

Quora Question Pairs

```
In [1]: import os
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
os.path
```

```
Out[1]: <module 'ntpath' from 'C:\\Users\\ankan\\Anaconda3\\lib\\ntpath.py'>
```

```
In [2]: os.chdir('C:\\Users\\ankan\\Desktop\\Quora')
```

```
In [3]: os.getcwd()
```

```
Out[3]: 'C:\\Users\\ankan\\Desktop\\Quora'
```

```
In [4]: df = pd.read_csv("train.csv")

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

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

df.head()
```

```
Number of data points: 404290
~> Total number of question pairs for training <-- :
    404290
```

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

```
In [5]: df.info()
```

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

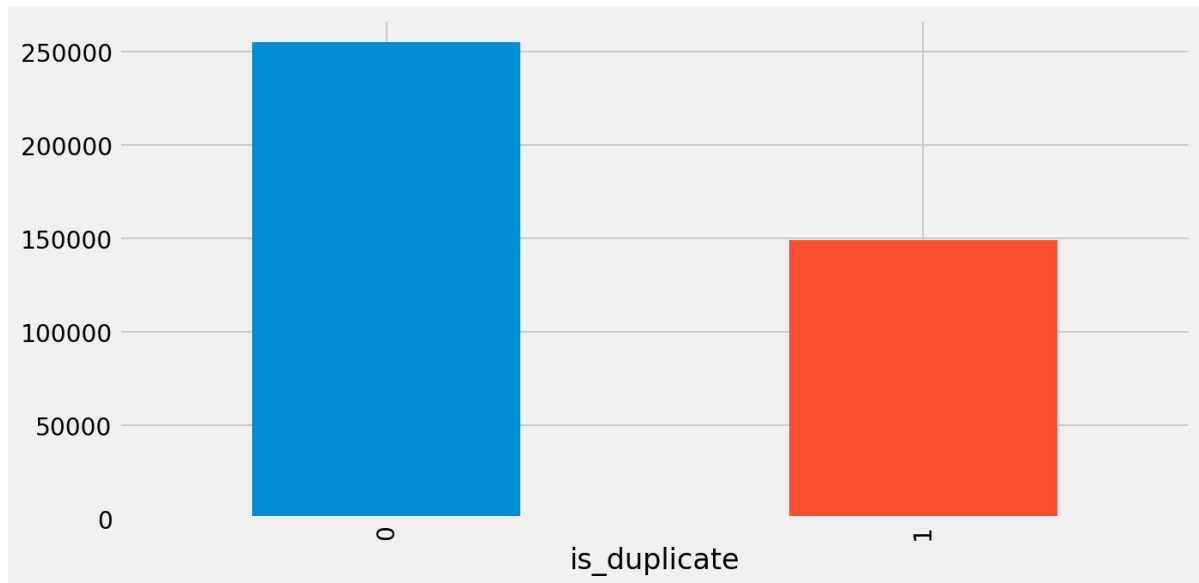
```
In [6]: import warnings
plt.style.use('fivethirtyeight')
plt.rcParams['figure.figsize'] = [10, 5]
warnings.filterwarnings("ignore", category=FutureWarning)
%config InlineBackend.figure_format = 'retina'
```

```
In [7]: df.groupby("is_duplicate")['id'].count().plot.bar()

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%
```



```
In [8]: qids = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())
unique_qs = len(np.unique(qids))
qs_morethan_onetime = np.sum(qids.value_counts() > 1)
print('Total number of Unique Questions are: {}'.format(unique_qs))
#print len(np.unique(qids))

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

print('Max number of times a single question is repeated: {}'.format(max(qids.va
lue_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 (20.7795394593
7505%)

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

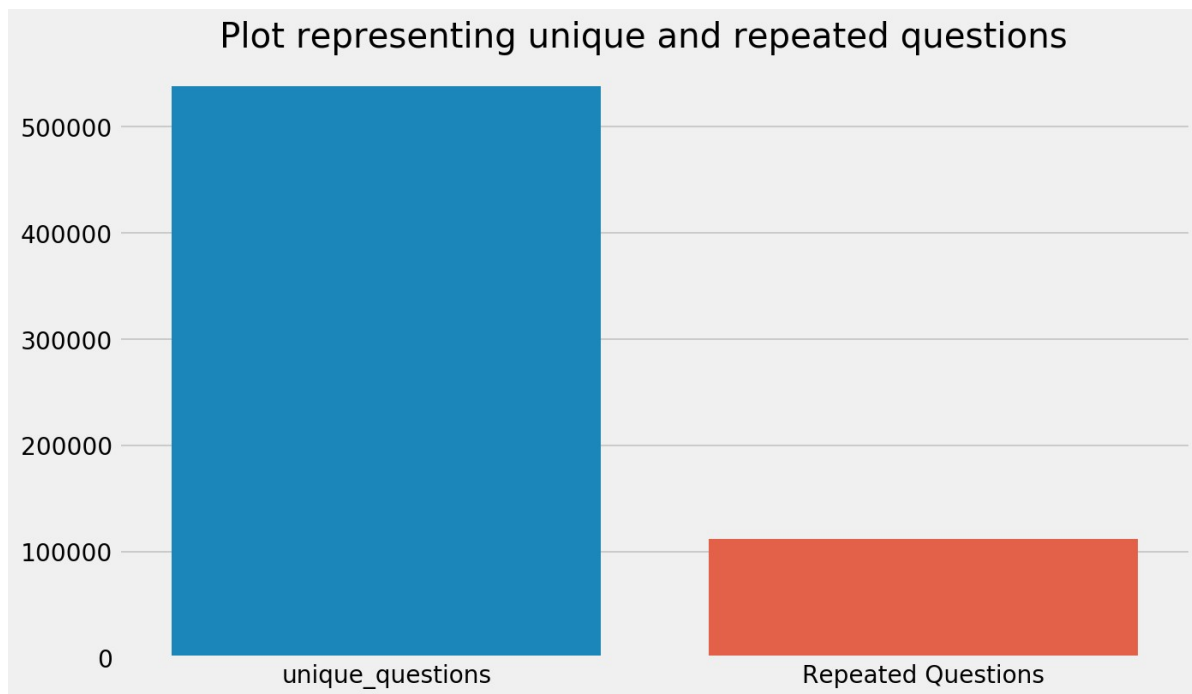
```
In [9]: x = ["unique_questions" , "Repeated Questions"]
        y = [unique_qs , qs_morethan_onetime]

