90.512967	9.149163	98.874355	33.330166	-	35.828580	46.673881	...	
1	124.949524	-40.560911	-	91.266866	110.465091	104.608168	34.715419	105.596472	16.691853	-	66.591777	73.211175	...

2 rows × 96 columns

In [128]:

```
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 train question1 w2v  dataframe : ", df3_q1_train.shape[1])
print("Number of features in train question2 w2v  dataframe : ", df3_q2_train.shape[1])
print("Number of features in test question1 w2v  dataframe : ", df4_q1_test.shape[1])
print("Number of features in test question2 w2v  dataframe : ", df4_q2_test.shape[1])
print("Number of features in final dataframe  : ", df1.shape[1]+df2.shape[1]+df3_q1_train.shape[1]+df3_q2_train.shape[1])
```

```
Number of features in nlp dataframe : 17
Number of features in preprocessed dataframe : 12
Number of features in train question1 w2v  dataframe : 96
Number of features in train question2 w2v  dataframe : 96
Number of features in test question1 w2v  dataframe : 96
Number of features in test question2 w2v  dataframe : 96
Number of features in final dataframe  : 221
```

Merging all the above features and creating a final Train data

In [135]:

```
# storing the final features to csv file
if not os.path.isfile('final_features_train.csv'):
    df3_q1_train['id']=df1_train['id']
    df3_q2_train['id']=df1_train['id']
    df3_q1_train['q1']=df1_train['q1']
    df3_q2_train['q2']=df1_train['q2']
```

```
result = df1_train.merge(df2_train, on='id', how='left')
result.to_csv('final_features_train.csv')
```

Merging all the above features and creating a final Test data

In [138]:

```
# storing the final features to csv file
if not os.path.isfile('final_features_test.csv'):
    df4_q1_test['id']=df1_test['id']
    df4_q2_test['id']=df1_test['id']
    df1_test = df1_test.merge(df2_test, on='id', how='left')
    df2_test = df4_q1_test.merge(df4_q2_test, on='id', how='left')
    result = df1_test.merge(df2_test, on='id', how='left')
    result.to_csv('final_features_test.csv')
```

In [136]:

```
result.head()
```

Out[136]:

	id	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	...	86_y
0	0	0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	...	25.909103
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	...	108.655479
2	2	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	...	126.138883
3	3	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	...	70.145525
4	4	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	...	18.598218

5 rows × 220 columns

In [10]:

```
import sqlite3
from sqlalchemy import create_engine
```

In [17]:

```
#Creating db file from csv
if not os.path.isfile('train.db'):
    disk_engine = create_engine('sqlite:///train.db')
    chunksize = 1
    j = 0
    index_start = 1
    for df in pd.read_csv('final_features.csv', names=a, chunksize=chunksize, iterator=True, encoding='utf-8'):
        df.index += index_start
        j+=1
        df.to_sql('data', disk_engine, if_exists='append')
        index_start = df.index[-1] + 1
```

In [18]:

```
#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)
```

```

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 [19]:

```

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

```

Tables in the database:
data

In [20]:

```

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

        # for selecting random points
        data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 100001;", conn_r)
    conn_r.commit()
    conn_r.close()

```

In [21]:

```

# 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 [129]:

```
data.head(2)
```

Out[129]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	mean_len	...
1	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0	12.5	...
2	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	4.0	12.0	...

2 rows × 218 columns

In [130]:

```
# after we read from sql table each entry was read it as a string
```

```
cols = list(result.columns)
for i in cols:
    result[i] = result[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
```

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
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

```
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
```

In [131]:

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

4.3 loading final train data and final test data created earlier

In [154]:

```
X_train=pd.read_csv('final_features_train.csv')
X_test=pd.read_csv('final_features_test.csv')

y_train=X_train['is_duplicate'].values
y_test=X_test['is_duplicate'].values
```

In [160]:

```
X_train.drop(['Unnamed: 0', 'id','is_duplicate'], axis=1, inplace=True)
X_test.drop(['Unnamed: 0', 'id','is_duplicate'], axis=1, inplace=True)
```

In [161]:

```
X_train.head(2)
```

Out[161]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	mean_len	...
0	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	2.0	13.0	...
1	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0	12.5	...

