

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
```

```
In [4]: C:\Users\SaquibRahman\OneDrive - TheMathCompany Private Limited\Desktop\Training\Python\NLP by campus X\quora-question-pairs_only-bow.ipynb at main · campusx-official_quora-question-pairs · GitHub_files\data
```

```
In [5]: 1 df = pd.read_csv("train.csv")
2 df.shape
```

Out[5]: (404290, 6)

```
In [6]: 1 df.sample(10)
```

Out[6]:

|        | <b>id</b> | <b>qid1</b> | <b>qid2</b> | <b>question1</b>                                  | <b>question2</b>                                  | <b>is_duplicate</b> |
|--------|-----------|-------------|-------------|---|---|---------------------|
| 116556 | 116556    | 116583      | 189843      | What are some of the best places in India to v... | What are some of the best places to visit in I... | 1                   |
| 344164 | 344164    | 472354      | 472355      | Why does Uranus not orbit anticlockwise like t... | How does the orbit of Uranus look like (from t... | 0                   |
| 341253 | 341253    | 469053      | 119568      | *>  <* 1800><251><4919 *>  <* Cisco Router@@Te... | What is cisco router technical support phone n... | 0                   |
| 284183 | 284183    | 404379      | 404380      | How do I install a custom Rom in Spice mettle ... | How do I install a custom ROM in my rooted phone? | 0                   |
| 120771 | 120771    | 195843      | 195844      | How does MeetMe define their target market?       | How does Twitter define their target market?      | 0                   |
| 46062  | 46062     | 18952       | 82436       | How do I gain weight?                             | What is the best way for underweight to gain w... | 1                   |
| 183260 | 183260    | 42810       | 31229       | What is actual meaning of life?                   | What is the meaning of this life?                 | 1                   |
| 393906 | 393906    | 526744      | 526745      | Is it illegal for a 19 dating a 15?               | I'm 36 year old man and in love with a 19 year... | 0                   |
| 60705  | 60705     | 106123      | 30782       | Porn Stars: Where is Shy Love now?                | How can you look at someone's private Instagra... | 0                   |
| 196121 | 196121    | 296759      | 296760      | What are the health benefits of darjeeling tea?   | What are the health benefits of drinking Darje... | 1                   |

```
In [7]: 1 df.info()
```

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

```
In [8]: 1 # missing values
2 df.isnull().sum()
```

```
Out[8]: id          0
qid1        0
qid2        0
question1   1
question2   2
is_duplicate 0
dtype: int64
```

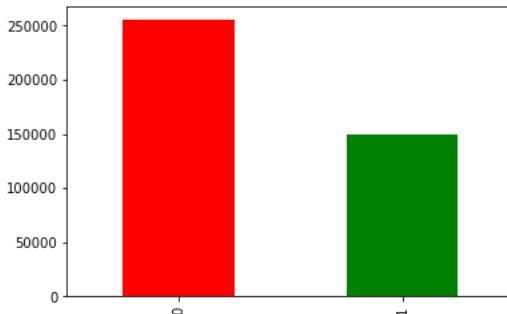
```
In [9]: 1 # duplicate rows
2 df.duplicated().sum()
```

Out[9]: 0

```
In [11]: 1 # Distribution of duplicate and non-duplicate questions
2 c= ['red', 'green']
3 print(df['is_duplicate'].value_counts())
4 print((df['is_duplicate'].value_counts()/df['is_duplicate'].count())*100)
5 df['is_duplicate'].value_counts().plot(kind='bar', color = c)

0    255027
1    149263
Name: is_duplicate, dtype: int64
0    63.080215
1    36.919785
Name: is_duplicate, dtype: float64
```

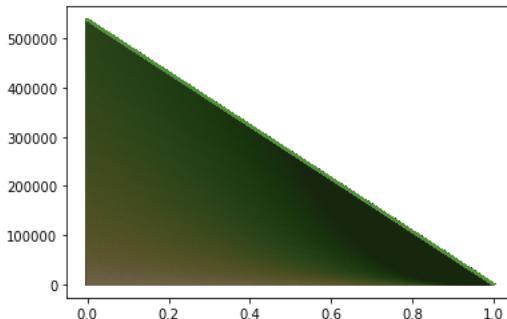
Out[11]: &lt;AxesSubplot:&gt;



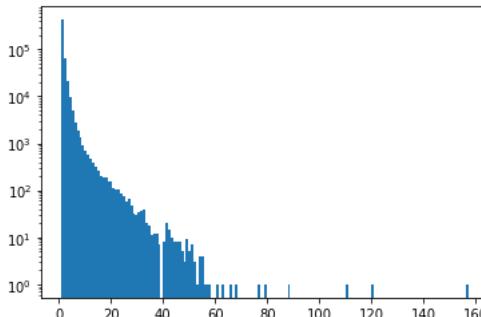
```
In [69]: 1 # Repeated questions
2
3 qid = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())
4 print('Number of unique questions',np.unique(qid).shape[0])
5 x = qid.value_counts()>1
6 print('Number of questions getting repeated',x[x].shape[0])
```

Number of unique questions 537933  
 Number of questions getting repeated 111780

