



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

__ Useful Links __

- Discussions : <https://www.kaggle.com/anokas/data-analysis-xgboost-starter-0-35460-lb/comments> (<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> (<https://www.dropbox.com/sh/93968nfnrzh8bp5/AACZdtsApc1QSTQc7X0H3QZ5a?dl=0>)
- Blog 1 : <https://engineering.quora.com/Semantic-Question-Matching-with-Deep-Learning> (<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> (<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> (<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

2.3 Train and Test Construction

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

3. Exploratory Data Analysis

In [1]:

```
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
19
20 import time
21 import warnings
22 from sklearn.preprocessing import normalize
23 from sklearn.feature_extraction.text import CountVectorizer
24 from sklearn.feature_extraction.text import TfidfVectorizer
25 warnings.filterwarnings("ignore")
26 import sys
27 from tqdm import tqdm
28
29 # extract word2vec vectors
30 # https://github.com/explosion/spaCy/issues/1721
31 # http://Landinghub.visualstudio.com/visual-cpp-build-tools
32 import spacy
33
34 # This package is used for finding Longest common subsequence between two strings
35 # you can write your own dp code for this
36 import distance
37 from fuzzywuzzy import fuzz
38 from MulticoreTSNE import MulticoreTSNE as TSNE
39 # Import the Required Lib packages for WORD-Cloud generation
40 # https://stackoverflow.com/questions/45625434/how-to-install-wordcloud-in-python3-6
41 from wordcloud import WordCloud, STOPWORDS
```

```
42 from os import path  
43 from PIL import Image
```

3.1 Reading data and basic stats

In [2]:

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

Number of data points: 404290

In [3]:

```
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 23^{24} is divided by 1... | 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.

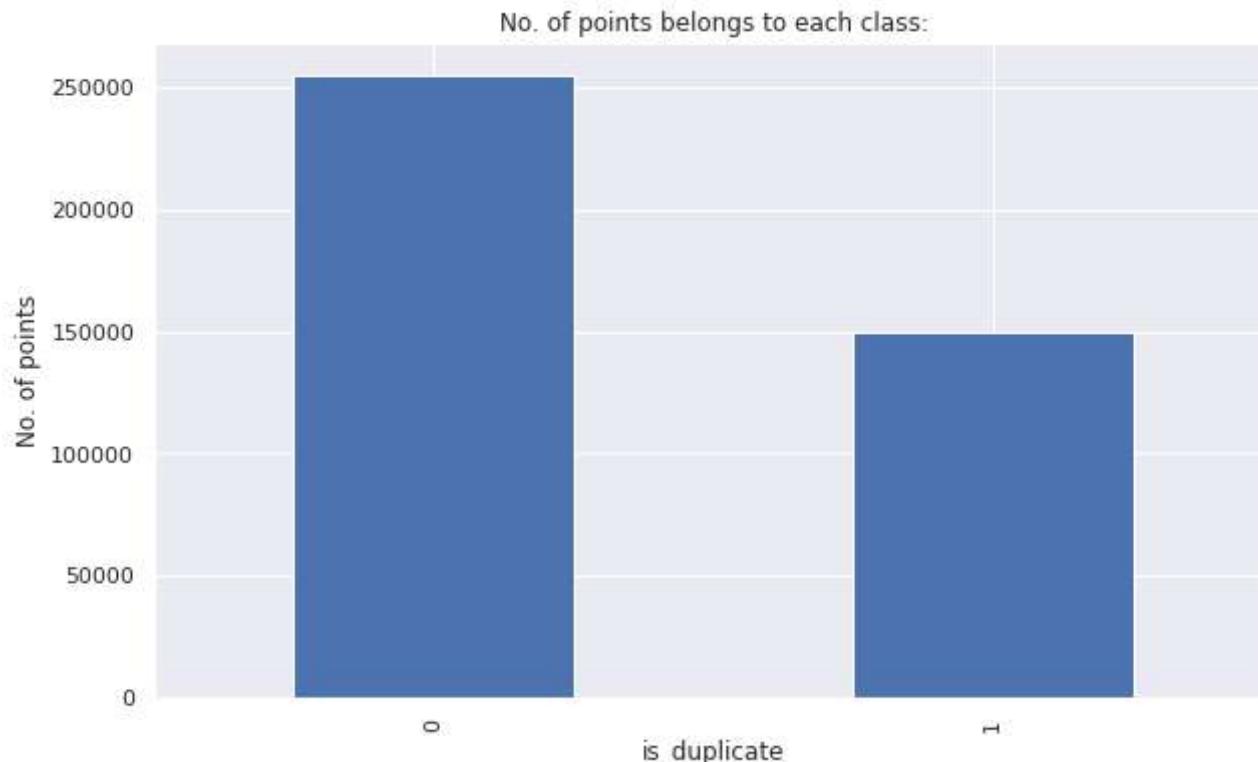
3.2.1 Distribution of data points among output classes

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

In [7]:

```
1 # GROUP BY BASED ON ='is_duplicate' return dict with key as columns and categorize based on unique values in is_dupe
2 sns.set_style('darkgrid')
3 sns.set(rc={'figure.figsize':(10,6)})
4 plt.title('No. of points belongs to each class:')
5 plt.ylabel('No. of points')
6 df.groupby("is_duplicate")['id'].count().plot.bar()
```

Out[7]: <matplotlib.axes._subplots.AxesSubplot at 0x7f3eb716c9b0>



In [8]:

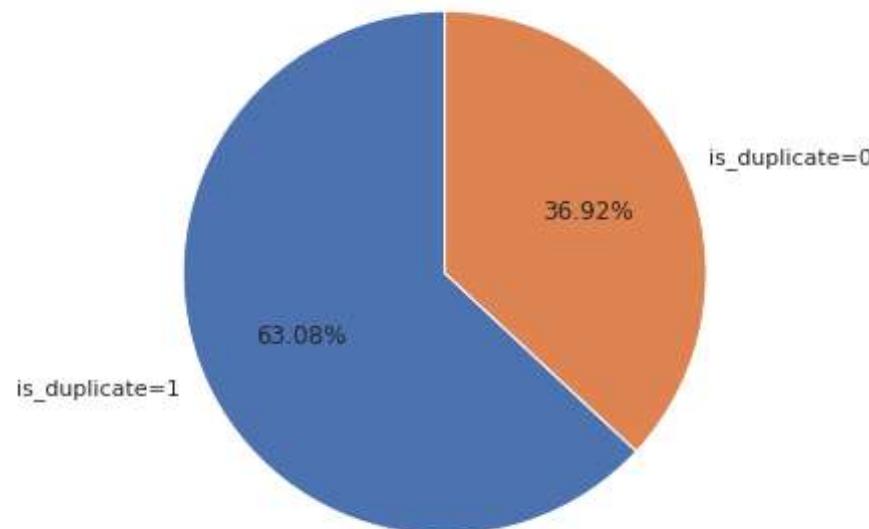
```
1 print('~> Total number of question pairs for training:\n    {}'.format(len(df)))
```

~> Total number of question pairs for training:
404290

In [19]:

```
1 size1=100 - round(df.is_duplicate.mean()*100, 2);
2 size2=round(df['is_duplicate'].mean()*100, 2)
3 sizes = [size1, size2]
4 plt.title('percentage of points belongs to each class',size=15)
5 plt.pie(sizes, labels=('is_duplicate=1', 'is_duplicate=0'), autopct='%1.2f%%',startangle=90)
6 plt.show()
```

percentage of points belongs to each class



3.2.2 Number of unique questions