2 rows × 218 columns

In [162]:

```
X_test.head(2)
```

Out[162]:

	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	mean_len	...
0	0.666644	0.666644	0.999967	0.999967	0.833319	0.833319	0.0	1.0	0.0	6.0	...
1	0.666656	0.444440	0.599988	0.599988	0.636358	0.466664	1.0	1.0	4.0	13.0	...

2 rows × 218 columns

In [164]:

```
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, 218)
Number of data points in test data : (30000, 218)
```

In [36]:

```
from collections import Counter
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)
```

```
----- Distribution of output variable in train data -----
Class 0: 0.6275434192534541 Class 1: 0.37245658074654586
----- Distribution of output variable in test data -----
Class 0: 0.37247619047619046 Class 1: 0.37247619047619046
```

In [109]:

```
# 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
    #diamensional array
    # C.sum(axis =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
    #diamensional array
    # C.sum(axis =0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                      [3/4, 4/6]]
    plt.figure(figsize=(20,4))

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

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

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

    plt.show()
```

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

In [40]:

```
# 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/18662466/4094039
```

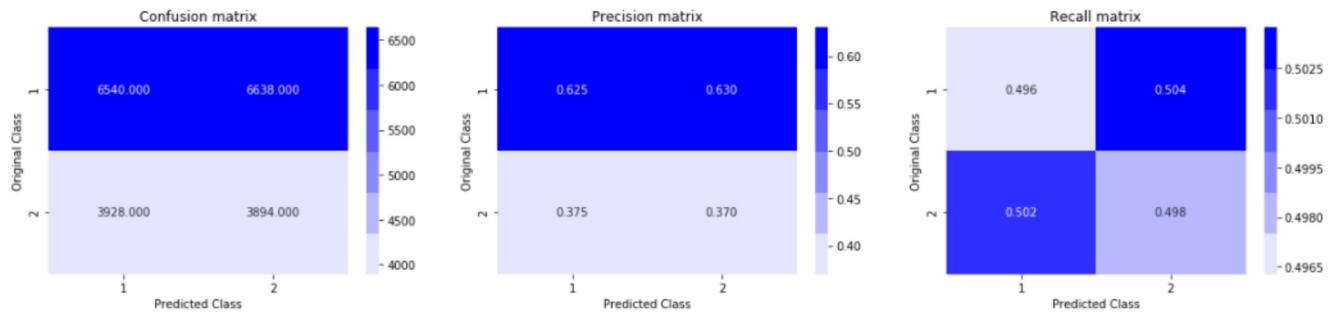
```

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 0.8847942690221885



4.4 Logistic Regression with hyperparameter tuning

In [46]:

```

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=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()

```

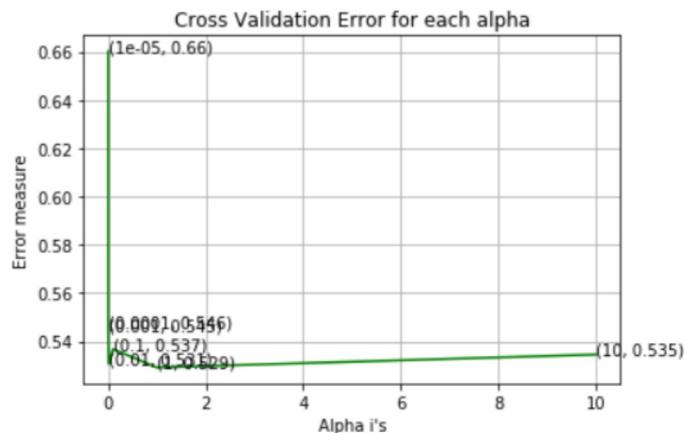
```

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)