```
In [85]: 1 # a = np.unique(qid)
2 # lst = list(np.unique(a, return_counts=True))
3 # plt.plot(lst)
4 # plt.show()
```



```
In [13]: 1 # Repeated questions histogram
2
3 plt.hist(qid.value_counts().values,bins=160)
4 plt.yscale('log')
5 plt.show()
```



### using bag of words

```
In [15]: 1 new_df = df.sample(30000)
```

```
In [16]: 1 new_df.isnull().sum()
```

```
Out[16]: id          0
          qid1       0
          qid2       0
          question1  0
          question2  0
          is_duplicate 0
          dtype: int64
```

```
In [17]: 1 new_df.duplicated().sum()
```

```
Out[17]: 0
```

```
In [18]: 1 ques_df = new_df[['question1','question2']]
2 ques_df.head()
```

```
Out[18]:
```

|        | question1  | question2                               |
|--------|--|---|
| 212317 | Why do some educated people in developed count...  | Why do people join ISIS?                |
| 376657 | Can you be an atheist but believe in ghosts? Is it justified for an atheist to believe in g... |   |
| 278723 | Why do some people have sweaty hands all the t...  | Why do some people have sweaty palms?   |
| 382943 | How do I approach a girl for sex? What is the best way to approach a girl?                     |   |
| 198249 | What are some affordable tax law firms in Kara...  | What are some affordable tax law firms? |

### count vectorizer!

```
In [20]: 1 from sklearn.feature_extraction.text import CountVectorizer
2 # merge texts
3 questions = list(ques_df['question1']) + list(ques_df['question2'])
4
5 cv = CountVectorizer(max_features=3000)
6 q1_arr, q2_arr = np.vsplit(cv.fit_transform(questions).toarray(), 2)
```

```
In [22]: 1 temp_df1 = pd.DataFrame(q1_arr, index= ques_df.index)
2 temp_df2 = pd.DataFrame(q2_arr, index= ques_df.index)
3 temp_df = pd.concat([temp_df1, temp_df2], axis=1)
4 temp_df.shape
```

```
Out[22]: (30000, 6000)
```

```
In [24]: 1 temp_df['is_duplicate'] = new_df['is_duplicate']
```

```
In [25]: 1 temp_df
```

```
Out[25]:
```

| 0      | 1   | 2   | 3   | 4   | 5   | 6   | 7   | 8   | 9   | ... | 2991 | 2992 | 2993 | 2994 | 2995 | 2996 | 2997 | 2998 | 2999 | is_duplicate |
|--------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|------|------|------|------|------|------|------|------|------|--------------|
| 212317 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |
| 376657 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 1            |
| 278723 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 1            |
| 382943 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |
| 198249 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |
| ...    | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ...  | ...  | ...  | ...  | ...  | ...  | ...  | ...  | ...  |              |
| 106526 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 1            |
| 91946  | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 1    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |
| 164002 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |
| 60972  | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 3    | 0    | 0    | 1    | 0    | 0    | 0    | 0    | 0    | 0            |
| 342272 | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | 0   | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0    | 0            |

30000 rows × 6001 columns

```
In [26]: 1 from sklearn.model_selection import train_test_split
2
3 X_train,X_test,y_train,y_test = train_test_split(temp_df.iloc[:,0:-1].values,temp_df.iloc[:, -1].values,test_size=0.2,random_
4
```

**random forest & xgboost and lightgbm**

```
In [27]: 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score
3 rf = RandomForestClassifier()
4 rf.fit(X_train,y_train)
5 y_pred = rf.predict(X_test)
6 accuracy_score(y_test,y_pred)
```

Out[27]: 0.736

```
In [29]: 1 from xgboost import XGBClassifier
2 xgb = XGBClassifier()
3 xgb.fit(X_train,y_train)
4 y_pred = xgb.predict(X_test)
5 accuracy_score(y_test,y_pred)
```

Out[29]: 0.7376666666666667

```
In [30]: 1 from lightgbm import LGBMClassifier
2 lgb = LGBMClassifier()
3 lgb.fit(X_train,y_train)
4 y_pred = lgb.predict(X_test)
5 accuracy_score(y_test,y_pred)
```

Out[30]: 0.7348333333333333

**BOW with basic feature engineering**