```
In [21]:  
1 #CREATING A LIST OF HAVING ALL QID FROM qid1 AND qid2 COLUMNS  
2 #value_counts() return count of unique values in a selected column  
3 qids = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())  
4 unique_qs = len(np.unique(qids))  
5 qs_morethan_onetime = np.sum(qids.value_counts() > 1)  
6 print ('Total number of Unique Questions are: {} \n'.format(unique_qs))  
7  
8 print ('Number of unique questions that appear more than one time: {} ({}%)\n'.\n9     format(qs_morethan_onetime,round(qs_morethan_onetime/unique_qs*100,2)))  
10  
11 print ('Max number of times a single question is repeated: {} \n'.format(max(qids.value_counts())))  
12  
13 q_vals=qids.value_counts()  
14 #convert to array  
15 q_vals=q_vals.values
```

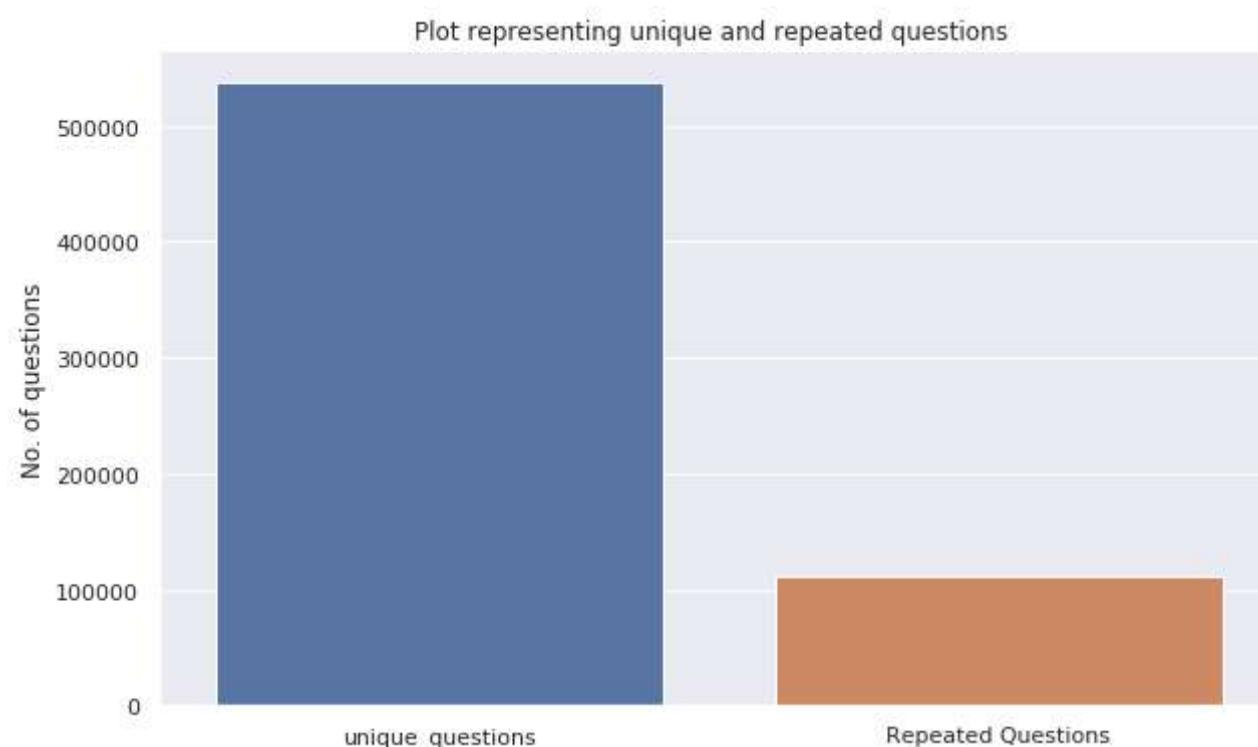
Total number of Unique Questions are: 537933

Number of unique questions that appear more than one time: 111780 (20.78%)

Max number of times a single question is repeated: 157

In [22]:

```
1 x = ["unique_questions" , "Repeated Questions"]
2 y = [unique_qs , qs_morethan_onetime]
3
4 plt.figure(figsize=(10, 6))
5 sns.set_style('darkgrid')
6 plt.ylabel('No. of questions')
7 plt.title ("Plot representing unique and repeated questions   ")
8 sns.barplot(x,y)
9 plt.show()
```



3.2.3 Checking for Duplicates

```
In [23]: 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

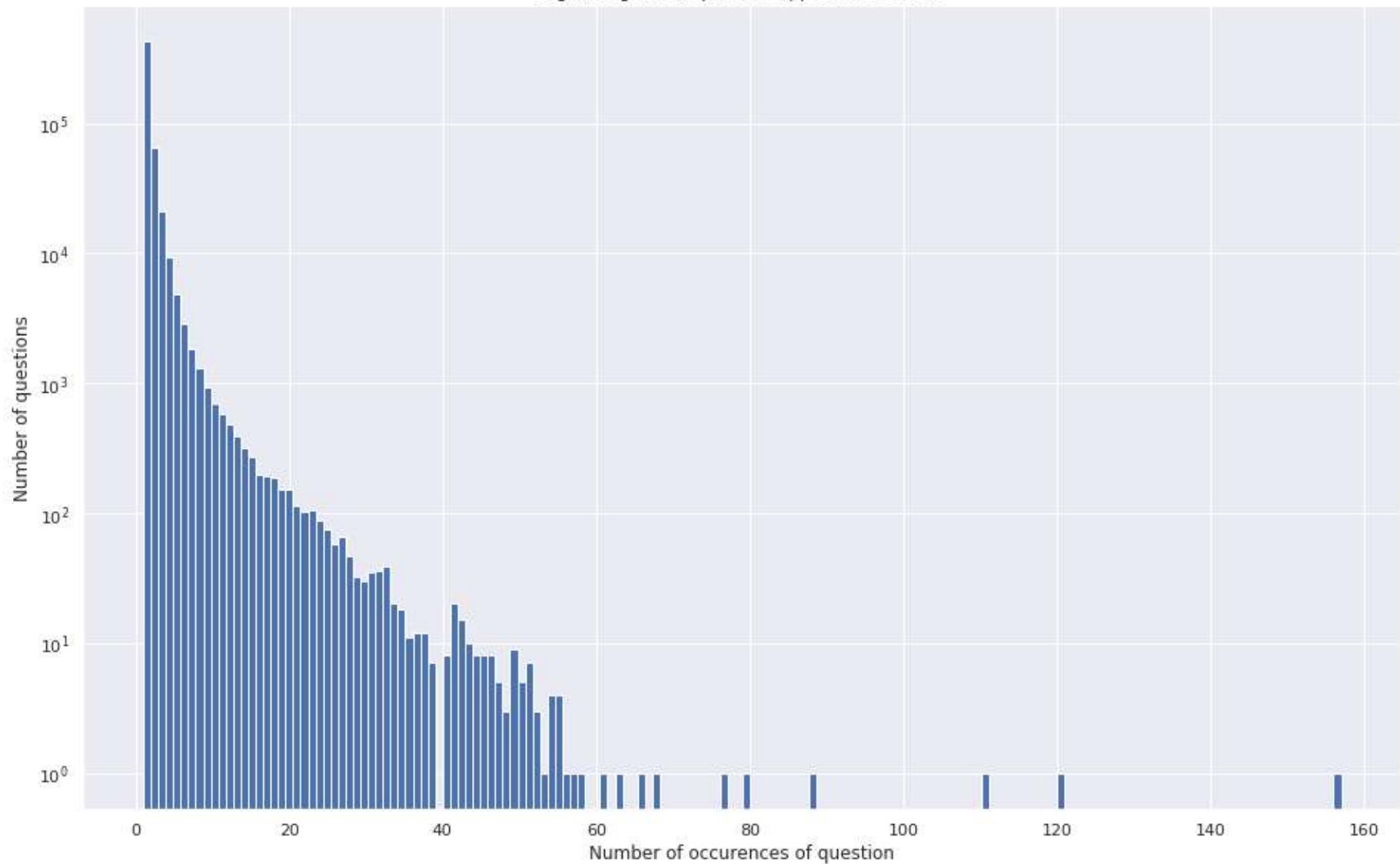
3.2.4 Number of occurrences of each question