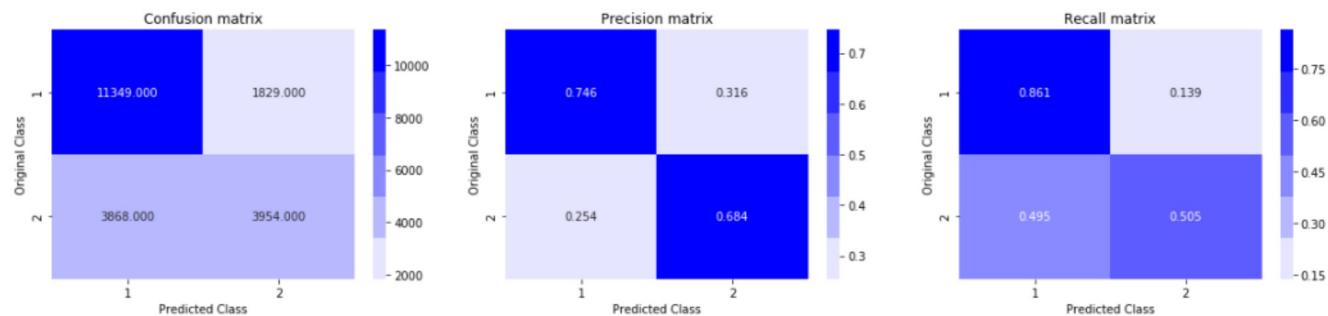
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test, 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: 0.6602604108349805
For values of alpha = 0.0001 The log loss is: 0.5459562711029176
For values of alpha = 0.001 The log loss is: 0.5445127199019395
For values of alpha = 0.01 The log loss is: 0.5305912183114352
For values of alpha = 0.1 The log loss is: 0.5367786929207705
For values of alpha = 1 The log loss is: 0.5291676629786991
For values of alpha = 10 The log loss is: 0.5345252046108739
```



For values of best alpha = 1 The train log loss is: 0.5306005424795032
For values of best alpha = 1 The test log loss is: 0.5291676629786991
Total number of data points : 30000



4.5 Linear SVM with hyperparameter tuning

In [48]:

```
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,
```

```

# class_weight=None, warm_start=False, average=False, n_iter=None)

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

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


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

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

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

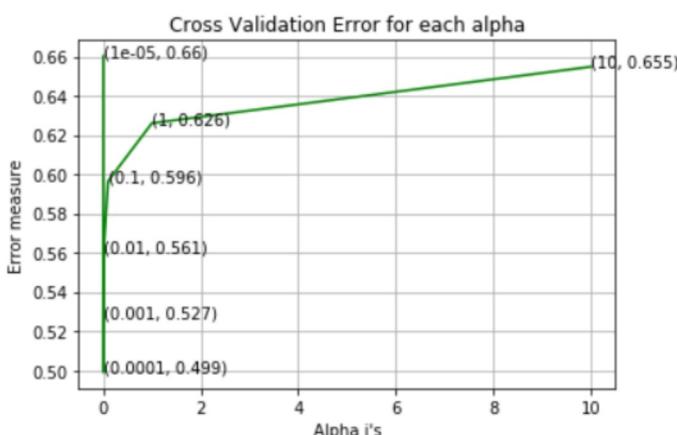
predict_y = sig_clf.predict_proba(X_train)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss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

```

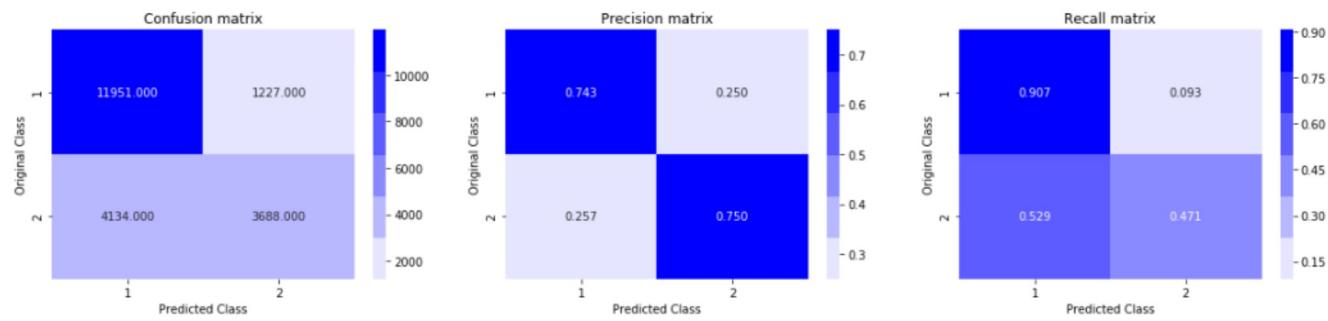
```

