

Quora Question Pairs

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.

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.

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.

2.1.1 Data Overview

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

2.1.2 Example Data point

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

Performance Metric

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

Metric(s):

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

Train and Test Construction

Data has been randomly split in the ratio of 70:30 or 80:20 whatever we choose as we have sufficient points to work with.

Exploratory Data Analysis

In [3]:

```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 from subprocess import check_output
6 %matplotlib inline
7 import plotly.offline as py
8 py.init_notebook_mode(connected=True)
9 import plotly.graph_objs as go
10 import plotly.tools as tls
11 import os
12 import gc
13
14 import re
15 from nltk.corpus import stopwords
16 import distance
17 from nltk.stem import PorterStemmer
18 from bs4 import BeautifulSoup
```

Reading data and basic stats

In [4]:

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

Number of data points: 404290

In [5]:

```
1 df.head()
```

Out[5]:

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

In [6]:

```
1 df.info()
```

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

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

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

Distribution of data points among output classes

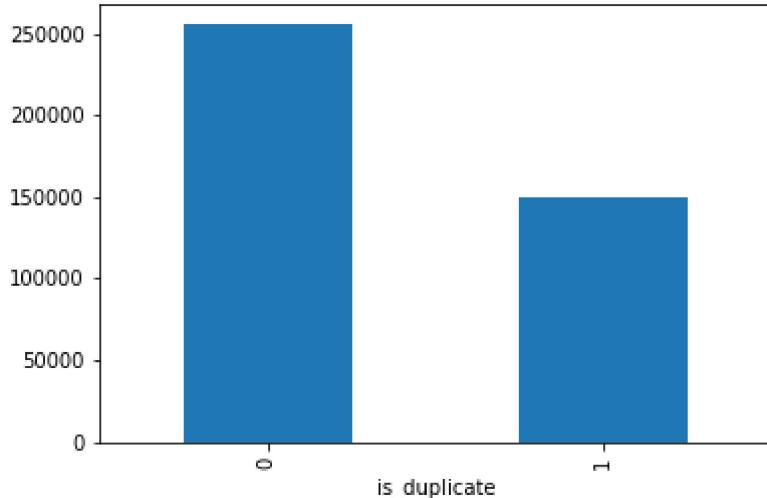
- Number of duplicate(similar) and non-duplicate(non similar) questions

In [7]:

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

Out[7]:

```
<matplotlib.axes._subplots.AxesSubplot at 0x1e706abe160>
```



In [0]:

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

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

In [11]:

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

Number of unique questions

In [0]:

```

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

```

Total num of Unique Questions are: 537933

Number of unique questions that appear more than one time: 111780 (20.779539
45937505%)

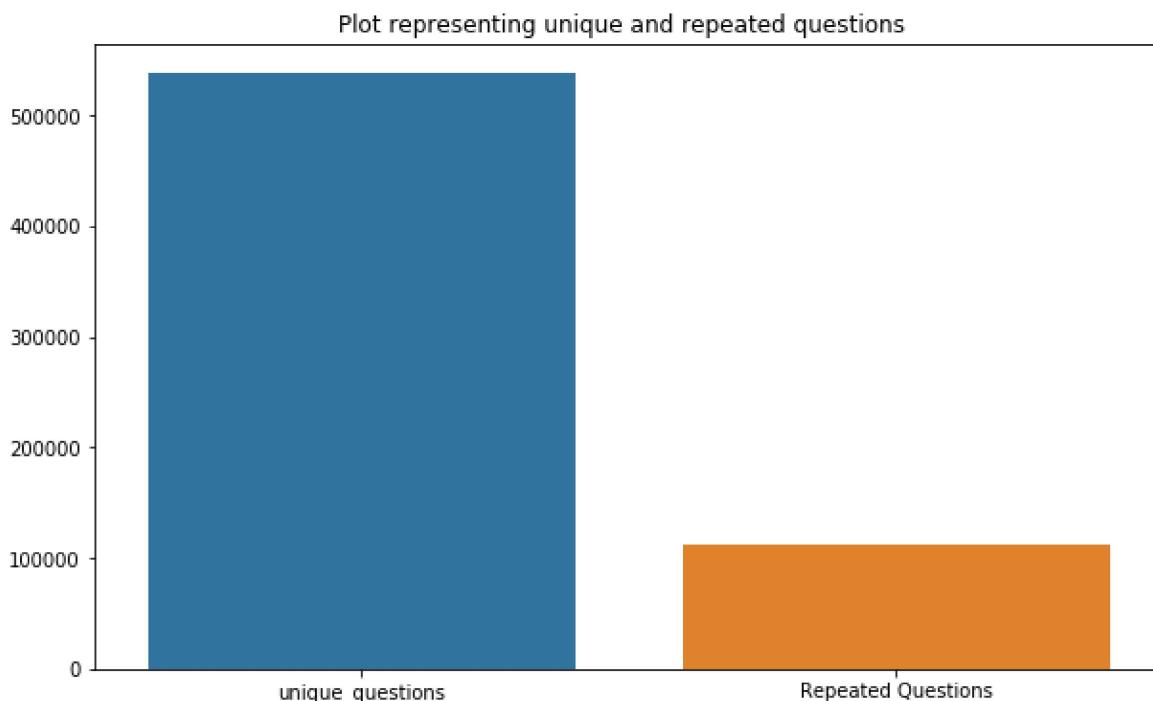
Max number of times a single question is repeated: 157

In [0]:

```

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

```



Checking for Duplicates

In [0]:

```

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

```

Number of duplicate questions 0

Number of occurrences of each question

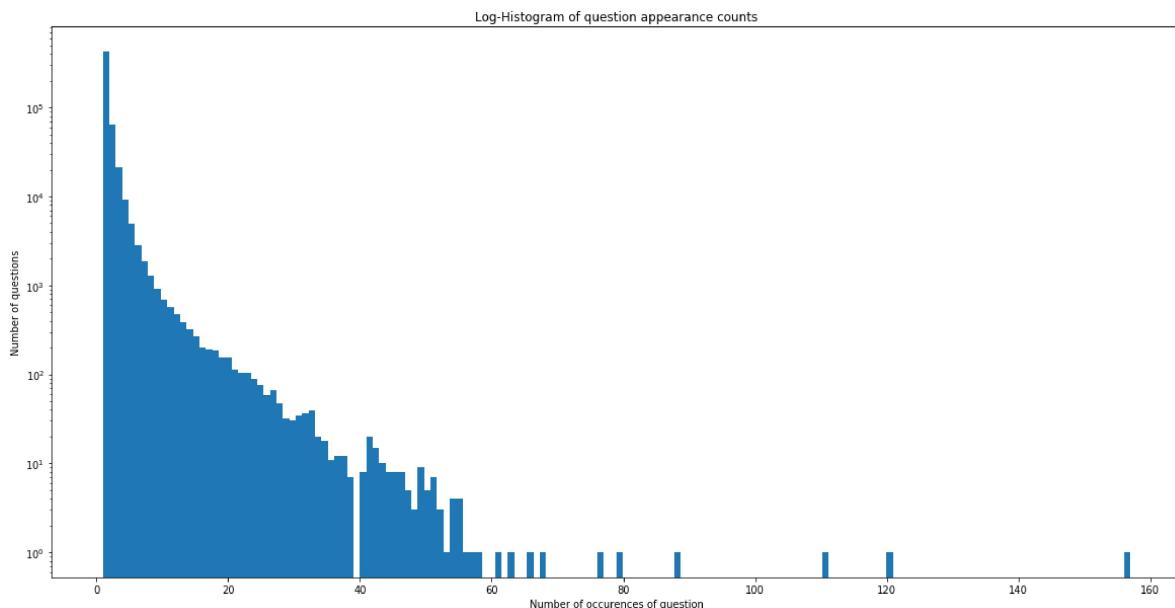
In [0]:

```

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

```

Maximum number of times a single question is repeated: 157



Checking for NULL values

In [0]:

```

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

```

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

- There are two rows with null values in question2

In [0]:

```

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

```

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

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 [0]:

```

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

```

Out[20]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words
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	
1	1	3	4	Kohinoor (Koh-i-Noor) Dia...	What is the story of Kohinoor (Koh-i-Noor) Dia...	0	4	1	51	88	8	
				How can I increase ..	How can Internet							

Analysis of some of the extracted features

- Here are some questions have only one single words.