In [24]:

```
1 plt.figure(figsize=(16, 10))
2 sns.set_style('darkgrid')
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: {}\\n'.format(max(qids.value_counts()))))
```

Maximum number of times a single question is repeated: 157

Log-Histogram of question appearance counts



3.2.5 Checking for NULL values

```
In [25]: 1 #Checking whether there are any rows with null values
          2 nan_rows = df[df.isnull().any(1)]
          3 nan_rows
```

Out[25]:

| | 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 My Chinese name is Haichao Yu. What English na... | 0 |

- There are two rows with null values in question2

```
In [26]: 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: []

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 [27]:
```

```

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.str.len() = return a df contains the length of each element in the Index of a df.
7     df['q1len'] = df['question1'].str.len()
8     df['q2len'] = df['question2'].str.len()
9     df['q1_n_words'] = df['question1'].apply(lambda row: len(row.split(" ")))
10    df['q2_n_words'] = df['question2'].apply(lambda row: len(row.split(" ")))
11
12    def normalized_word_Common(row):
13        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
14        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
15        return 1.0 * len(w1 & w2)
16    df['word_Common'] = df.apply(normalized_word_Common, axis=1)
17
18    def normalized_word_Total(row):
19        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
20        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
21        return 1.0 * (len(w1) + len(w2))
22    df['word_Total'] = df.apply(normalized_word_Total, axis=1)
23
24    def normalized_word_share(row):
25        w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
26        w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
27        return 1.0 * len(w1 & w2)/(len(w1) + len(w2))
28    df['word_share'] = df.apply(normalized_word_share, axis=1)
29
30    df['freq_q1+q2'] = df['freq_qid1']+df['freq_qid2']
31    df['freq_q1-q2'] = abs(df['freq_qid1']-df['freq_qid2'])
32
33    df.to_csv("df_fe_without_preprocessing_train.csv", index=False)
34
35    df.head()

```

Out[27]:

| id | qid1 | qid2 | question1 | question2 | is_duplicate | freq_qid1 | freq_qid2 | q1len | q2len | q1_n_words | q2_n_words | word_Common | word_Total |
|----|------|------|-----------|-----------|--------------|-----------|-----------|-------|-------|------------|------------|-------------|------------|
|----|------|------|-----------|-----------|--------------|-----------|-----------|-------|-------|------------|------------|-------------|------------|

| | id | qid1 | qid2 | question1 | question2 | is_duplicate | freq_qid1 | freq_qid2 | q1len | q2len | q1_n_words | q2_n_words | word_Common | word_Total |
|---|-----------|-------------|-------------|--|--|---------------------|------------------|------------------|--------------|--------------|-------------------|-------------------|--------------------|-------------------|
| 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 | 23.0 |
| 1 | 1 | 3 | 4 | Kohinoor (Koh-i-Noor) Dia... | What would happen if the Indian government sto... | 0 | 4 | 1 | 51 | 88 | 8 | 13 | 4.0 | 20.0 |
| 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 | 24.0 |
| 3 | 3 | 7 | 8 | Why am I mentally very lonely? How can I solve... | Find the remainder when $[math]23^{24} [/math] i...$ | 0 | 1 | 1 | 50 | 65 | 11 | 9 | 0.0 | 19.0 |
| 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 | 20.0 |

3.3.1 Analysis of some of the extracted features

- Here are some questions have only one single words.

```
In [28]: 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

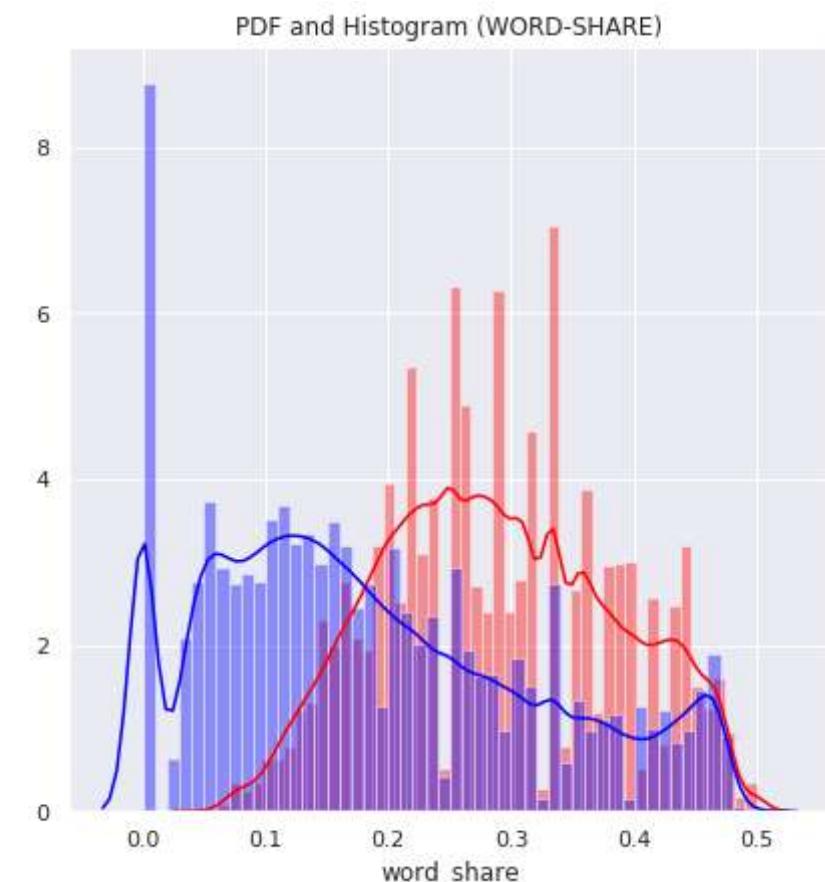
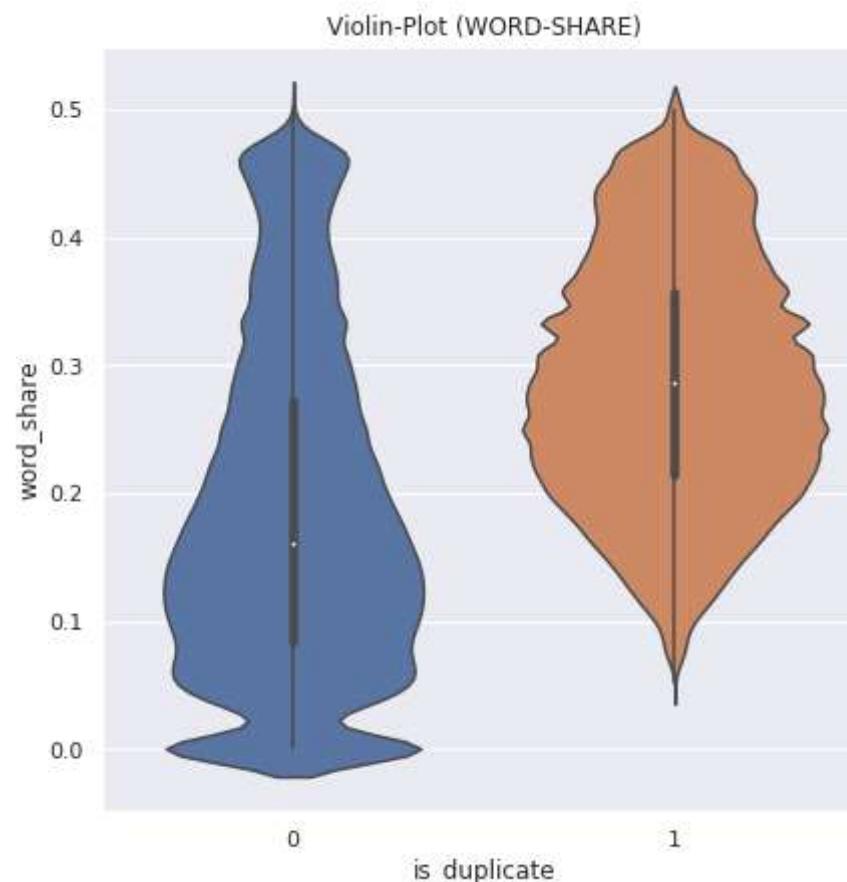
3.3.1.1 Feature: word_share

In [31]:

```

1 plt.figure(figsize=(15, 7))
2 sns.set_style('darkgrid')
3 plt.subplot(1,2,1)
4 plt.title('Violin-Plot (WORD-SHARE)')
5 sns.violinplot(x = 'is_duplicate', y = 'word_share', data = df) #df[0:]
6
7 plt.subplot(1,2,2)
8 plt.title('PDF and Histogram (WORD-SHARE)')
9 sns.distplot(df[df['is_duplicate'] == 1.0]['word_share'][0:] , label = "1", color = 'red')
10 sns.distplot(df[df['is_duplicate'] == 0.0]['word_share'][0:] , label = "0" , color = 'blue' )
11 plt.show()

```



- The distributions for normalized word_share have some overlap on the far right-hand side, i.e., there are quite a lot of questions with high word

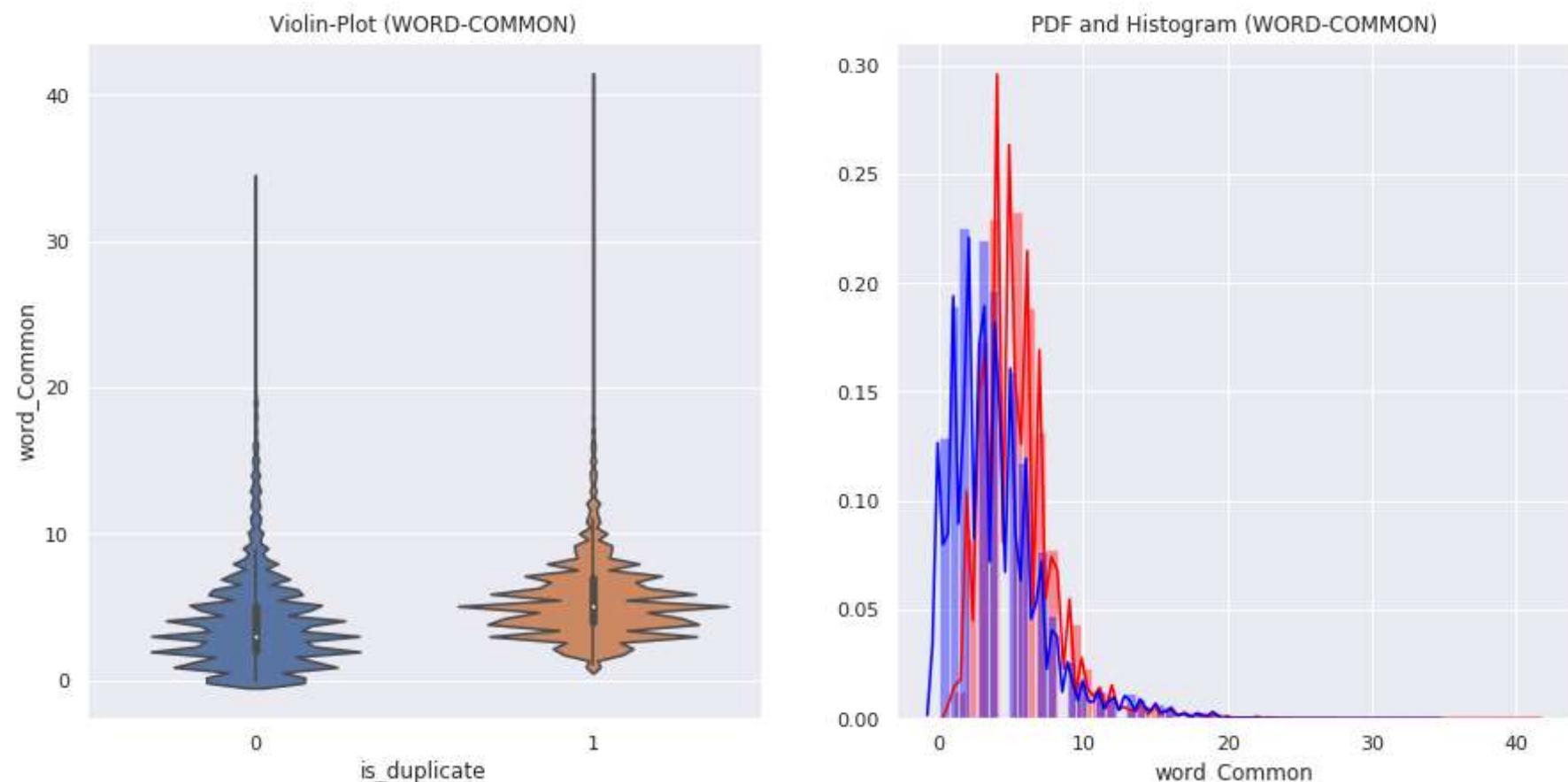
similarity

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

3.3.1.2 Feature: **word_Common**

In [32]:

```
1 plt.figure(figsize=(15, 7))
2 sns.set_style('darkgrid')
3 plt.subplot(1,2,1)
4 plt.title('Violin-Plot (WORD-COMMON)')
5 sns.violinplot(x = 'is_duplicate', y = 'word_Common', data = df[0:])
6
7 plt.subplot(1,2,2)
8 plt.title('PDF and Histogram (WORD-COMMON)')
9 sns.distplot(df[df['is_duplicate'] == 1.0]['word_Common'][0:] , label = "1", color = 'red')
10 sns.distplot(df[df['is_duplicate'] == 0.0]['word_Common'][0:] , label = "0" , color = 'blue' )
11 plt.show()
```



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

1.2.1 : EDA: Advanced Feature Extraction.

```
In [2]: 1 #https://stackoverflow.com/questions/12468179/unicodedecodeerror-utf8-codec-cant-decode-byte-0x9c
2 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
3     df = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
4     df = df.fillna('')
5     df.head()
6 else:
7     print("get df_fe_without_preprocessing_train.csv from drive or run the previous notebook")
```

```
In [3]: 1 df.head(2)
```

Out[3]:

| | id | qid1 | qid2 | question1 | question2 | is_duplicate | freq_qid1 | freq_qid2 | q1len | q2len | q1_n_words | q2_n_words | word_Common | word_Total | word_Diff | |
|---|----|------|------|---|---|---|-----------|-----------|-------|-------|------------|------------|-------------|------------|-----------|------|
| 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 | 23.0 | 13.0 | |
| 1 | 1 | 3 | 4 | Kohinoor (Koh-i-Noor) Dia... | 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 | 20.0 | 16.0 |

3.4 Preprocessing of Text

- Preprocessing:
 - Removing html tags
 - Removing Punctuations
 - Performing stemming
 - Removing Stopwords

- Expanding contractions etc.