For values of alpha =  1e-05 The log loss is: 0.6602604108349805
For values of alpha =  0.0001 The log loss is: 0.49899986454926387
For values of alpha =  0.001 The log loss is: 0.5272283694459158
For values of alpha =  0.01 The log loss is: 0.5606845355119315
For values of alpha =  0.1 The log loss is: 0.5964203604091971
For values of alpha =  1 The log loss is: 0.6260445118427591
For values of alpha =  10 The log loss is: 0.6547962018247225

```



```
for values or best alpha = 0.0001 the test log loss is: 0.4989998645492638 /  
Total number of data points : 30000
```



4.6 XGBoost

```
In [119]:
```

```
import xgboost as xgb
from sklearn.model_selection import RandomizedSearchCV
XG = xgb.XGBClassifier(subsample=0.7, colsample_bytree=0.7, random_state=0, reg_alpha=1, reg_lambda=0)

parameters = {'max_depth':[1, 5, 10, 50], 'n_estimators':[50, 100, 200]}
clf = RandomizedSearchCV(XG, parameters, cv=2, scoring='r2')
clf.fit(X_train, y_train)
print(clf.best_params_)

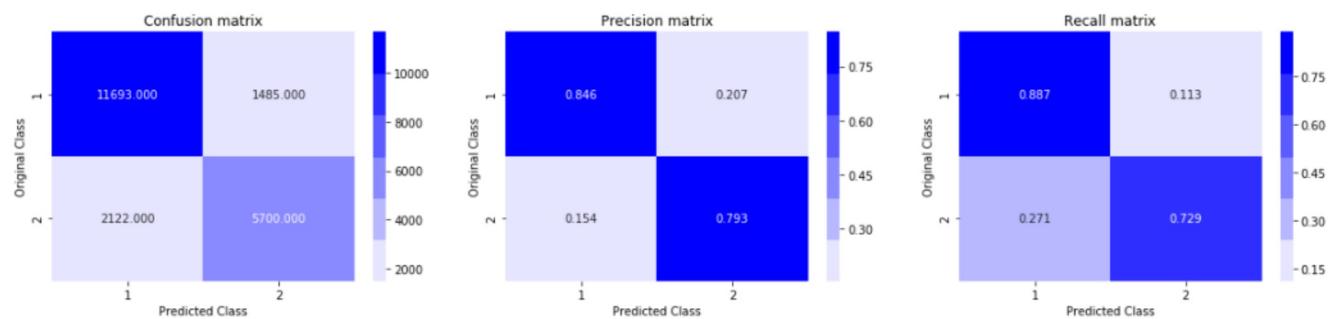
{'n_estimators': 100, 'max_depth': 5}
```

```
In [144]:
```

```
clf = xgb.XGBClassifier(n_estimators=100, max_depth=5, subsample=0.7, colsample_bytree=0.7, random_state=0, reg_alpha=1, reg_lambda=0)
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 n_estimator = ', 100, "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 n_estimator = ', 100, "The test log loss is:", log_loss(y_test, predict_y, labels=clf.classes_, eps=1e-15))
predicted_y = np.argmax(predict_y, axis=1)
print("Total number of data points :", len(predicted_y))
plot_confusion_matrix(y_test, predicted_y)
```

```
For n_estimator = 100 The train log loss is: 0.3083573356061459
For n_estimator = 100 The test log loss is: 0.35870354120629905
Total number of data points : 30000
```



4.4 Logistic Regression with hyperparameter tuning and simple TFIDF vecotrs

