

Quora Question Pairs

1. Business Problem

1.1 Description

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

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

Credits: Kaggle

Problem Statement

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

1.2 Sources/Useful Links

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

Useful Links

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

1.3 Real world/Business Objectives and Constraints

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

2. Machine Learning Problem

2.1 Data

2.1.1 Data Overview

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

2.1.2 Example Data point

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

2.2 Mapping the real world problem to an ML problem

2.2.1 Type of Machine Learning Problem

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

2.2.2 Performance Metric

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

Metric(s):

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

2.3 Train and Test Construction

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

3. Exploratory Data Analysis

In [3]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import sqlite3
from sqlalchemy import create_engine # database connection
import csv
import os
warnings.filterwarnings("ignore")
import datetime as dt
import numpy as np
from nltk.corpus import stopwords
from sklearn.decomposition import TruncatedSVD
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
```

```

from sklearn.manifold import TSNE
import seaborn as sns
from sklearn.neighbors import KNeighborsClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics.classification import accuracy_score, log_loss
from sklearn.feature_extraction.text import TfidfVectorizer
from collections import Counter
from scipy.sparse import hstack
from sklearn.multiclass import OneVsRestClassifier
from sklearn.svm import SVC
from sklearn.model_selection import StratifiedKFold
from collections import Counter, defaultdict
from sklearn.calibration import CalibratedClassifierCV
from sklearn.naive_bayes import MultinomialNB
from sklearn.naive_bayes import GaussianNB
from sklearn.model_selection import train_test_split
from sklearn.model_selection import GridSearchCV
import math
from sklearn.metrics import normalized_mutual_info_score
from sklearn.ensemble import RandomForestClassifier

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

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

```

3.1 Reading data and basic stats

In [4]:

```

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

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

```

Number of data points: 404290

In [5]:

```
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 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 [6]:

```
df.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 404290 entries, 0 to 404289
Data columns (total 6 columns):
id                404290 non-null int64
qid1              404290 non-null int64
qid2              404290 non-null int64
question1         404289 non-null object

```

```
question2      404288 non-null object
is_duplicate    404290 non-null int64
dtypes: int64(4), object(2)
memory usage: 18.5+ MB
```

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

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

3.2.1 Distribution of data points among output classes

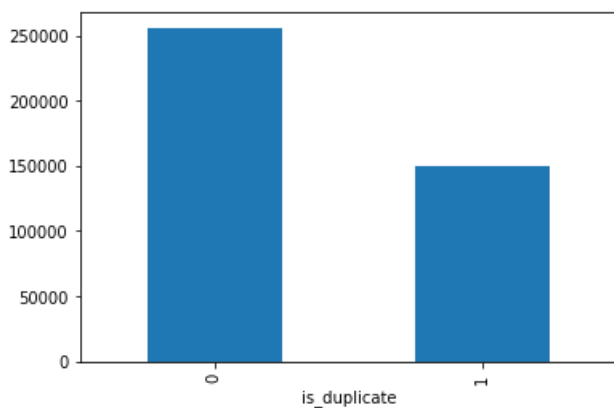
- Number of duplicate(similar) and non-duplicate(non similar) questions

In [7]:

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

Out[7]:

<matplotlib.axes._subplots.AxesSubplot at 0x179c0350fd0>



In [8]:

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

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

In [9]:

```
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 [10]:

```
qids = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())
```

```

qids = pd.Series([q1 , q2 ], dtype='<U', index=[0,1])
unique_qs = len(np.unique(qids))
qs_morethan_onetime = np.sum(qids.value_counts() > 1)
print ('Total number of Unique Questions are: {}'.format(unique_qs))
#print len(np.unique(qids))

print ('Number of unique questions that appear more than one time: {}
({}%)\n'.format(qs_morethan_onetime,qs_morethan_onetime/unique_qs*100))

print ('Max number of times a single question is repeated: {}'.format(max(qids.value_counts()))))

q_vals=qids.value_counts()

q_vals=q_vals.values

```

Total number of Unique Questions are: 537933

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

Max number of times a single question is repeated: 157

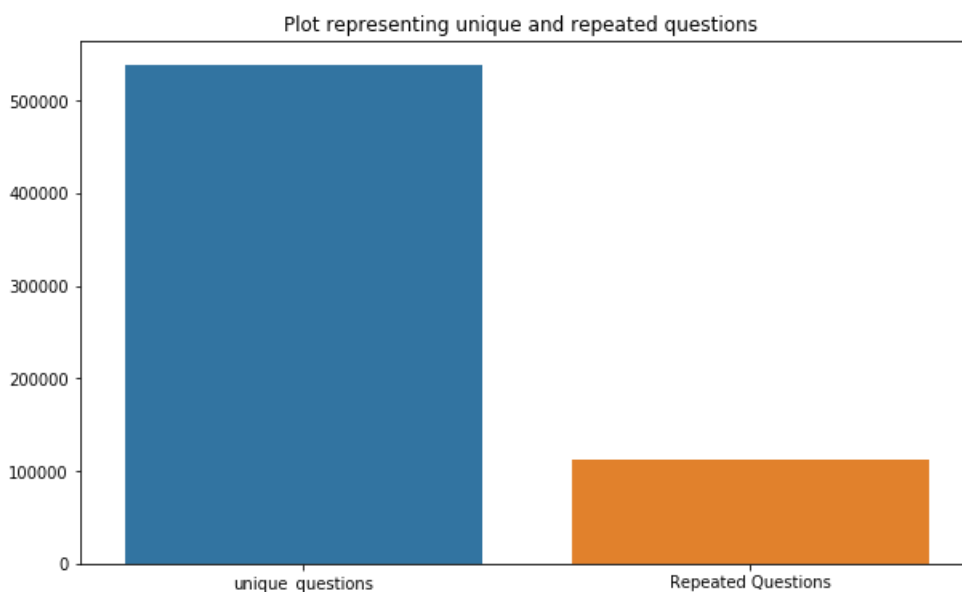
In [11]:

```

x = ["unique_questions" , "Repeated Questions"]
y = [unique_qs , qs_morethan_onetime]

plt.figure(figsize=(10, 6))
plt.title ("Plot representing unique and repeated questions ")
sns.barplot(x,y)
plt.show()

```



3.2.3 Checking for Duplicates

In [12]:

```

#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 [13]:

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

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

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

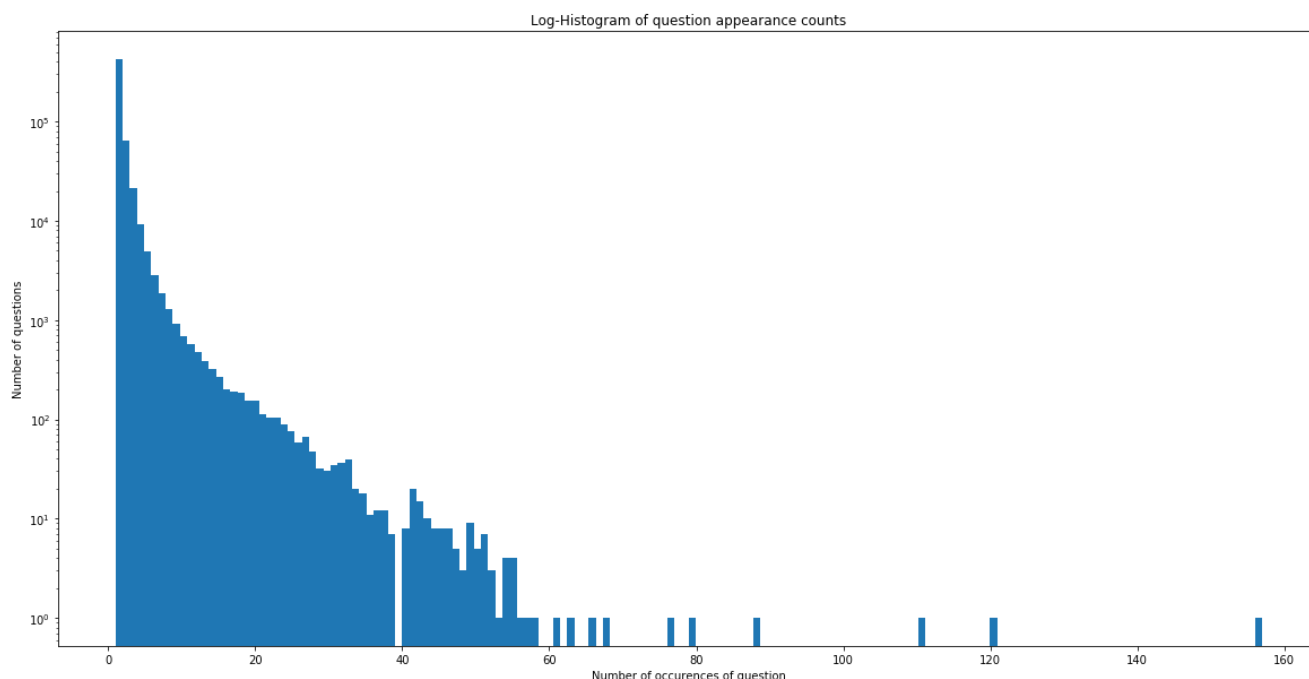
plt.title('Log-Histogram of question appearance counts')

plt.xlabel('Number of occurrences of question')

plt.ylabel('Number of questions')

print ('Maximum number of times a single question is repeated: {}'.format(max(qids.value_counts(
))))
```

Maximum number of times a single question is repeated: 157



3.2.5 Checking for NULL values

In [14]:

```
#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		
363362	My Chinese name is Haichao Yu. What English na...					0

- There are two rows with null values in question2

In [15]:

```
# Filling the null values with ' '
df = df.fillna(' ')
```

```

q1 = q1.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 [16]:

```

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[16]:

	id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common
0	0	1	2	What is the step by step guide to invest in eh	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0

id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	
1	1	3	4	What is the story of Kohinoor (Koh-i-Noor) Dia...	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0
2	2	5	6	How can I increase the speed of my internet co...	How can Internet speed be increased by hacking...	0	1	1	73	59	14	10	4.0
3	3	7	8	Why am I mentally very lonely? How can I solve...	Find the remainder when 23^{24} is divided by 24	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 [17]:

```
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 [18]:

```
plt.figure(figsize=(12, 8))

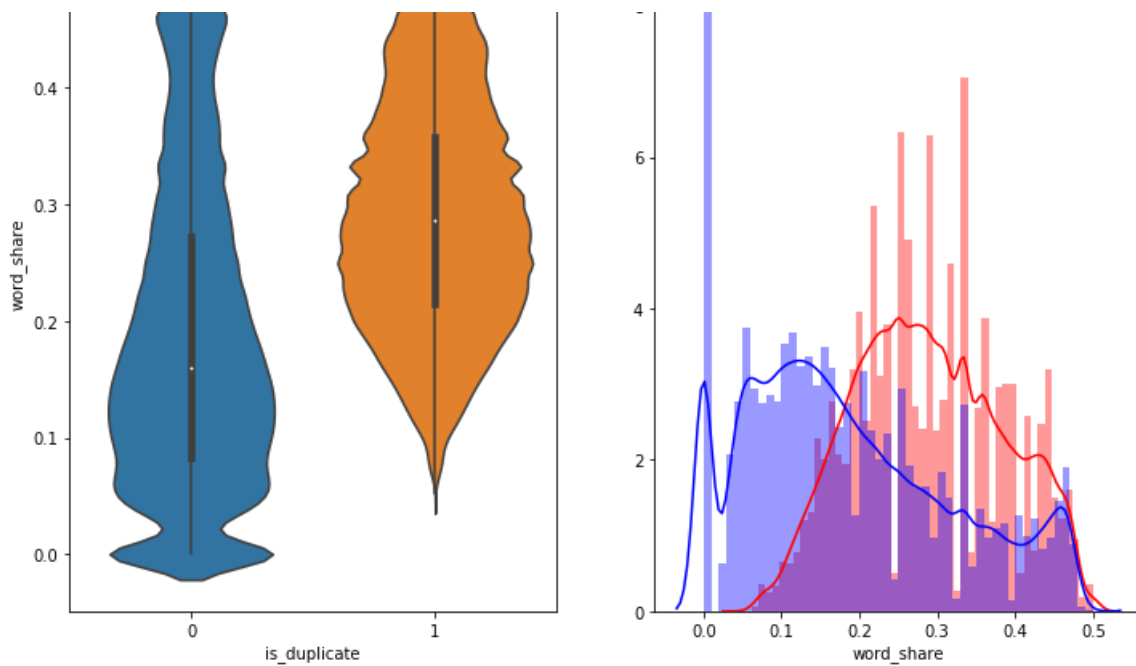
plt.subplot(1,2,1)
sns.violinplot(x = 'is_duplicate', y = 'word_share', data = df[0:])

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()
```

C:\Users\rdbz3b\AppData\Local\Continuum\anaconda3\lib\site-packages\scipy\stats\stats.py:1713: FutureWarning:

Using a non-tuple sequence for multidimensional indexing is deprecated; use `arr[tuple(seq)]` instead of `arr[seq]`. In the future this will be interpreted as an array index, `arr[np.array(seq)]`, which will result either in an error or a different result.





- 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)

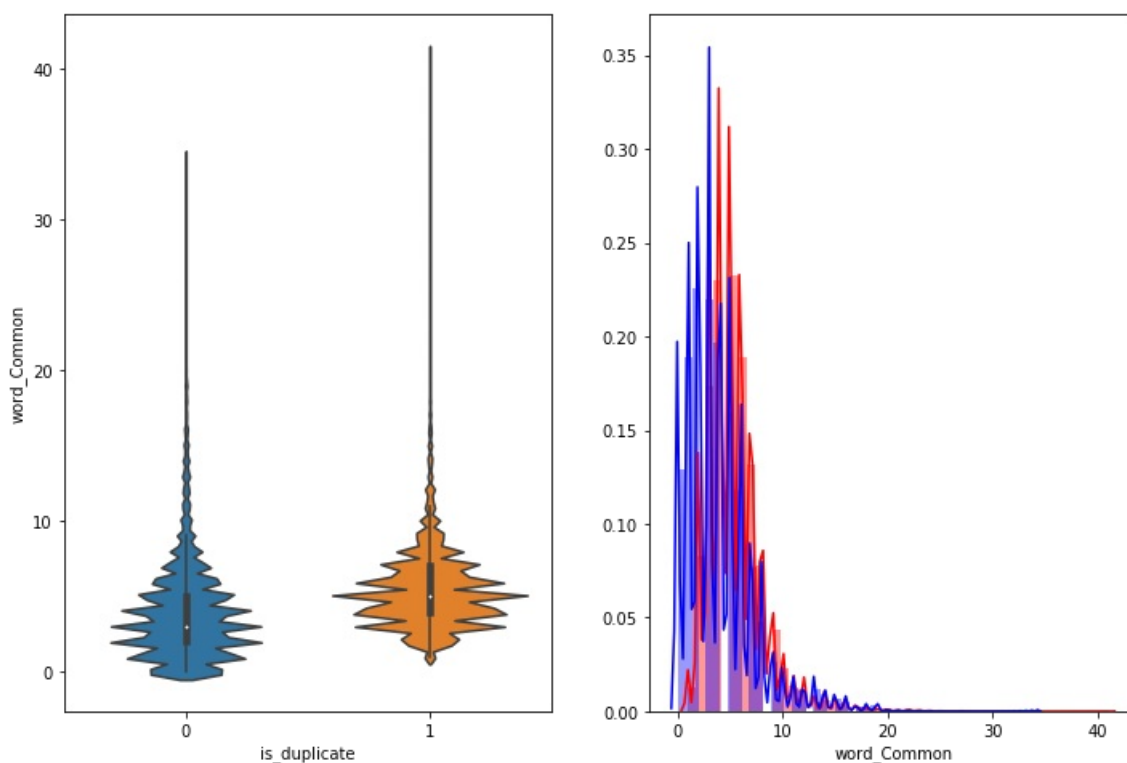
3.3.1.2 Feature: word_Common

In [19]:

```
plt.figure(figsize=(12, 8))

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

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



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

3.4 EDA: Advanced Feature Extraction.

In [20]:

```
import warnings
warnings.filterwarnings("ignore")
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
from subprocess import check_output
%matplotlib inline
import plotly.offline as py
py.init_notebook_mode(connected=True)
import plotly.graph_objs as go
import plotly.tools as tls
import os
import gc

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

In [21]:

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

In [22]:

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

Out[22]:

id	qid1	qid2	question1	question2	is_duplicate	freq_qid1	freq_qid2	q1len	q2len	q1_n_words	q2_n_words	word_Common	v
0	0	1	2	What is the step by step guide to invest in sh...	What is the step by step guide to invest in sh...	0	1	1	66	57	14	12	10.0

1	1	3	4	What is the story of Kohinoor (Koh-i-Noor)	What would happen if the Indian government sto...	0	4	1	51	88	8	13	4.0
---	---	---	---	--	---	---	---	---	----	----	---	----	-----

3.5 Preprocessing of Text

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

In [23]:

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

STOP_WORDS = stopwords.words("english")

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

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

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

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

    return x
```

3.6 Advanced Feature Extraction (NLP and Fuzzy Features)

In [24]:

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

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

    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
        return token_features

    # Get the non-stopwords in Questions
    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])
```

```

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

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

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

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

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

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

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

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

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

# get the Longest Common sub string

def get_longest_substr_ratio(a, b):
    strs = list(distance.lcs substrings(a, b))
    if len(strs) == 0:
        return 0
    else:
        return len(strs[0]) / (min(len(a), len(b)) + 1)

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

    print("token features...")