```
In [4]: 1 # To get the results in 4 decimal points
2 SAFE_DIV = 0.0001
3
4 STOP_WORDS = stopwords.words("english")
5
6
7 def preprocess(x):
8     x = str(x).lower()
9     #it will apply 1st replace to wholestring then 2nd and so on.
10    x = x.replace(",000,000", "m").replace(",000", "k").replace("//", "").replace("'", "")\
11        .replace("won't", "will not").replace("cannot", "can not").replace("can't", "can not")\
12        .replace("n't", " not").replace("what's", "what is").replace("it's", "it is")\
13        .replace("ve", " have").replace("i'm", "i am").replace('re", " are")\
14        .replace("he's", "he is").replace("she's", "she is").replace("'s", " own")\
15        .replace("%", " percent ").replace("₹", " rupee ").replace("$", " dollar ")\
16        .replace("€", " euro ").replace('ll", " will")
17    x = re.sub(r"([0-9]+)000000", r"\1m", x)
18    x = re.sub(r"([0-9]+)000", r"\1k", x)
19
20
21    porter = PorterStemmer()
22    pattern = re.compile('\W')
23
24    if type(x) == type(''):
25        x = re.sub(pattern, ' ', x)
26
27
28    if type(x) == type(''):
29        x = porter.stem(x)
30        example1 = BeautifulSoup(x)
31        x = example1.get_text()
32
33
34    return x
35
```

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

3.5 Advanced Feature Extraction (NLP and Fuzzy Features)

Definition:

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

Features:

- **cwc_min** : Ratio of common_word_count to min length of word count of Q1 and Q2
 $cwc_min = \text{common_word_count} / (\min(\text{len}(q1_words), \text{len}(q2_words)))$
- **cwc_max** : Ratio of common_word_count to max length of word count of Q1 and Q2
 $cwc_max = \text{common_word_count} / (\max(\text{len}(q1_words), \text{len}(q2_words)))$
- **csc_min** : Ratio of common_stop_count to min length of stop count of Q1 and Q2
 $csc_min = \text{common_stop_count} / (\min(\text{len}(q1_stops), \text{len}(q2_stops)))$
- **csc_max** : Ratio of common_stop_count to max length of stop count of Q1 and Q2
 $csc_max = \text{common_stop_count} / (\max(\text{len}(q1_stops), \text{len}(q2_stops)))$
- **ctc_min** : Ratio of common_token_count to min length of token count of Q1 and Q2
 $ctc_min = \text{common_token_count} / (\min(\text{len}(q1_tokens), \text{len}(q2_tokens)))$
- **ctc_max** : Ratio of common_token_count to max length of token count of Q1 and Q2
 $ctc_max = \text{common_token_count} / (\max(\text{len}(q1_tokens), \text{len}(q2_tokens)))$
- **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

```
first_word_eq = int(q1_tokens[0] == q2_tokens[0])
```

- **abs_len_diff** : Abs. length difference

```
abs_len_diff = abs(len(q1_tokens) - len(q2_tokens))
```

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

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

- **fuzz_ratio** : <https://github.com/seatgeek/fuzzywuzzy#usage> (<https://github.com/seatgeek/fuzzywuzzy#usage>)

<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)

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

<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)

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

<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/>)

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

<http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/> (<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 [5]:

```
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"]), axis=1)
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 # do read this blog: http://chairnerd.seatgeek.com/fuzzywuzzy-fuzzy-string-matching-in-python/
81 # https://stackoverflow.com/questions/31806695/when-to-use-which-fuzz-function-to-compare-2-strings
82 # https://github.com/seatgeek/fuzzywuzzy
83     print("fuzzy features..")

```

```

84
85 df["token_set_ratio"]      = df.apply(lambda x: fuzz.token_set_ratio(x["question1"], x["question2"]), axis=1)
86 # The token sort approach involves tokenizing the string in question, sorting the tokens alphabetically, and
87 # then joining them back into a string We then compare the transformed strings with a simple ratio().
88 df["token_sort_ratio"]    = df.apply(lambda x: fuzz.token_sort_ratio(x["question1"], x["question2"]), axis=1)
89 df["fuzz_ratio"]          = df.apply(lambda x: fuzz.QRatio(x["question1"], x["question2"]), axis=1)
90 df["fuzz_partial_ratio"]  = df.apply(lambda x: fuzz.partial_ratio(x["question1"], x["question2"]), axis=1)
91 df["longest_substr_ratio"] = df.apply(lambda x: get_longest_substr_ratio(x["question1"], x["question2"]), axis=1)
92
93 return df

```

In [6]:

```

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

```

CPU times: user 2.46 s, sys: 204 ms, total: 2.67 s
Wall time: 2.8 s

3.5.1 Analysis of extracted features

3.5.1.1 Plotting Word clouds

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

In [7]:

```

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 {1,2,3,4}
5 #print(np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).shape) #(1, 149263, 2)
6 p = np.dstack([df_duplicate["question1"], df_duplicate["question2"]]).flatten()
7 n = np.dstack([dfp_nonduplicate["question1"], dfp_nonduplicate["question2"]]).flatten()
8
9 print ("Number of data points in class 1 (duplicate pairs) :",len(p))
10 print ("Number of data points in class 0 (non duplicate pairs) :",len(n))
11
12 #Saving the np array into a text file
13 np.savetxt('train_p.txt', p, delimiter=' ', fmt='%s')
14 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 [8]:

```

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))

```

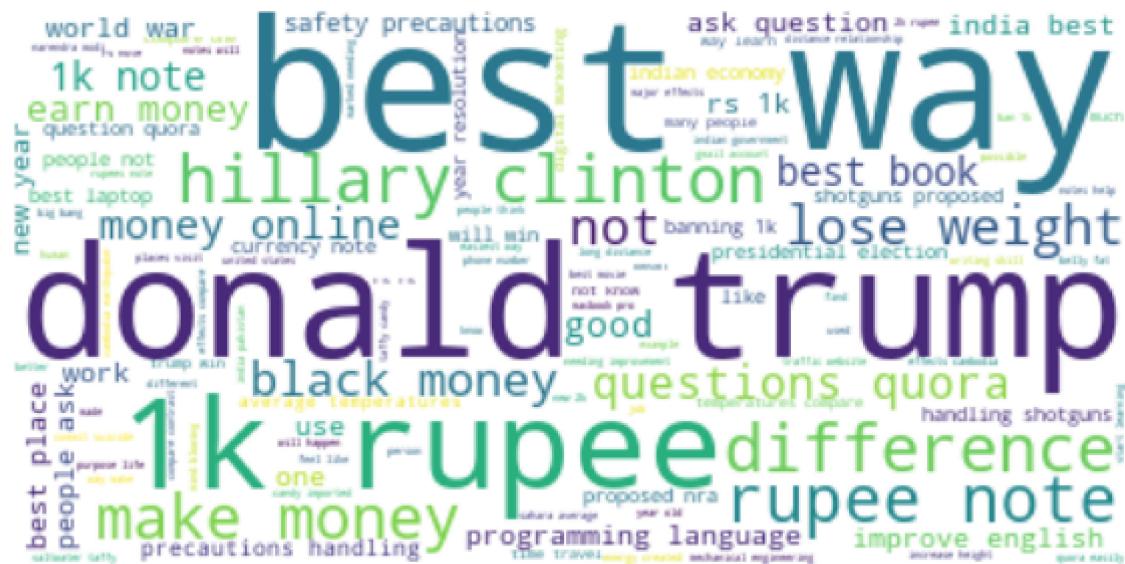
Total number of words in duplicate pair questions : 16109886