```
In [31]: 1 import warnings
2 warnings.filterwarnings('ignore')
```

```
In [37]: 1 new_df = df.sample(30000,random_state=2)
2 new_df
```

Out[37]:

|        | id     | qid1   | qid2   | question1   | question2   | is_duplicate |
|--------|--------|--------|--------|---|---|--------------|
| 398782 | 398782 | 496695 | 532029 | What is the best marketing automation tool for... | What is the best marketing automation tool for... | 1            |
| 115086 | 115086 | 187729 | 187730 | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh... | 0            |
| 327711 | 327711 | 454161 | 454162 | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers... | 0            |
| 367788 | 367788 | 498109 | 491396 | Why do so many people in the U.S. hate the sou... | My boyfriend doesnt feel guilty when he hurts ... | 0            |
| 151235 | 151235 | 237843 | 50930  | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag... | 0            |
| ...    | ...    | ...    | ...    | ...   | ...   | ...          |
| 243932 | 243932 | 26193  | 356455 | What are some good web scraping tutorials?        | What are some good web scraping programs?         | 1            |
| 91980  | 91980  | 154063 | 154064 | Can I apply for internet banking in SBI withou... | I have internet banking kit of SBI but it's no... | 0            |
| 266955 | 266955 | 133017 | 384210 | How much HE laundry detergent do you use in a ... | Can I use regular Dawn dishsoap in my dishwash... | 0            |
| 71112  | 71112  | 122427 | 122428 | What is the best way to understand and learn m... | What are some of the best ways to learn math?     | 1            |
| 312470 | 312470 | 436915 | 436916 | What would the Modi-led government do in case ... | If Pakistan mounts a 26/11 type attack again, ... | 1            |

30000 rows × 6 columns

```
In [35]: 1 new_df.isnull().sum()
```

```
Out[35]: id          0
qid1        0
qid2        0
question1    0
question2    0
is_duplicate 0
dtype: int64
```

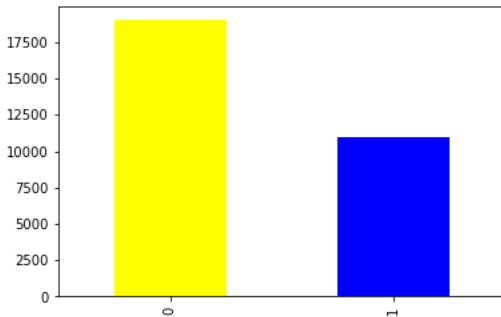
```
In [38]: 1 new_df.duplicated().sum()
```

Out[38]: 0

```
In [40]: 1 # Distribution of duplicate and non-duplicate questions
2 c= ['yellow', 'blue']
3 print(new_df['is_duplicate'].value_counts())
4 print((new_df['is_duplicate'].value_counts()/new_df['is_duplicate'].count())*100)
5 new_df['is_duplicate'].value_counts().plot(kind='bar', color =c)

0    19013
1   10987
Name: is_duplicate, dtype: int64
0    63.376667
1    36.623333
Name: is_duplicate, dtype: float64
```

Out[40]: &lt;AxesSubplot:&gt;



```
In [41]: 1 # Repeated questions
2
3 qid = pd.Series(new_df['qid1'].tolist() + new_df['qid2'].tolist())
4 print('Number of unique questions',np.unique(qid).shape[0])
5 x = qid.value_counts()>1
6 print('Number of questions getting repeated',x[x].shape[0])
```

Number of unique questions 55299  
 Number of questions getting repeated 3480

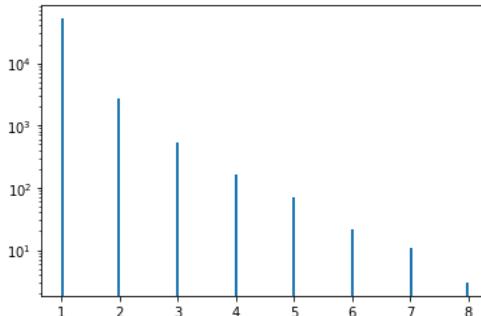
In [88]: 1 qid

```
Out[88]: 0      1
1      3
2      5
3      7
4      9
...
808575  379845
808576  155606
808577  537929
808578  537931
808579  537933
Length: 808580, dtype: int64
```

In [87]: 1 55299+3480

Out[87]: 58779

```
In [42]: 1 # Repeated questions histogram
2
3 plt.hist(qid.value_counts().values,bins=160)
4 plt.yscale('log')
5 plt.show()
```



```
In [43]: 1 # Feature Engineering
2
3 new_df['q1_len'] = new_df['question1'].str.len()
4 new_df['q2_len'] = new_df['question2'].str.len()
```