    # Merging Features with dataset

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

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

    #Computing Fuzzy Features and Merging with Dataset

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

    df["token_set_ratio"] = df.apply(lambda x: fuzz.token_set_ratio(x["question1"],
x["question2"]), axis=1)
    # The token sort approach involves tokenizing the string in question, sorting the tokens alpha
    betically, and
    # then joining them back into a string We then compare the transformed strings with a simple r

```

```

atio()).
    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 [25]:

```

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

```

Out[25]:

	id	qid1	qid2	question1	question2	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	...	ctc_max	last_word_eq	first_word
0	0	1	2	what is the step by step guide to invest in sh...	what is the step by step guide to invest in sh...	0	0.999980	0.833319	0.999983	0.999983	...	0.785709	0.0	
1	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	

2 rows × 21 columns

3.6.1 Analysis of extracted features

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

In [27]:

```

df_duplicate = df[df['is_duplicate'] == 1]
dfp_nonduplicate = df[df['is_duplicate'] == 0]

# Converting 2d array of q1 and q2 and flatten the array: like {{1,2},{3,4}} to {1,2,3,4}
p = np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
n = np.dstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).flatten()

print ("Number of data points in class 1 (duplicate pairs) :",len(p))
print ("Number of data points in class 0 (non duplicate pairs) :",len(n))

#Saving the np array into a text file
np.savetxt('train_p.txt', p, delimiter=' ', fmt='%s',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 [28]:

```

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

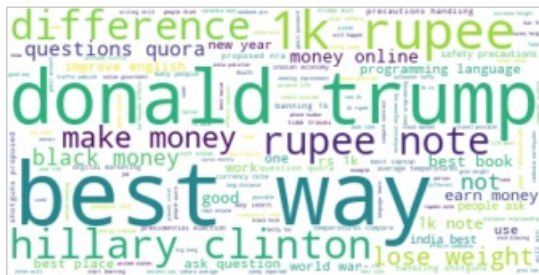
textp_w = open(path.join(d, 'train_p.txt')).read()
textn_w = open(path.join(d, 'train_n.txt')).read()

```

Total number of words in duplicate pair questions : 16110303
Total number of words in non duplicate pair questions : 33194892

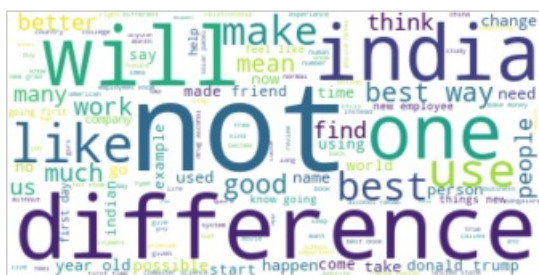
In [29]:

Word Cloud for Duplicate Question pairs



In [30]:

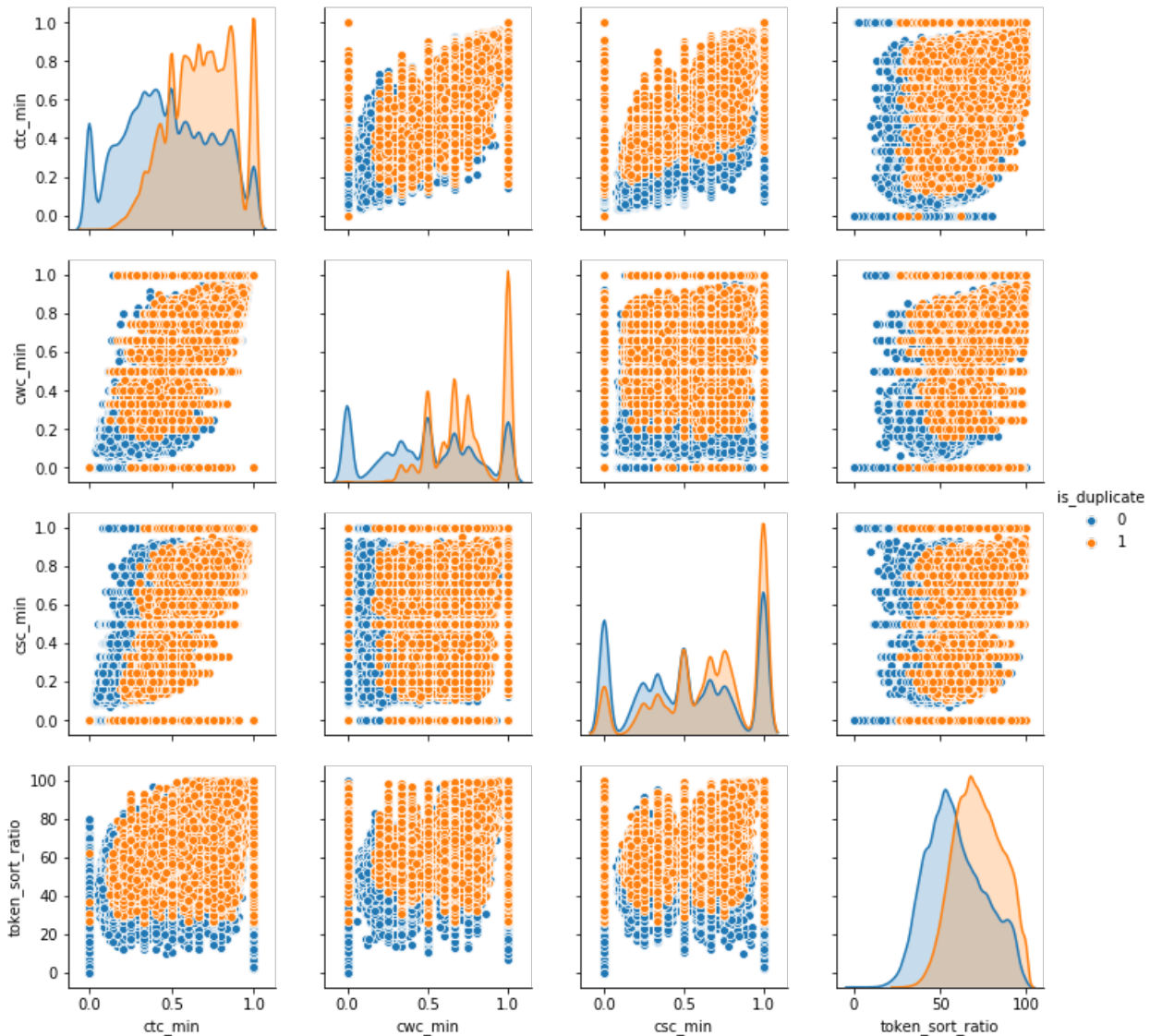
Word Cloud for non-Duplicate Question pairs:



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

In [31]:

```
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()
```

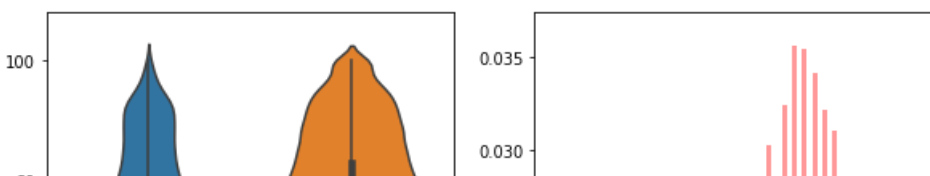


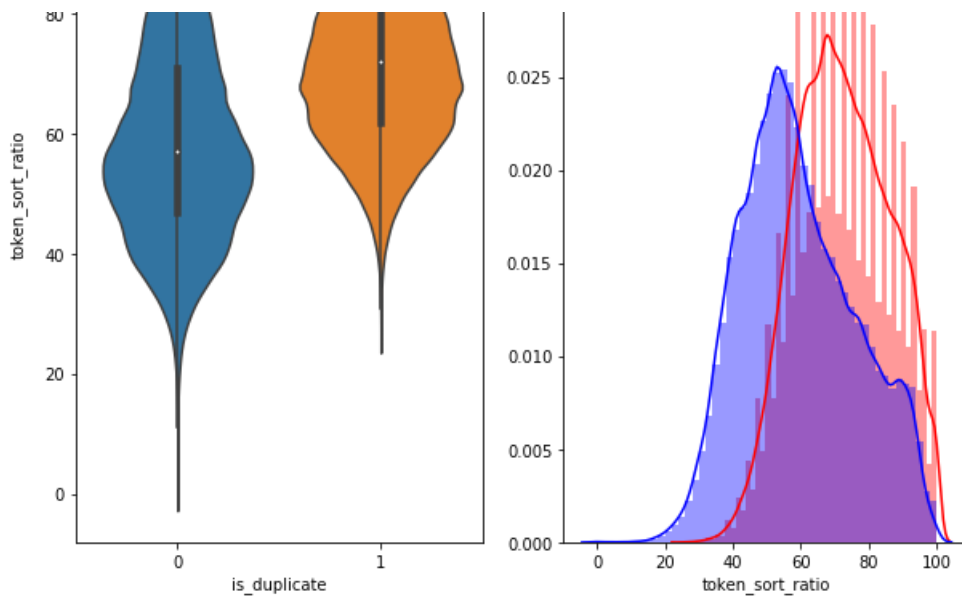
In [32]:

```
# 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()
```



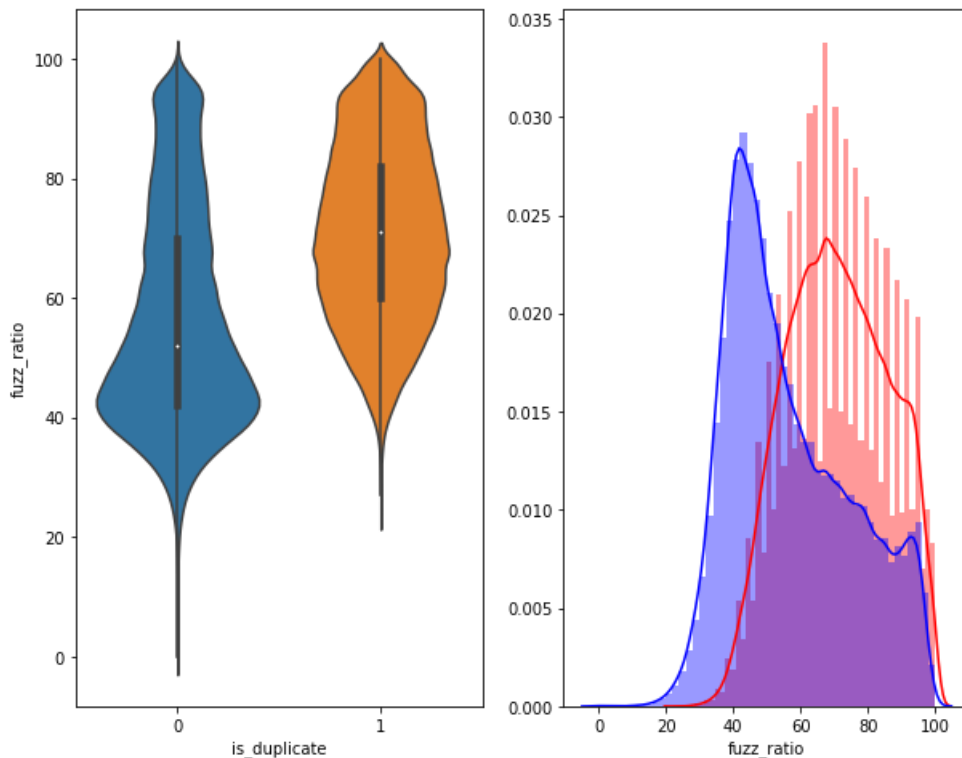


In [33]:

```
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()
```



3.6.2 Visualization

In [34]:

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

from sklearn.preprocessing import MinMaxScaler
```



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

In [35]:

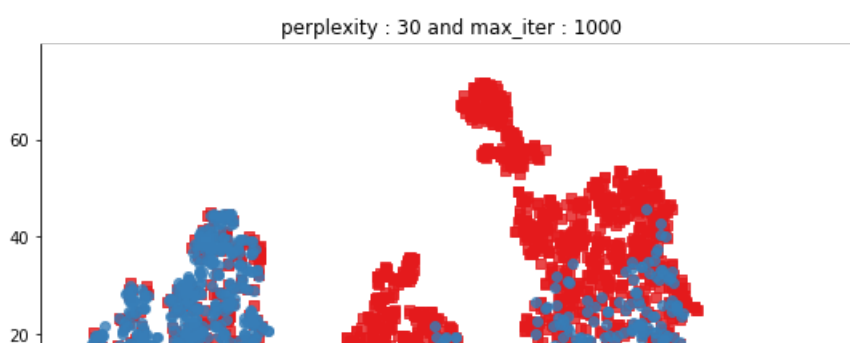
```
tsne2d = TSNE(
    n_components=2,
    init='random', # pca
    random_state=101,
    method='barnes_hut',
    n_iter=1000,
    verbose=2,
    angle=0.5
).fit_transform(X)
```

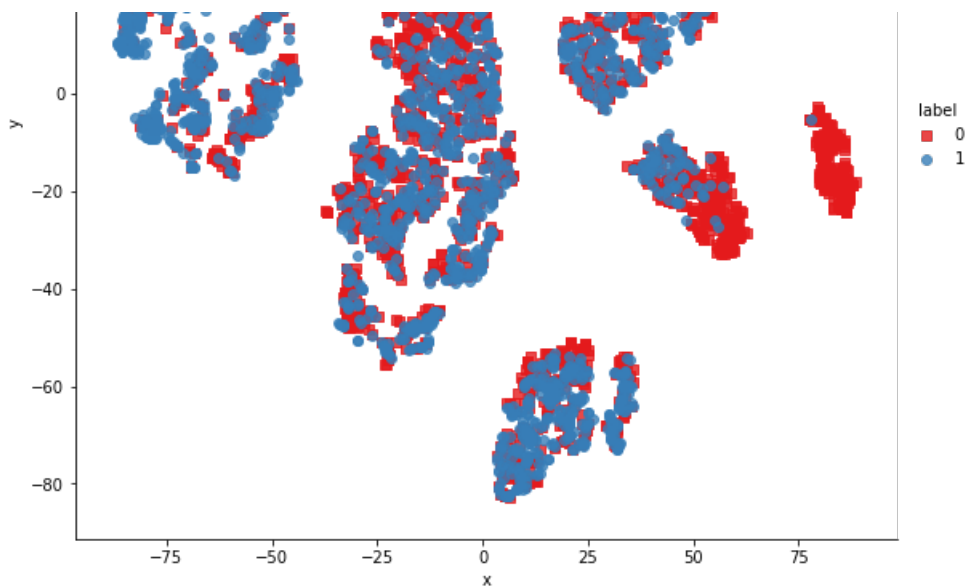
```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.010s...
[t-SNE] Computed neighbors for 5000 samples in 0.423s...
[t-SNE] Computed conditional probabilities for sample 1000 / 5000
[t-SNE] Computed conditional probabilities for sample 2000 / 5000
[t-SNE] Computed conditional 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 0.264s
[t-SNE] Iteration 50: error = 80.8968964, gradient norm = 0.0430571 (50 iterations in 6.918s)
[t-SNE] Iteration 100: error = 70.3833160, gradient norm = 0.0099593 (50 iterations in 5.301s)
[t-SNE] Iteration 150: error = 68.6159134, gradient norm = 0.0056708 (50 iterations in 5.252s)
[t-SNE] Iteration 200: error = 67.7694321, gradient norm = 0.0040581 (50 iterations in 5.470s)
[t-SNE] Iteration 250: error = 67.2746048, gradient norm = 0.0033067 (50 iterations in 5.482s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 67.274605
[t-SNE] Iteration 300: error = 1.7729300, gradient norm = 0.0011900 (50 iterations in 5.899s)
[t-SNE] Iteration 350: error = 1.3714967, gradient norm = 0.0004818 (50 iterations in 5.764s)
[t-SNE] Iteration 400: error = 1.2036748, gradient norm = 0.0002779 (50 iterations in 5.558s)
[t-SNE] Iteration 450: error = 1.1132656, gradient norm = 0.0001889 (50 iterations in 5.603s)
[t-SNE] Iteration 500: error = 1.0582460, gradient norm = 0.0001434 (50 iterations in 5.608s)
[t-SNE] Iteration 550: error = 1.0222589, gradient norm = 0.0001180 (50 iterations in 5.622s)
[t-SNE] Iteration 600: error = 0.9984865, gradient norm = 0.0001015 (50 iterations in 5.631s)
[t-SNE] Iteration 650: error = 0.9830498, gradient norm = 0.0000958 (50 iterations in 5.614s)
[t-SNE] Iteration 700: error = 0.9726909, gradient norm = 0.0000877 (50 iterations in 5.566s)
[t-SNE] Iteration 750: error = 0.9647216, gradient norm = 0.0000823 (50 iterations in 5.584s)
[t-SNE] Iteration 800: error = 0.9582971, gradient norm = 0.0000755 (50 iterations in 5.598s)
[t-SNE] Iteration 850: error = 0.9531373, gradient norm = 0.0000697 (50 iterations in 5.579s)
[t-SNE] Iteration 900: error = 0.9484153, gradient norm = 0.0000696 (50 iterations in 5.600s)
[t-SNE] Iteration 950: error = 0.9445393, gradient norm = 0.0000659 (50 iterations in 5.593s)
[t-SNE] Iteration 1000: error = 0.9412127, gradient norm = 0.0000674 (50 iterations in 5.579s)
[t-SNE] Error after 1000 iterations: 0.941213
```