Total number of words in non duplicate pair questions : 33193130

Word Clouds generated from duplicate pair question's text

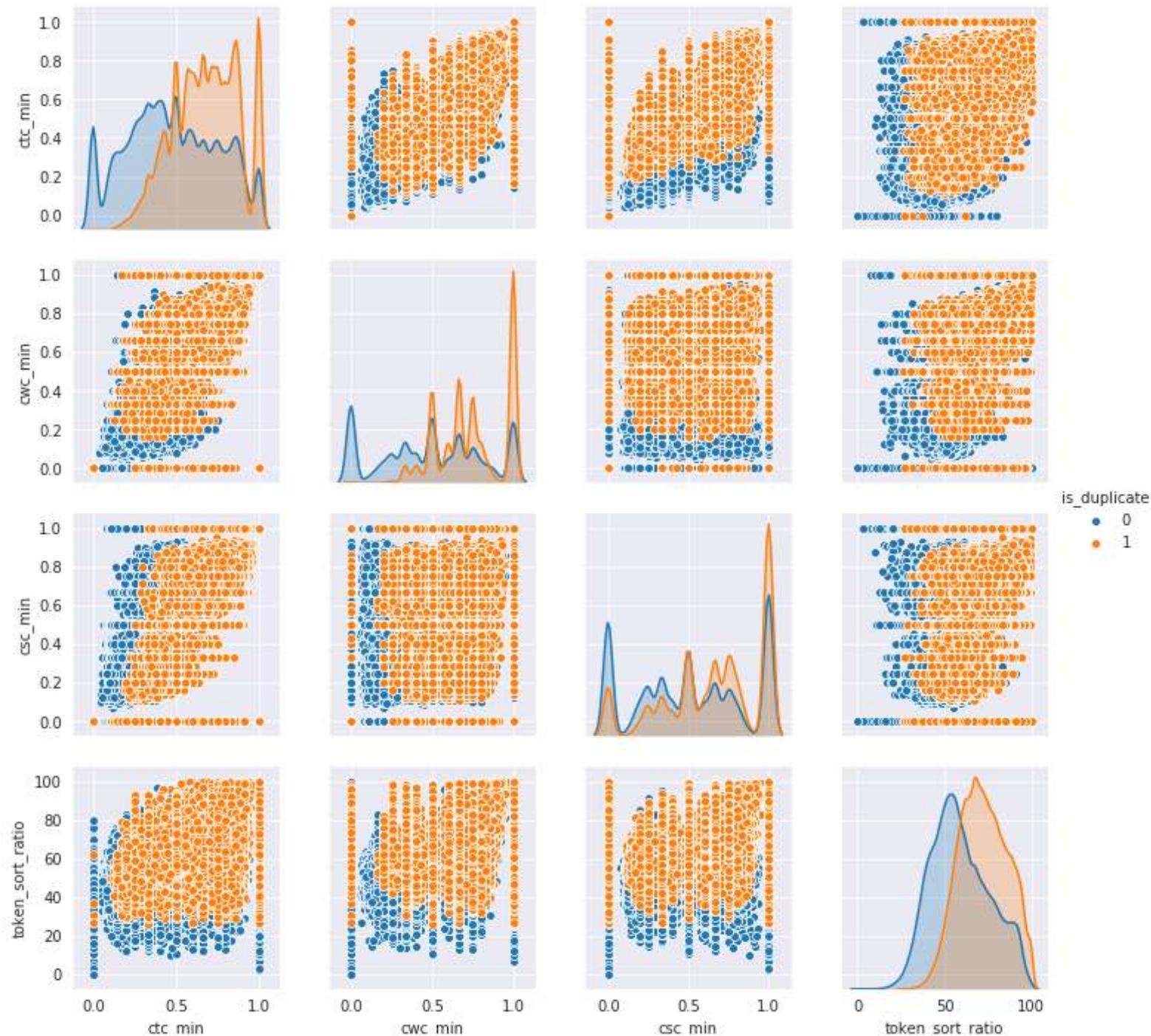
```
In [9]: 1 plt.figure(1,figsize=(10,7))
2 wc = WordCloud(background_color="white", max_words=len(textp_w), stopwords=stopwords)
3 wc.generate(textp_w)
4 print ("Word Cloud for Duplicate Question pairs")
5 plt.imshow(wc, interpolation='bilinear')
6 plt.axis("off")
7 plt.show()
```

Word Cloud for Duplicate Question pairs

Word Clouds generated from non duplicate pair question's text

In [12]:

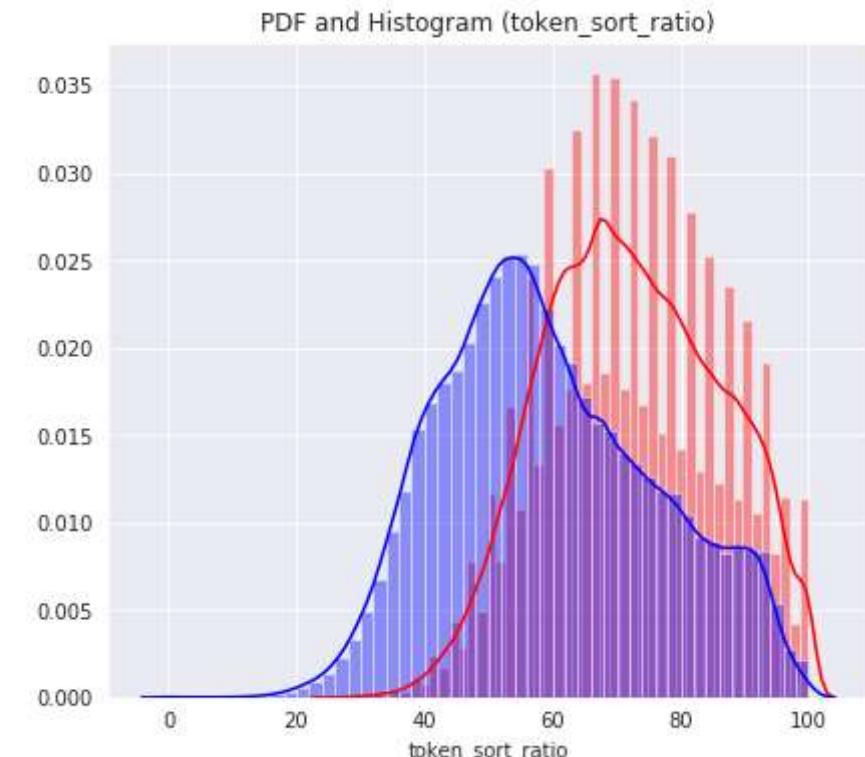
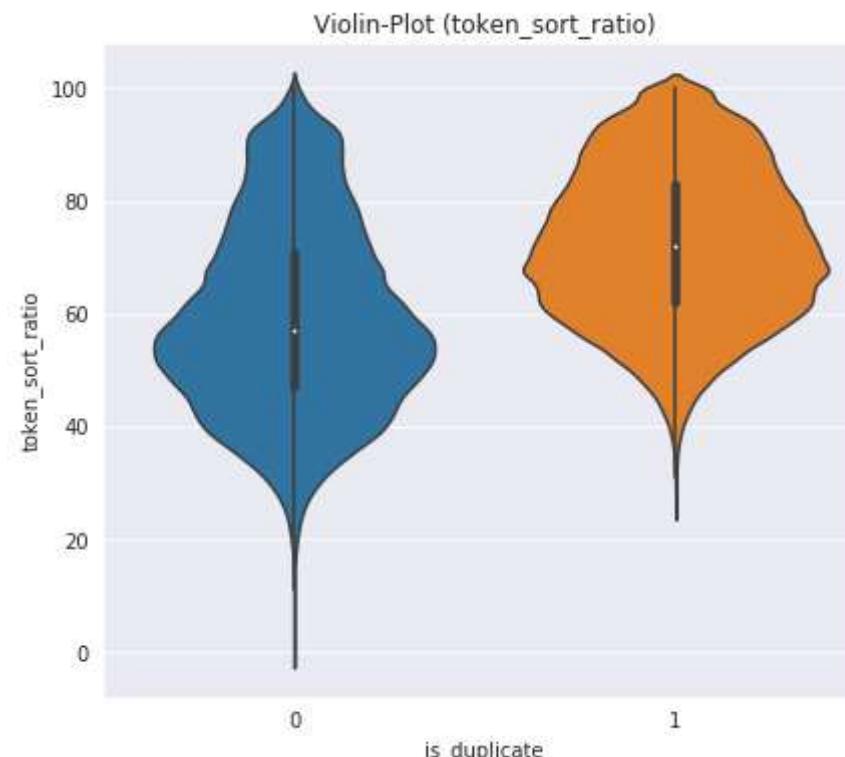
```
1 %%time
2 n = df.shape[0]
3 sns.set_style('darkgrid')
4 sns.pairplot(df[['ctc_min', 'cwc_min', 'csc_min', 'token_sort_ratio', 'is_duplicate']][0:n], hue='is_duplicate', var
```