```
in [96]:  
  
#prepro_features_train.csv (Simple Preprocessing Features)  
#nlp_features_train.csv (NLP Features)  
if os.path.isfile('nlp_features_train.csv'):  
    dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')  
else:  
    print("download nlp_features_train.csv from drive or run previous notebook")  
  
if os.path.isfile('df_fe_without.preprocessing_train.csv'):  
    dfppro = pd.read_csv("df_fe_without.preprocessing_train.csv",encoding='latin-1')  
else:  
    print("download df_fe_without.preprocessing_train.csv from drive or run previous notebook")
```

In [97]:

```
df1 = dfnlp.drop(['qid1','qid2','question1','question2'],axis=1)
df2 = dfppro.drop(['qid1','qid2','question1','question2','is duplicate'],axis=1)
```

In [98]:

```
df = pd.read_csv("train.csv")
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

In [99]:

```
df.drop(['qid1', 'qid2', 'is duplicate'], axis=1, inplace=True)
```

In [100]:

```
warnings.filterwarnings("ignore")
df["question1"] = df["question1"].fillna("").apply(preprocess)
df["question2"] = df["question2"].fillna("").apply(preprocess)
```

In [11]:

```
from functools import reduce
```

In [101]:

```
dfs = [df, df1, df2]
df_final = reduce(lambda left,right: pd.merge(left,right,on='id')), dfs)
```

DataFrame consisting NLP, fuzzy features and Basic features

In [80]:

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

Out[80]:

		question1	sto... question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	...	freq_qid2	q1len
--	--	-----------	---------------------	--------------	---------	---------	---------	---------	---------	---------	-----	-----------	-------

2 rows × 30 columns

In [102]:

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

In [103]:

```
df_final.head(2)
```

Out[103]:

	question1	question2	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	...	freq
0	what is the step by step guide to invest in sh...	what is the step by step guide to invest in sh...	0.999980	0.833319	0.999983	0.999983	0.916659	0.785709	0.0	1.0	...	1
1	what is the story of kohinoor koh i noor dia...	what would happen if the indian government sto...	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	...	1

2 rows × 28 columns

In [104]:

```
from sklearn.model_selection import train_test_split
X_train2,X_test2, y_train2, y_test2 = train_test_split(df_final, y_true, stratify=y_true, test_size=0.3, random_state=0)
```

In [105]:

```
print("Number of data points in train data :",X_train2.shape)
print("Number of data points in test data :",X_test2.shape)
```

Number of data points in train data :(70000, 28)
Number of data points in test data :(30000, 28)

In [106]:

```
X_train2=X_train2.reset_index()
X_test2=X_test2.reset_index()
```

In [107]:

```
X_train2.drop(['index'], axis=1, inplace=True)
X_test2.drop(['index'], axis=1, inplace=True)
```

In [108]:

```
X_train2.head()
```

Out[108]:

	question1	question2	ewe_min	ewe_max	ese_min	ese_max	ete_min	ete_max	last_word_eq	first_word_eq	...	fre
0	what are some best lines said by any celebrities	what is the best line said by your x	0.499988	0.499988	0.399992	0.399992	0.444440	0.444440	0.0	1.0	...	1
1	how can i improve my pronunciation of english ...	how can i improve my speaking	0.499975	0.249994	0.999975	0.799984	0.833319	0.555549	0.0	1.0	...	23
2	which is the stronger acid propionic acid ace...	why is benzoic acid stronger than acetic acid ...	0.799984	0.666656	0.249994	0.249994	0.416663	0.384612	1.0	0.0	...	1
3	if evidence shows there own only a 49 percent ...	is there any evidence that spanking causes psy...	0.333328	0.117646	0.833319	0.357140	0.583328	0.205882	0.0	0.0	...	1
4	any ideas for a grade 10 biology research proj...	what are some good biological research topics	0.249994	0.166664	0.000000	0.000000	0.142855	0.111110	0.0	0.0	...	2