In [36]:

```
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 [37]:

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

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 5000 samples in 0.016s...
[t-SNE] Computed neighbors for 5000 samples in 0.383s...
[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 0.280s
[t-SNE] Iteration 50: error = 80.3592682, gradient norm = 0.0335202 (50 iterations in 12.571s)
[t-SNE] Iteration 100: error = 69.1112671, gradient norm = 0.0036575 (50 iterations in 6.665s)
[t-SNE] Iteration 150: error = 67.6171112, gradient norm = 0.0017708 (50 iterations in 6.081s)
[t-SNE] Iteration 200: error = 67.0565109, gradient norm = 0.0011567 (50 iterations in 6.063s)
[t-SNE] Iteration 250: error = 66.7296524, gradient norm = 0.0009161 (50 iterations in 6.056s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 66.729652
[t-SNE] Iteration 300: error = 1.4983541, gradient norm = 0.0006807 (50 iterations in 7.522s)
[t-SNE] Iteration 350: error = 1.1549147, gradient norm = 0.0001922 (50 iterations in 9.185s)
[t-SNE] Iteration 400: error = 1.0101781, gradient norm = 0.0000912 (50 iterations in 9.277s)
[t-SNE] Iteration 450: error = 0.9388669, gradient norm = 0.0000628 (50 iterations in 9.194s)
[t-SNE] Iteration 500: error = 0.9029322, gradient norm = 0.0000524 (50 iterations in 9.059s)
[t-SNE] Iteration 550: error = 0.8841860, gradient norm = 0.0000482 (50 iterations in 8.988s)
[t-SNE] Iteration 600: error = 0.8722453, gradient norm = 0.0000365 (50 iterations in 8.844s)
[t-SNE] Iteration 650: error = 0.8627461, gradient norm = 0.0000347 (50 iterations in 8.669s)
[t-SNE] Iteration 700: error = 0.8549610, gradient norm = 0.0000312 (50 iterations in 8.670s)
[t-SNE] Iteration 750: error = 0.8487639, gradient norm = 0.0000311 (50 iterations in 8.683s)
[t-SNE] Iteration 800: error = 0.8440317, gradient norm = 0.0000281 (50 iterations in 8.681s)
[t-SNE] Iteration 850: error = 0.8396705, gradient norm = 0.0000250 (50 iterations in 8.699s)
[t-SNE] Iteration 900: error = 0.8354425, gradient norm = 0.0000242 (50 iterations in 8.643s)
[t-SNE] Iteration 950: error = 0.8317489, gradient norm = 0.0000233 (50 iterations in 8.678s)
[t-SNE] Iteration 1000: error = 0.8288577, gradient norm = 0.0000257 (50 iterations in 8.636s)
[t-SNE] Error after 1000 iterations: 0.828858
```

In [38]:

```
tracel = go.Scatter3d(
    x=tsne3d[:,0],
    y=tsne3d[:,1],
    z=tsne3d[:,2],
```

```

mode='markers',
marker=dict(
    sizemode='diameter',
    color = y,
    colorscale = 'Portland',
    colorbar = dict(title = 'duplicate'),
    line=dict(color='rgb(255, 255, 255)'),
    opacity=0.75
)
)

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

```

3.6 Featurizing text data with tfidf vectors

In [40]:

```

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

# encode questions to unicode
# https://stackoverflow.com/a/6812069
# ----- python 2 -----

```

```

# df['question1'] = df['question1'].apply(lambda x: unicode(str(x), "utf-8"))
# df['question2'] = df['question2'].apply(lambda x: unicode(str(x), "utf-8"))
# ----- python 3 -----
df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))

```

In [41]:

```
df.head()
```

Out[41]:

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

In [42]:

```

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import CountVectorizer
# merge texts
questions = list(df['question1']) + list(df['question2'])

tfidf = TfidfVectorizer(lowercase=False)
tfidf.fit(questions)

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

```

In [45]:

```

#prepro_features_train.csv (Simple Preprocessing Featues)
#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 [46]:

```
dfnlp.columns
```

Out[46]:

```

Index(['id', 'qid1', 'qid2', 'question1', 'question2', 'is_duplicate',
      'cwc_min', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max',
      'last_word_eq', 'first_word_eq', 'abs_len_diff', 'mean_len',
      'token_set_ratio', 'token_sort_ratio', 'fuzz_ratio',
      'fuzz_partial_ratio', 'longest_substr_ratio'],
      dtype='object')

```

In [47]:

```
df.columns
```

Out[47]:

```

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

```

```
Index(['id', 'qid1', 'qid2', 'question1', 'question2', 'is_duplicate'], dtype='object')
```

In [48]:

```
dfppro.columns
```

Out[48]:

```
Index(['id', 'qid1', 'qid2', 'question1', 'question2', 'is_duplicate',
      'freq_qid1', 'freq_qid2', 'qlen', 'q2len', 'q1_n_words', 'q2_n_words',
      'word_Common', 'word_Total', 'word_share', 'freq_q1+q2', 'freq_q1-q2'],
      dtype='object')
```

In [49]:

```
df_1 = dfppro.drop(['qid1', 'qid2', 'question1', 'question2', 'is_duplicate'],axis=1)
df_2 = dfnlp.drop(['qid1', 'qid2', 'question1', 'question2', 'is_duplicate'],axis=1)
df_3 = df.drop(['qid1', 'qid2'],axis=1)
```

```
print("df_1\n",df_1.columns)
print("df_2\n",df_2.columns)
print("df_3\n",df_3.columns)
```

```
df_1
Index(['id', 'freq_qid1', 'freq_qid2', 'qlen', 'q2len', 'q1_n_words',
      'q2_n_words', 'word_Common', 'word_Total', 'word_share', 'freq_q1+q2',
      'freq_q1-q2'],
      dtype='object')
df_2
Index(['id', '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'],
      dtype='object')
df_3
Index(['id', 'question1', 'question2', 'is_duplicate'], dtype='object')
```

In [50]:

```
df_final = df_1.merge(df_2,on="id")
print(df_final.columns)
```

```
Index(['id', 'freq_qid1', 'freq_qid2', 'qlen', 'q2len', 'q1_n_words',
      'q2_n_words', 'word_Common', 'word_Total', 'word_share', 'freq_q1+q2',
      'freq_q1-q2', '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'],
      dtype='object')
```

In [51]:

```
df_final = df_final.merge(df_3,on="id")
print(df_final.shape)
print(df_final.columns)
```

```
(404290, 30)
Index(['id', 'freq_qid1', 'freq_qid2', 'qlen', 'q2len', 'q1_n_words',
      'q2_n_words', 'word_Common', 'word_Total', 'word_share', 'freq_q1+q2',
      'freq_q1-q2', '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', 'question1', 'question2',
      'is_duplicate'],
      dtype='object')
```

In [52]:

```
from sklearn.utils import resample
df_final=resample(df_final, n_samples=100000, random_state=30)
```

In [53]:

```
df_final["is_duplicate"].value_counts()
```

Out[53]:

```
0    62864
1    37136
Name: is_duplicate, dtype: int64
```

Splitting data into train and test to avoid data leakage

In [54]:

```
x = df_final.drop(["is_duplicate", "id"], axis=1)
y = df_final["is_duplicate"]
```

Random train test split(70:30)

In [55]:

```
from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(x, y, stratify=y, test_size=0.3)
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, 28)
Number of data points in test data : (30000, 28)
```

In [56]:

```
"is_duplicate" in X_train.columns
```

Out[56]:

```
False
```

In [57]:

```
questions = list(X_train['question1']) + list(X_train['question2'])
questions1 = list(X_train['question1'])
questions2 = list(X_train['question2'])

questions = np.asarray(questions)
questions1 = np.asarray(questions1)
questions2 = np.asarray(questions2)
```

In [58]:

```
tfidf = TfidfVectorizer(lowercase=False)

# fitting on train data
tfidf.fit(questions.astype('U'))

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

questions1 = tfidf.transform(questions1.astype('U'))
questions2 = tfidf.transform(questions2.astype('U'))

test_questions1 = tfidf.transform(np.asarray(list(X_test["question1"]))).astype('U')
test_questions2 = tfidf.transform(np.asarray(list(X_test["question2"]))).astype('U')

print("*"*50)
print(questions1.shape)
print(questions2.shape)
print("*"*50)
```

```
print(test_questions1.shape)
print(test_questions2.shape)
```

```
*****
(70000, 46078)
(70000, 46078)
*****
(30000, 46078)
(30000, 46078)
```

In [59]:

```
X_train.columns
```

Out[59]:

```
Index(['freq_qid1', 'freq_qid2', 'q1len', 'q2len', 'q1_n_words', 'q2_n_words',
      'word_Common', 'word_Total', 'word_share', 'freq_q1+q2', 'freq_q1-q2',
      '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', 'question1', 'question2'],
      dtype='object')
```

In [60]:

```
train_values = X_train.drop(['question1', 'question2'],axis=1).values
test_values = X_test.drop(['question1', 'question2'],axis=1).values
```

In [61]:

```
for i in (train_values,test_values,questions1,questions2):
    print(type(i),i.shape)
```

```
<class 'numpy.ndarray'> (70000, 26)
<class 'numpy.ndarray'> (30000, 26)
<class 'scipy.sparse.csr.csr_matrix'> (70000, 46078)
<class 'scipy.sparse.csr.csr_matrix'> (70000, 46078)
```

In [62]:

```
from scipy.sparse import csr_matrix,hstack
```

In [63]:

```
# combining all features into 1 set
# here we are stacking csr matrix with dense array
# so 1st convert dense array to csr matrix.

train_values = csr_matrix(train_values)
print(train_values.shape)

test_values = csr_matrix(test_values)
print(test_values.shape)

set_ = hstack((train_values,questions1,questions2))
set_t = hstack((test_values,test_questions1,test_questions2))
print(set_t.shape)
```

```
(70000, 26)
(30000, 26)
(30000, 92182)
```

Model Making

In [64]:

```
from collections import Counter
```

```

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

```

```

----- Distribution of output variable in train data -----
Class 0:  0.6286428571428572 Class 1:  0.37135714285714283
----- Distribution of output variable in train data -----
Class 0:  0.3713666666666667 Class 1:  0.3713666666666667

```

In [65]:

```

# This function plots the confusion matrices given y_i, y_i_hat.
def plot_confusion_matrix(test_y, predict_y):
    C = confusion_matrix(test_y, predict_y)
    # C = 9,9 matrix, each cell (i,j) represents number of points of class i are predicted class j

    A = ((C.T)/(C.sum(axis=1))).T
    #divid each element of the confusion matrix with the sum of elements in that column

    # C = [[1, 2],
    #      [3, 4]]
    # C.T = [[1, 3],
    #        [2, 4]]
    # C.sum(axis = 1)  axis=0 corresonds to columns and axis=1 corresponds to rows in two
    dimensional array
    # C.sum(axix =1) = [[3, 7]]
    # ((C.T)/(C.sum(axis=1))) = [[1/3, 3/7]
    #                             [2/3, 4/7]]

    # ((C.T)/(C.sum(axis=1))).T = [[1/3, 2/3]
    #                               [3/7, 4/7]]
    # sum of row elements = 1

    B = (C/C.sum(axis=0))
    #divid each element of the confusion matrix with the sum of elements in that row
    # C = [[1, 2],
    #      [3, 4]]
    # C.sum(axis = 0)  axis=0 corresonds to columns and axis=1 corresponds to rows in two
    dimensional array
    # C.sum(axix =0) = [[4, 6]]
    # (C/C.sum(axis=0)) = [[1/4, 2/6],
    #                       [3/4, 4/6]]
    plt.figure(figsize=(20,4))

    labels = [1,2]
    # representing A in heatmap format
    cmap=sns.light_palette("blue")
    plt.subplot(1, 3, 1)
    sns.heatmap(C, annot=True, 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()

```

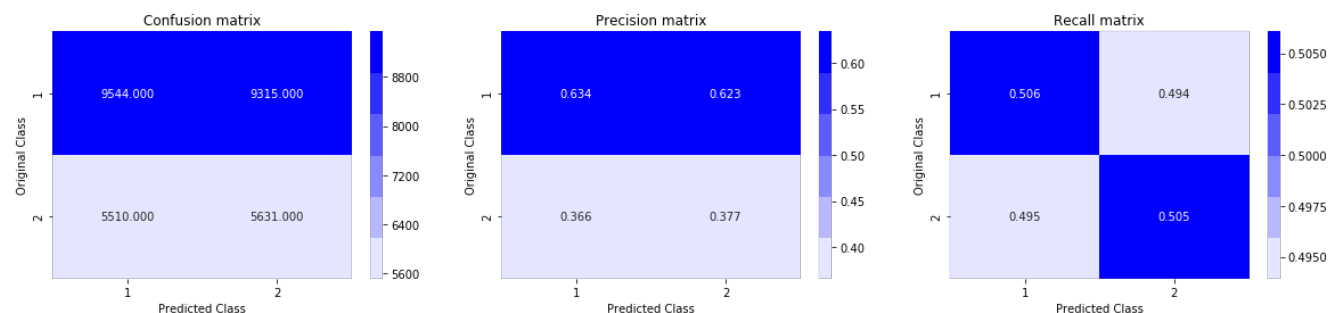

Random Model

In [66]:

```
# we need to generate 9 numbers and the sum of numbers should be 1
# one solution is to generate 9 numbers and divide each of the numbers by their sum
# ref: https://stackoverflow.com/a/18662466/4084039
# we create a output array that has exactly same size as the CV data
predicted_y = np.zeros((test_len,2))
for i in range(test_len):
    rand_probs = np.random.rand(1,2)
    predicted_y[i] = ((rand_probs/sum(sum(rand_probs)))[0])
print("Log loss on Test Data using Random Model",log_loss(y_test, predicted_y, eps=1e-15))