CPU times: user 1min 34s, sys: 26.9 s, total: 2min 1s
Wall time: 1min 21s

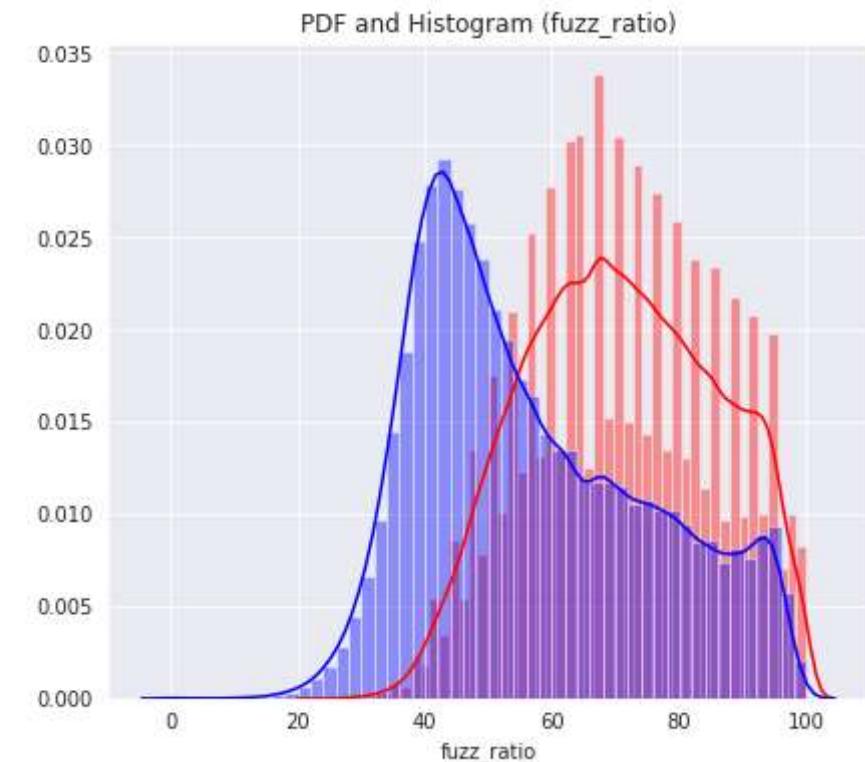
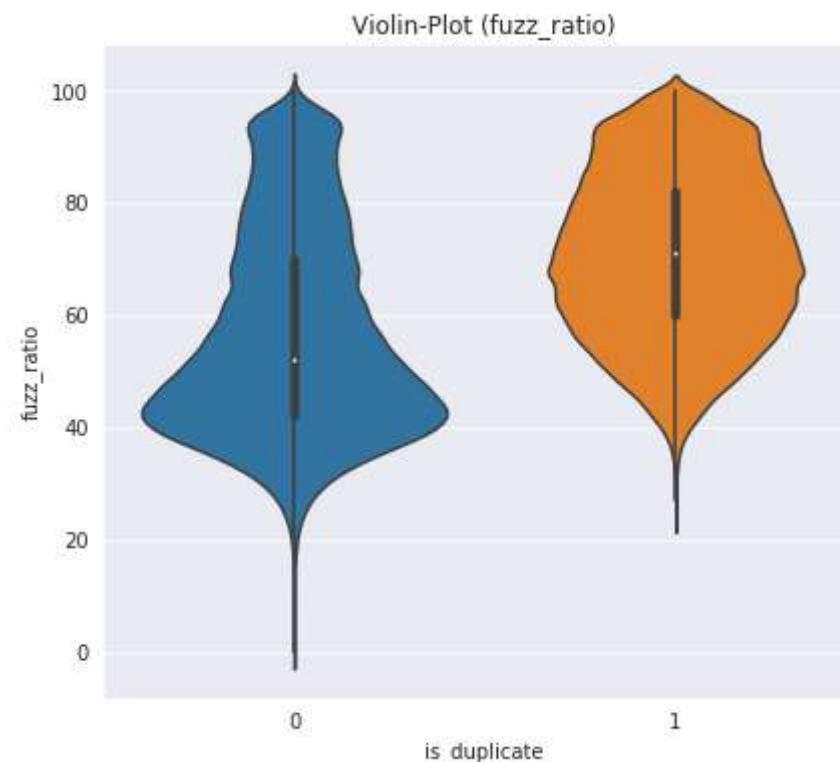
In [13]:

```
1 # Distribution of the token_sort_ratio
2 plt.figure(figsize=(15, 6))
3
4 sns.set_style('darkgrid')
5 plt.subplot(1,2,1)
6 plt.title('Violin-Plot (token_sort_ratio)')
7 sns.violinplot(x = 'is_duplicate', y = 'token_sort_ratio', data = df[0:] , )
8
9 plt.subplot(1,2,2)
10 plt.title('PDF and Histogram (token_sort_ratio)')
11 sns.distplot(df[df['is_duplicate'] == 1.0]['token_sort_ratio'][0:] , label = "1", color = 'red')
12 sns.distplot(df[df['is_duplicate'] == 0.0]['token_sort_ratio'][0:] , label = "0" , color = 'blue' )
13 plt.show()
```



In [14]:

```
1 plt.figure(figsize=(15, 6))
2
3 sns.set_style('darkgrid')
4 plt.subplot(1,2,1)
5 plt.title('Violin-Plot (fuzz_ratio)')
6 sns.violinplot(x = 'is_duplicate', y = 'fuzz_ratio', data = df[0:] , )
7
8 plt.subplot(1,2,2)
9 plt.title('PDF and Histogram (fuzz_ratio)')
10 sns.distplot(df[df['is_duplicate'] == 1.0]['fuzz_ratio'][0:] , label = "1", color = 'red')
11 sns.distplot(df[df['is_duplicate'] == 0.0]['fuzz_ratio'][0:] , label = "0" , color = 'blue' )
12 plt.show()
```



3.6 Featurizing text data with tfidf weighted word-vectors

In [2]:

```

1 from sklearn.externals import joblib
2 def saveModelToFile(obj,filename):
3     joblib.dump(obj,open(filename+".pkl","wb"))
4 def openModelFromFile(filename):
5     temp = joblib.load(open(filename+".pkl","rb"))
6     return temp

```

In [3]:

```

1 # avoid decoding problems
2 df = pd.read_csv("train.csv")
3
4 # encode questions to unicode
5 # https://stackoverflow.com/a/6812069
6 # ----- python 2 -----
7 # df['question1'] = df['question1'].apply(lambda x: unicode(str(x), "utf-8"))
8 # df['question2'] = df['question2'].apply(lambda x: unicode(str(x), "utf-8"))
9 # ----- python 3 -----
10 df['question1'] = df['question1'].apply(lambda x: str(x))
11 df['question2'] = df['question2'].apply(lambda x: str(x))

```

In [4]:

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

Out[4]:

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

In [5]:

```
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_))
```

- After we find TF-IDF scores, we convert each question to a weighted average of word2vec vectors by these scores.
- here we use a pre-trained GLOVE model which comes free with "Spacy". [\(https://spacy.io/usage/vectors-similarity\)](https://spacy.io/usage/vectors-similarity)
- It is trained on Wikipedia and therefore, it is stronger in terms of word semantics.