In [0]:

```

1 print ("Minimum length of the questions in question1 : " , min(df['q1_n_words']))
2
3 print ("Minimum length of the questions in question2 : " , min(df['q2_n_words']))
4
5 print ("Number of Questions with minimum length [question1] :", df[df['q1_n_words']==1].shape[0])
6 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

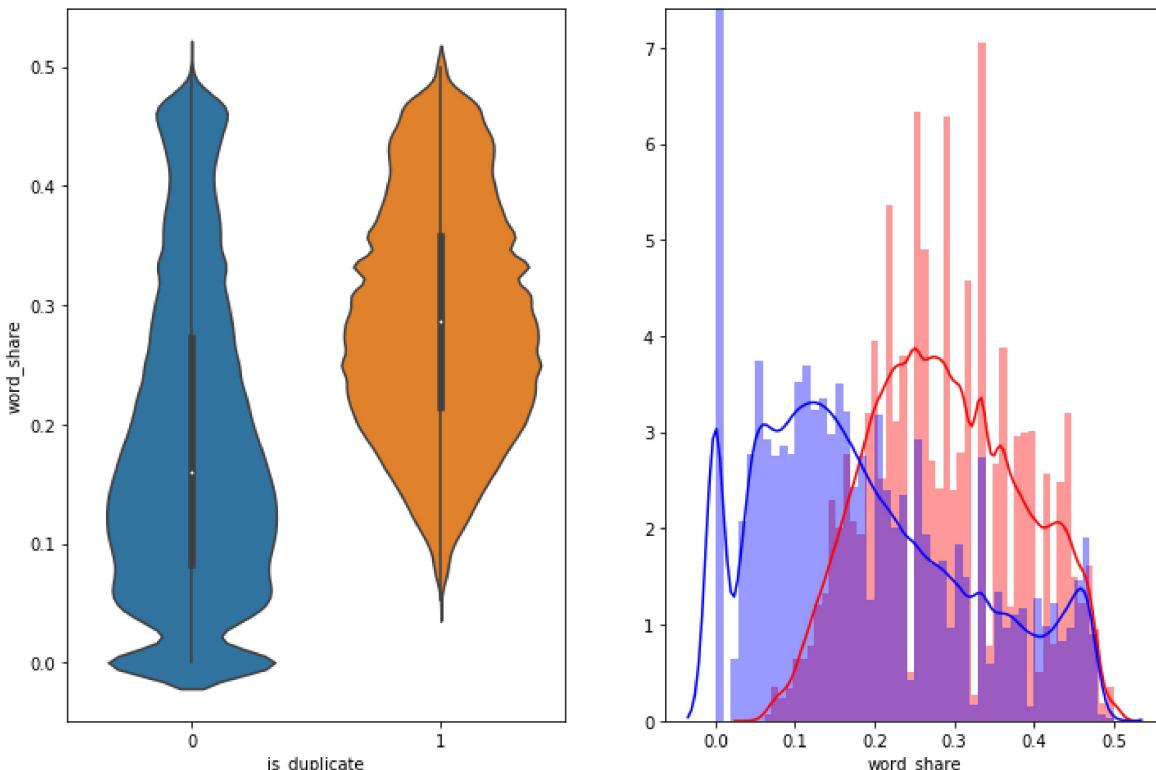
Feature: word_share

In [0]:

```

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

```

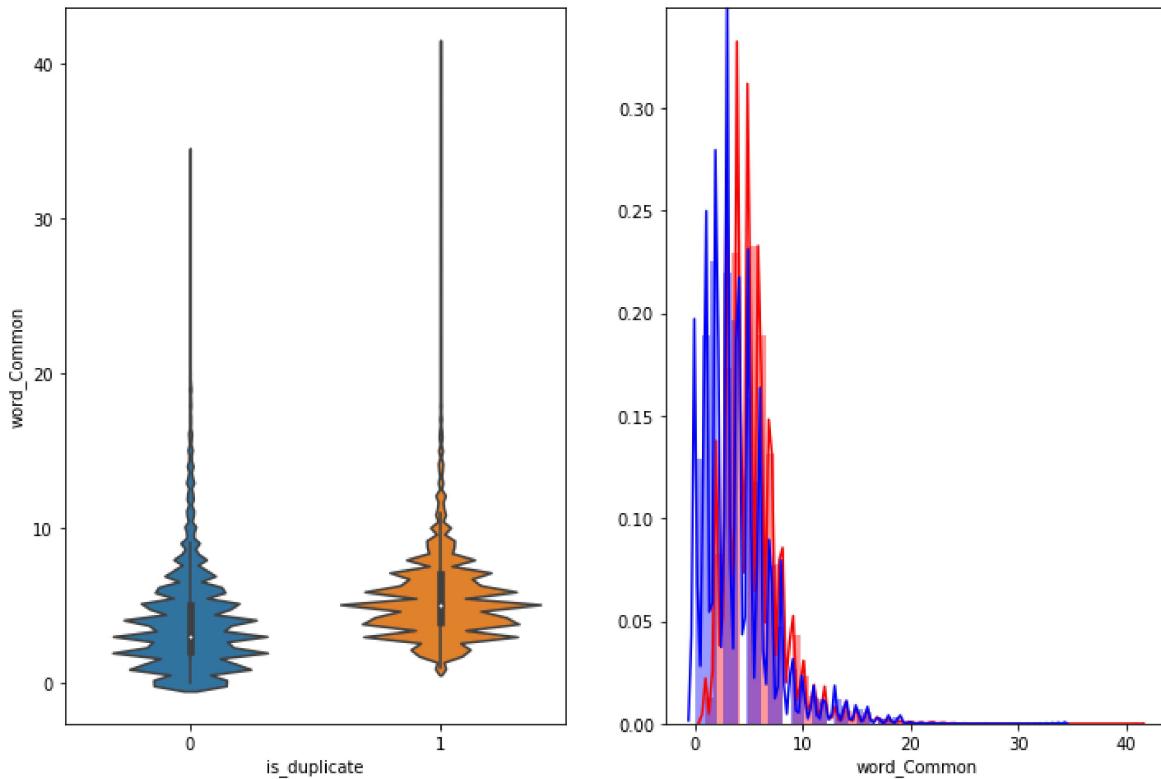


- The distributions for normalized word_share have some overlap on the far right-hand side, i.e., there are quite a lot of questions with high word similarity
- The average word share and Common no. of words of qid1 and qid2 is more when they are duplicate(Similar)

Feature: word_Common

In [0]:

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



The distributions of the word_Common feature in similar and non-similar questions are highly overlapping

EDA: Advanced Feature Extraction.

In [0]:

```
1 import warnings
2 warnings.filterwarnings("ignore")
3 import numpy as np
4 import pandas as pd
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 from subprocess import check_output
8 %matplotlib inline
9 import plotly.offline as py
10 py.init_notebook_mode(connected=True)
11 import plotly.graph_objs as go
12 import plotly.tools as tls
13 import os
14 import gc
15
16 import re
17 from nltk.corpus import stopwords
18 import distance
19 from nltk.stem import PorterStemmer
20 from bs4 import BeautifulSoup
21 import re
22 from nltk.corpus import stopwords
23 import distance
24 from nltk.stem import PorterStemmer
25 from bs4 import BeautifulSoup
26 from fuzzywuzzy import fuzz
27 from sklearn.manifold import TSNE
28 from wordcloud import WordCloud, STOPWORDS
29 from os import path
30 from PIL import Image
```

In [0]:

```
1 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
2     df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
3     df = df.fillna('')
4     df.head()
5 else:
6     print("get df_fe_without_preprocessing_train.csv from drive or run the previous not
```

In [0]:

1 df.head(2)

Out[8]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	

Preprocessing of Text

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

In [0]:

```

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

```

- 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 lenght 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 lenght 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>

$$(\text{https://github.com/seatgeek/fuzzywuzzy\#usage}) \text{ http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/ } (\text{http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/})$$

- **fuzz_partial_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>

$$(\text{https://github.com/seatgeek/fuzzywuzzy\#usage}) \text{ http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/ } (\text{http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/})$$

- **token_sort_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>

$$(\text{https://github.com/seatgeek/fuzzywuzzy\#usage}) \text{ http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/ } (\text{http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/})$$

- **token_set_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage>

$$(\text{https://github.com/seatgeek/fuzzywuzzy\#usage}) \text{ http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/ } (\text{http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/})$$

- **longest_substr_ratio** : Ratio of length longest common substring to min length of token count of Q1 and Q2

```
longest_substr_ratio = len(longest common substring) / (min(len(q1_tokens), len(q2_tokens)))
```

In [0]:

```

1 def get_token_features(q1, q2):
2     token_features = [0.0]*10
3
4     # Converting the Sentence into Tokens:
5     q1_tokens = q1.split()
6     q2_tokens = q2.split()
7
8     if len(q1_tokens) == 0 or len(q2_tokens) == 0:
9         return token_features
10    # Get the non-stopwords in Questions
11    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
12    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])
13
14    #Get the stopwords in Questions
15    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
16    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])
17
18    # Get the common non-stopwords from Question pair
19    common_word_count = len(q1_words.intersection(q2_words))
20
21    # Get the common stopwords from Question pair
22    common_stop_count = len(q1_stops.intersection(q2_stops))
23
24    # Get the common Tokens from Question pair
25    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))
26
27
28    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)) + SAFE_DIV)
29    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)) + SAFE_DIV)
30    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)) + SAFE_DIV)
31    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)) + SAFE_DIV)
32    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
33    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
34
35    # Last word of both question is same or not
36    token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])
37
38    # First word of both question is same or not
39    token_features[7] = int(q1_tokens[0] == q2_tokens[0])
40
41    token_features[8] = abs(len(q1_tokens) - len(q2_tokens))
42
43    #Average Token Length of both Questions
44    token_features[9] = (len(q1_tokens) + len(q2_tokens))/2
45    return token_features
46
47 # get the Longest Common sub string
48
49 def get_longest_substr_ratio(a, b):
50     strs = list(distance.lcsubstrings(a, b))
51     if len(strs) == 0:
52         return 0
53     else:
54         return len(strs[0]) / (min(len(a), len(b)) + 1)
55
56 def extract_features(df):
57     # preprocessing each question
58     df["question1"] = df["question1"].fillna("").apply(preprocess)
59     df["question2"] = df["question2"].fillna("").apply(preprocess)

```

```

60
61     print("token features...")
62
63     # Merging Features with dataset
64
65     token_features = df.apply(lambda x: get_token_features(x["question1"], x["question2"]))
66
67     df[ "cwc_min" ]      = list(map(lambda x: x[0], token_features))
68     df[ "cwc_max" ]      = list(map(lambda x: x[1], token_features))
69     df[ "csc_min" ]      = list(map(lambda x: x[2], token_features))
70     df[ "csc_max" ]      = list(map(lambda x: x[3], token_features))
71     df[ "ctc_min" ]      = list(map(lambda x: x[4], token_features))
72     df[ "ctc_max" ]      = list(map(lambda x: x[5], token_features))
73     df[ "last_word_eq" ] = list(map(lambda x: x[6], token_features))
74     df[ "first_word_eq" ]= list(map(lambda x: x[7], token_features))
75     df[ "abs_len_diff" ] = list(map(lambda x: x[8], token_features))
76     df[ "mean_len" ]      = list(map(lambda x: x[9], token_features))
77
78     #Computing Fuzzy Features and Merging with Dataset
79
80     print("fuzzy features...")
81
82     df[ "token_set_ratio" ]      = df.apply(lambda x: fuzz.token_set_ratio(x["question1"], x["question2"]))
83     # The token sort approach involves tokenizing the string in question, sorting the tokens
84     # then joining them back into a string We then compare the transformed strings with
85     df[ "token_sort_ratio" ]      = df.apply(lambda x: fuzz.token_sort_ratio(x["question1"], x["question2"]))
86     df[ "fuzz_ratio" ]           = df.apply(lambda x: fuzz.QRatio(x["question1"], x["question2"]))
87     df[ "fuzz_partial_ratio" ]    = df.apply(lambda x: fuzz.partial_ratio(x["question1"], x["question2"]))
88     df[ "longest_substr_ratio" ]  = df.apply(lambda x: get_longest_substr_ratio(x["question1"], x["question2"]))
89
90     return df

```

In [0]:

```

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

```

Out[12]:

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

2 rows × 21 columns

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 [0]:

```

1 df_duplicate = df[df['is_duplicate'] == 1]
2 dfp_nonduplicate = df[df['is_duplicate'] == 0]
3
4 # Converting 2d array of q1 and q2 and flatten the array: like {{1,2},{3,4}} to {1,2,3,
5 p = np.vstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
6 n = np.vstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).flatten()
7
8 print ("Number of data points in class 1 (duplicate pairs) :",len(p))
9 print ("Number of data points in class 0 (non duplicate pairs) :",len(n))
10
11 #Saving the np array into a text file
12 np.savetxt('train_p.txt', p, delimiter=' ', fmt='%s')
13 np.savetxt('train_n.txt', n, delimiter=' ', fmt='%s')
```

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

In [0]:

```

1 # reading the text files and removing the Stop Words:
2 d = path.dirname('.')
3
4 textp_w = open(path.join(d, 'train_p.txt')).read()
5 textn_w = open(path.join(d, 'train_n.txt')).read()
6 stopwords = set(STOPWORDS)
7 stopwords.add("said")
8 stopwords.add("br")
9 stopwords.add(" ")
10 stopwords.remove("not")
11
12 stopwords.remove("no")
13 #stopwords.remove("good")
14 #stopwords.remove("Love")
15 stopwords.remove("like")
16 #stopwords.remove("best")
17 #stopwords.remove("!")
18 print ("Total number of words in duplicate pair questions :",len(textp_w))
19 print ("Total number of words in non duplicate pair questions :",len(textn_w))
```

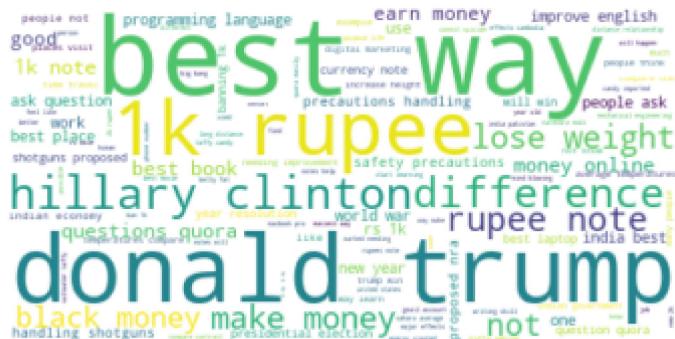
Total number of words in duplicate pair questions : 16109886
Total number of words in non duplicate pair questions : 33193130

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

In [0]:

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

Word Cloud for Duplicate Question pairs

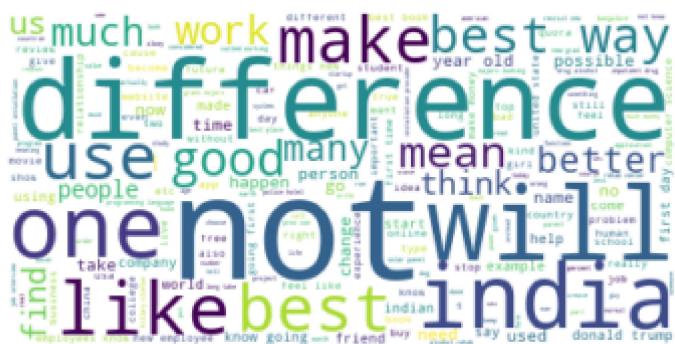


Word Clouds generated from non duplicate pair question's text

In [0]:

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

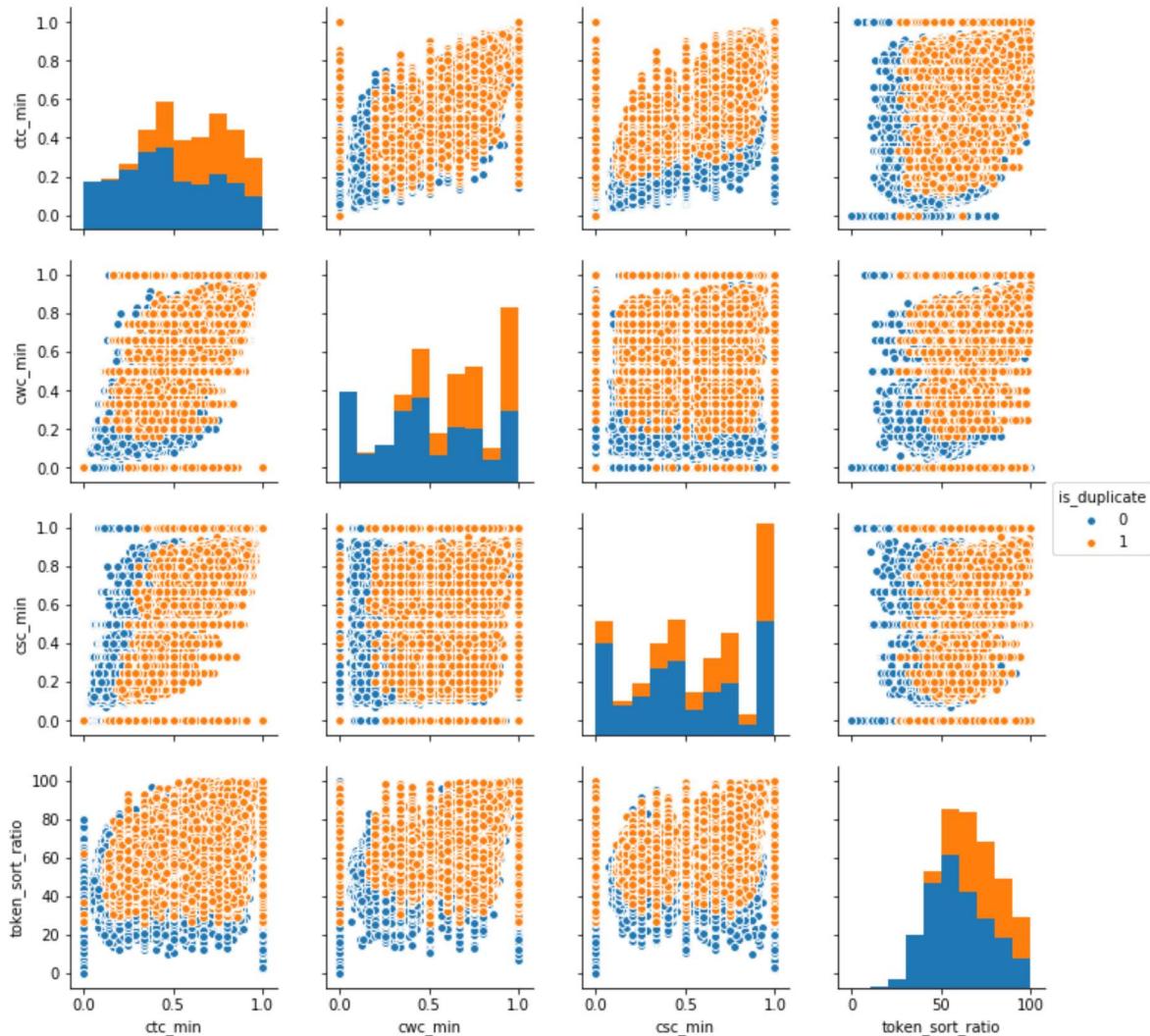
Word Cloud for non-Duplicate Question pairs:



Pair plot of features ['ctc min', 'cwc min', 'csc min', 'token sort ratio']

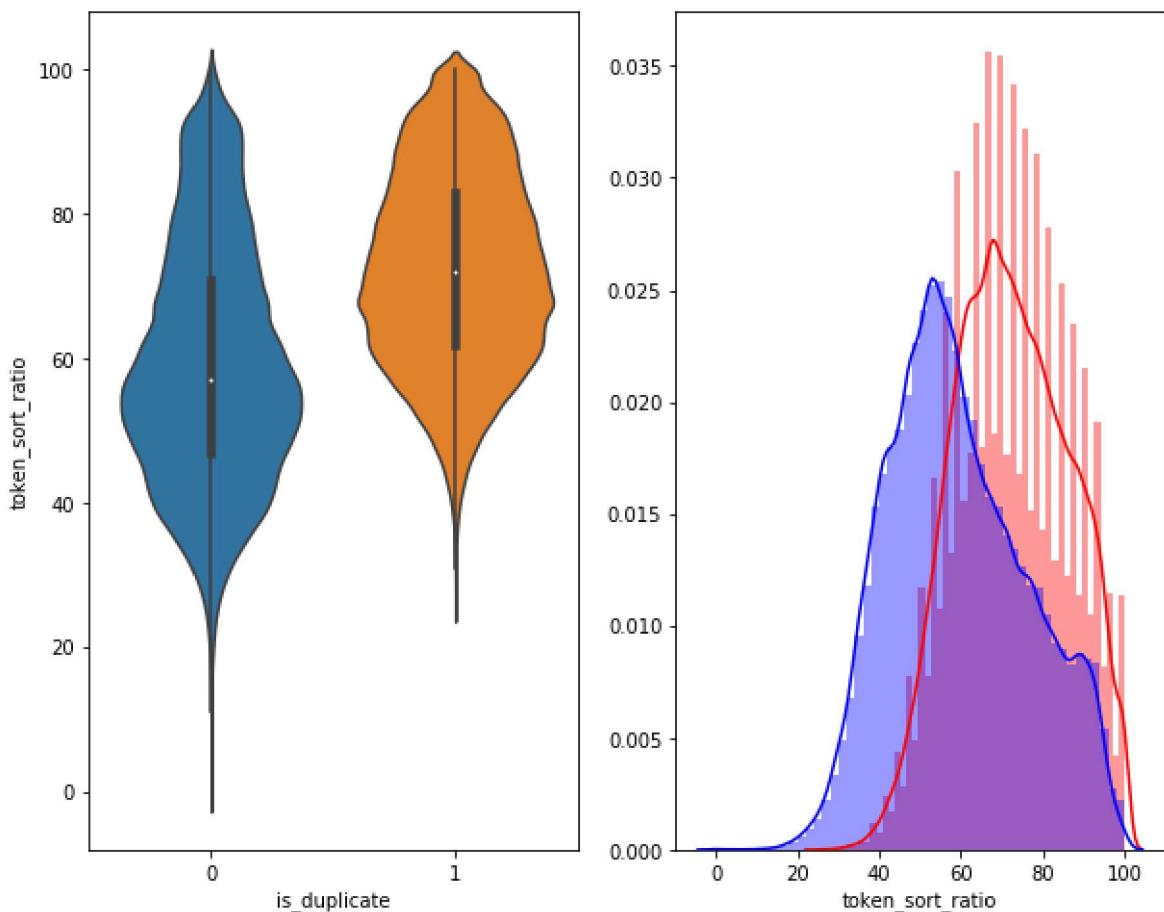
In [0]:

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



In [0]:

```
1 # Distribution of the token_sort_ratio
2 plt.figure(figsize=(10, 8))
3
4 plt.subplot(1,2,1)
5 sns.violinplot(x = 'is_duplicate', y = 'token_sort_ratio', data = df[0:] , )
6
7 plt.subplot(1,2,2)
8 sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = "1", color = "#E69138")
9 sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = "0" , color = "#4F81BD")
10 plt.show()
```

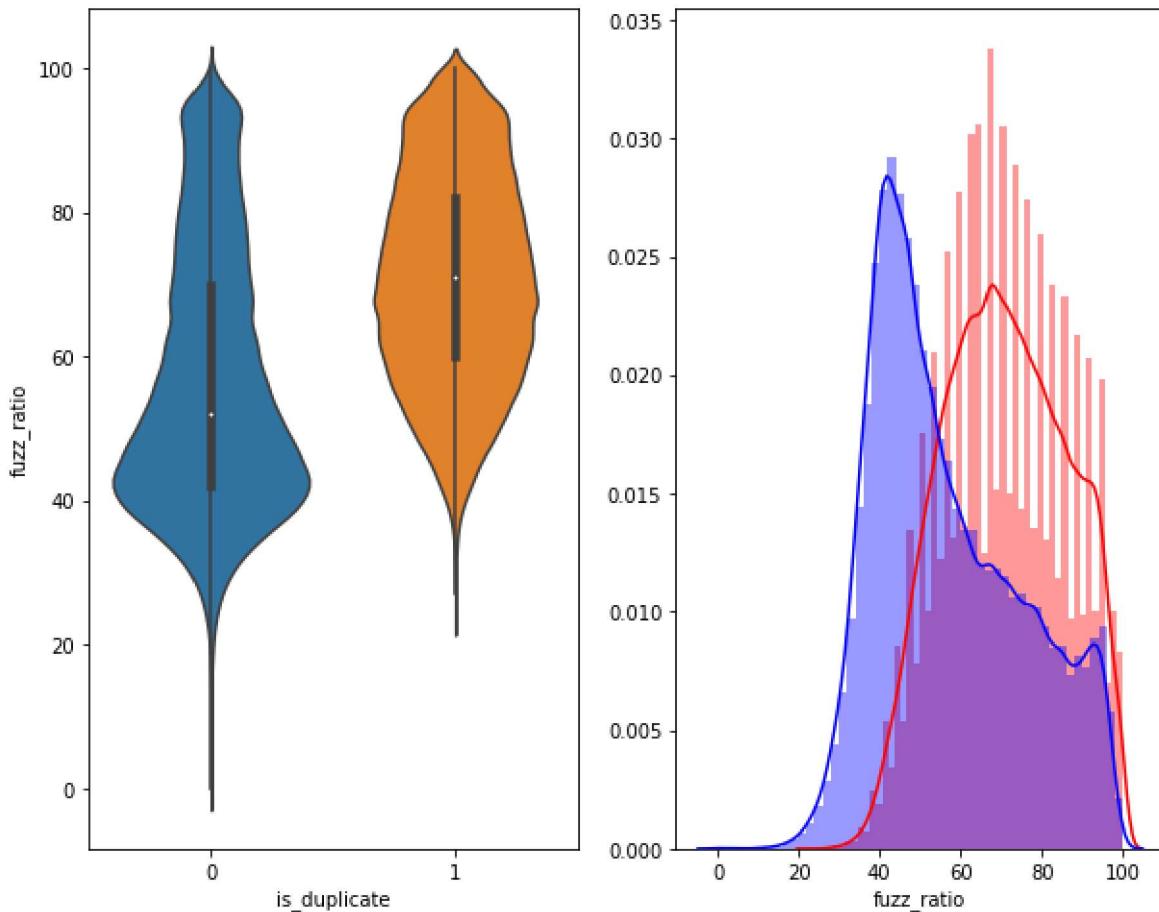


In [0]:

```

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

```



Visualization

In [0]:

```

1 # Using TSNE for Dimensionality reduction for 15 Features(Generated after cleaning the
2
3 from sklearn.preprocessing import MinMaxScaler
4
5 dfp_subsampled = df[0:5000]
6 X = MinMaxScaler().fit_transform(dfp_subsampled[['cwc_min', 'cwc_max', 'csc_min', 'csc_
7 y = dfp_subsampled['is_duplicate'].values

```

In [0]:

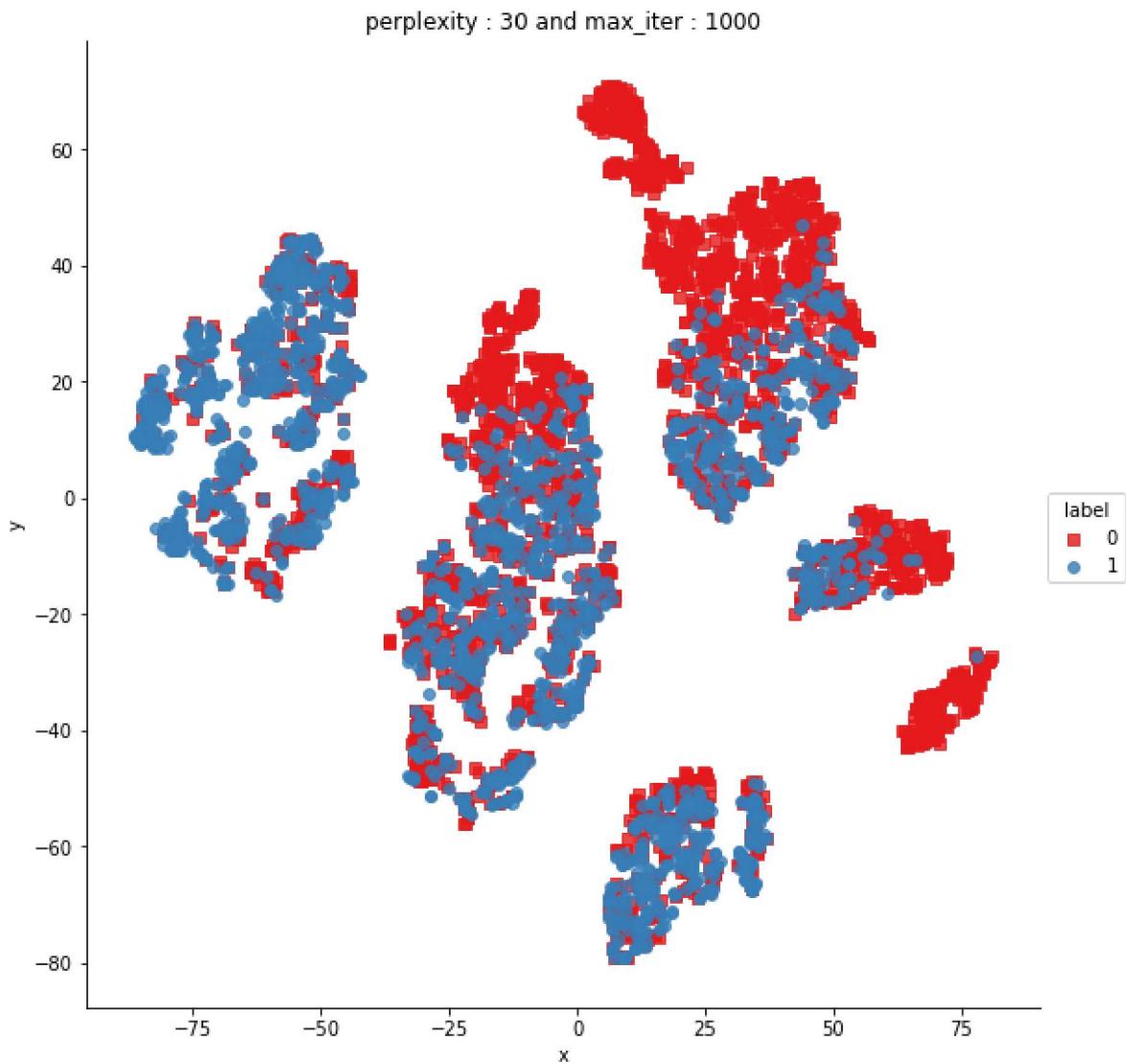
```
1 tsne2d = TSNE(  
2     n_components=2,  
3     init='random', # pca  
4     random_state=101,  
5     method='barnes_hut',  
6     n_iter=1000,  
7     verbose=2,  
8     angle=0.5  
9 ).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...  
[t-SNE] Indexed 5000 samples in 0.011s...  
[t-SNE] Computed neighbors for 5000 samples in 0.912s...  
[t-SNE] Computed conditional probabilities for sample 1000 / 5000  
[t-SNE] Computed conditional probabilities for sample 2000 / 5000  
[t-SNE] Computed conditional probabilities for sample 3000 / 5000  
[t-SNE] Computed conditional probabilities for sample 4000 / 5000  
[t-SNE] Computed conditional probabilities for sample 5000 / 5000  
[t-SNE] Mean sigma: 0.116557  
[t-SNE] Computed conditional probabilities in 0.433s  
[t-SNE] Iteration 50: error = 80.9244080, gradient norm = 0.0428133 (50 iterations in 13.099s)  
[t-SNE] Iteration 100: error = 70.3858795, gradient norm = 0.0100968 (50 iterations in 9.067s)  
[t-SNE] Iteration 150: error = 68.6138382, gradient norm = 0.0058392 (50 iterations in 9.602s)  
[t-SNE] Iteration 200: error = 67.7700119, gradient norm = 0.0036596 (50 iterations in 9.121s)  
[t-SNE] Iteration 250: error = 67.2725067, gradient norm = 0.0034962 (50 iterations in 11.305s)  
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.272507  
[t-SNE] Iteration 300: error = 1.7737305, gradient norm = 0.0011918 (50 iterations in 8.289s)  
[t-SNE] Iteration 350: error = 1.3720417, gradient norm = 0.0004822 (50 iterations in 10.526s)  
[t-SNE] Iteration 400: error = 1.2039998, gradient norm = 0.0002768 (50 iterations in 9.600s)  
[t-SNE] Iteration 450: error = 1.1133438, gradient norm = 0.0001881 (50 iterations in 11.827s)  
[t-SNE] Iteration 500: error = 1.0579143, gradient norm = 0.0001434 (50 iterations in 8.941s)  
[t-SNE] Iteration 550: error = 1.0221983, gradient norm = 0.0001164 (50 iterations in 11.092s)  
[t-SNE] Iteration 600: error = 0.9987167, gradient norm = 0.0001039 (50 iterations in 11.467s)  
[t-SNE] Iteration 650: error = 0.9831534, gradient norm = 0.0000938 (50 iterations in 11.799s)  
[t-SNE] Iteration 700: error = 0.9722011, gradient norm = 0.0000858 (50 iterations in 12.028s)  
[t-SNE] Iteration 750: error = 0.9643636, gradient norm = 0.0000799 (50 iterations in 12.120s)  
[t-SNE] Iteration 800: error = 0.9584482, gradient norm = 0.0000785 (50 iterations in 11.867s)  
[t-SNE] Iteration 850: error = 0.9538348, gradient norm = 0.0000739 (50 iterations in 11.461s)  
[t-SNE] Iteration 900: error = 0.9496906, gradient norm = 0.0000712 (50 iterations in 11.023s)
```

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

In [0]:

```
1 df = pd.DataFrame({'x':tsne2d[:,0], 'y':tsne2d[:,1] , 'label':y})
2
3 # draw the plot in appropriate place in the grid
4 sns.lmplot(data=df, x='x', y='y', hue='label', fit_reg=False, size=8, palette="Set1", mar
5 plt.title("perplexity : {} and max_iter : {}".format(30, 1000))
6 plt.show()
```



In [0]:

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

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.010s...
[t-SNE] Computed neighbors for 5000 samples in 0.935s...
[t-SNE] Computed conditional probabilities for sample 1000 / 5000
[t-SNE] Computed conditional probabilities for sample 2000 / 5000
[t-SNE] Computed conditional probabilities for sample 3000 / 5000
[t-SNE] Computed conditional probabilities for sample 4000 / 5000
[t-SNE] Computed conditional probabilities for sample 5000 / 5000
[t-SNE] Mean sigma: 0.116557
[t-SNE] Computed conditional probabilities in 0.363s
[t-SNE] Iteration 50: error = 77.7944183, gradient norm = 0.1014017 (50 iterations in 34.931s)
[t-SNE] Iteration 100: error = 69.2682266, gradient norm = 0.0248657 (50 iterations in 15.147s)
[t-SNE] Iteration 150: error = 67.7877655, gradient norm = 0.0150941 (50 iterations in 13.761s)
[t-SNE] Iteration 200: error = 67.1991119, gradient norm = 0.0126559 (50 iterations in 13.425s)
[t-SNE] Iteration 250: error = 66.8560715, gradient norm = 0.0074975 (50 iterations in 12.904s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 66.85607
1
[t-SNE] Iteration 300: error = 1.2356015, gradient norm = 0.0007033 (50 iterations in 13.302s)
[t-SNE] Iteration 350: error = 0.9948602, gradient norm = 0.0001997 (50 iterations in 18.898s)
[t-SNE] Iteration 400: error = 0.9168936, gradient norm = 0.0001430 (50 iterations in 13.397s)
[t-SNE] Iteration 450: error = 0.8863022, gradient norm = 0.0000975 (50 iterations in 16.379s)
[t-SNE] Iteration 500: error = 0.8681002, gradient norm = 0.0000854 (50 iterations in 17.791s)
[t-SNE] Iteration 550: error = 0.8564141, gradient norm = 0.0000694 (50 iterations in 17.060s)
[t-SNE] Iteration 600: error = 0.8470711, gradient norm = 0.0000640 (50 iterations in 15.454s)
[t-SNE] Iteration 650: error = 0.8389117, gradient norm = 0.0000561 (50 iterations in 17.562s)
[t-SNE] Iteration 700: error = 0.8325295, gradient norm = 0.0000529 (50 iterations in 13.443s)
[t-SNE] Iteration 750: error = 0.8268463, gradient norm = 0.0000528 (50 iterations in 17.981s)
[t-SNE] Iteration 800: error = 0.8219477, gradient norm = 0.0000477 (50 iterations in 17.448s)
[t-SNE] Iteration 850: error = 0.8180174, gradient norm = 0.0000490 (50 iterations in 18.376s)
[t-SNE] Iteration 900: error = 0.8150476, gradient norm = 0.0000456 (50 iterations in 17.778s)
```

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

In []:

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

Featurizing text data with tfidf weighted word-vectors

In [0]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import re
4 import time
5 import warnings
6 import numpy as np
7 from nltk.corpus import stopwords
8 from sklearn.preprocessing import normalize
9 from sklearn.feature_extraction.text import CountVectorizer
10 from sklearn.feature_extraction.text import TfidfVectorizer
11 warnings.filterwarnings("ignore")
12 import sys
13 import os
14 import pandas as pd
15 import numpy as np
16 from tqdm import tqdm
17
18 # extract word2vec vectors
19 # https://github.com/explosion/spaCy/issues/1721
20 # http://landinghub.visualstudio.com/visual-cpp-build-tools
21 import spacy

```

C:\Users\brahm\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.

"This module will be removed in 0.20.", DeprecationWarning)

In [0]:

```

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

```

In [0]:

```
1 df.head()
```

Out[3]:

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

In [0]:

```

1 from sklearn.feature_extraction.text import TfidfVectorizer
2 from sklearn.feature_extraction.text import CountVectorizer
3 # merge texts
4 questions = list(df['question1']) + list(df['question2'])
5
6 tfidf = TfidfVectorizer(lowercase=False, )
7 tfidf.fit_transform(questions)
8
9 # dict key:word and value:tf-idf score
10 word2tfidf = dict(zip(tfidf.get_feature_names(), tfidf.idf_))

```

In [0]:

```

1 nlp = spacy.load('en_core_web_sm')
2
3 vecs1 = []
4 for qu1 in tqdm(list(df['question1'])):
5     doc1 = nlp(qu1)
6     # 384 is the number of dimensions of vectors
7     mean_vec1 = np.zeros([len(doc1), len(doc1[0].vector)])
8     for word1 in doc1:
9         # word2vec
10        vec1 = word1.vector
11        # fetch df score
12        try:
13            idf = word2tfidf[str(word1)]
14        except:
15            idf = 0
16        # compute final vec
17        mean_vec1 += vec1 * idf
18    mean_vec1 = mean_vec1.mean(axis=0)
19    vecs1.append(mean_vec1)
20 df['q1_feats_m'] = list(vecs1)
21

```

100% |██████████| 404290/404290 [2:13:51<00:00, 50.34it/s]

In [0]:

```

1 vecs2 = []
2 for qu2 in tqdm(list(df['question2'])):
3     doc2 = nlp(qu2)
4     mean_vec1 = np.zeros([len(doc1), len(doc2[0].vector)])
5     for word2 in doc2:
6         # word2vec
7         vec2 = word2.vector
8         # fetch df score
9         try:
10             idf = word2tfidf[str(word2)]
11         except:
12             #print word
13             idf = 0
14         # compute final vec
15         mean_vec2 += vec2 * idf
16     mean_vec2 = mean_vec2.mean(axis=0)
17     vecs2.append(mean_vec2)
18 df['q2_feats_m'] = list(vecs2)

```

100% |██████████| 404290/404290 [1:47:52<00:00, 62.46it/s]

In [0]:

```

1 if os.path.isfile('nlp_features_train.csv'):
2     dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
3 else:
4     print("download nlp_features_train.csv from drive or run previous notebook")
5
6 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
7     dfppro = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
8 else:
9     print("download df_fe_without_preprocessing_train.csv from drive or run previous no")

```

In [0]:

```

1 df1 = dfnlp.drop(['qid1','qid2','question1','question2'],axis=1)
2 df2 = dfppro.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
3 df3 = df.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
4 df3_q1 = pd.DataFrame(df3.q1_feats_m.values.tolist(), index=df3.index)
5 df3_q2 = pd.DataFrame(df3.q2_feats_m.values.tolist(), index=df3.index)

```

In [0]:

```

1 # dataframe of nlp features
2 df1.head()

```

Out[9]:

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

In [0]:

```

1 # data before preprocessing
2 df2.head()

```

Out[10]:

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

In [0]:

```

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

```

Out[11]:

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

5 rows × 384 columns

In [0]:

```

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

```

Out[12]:

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

5 rows × 384 columns

In [0]:

```

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

```

Number of features in nlp dataframe : 17
Number of features in preprocessed dataframe : 12
Number of features in question1 w2v dataframe : 384
Number of features in question2 w2v dataframe : 384
Number of features in final dataframe : 794

In [0]:

```

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

```

In [0]:

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import re
4 import time
5 import warnings
6 import sqlite3
7 from sqlalchemy import create_engine # database connection
8 import csv
9 import os
10 warnings.filterwarnings("ignore")
11 import datetime as dt
12 import numpy as np
13 from nltk.corpus import stopwords
14 from sklearn.decomposition import TruncatedSVD
15 from sklearn.preprocessing import normalize
16 from sklearn.feature_extraction.text import CountVectorizer
17 from sklearn.manifold import TSNE
18 import seaborn as sns
19 from sklearn.neighbors import KNeighborsClassifier
20 from sklearn.metrics import confusion_matrix
21 from sklearn.metrics.classification import accuracy_score, log_loss
22 from sklearn.feature_extraction.text import TfidfVectorizer
23 from collections import Counter
24 from scipy.sparse import hstack
25 from sklearn.multiclass import OneVsRestClassifier
26 from sklearn.svm import SVC
27 from sklearn.cross_validation import StratifiedKFold
28 from collections import Counter, defaultdict
29 from sklearn.calibration import CalibratedClassifierCV
30 from sklearn.naive_bayes import MultinomialNB
31 from sklearn.naive_bayes import GaussianNB
32 from sklearn.model_selection import train_test_split
33 from sklearn.model_selection import GridSearchCV
34 import math
35 from sklearn.metrics import normalized_mutual_info_score
36 from sklearn.ensemble import RandomForestClassifier
37
38
39
40 from sklearn.model_selection import cross_val_score
41 from sklearn.linear_model import SGDClassifier
42 from mlxtend.classifier import StackingClassifier
43
44 from sklearn import model_selection
45 from sklearn.linear_model import LogisticRegression
46 from sklearn.metrics import precision_recall_curve, auc, roc_curve

```

C:\Users\brahm\Anaconda3\lib\site-packages\sklearn\cross_validation.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favor of the model_selection module into which all the refactored classes and functions are moved. Also note that the interface of the new CV iterators are different from that of this module. This module will be removed in 0.20.
 "This module will be removed in 0.20.", DeprecationWarning)

Machine Learning Models

Reading data from file and storing into sql table

In [0]:

```

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

```

In [0]:

```

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

```

In [0]:

```

1 read_db = 'train.db'
2 conn_r = create_connection(read_db)
3 checkTableExists(conn_r)
4 conn_r.close()

```

Tables in the databse:
data

In [0]:

```

1 if os.path.isfile(read_db):
2     conn_r = create_connection(read_db)
3     if conn_r is not None:
4         # for selecting first 1M rows
5         # data = pd.read_sql_query("""SELECT * FROM data LIMIT 100001;""", conn_r)
6
7         data = pd.read_sql_query("SELECT * From data ORDER BY RANDOM() LIMIT 100001;", conn_r.commit())
8         conn_r.close()
9

```

In [0]:

```

1 data.drop(data.index[0], inplace=True)
2 y_true = data['is_duplicate']
3 data.drop(['Unnamed: 0', 'id','index','is_duplicate'], axis=1, inplace=True)

```

In [0]:

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

Out[9]:

	cwc_min	cwc_max	csc_min	csc_max	ctc
1	0.199996000079998	0.166663888935184	0.0	0.0	0.1428551020
2	0.399992000159997	0.399992000159997	0.499987500312492	0.499987500312492	0.4444395062
3	0.833319444675922	0.714275510349852	0.999983333611106	0.857130612419823	0.6874957031
4	0.0	0.0	0.599988000239995	0.499991666805553	0.2499979166
5	0.749981250468738	0.749981250468738	0.499987500312492	0.499987500312492	0.6249921875

5 rows × 794 columns

Converting strings to numerics

In [0]:

```

1 cols = list(data.columns)
2 for i in cols:
3     data[i] = data[i].apply(pd.to_numeric)
4 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
    .

```

In [0]:

```

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

```

Random train test split(70:30)

In [0]:

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

In [0]:

```

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

```

Number of data points in train data : (70000, 794)
Number of data points in test data : (30000, 794)

In [0]:

```

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

----- Distribution of output variable in train data -----
Class 0: 0.6324857142857143 Class 1: 0.36751428571428574
----- Distribution of output variable in train data -----
Class 0: 0.3675 Class 1: 0.3675

```

In [0]:

```

1 # This function plots the confusion matrices given y_i, y_i_hat.
2 def plot_confusion_matrix(test_y, predict_y):
3     C = confusion_matrix(test_y, predict_y)
4     # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted
5
6     A = (((C.T)/(C.sum(axis=1))).T)
7
8     B =(C/C.sum(axis=0))
9
10    plt.figure(figsize=(20,4))
11
12    labels = [1,2]
13    # representing A in heatmap format
14    cmap=sns.light_palette("blue")
15    plt.subplot(1, 3, 1)
16    sns.heatmap(C, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
17    plt.xlabel('Predicted Class')
18    plt.ylabel('Original Class')
19    plt.title("Confusion matrix")
20
21    plt.subplot(1, 3, 2)
22    sns.heatmap(B, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
23    plt.xlabel('Predicted Class')
24    plt.ylabel('Original Class')
25    plt.title("Precision matrix")
26
27    plt.subplot(1, 3, 3)
28    # representing B in heatmap format
29    sns.heatmap(A, annot=True, cmap=cmap, fmt=".3f", xticklabels=labels, yticklabels=la
30    plt.xlabel('Predicted Class')
31    plt.ylabel('Original Class')
32    plt.title("Recall matrix")
33
34    plt.show()

```

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

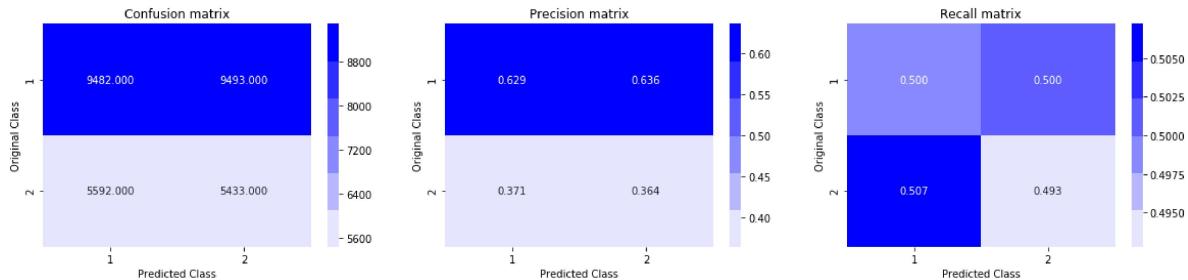
In [0]:

```

1 predicted_y = np.zeros((test_len,2))
2 for i in range(test_len):
3     rand_probs = np.random.rand(1,2)
4     predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
5 print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-1
6
7 predicted_y =np.argmax(predicted_y, axis=1)
8 plot_confusion_matrix(y_test, predicted_y)

```

Log loss on Test Data using Random Model 0.887242646958



Logistic Regression with hyperparameter tuning

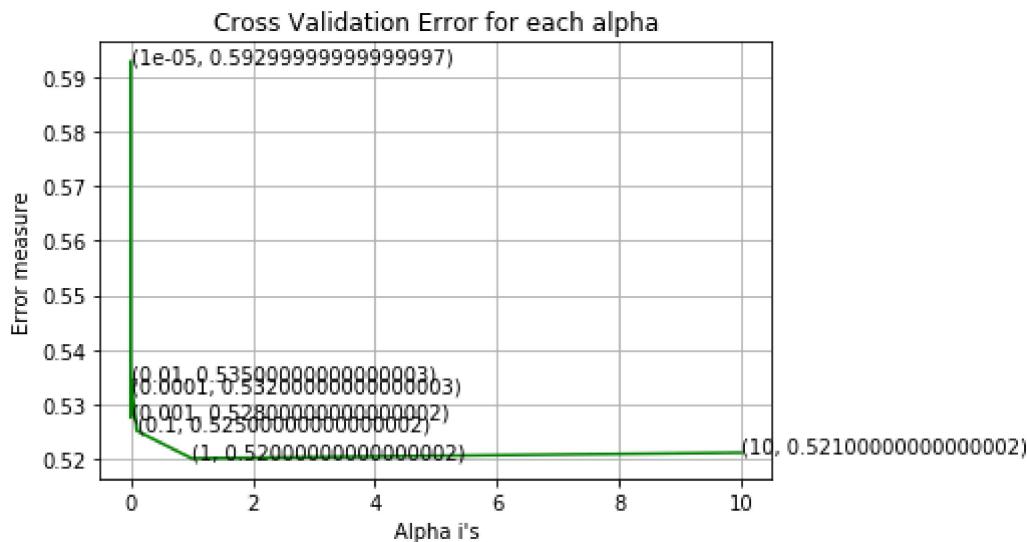
In [0]:

```

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