5 rows × 28 columns



TFIDF Vectorization of question1

In [109]:

```
from sklearn.feature_extraction.text import TfidfVectorizer
vectorizer = TfidfVectorizer(min_df=2, ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train2['question1'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_q1_tfidf = vectorizer.transform(X_train2['question1'].values)
X_test_q1_tfidf = vectorizer.transform(X_test2['question1'].values)
```

TFIDF Vectorization of question2

In [111]:

```
vectorizer = TfidfVectorizer(min_df=2, ngram_range=(1,4), max_features=5000)
vectorizer.fit(X_train2['question2'].values) # fit has to happen only on train data

# we use the fitted CountVectorizer to convert the text to vector
X_train_q2_tfidf = vectorizer.transform(X_train2['question2'].values)
X_test_q2_tfidf = vectorizer.transform(X_test2['question2'].values)
```

In [112]:

```
X_train3=X_train2.drop(['question1','question2'], axis=1)
```

```
ValueError: Input contains NaN, infinity or a value too large for an int64.
```

```
In [115]:
```

```
q1_train = pd.DataFrame(X_train_q1_tfidf.toarray())
q1_test= pd.DataFrame(X_test_q1_tfidf.toarray())

q2_train= pd.DataFrame(X_train_q2_tfidf.toarray())
q2_test= pd.DataFrame(X_test_q2_tfidf.toarray())
```

```
In [117]:
```

```
X_train_final=pd.concat([X_train3,q1_train,q2_train], axis=1)
```

X_train_final dataframe consist of simple TFIDF feature, NLP, Fuzzy feature and Basic Feature for each question

```
In [107]:
```

```
print(X_train_final.shape)
```

```
(70000, 10026)
```

```
In [119]:
```

```
X_test_final=pd.concat([X_test3,q1_test,q2_test], axis=1)
```

```
In [108]:
```

```
print(X_test_final.shape)
```

```
(30000, 10026)
```

```
In [91]:
```

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

```
In [122]:
```

```
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=0)
    clf.fit(X_train_final, y_train2)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train_final, y_train2)
```

```

    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test2, 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=0)
clf.fit(X_train_final, y_train2)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train_final, y_train2)

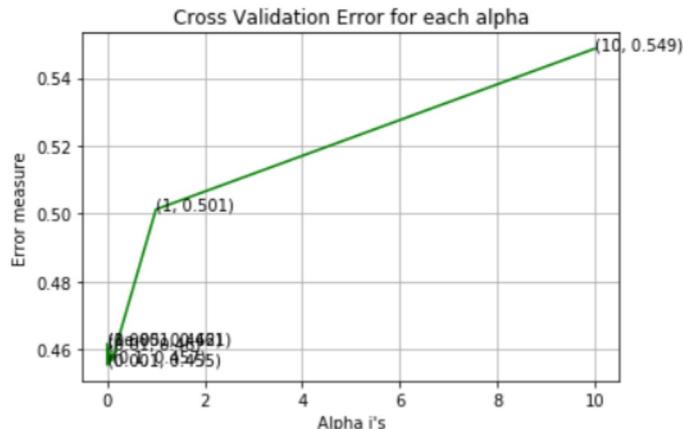
predict_y = sig_clf.predict_proba(X_train_final)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train2, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test_final)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test2, 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_test2, predicted_y)