In [11]:

```
1 # en_vectors_web_lg, which includes over 1 million unique vectors.
2 nlp = spacy.load('en_core_web_sm')
3
4 vecs1 = []
5 # https://github.com/noamraph/tqdm
6 # tqdm is used to print the progress bar
7 for qu1 in tqdm(list(df['question1'])):
8     doc1 = nlp(qu1)
9     # 384 is the number of dimensions of vectors of spacy model(predefined)
10    #our quest/doc1 contain many words nd we want a single vector for doc1(qus1),
11    #so we are creating (#words in qus1 * 384) dim array nd taking mean of that array as final vector for qus1
12    mean_vec1 = np.zeros([len(doc1), len(doc1[0].vector)])
13    for word1 in doc1:
14        # word2vec
15        vec1 = word1.vector
16        # fetch df score
17        try:
18            idf = word2tfidf[str(word1)]
19        except:
20            idf = 0
21        # compute final vec
22        mean_vec1 += vec1 * idf
23        mean_vec1 = mean_vec1.mean(axis=0)
24    vecs1.append(mean_vec1)
25 df['q1_feats_m'] = list(vecs1)
26 #saveModeltofile(obj,filename)
```

100% |██████████| 404290/404290 [1:09:29<00:00, 96.97it/s]

NameError Traceback (most recent call last)

<ipython-input-11-26cf48a8345f> in <module>

```
24     vecs1.append(mean_vec1)
25 df['q1_feats_m'] = list(vecs1)
---> 26 saveModeltofile(obj,filename)
```

NameError: name 'obj' is not defined

In [12]:

```

1 vecs2 = []
2 for qu2 in tqdm(list(df['question2'])):
3     doc2 = nlp(qu2)
4     mean_vec2 = np.zeros([len(doc2), 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 [1:24:02<00:00, 80.17it/s]

In [13]:

```

1 #prepro_features_train.csv (Simple Preprocessing Features)
2 #nlp_features_train.csv (NLP Features)
3 if os.path.isfile('nlp_features_train.csv'):
4     dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
5 else:
6     print("download nlp_features_train.csv from drive or run previous notebook")
7
8 if os.path.isfile('df_fe_without_preprocessing_train.csv'):
9     dfppro = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
10 else:
11     print("download df_fe_without_preprocessing_train.csv from drive or run previous notebook")

```

In [14]:

```

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 [15]: 1 df1_tf = dfnlp.drop(['qid1','qid2'],axis=1)
2 df2_tf = dfppro.drop(['qid1','qid2','question1','question2','is_duplicate'],axis=1)
```

```
In [16]: 1 # dataframe of nlp features
2 df1.head()
```

Out[16]:

| | id | is_duplicate | cwc_min | cwc_max | csc_min | csc_max | ctc_min | ctc_max | last_word_eq | first_word_eq | abs_len_diff | mean_len | token_set_ratio |
|---|-----------|---------------------|----------------|----------------|----------------|----------------|----------------|----------------|---------------------|----------------------|---------------------|-----------------|------------------------|
| 0 | 0 | 0 | 0.999980 | 0.833319 | 0.999983 | 0.999983 | 0.916659 | 0.785709 | 0.0 | 1.0 | 2.0 | 13.0 | 100.0 |
| 1 | 1 | 0 | 0.799984 | 0.399996 | 0.749981 | 0.599988 | 0.699993 | 0.466664 | 0.0 | 1.0 | 5.0 | 12.5 | 80.0 |
| 2 | 2 | 0 | 0.399992 | 0.333328 | 0.399992 | 0.249997 | 0.399996 | 0.285712 | 0.0 | 1.0 | 4.0 | 12.0 | 60.0 |
| 3 | 3 | 0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 | 2.0 | 12.0 | 30.0 |
| 4 | 4 | 0 | 0.399992 | 0.199998 | 0.999950 | 0.666644 | 0.571420 | 0.307690 | 0.0 | 1.0 | 6.0 | 10.0 | 60.0 |

```
In [17]: 1 # data before preprocessing
2 df2.head()
```

Out[17]:

| | id | freq_qid1 | freq_qid2 | q1len | q2len | q1_n_words | q2_n_words | word_Common | word_Total | word_share | freq_q1+q2 | freq_q1-q2 |
|---|-----------|------------------|------------------|--------------|--------------|-------------------|-------------------|--------------------|-------------------|-------------------|-------------------|-------------------|
| 0 | 0 | 1 | 1 | 66 | 57 | 14 | 12 | 10.0 | 23.0 | 0.434783 | 2 | 0 |
| 1 | 1 | 4 | 1 | 51 | 88 | 8 | 13 | 4.0 | 20.0 | 0.200000 | 5 | 3 |
| 2 | 2 | 1 | 1 | 73 | 59 | 14 | 10 | 4.0 | 24.0 | 0.166667 | 2 | 0 |
| 3 | 3 | 1 | 1 | 50 | 65 | 11 | 9 | 0.0 | 19.0 | 0.000000 | 2 | 0 |
| 4 | 4 | 3 | 1 | 76 | 39 | 13 | 7 | 2.0 | 20.0 | 0.100000 | 4 | 2 |

In [18]:

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

Out[18]:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | 86 |
|---|-------------|-------------|-------------|-------------|-------------|------------|------------|------------|-------------|------------|-----|----------------|
| 0 | 211.129853 | -144.683060 | -68.811236 | -153.662138 | -89.931593 | 2.311301 | 136.743756 | 50.449102 | -64.150964 | 56.627520 | ... | 33.176595 -102 |
| 1 | 144.124684 | -114.012482 | -111.716699 | -104.885034 | -88.238481 | 16.441832 | 58.238008 | 102.095136 | 6.026963 | 178.498501 | ... | 67.386301 22 |
| 2 | 81.757897 | -142.184505 | 0.559867 | -104.660072 | -84.156624 | 22.515127 | 115.521664 | 50.436939 | -111.740921 | 51.713305 | ... | 61.731844 -49 |
| 3 | -126.651910 | -59.747162 | -67.763200 | -138.114751 | -101.038708 | 88.148514 | -22.912262 | 85.941421 | 27.784231 | 50.810640 | ... | 117.847461 -13 |
| 4 | 299.444019 | -188.632012 | -22.946285 | -273.683348 | -188.480383 | 107.123054 | 174.946303 | -72.042320 | -98.290509 | 137.439968 | ... | 62.256708 -173 |

5 rows × 96 columns

In [19]:

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

Out[19]:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | 86 |
|---|------------|-------------|-------------|-------------|-------------|-----------|------------|------------|------------|------------|-----|-------------------|
| 0 | 151.268527 | -127.013161 | -31.546284 | -142.905810 | -97.249094 | 9.485752 | 106.682272 | 36.754195 | -36.541906 | 53.162199 | ... | 28.763399 -108.30 |
| 1 | 152.023101 | -44.955390 | -103.559243 | -128.467608 | -118.567603 | 44.577921 | 137.906144 | 26.984743 | -78.328346 | 86.576890 | ... | 117.588063 5.23 |
| 2 | 4.930219 | -29.029577 | -117.808814 | -98.332280 | -19.064110 | -9.867806 | 141.808200 | 91.269548 | 50.727197 | 12.816846 | ... | 128.699951 17.43 |
| 3 | -6.951930 | -44.951732 | -17.343078 | -61.444447 | -7.469149 | 16.942010 | 95.049246 | -2.631597 | -13.050916 | -28.038383 | ... | 70.475260 -31.17 |
| 4 | 96.174522 | -71.613959 | 21.584889 | -92.742468 | -106.643133 | 10.646789 | 92.190156 | -40.565987 | -34.739524 | 56.340516 | ... | 18.852812 -70.41 |

5 rows × 96 columns

In [20]:

```
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]+df3_q2.shape[1])
```

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

In [21]:

```
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 [22]:

```
1 if not os.path.isfile('final_feat_with_ques.csv'):
2     result=df1_tf.merge(df2_tf, how='left', on='id')
3     result.to_csv('final_feat_with_ques.csv')
```

In []:

```
1
```