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

Log loss on Test Data using Random Model 0.8787922657236616



Assignment

Task 1 Logistic Regression

In [68]:

```
alpha = [0.0001,0.0005,0.0009,0.001,0.005,0.01,0.02,0.03,0.5,0.4,0.1,1,10] # 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(set_, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(set_, y_train)
    predict_y = sig_clf.predict_proba(set_t)
    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()
```

```

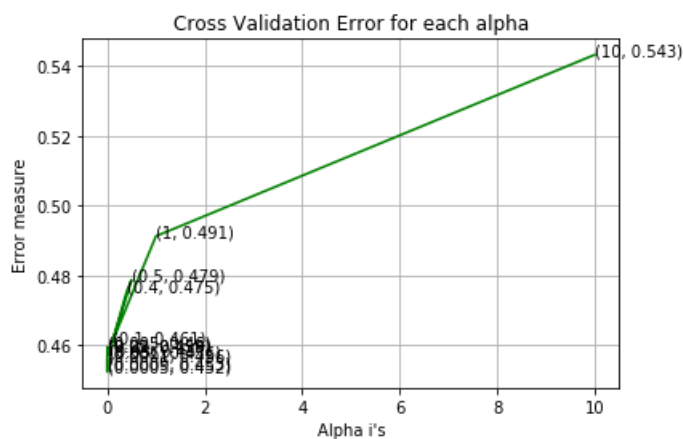
ax.plot(alpha, log_error_array,c='g')
for i, txt in enumerate(np.round(log_error_array,3)):
    ax.annotate((alpha[i],np.round(txt,3)), (alpha[i],log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error measure")
plt.show()

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

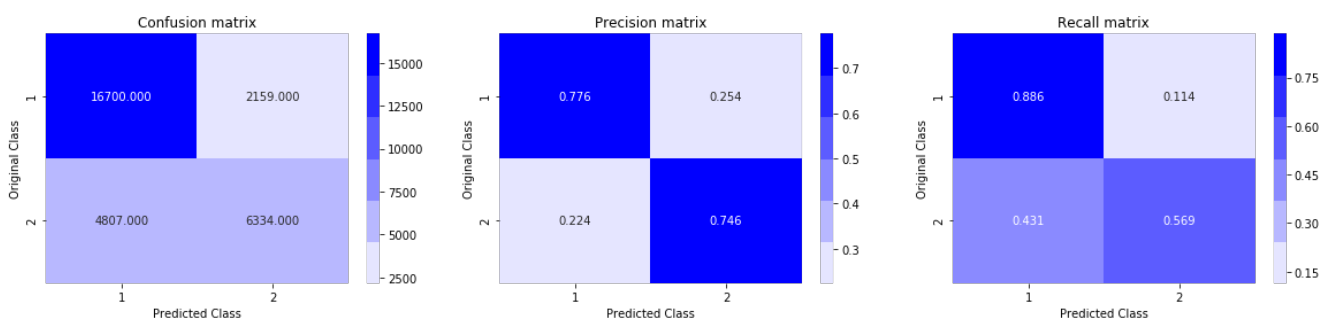
predict_y = sig_clf.predict_proba(set_)
print('For values of best alpha = ', alpha[best_alpha], "The train log loss is:",log_loss(y_train,
predict_y, labels=clf.classes_, eps=1e-15))
predict_y = sig_clf.predict_proba(set_t)
print('For values of best alpha = ', alpha[best_alpha], "The test log loss is:",log_loss(y_test, 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 = 0.0001 The log loss is: 0.45589711036717817
 For values of alpha = 0.0005 The log loss is: 0.45244603475267575
 For values of alpha = 0.0009 The log loss is: 0.4534203188525349
 For values of alpha = 0.001 The log loss is: 0.4564701809298237
 For values of alpha = 0.005 The log loss is: 0.4596113168814216
 For values of alpha = 0.01 The log loss is: 0.4587253487687639
 For values of alpha = 0.02 The log loss is: 0.45848253711796494
 For values of alpha = 0.03 The log loss is: 0.45735780002491044
 For values of alpha = 0.5 The log loss is: 0.47872487918451934
 For values of alpha = 0.4 The log loss is: 0.4754148085684936
 For values of alpha = 0.1 The log loss is: 0.46088166768983574
 For values of alpha = 1 The log loss is: 0.4912790449302387
 For values of alpha = 10 The log loss is: 0.5432193415271728



For values of best alpha = 0.0005 The train log loss is: 0.4513690350712787
 For values of best alpha = 0.0005 The test log loss is: 0.45244603475267575
 Total number of data points : 30000



Task2 Linear SVM

In [69]:

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

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

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

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

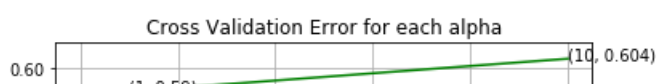
log_error_array=[]
for i in alpha:
    clf = SGDClassifier(alpha=i, penalty='l1', loss='hinge', random_state=42)
    clf.fit(set_, y_train)
    sig_clf = CalibratedClassifierCV(clf, method="sigmoid")
    sig_clf.fit(set_, y_train)
    predict_y = sig_clf.predict_proba(set_t)
    log_error_array.append(log_loss(y_test, 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))

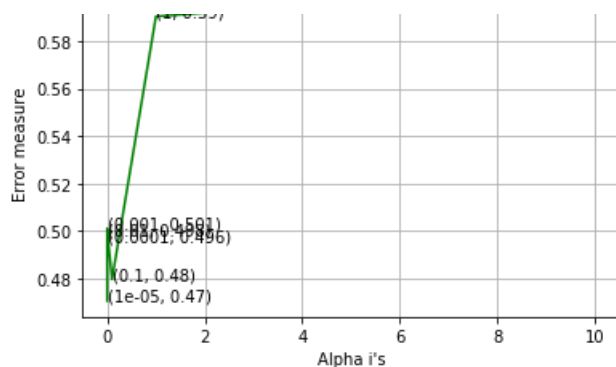
fig, ax = plt.subplots()
ax.plot(alpha, log_error_array, c='g')
for i, txt in enumerate(np.round(log_error_array, 3)):
    ax.annotate((alpha[i], np.round(txt, 3)), (alpha[i], log_error_array[i]))
plt.grid()
plt.title("Cross Validation Error for each alpha")
plt.xlabel("Alpha i's")
plt.ylabel("Error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

predict_y = sig_clf.predict_proba(set_)
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(set_t)
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: 0.47049602075990876
For values of alpha = 0.0001 The log loss is: 0.495906483627825
For values of alpha = 0.001 The log loss is: 0.5009965701080631
For values of alpha = 0.01 The log loss is: 0.49810173365890414
For values of alpha = 0.1 The log loss is: 0.47961150952361264
For values of alpha = 1 The log loss is: 0.5903394603865542
For values of alpha = 10 The log loss is: 0.6041327163127198
```

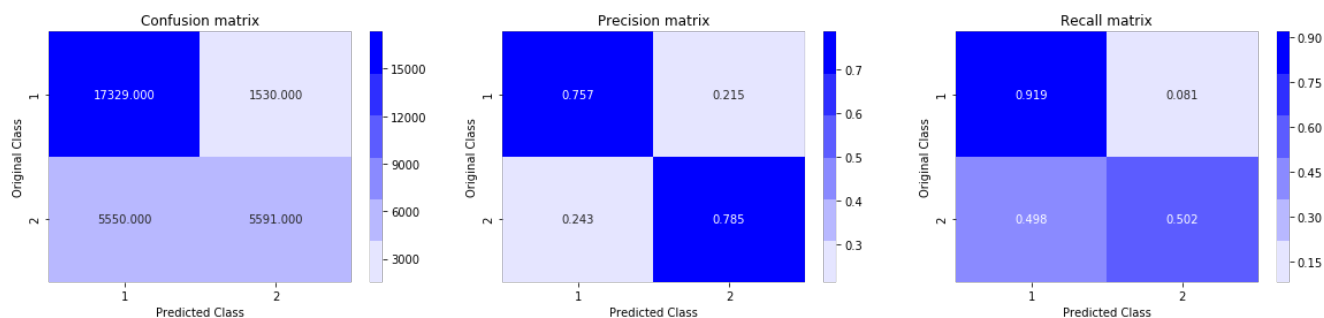




For values of best alpha = 1e-05 The train log loss is: 0.4665738151994522

For values of best alpha = 1e-05 The test log loss is: 0.47049602075990876

Total number of data points : 30000



Task3 XGBoost Tuning

In [70]:

```
# importing necessary Libraries
from sklearn.model_selection import RandomizedSearchCV
from scipy import stats
from xgboost import XGBClassifier
```

In [71]:

```
xgb = XGBClassifier(n_jobs=-1)
print(xgb.get_params())
```

```
<bound method XGBModel.get_params of XGBClassifier(base_score=0.5, booster='gbtree',
colsample_bylevel=1,
colsample_bytree=1, gamma=0, learning_rate=0.1, max_delta_step=0,
max_depth=3, min_child_weight=1, missing=None, n_estimators=100,
n_jobs=-1, nthread=None, objective='binary:logistic',
random_state=0, reg_alpha=0, reg_lambda=1, scale_pos_weight=1,
seed=None, silent=True, subsample=1)>
```