```

```

For values of alpha =  1e-05 The log loss is: 0.46156390947182524
For values of alpha =  0.0001 The log loss is: 0.4614824384091323
For values of alpha =  0.001 The log loss is: 0.45542944984663486
For values of alpha =  0.01 The log loss is: 0.46034160194614343
For values of alpha =  0.1 The log loss is: 0.4567022667491131
For values of alpha =  1 The log loss is: 0.501278367807922
For values of alpha =  10 The log loss is: 0.5486800223062623

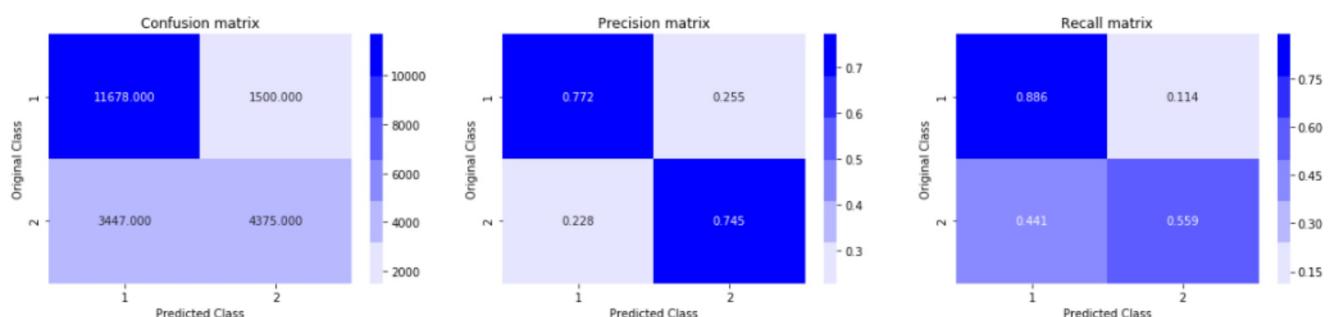
```



```

For values of best alpha =  0.001 The train log loss is: 0.44386085454825824
For values of best alpha =  0.001 The test log loss is: 0.45542944984663486
Total number of data points : 30000

```



4.4 Linear SVM with hyperparameter tuning and simple Tfidf vectors

In [124]:

```
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='hinge', random_state=0)
    clf.fit(X_train_final, y_train2)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(X_train_final, y_train2)
    predict_y = sig_clf.predict_proba(X_test_final)
    log_error_array.append(log_loss(y_test2, predict_y, labels=clf.classes_, eps=1e-15))
    print('For values of alpha = ', i, "The log loss is:", log_loss(y_test2, 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='hinge', random_state=0)
clf.fit(X_train_final, y_train2)
sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
sig_clf.fit(X_train_final, y_train2)

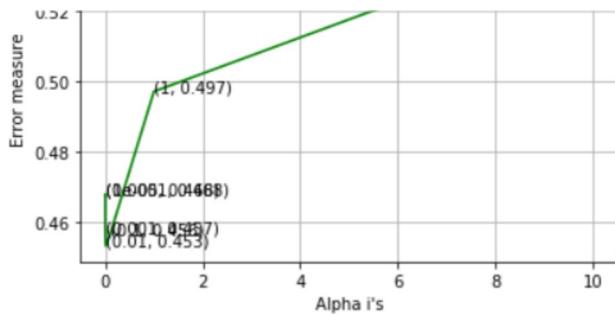
predict_y = sig_clf.predict_proba(X_train_final)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:", log_loss(y_train2, predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(X_test_final)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:", log_loss(y_test2, 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_test2, predicted_y)
```



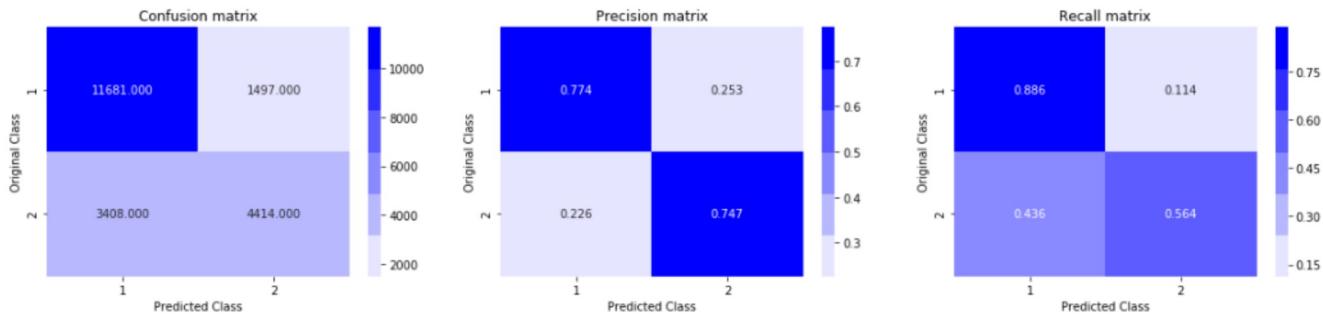
```
For values of alpha = 1e-05 The log loss is: 0.46777685745543085
For values of alpha = 0.0001 The log loss is: 0.46796243576412566
For values of alpha = 0.001 The log loss is: 0.4571789163962353
For values of alpha = 0.01 The log loss is: 0.45329059955068113
For values of alpha = 0.1 The log loss is: 0.4564920747109389
For values of alpha = 1 The log loss is: 0.49717029468270996
For values of alpha = 10 The log loss is: 0.5431967865726226
```