```

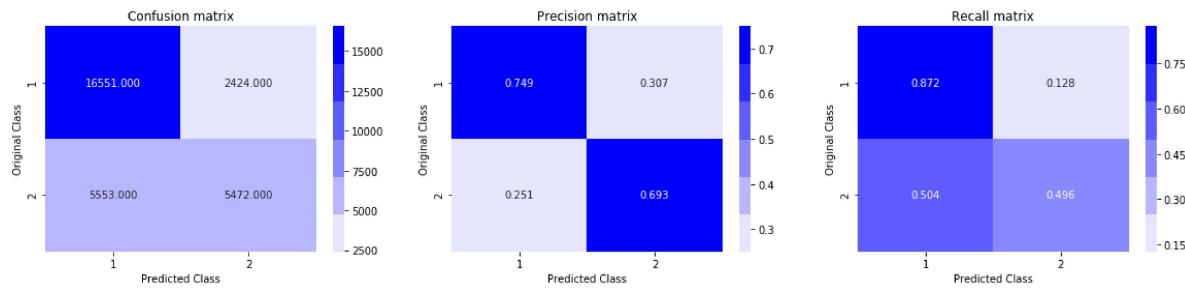
For values of alpha = 1e-05 The log loss is: 0.592800211149
 For values of alpha = 0.0001 The log loss is: 0.532351700629
 For values of alpha = 0.001 The log loss is: 0.527562275995
 For values of alpha = 0.01 The log loss is: 0.534535408885
 For values of alpha = 0.1 The log loss is: 0.525117052926
 For values of alpha = 1 The log loss is: 0.520035530431
 For values of alpha = 10 The log loss is: 0.521097925307



For values of best alpha = 1 The train log loss is: 0.513842874233

For values of best alpha = 1 The test log loss is: 0.520035530431

Total number of data points : 30000



Linear SVM with hyperparameter tuning

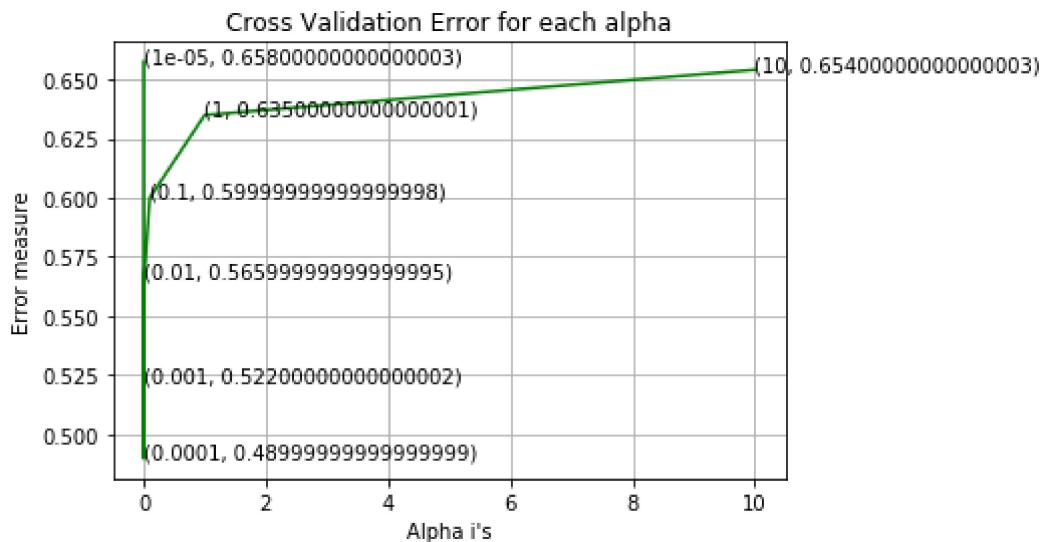
In [0]:

```

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

```

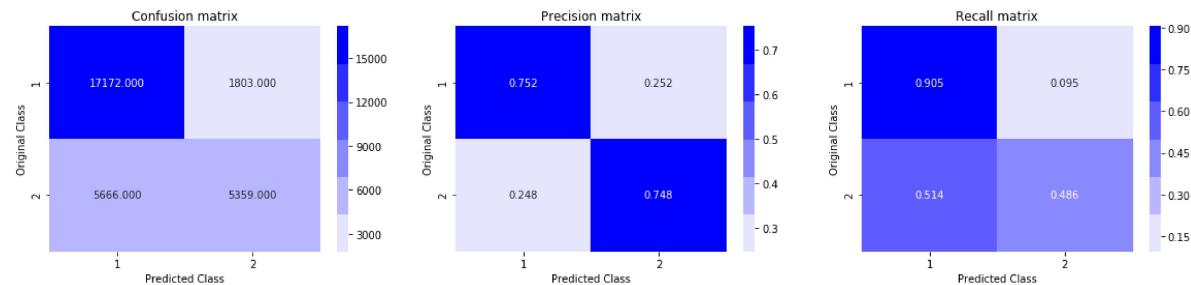
For values of alpha = 1e-05 The log loss is: 0.657611721261
 For values of alpha = 0.0001 The log loss is: 0.489669093534
 For values of alpha = 0.001 The log loss is: 0.521829068562
 For values of alpha = 0.01 The log loss is: 0.566295616914
 For values of alpha = 0.1 The log loss is: 0.599957866217
 For values of alpha = 1 The log loss is: 0.635059427016
 For values of alpha = 10 The log loss is: 0.654159467907



For values of best alpha = 0.0001 The train log loss is: 0.478054677285

For values of best alpha = 0.0001 The test log loss is: 0.489669093534

Total number of data points : 30000



XGBoost

In [0]:

```

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

```

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

Multiple eval metrics have been passed: 'valid-logloss' will be used for early stopping.

Will train until valid-logloss hasn't improved in 20 rounds.

```

[10] train-logloss:0.61583  valid-logloss:0.616104
[20] train-logloss:0.564616  valid-logloss:0.565273
[30] train-logloss:0.525758  valid-logloss:0.52679
[40] train-logloss:0.496661  valid-logloss:0.498021
[50] train-logloss:0.473563  valid-logloss:0.475182
[60] train-logloss:0.455315  valid-logloss:0.457186
[70] train-logloss:0.440442  valid-logloss:0.442482
[80] train-logloss:0.428424  valid-logloss:0.430795
[90] train-logloss:0.418803  valid-logloss:0.421447
[100] train-logloss:0.41069  valid-logloss:0.413583
[110] train-logloss:0.403831  valid-logloss:0.40693
[120] train-logloss:0.398076  valid-logloss:0.401402
[130] train-logloss:0.393305  valid-logloss:0.396851
[140] train-logloss:0.38913  valid-logloss:0.392952
[150] train-logloss:0.385469  valid-logloss:0.389521
[160] train-logloss:0.382327  valid-logloss:0.386667
[170] train-logloss:0.379541  valid-logloss:0.384148
[180] train-logloss:0.377014  valid-logloss:0.381932
[190] train-logloss:0.374687  valid-logloss:0.379883
[200] train-logloss:0.372585  valid-logloss:0.378068
[210] train-logloss:0.370615  valid-logloss:0.376367
[220] train-logloss:0.368559  valid-logloss:0.374595
[230] train-logloss:0.366545  valid-logloss:0.372847
[240] train-logloss:0.364708  valid-logloss:0.371311
[250] train-logloss:0.363021  valid-logloss:0.369886
[260] train-logloss:0.36144  valid-logloss:0.368673
[270] train-logloss:0.359899  valid-logloss:0.367421
[280] train-logloss:0.358465  valid-logloss:0.366395
[290] train-logloss:0.357128  valid-logloss:0.365361
[300] train-logloss:0.355716  valid-logloss:0.364315
[310] train-logloss:0.354425  valid-logloss:0.363403
[320] train-logloss:0.353276  valid-logloss:0.362595
[330] train-logloss:0.352084  valid-logloss:0.361823
[340] train-logloss:0.351051  valid-logloss:0.361167
[350] train-logloss:0.349867  valid-logloss:0.36043
[360] train-logloss:0.348829  valid-logloss:0.359773

```

```
[370]    train-logloss:0.347689  valid-logloss:0.359019
[380]    train-logloss:0.346607  valid-logloss:0.358311
[390]    train-logloss:0.345568  valid-logloss:0.357674
The test log loss is: 0.357054433715
```

In [0]:

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

Total number of data points : 30000