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

#checking 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

```
In [10]: 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()))))

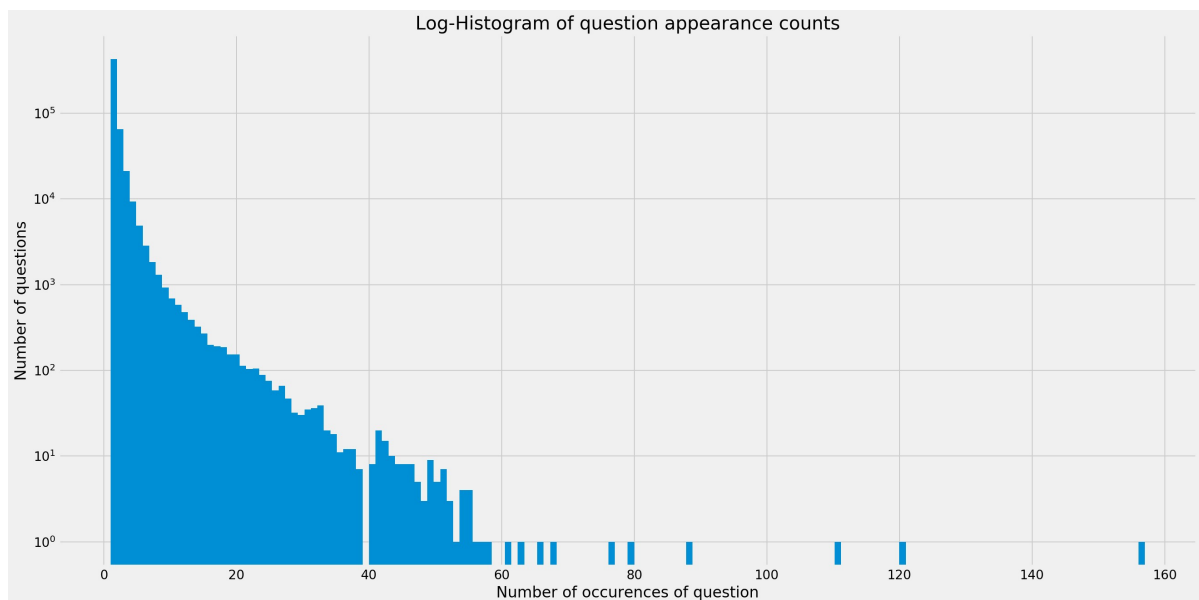
#Checking whether there are any rows with null values
nan_rows = df[df.isnull().any(1)]
print (nan_rows)

print('----->There are two rows with null values in question2<----- ')
```

Maximum number of times a single question is repeated: 157

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

----->There are two rows with null values in question2<-----



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

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

```
In [12]: 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[12]:

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