Cross Validation Error for each alpha





For values of best alpha = 0.01 The train log loss is: 0.4415680254182613
 For values of best alpha = 0.01 The test log loss is: 0.45329059955068113
 Total number of data points : 30000



In [123]:

```
from prettytable import PrettyTable
x = PrettyTable()
x.field_names = ["Vectorizer", "Train log loss", "Test log loss"]
x.add_row(["Logistic(TFIDF-weighted W2V)", 0.5306, 0.5291])
x.add_row(["Linear(TFIDF-weighted W2V)", 0.4996, 0.4989])
x.add_row(["Logistic(TFIDF)", 0.4438, 0.4554])
x.add_row(["Linear(TFIDF)", 0.4415, 0.4532])
x.add_row(["XGBoost", 0.3083, 0.3487])
print(x)
```

Vectorizer	Train log loss	Test log loss
Logistic(TFIDF-weighted W2V)	0.5306	0.5291
Linear(TFIDF-weighted W2V)	0.4996	0.4989
Logistic(TFIDF)	0.4438	0.4554
Linear(TFIDF)	0.4415	0.4532
XGBoost	0.3083	0.3487

Step By Step Procedure

1. Problem statement is Identify which questions asked on Quora are duplicates of questions that have already been asked. We are tasked with predicting whether a pair of questions are duplicates or not.
2. Data can be found at <https://www.kaggle.com/c/quora-question-pairs>
3. It is a binary classification problem
4. Our performance metric will be log loss and binary confusion matrix as we want a probability of a pair of questions to be duplicates so that you can choose any threshold of choice.
5. Total number of question pairs for training: 404290
6. Question pairs are not Similar (is_duplicate = 0): 63.08%
7. Question pairs are Similar (is_duplicate = 1): 36.92%
8. In our dataset classes are imbalanced
9. There are no duplicate questions in our dataset
10. Maximum number of times a single question is repeated: 157
11. I will create some basic feature before cleaning of text questions like question1 length , question2 length, frequency of question1 and question2 ID, Number of common words in question1 and question2 etc.
12. Minimum length of the question is 1 and maximum is 2 in 1

qid1 and qid2 is more when they are duplicate(Similar)

14. I perform the univariate analysis on word common and found out that The distributions of the word Common feature in similar and non-similar questions are highly overlapping
15. Now I performed the Advance feature extraction i.e NLP and fuzzy features
16. NLP and fuzzy feature includes cwc_min, csc_max, first_word_eq, mean_len, fuzz_ratio, fuzz_partial_ratio, token_set_ratio, token_sort_ratio etc.
17. I generated the word cloud for duplicate and non duplicate questions, In duplicate questions words like Donald Trump, Best Way occurs maximum no of times
18. I performed the univariate analysis on fuzz ratio and found out that it is higher for duplicate questions
19. I plot the 2D and 3D TSNE visualizations for the duplicate and non duplicate questions
20. I perform the featurization of question 1 and question 2 and transform those questions into Tfifd weighted word to vectors
21. I stored the final featurization of my dataset into a database named train.db
22. I performed the random splitting of my dataset into train and test
23. I build a random model and found out that log loss for a random model is 0.88
24. I applied various models such as Logistic Regression, Linear SVM and XGBoost
25. Now I tried Logistic Regression and Linear SVM with simple TF-IDF vectors instead of TF-IDF weighted word to vectors