XGBOOST Tuning1

In [72]:

```
param_dist = {'n_estimators': stats.randint(10, 100),
              'learning_rate': stats.uniform(0.01, 0.05),
              'subsample': stats.uniform(0.3, 0.5),
              'max_depth': stats.randint(5, 30),
              'colsample_bytree': stats.uniform(0.3, 0.4),
              'min_child_weight': [1, 2]
            }
```

In [73]:

```
model=RandomizedSearchCV(XGBClassifier(n_jobs=-1), param_distributions=param_dist, scoring = 'neg_l
og_loss', cv=3)
```

```

sg_1000, cv=3)
model.fit(set_, y_train)
print(model.best_estimator_)
print(model.score(set_t, y_test))

```

```

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
               colsample_bytree=0.6606035469632427, gamma=0,
               learning_rate=0.05406574619094252, max_delta_step=0, max_depth=22,
               min_child_weight=2, missing=None, n_estimators=98, n_jobs=-1,
               nthread=None, objective='binary:logistic', random_state=0,
               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
               silent=True, subsample=0.5430112800235185)
-0.3090359784790159

```

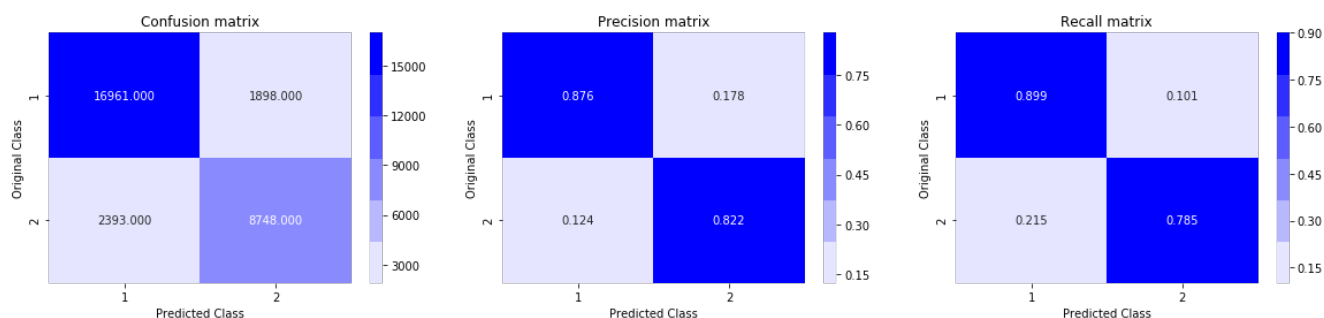
In [74]:

```

predict_y = model.predict_proba(set_)
print('For values of best parameters', "The train log loss is:", log_loss(y_train, predict_y, labels=
odel.classes_, eps=1e-15))
predict_y = model.predict_proba(set_t)
print('For values of best Parameters', "The test log loss is:", log_loss(y_test, predict_y, labels=m
odel.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 best parameters The train log loss is: 0.18769476308212615
For values of best Parameters The test log loss is: 0.3090359784790159
Total number of data points : 30000



XGBOOST Tuning2

In [75]:

```

param_dist = {'n_estimators': stats.randint(200, 400),
              'learning_rate': stats.uniform(0.01, 0.07),
              'subsample': stats.uniform(0.3, 0.5),
              'max_depth': stats.randint(5, 20),
              'colsample_bytree': stats.uniform(0.3, 0.4),
              'reg_lambda': [0.01, 0.1, 1]}

model = RandomizedSearchCV(XGBClassifier(n_jobs=-1), param_distributions=param_dist, scoring='neg_log_loss', cv=3)
model.fit(set_, y_train)
print(model.best_estimator_)
print(model.score(set_t, y_test))

```

```

XGBClassifier(base_score=0.5, booster='gbtree', colsample_bylevel=1,
               colsample_bytree=0.492537326837341, gamma=0,
               learning_rate=0.06408531601992787, max_delta_step=0, max_depth=14,
               min_child_weight=1, missing=None, n_estimators=323, n_jobs=-1,
               nthread=None, objective='binary:logistic', random_state=0,
               reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
               silent=True, subsample=0.7706996362058812)
-0.29418761841589947

```

In [76]:

```

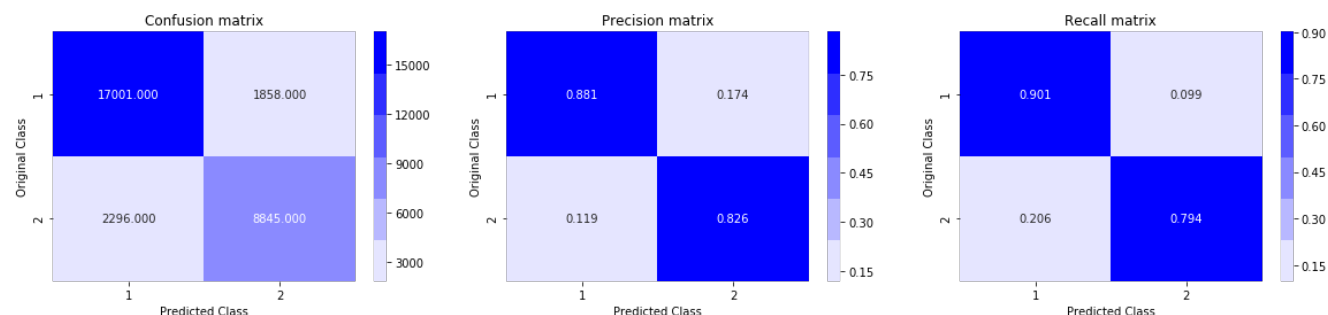
predict_y = model.predict_proba(set_)
print('For values of best parameters', "The train log loss is:",log_loss(y_train, predict_y, label
s=model.classes_, eps=1e-15))
predict_y = model.predict_proba(set_t)
print('For values of best Parameters',"The test log loss is:",log_loss(y_test, predict_y, labels=m
odel.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 best parameters The train log loss is: 0.15849154974980703

For values of best Parameters The test log loss is: 0.29418761841589947

Total number of data points : 30000



Summmary

In [83]:

```

from prettytable import PrettyTable
summary = PrettyTable()

```

In [84]:

```
summary.field_names = ["Task", "Vectorizer", "Train","Test"]
```

In [85]:

```

summary.add_row(["Logistic Regression","TFIDF","0.451","0.452"])
summary.add_row(["Linear SVM","TFIDF","0.466","0.470"])
summary.add_row(["Xgboost Tuning 1","TFIDF","0.187","0.309"])
summary.add_row(["Xgboost Tuning 2","TFIDF","0.158","0.294"])

```

In [86]:

```

print("Model Performance summary usning TFIDF")
print("*"*60)
print(summary)

```

Model Performance summary...

Task	Vectorizer	Train	Test
Logistic Regression	TFIDF	0.451	0.452
Linear SVM	TFIDF	0.466	0.470
Xgboost Tuning 1	TFIDF	0.187	0.309
Xgboost Tuning 2	TFIDF	0.158	0.294

Performance of models with w2v vectorizer

In [87]:

```

summary = PrettyTable()
summary.field_names = ["Task", "Vectorizer", "Train", "Test"]

```

```
summary.field_names = [ task , vectorizer , train , test ]
```

In [88]:

```
summary.add_row(["Logistic Regression","TFIDF W2V","0.513","0.520"])
summary.add_row(["Linear SVM","TFIDF W2V","0.478","0.489"])
summary.add_row(["Xgboost No tuning","TFIDF W2V","0.684","0.357"])
```

In [89]:

```
print("Model Performance summary usning TFIDF-W2V")
print("*"*60)
print(summary)
```

Model Performance summary usning TFIDF-W2V

Task	Vectorizer	Train	Test
Logistic Regression	TFIDF W2V	0.513	0.520
Linear SVM	TFIDF W2V	0.478	0.489
Xgboost No tuning	TFIDF W2V	0.684	0.357

Observation

- Our aim is to get loss less than **Random Model** which is 0.8872
- We Used TFIDF vectorizer to get better results.
- TFIDF vector gave better result in each case as compared to TFIDF-W2V
- The models trained on TF-IDF are not included here
- XGBOOST tuning 2 gave the best performance with **minimum training of 0.158**