```
In [44]: 1 new_df.head()
```

Out[44]:

|        | <b>id</b> | <b>qid1</b> | <b>qid2</b> | <b>question1</b>                                  | <b>question2</b>                                   | <b>is_duplicate</b> | <b>q1_len</b> | <b>q2_len</b> |
|--------|-----------|-------------|-------------|---|--|---------------------|---------------|---------------|
| 398782 | 398782    | 496695      | 532029      | What is the best marketing automation tool for... | What is the best marketing automation tool for...  | 1                   | 76            | 77            |
| 115086 | 115086    | 187729      | 187730      | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh...  | 0                   | 49            | 57            |
| 327711 | 327711    | 454161      | 454162      | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers...  | 0                   | 105           | 120           |
| 367788 | 367788    | 498109      | 491396      | Why do so many people in the U.S. hate the sou... | My boyfriend doesn't feel guilty when he hurts ... | 0                   | 59            | 146           |
| 151235 | 151235    | 237843      | 50930       | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag...  | 0                   | 35            | 50            |

```
In [45]: 1 new_df['q1_num_words'] = new_df['question1'].apply(lambda row: len(row.split(" ")))
2 new_df['q2_num_words'] = new_df['question2'].apply(lambda row: len(row.split(" ")))
3 new_df.head()
```

Out[45]:

|        | <b>id</b> | <b>qid1</b> | <b>qid2</b> | <b>question1</b>                                  | <b>question2</b>                                   | <b>is_duplicate</b> | <b>q1_len</b> | <b>q2_len</b> | <b>q1_num_words</b> | <b>q2_num_words</b> |
|--------|-----------|-------------|-------------|---|--|---------------------|---------------|---------------|---------------------|---------------------|
| 398782 | 398782    | 496695      | 532029      | What is the best marketing automation tool for... | What is the best marketing automation tool for...  | 1                   | 76            | 77            | 12                  | 12                  |
| 115086 | 115086    | 187729      | 187730      | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh...  | 0                   | 49            | 57            | 12                  | 15                  |
| 327711 | 327711    | 454161      | 454162      | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers...  | 0                   | 105           | 120           | 25                  | 17                  |
| 367788 | 367788    | 498109      | 491396      | Why do so many people in the U.S. hate the sou... | My boyfriend doesn't feel guilty when he hurts ... | 0                   | 59            | 146           | 12                  | 30                  |
| 151235 | 151235    | 237843      | 50930       | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag...  | 0                   | 35            | 50            | 5                   | 9                   |

```
In [46]: 1 def common_words(row):
2     w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
3     w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
4     return len(w1 & w2)
```

```
In [47]: 1 new_df['word_common'] = new_df.apply(common_words, axis=1)
2 new_df.head()
```

Out[47]:

|        | <b>id</b> | <b>qid1</b> | <b>qid2</b> | <b>question1</b>                                  | <b>question2</b>                                   | <b>is_duplicate</b> | <b>q1_len</b> | <b>q2_len</b> | <b>q1_num_words</b> | <b>q2_num_words</b> | <b>word_common</b> |
|--------|-----------|-------------|-------------|---|--|---------------------|---------------|---------------|---------------------|---------------------|--------------------|
| 398782 | 398782    | 496695      | 532029      | What is the best marketing automation tool for... | What is the best marketing automation tool for...  | 1                   | 76            | 77            | 12                  | 12                  | 11                 |
| 115086 | 115086    | 187729      | 187730      | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh...  | 0                   | 49            | 57            | 12                  | 15                  | 7                  |
| 327711 | 327711    | 454161      | 454162      | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers...  | 0                   | 105           | 120           | 25                  | 17                  | 2                  |
| 367788 | 367788    | 498109      | 491396      | Why do so many people in the U.S. hate the sou... | My boyfriend doesn't feel guilty when he hurts ... | 0                   | 59            | 146           | 12                  | 30                  | 0                  |
| 151235 | 151235    | 237843      | 50930       | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag...  | 0                   | 35            | 50            | 5                   | 9                   | 3                  |

```
In [93]: 1 def total_words(row):
2     w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))
3     w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))
4     return (len(w1) + len(w2))
5     print(len(w1))
```

```
In [49]: 1 new_df['word_total'] = new_df.apply(total_words, axis=1)
2 new_df.head()
```

Out[49]:

|        |        |        |        |  | question1   | question2  | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total |
|--------|--------|--------|--------|--|---|--|--------------|--------|--------|--------------|--------------|-------------|------------|
| 398782 | 398782 | 496695 | 532029 |  | What is the best marketing automation tool for... | What is the best marketing automation tool for...  | 1            | 76     | 77     | 12           | 12           | 11          | 24         |
| 115086 | 115086 | 187729 | 187730 |  | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh...  | 0            | 49     | 57     | 12           | 15           | 7           | 23         |
| 327711 | 327711 | 454161 | 454162 |  | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers...  | 0            | 105    | 120    | 25           | 17           | 2           | 34         |
| 367788 | 367788 | 498109 | 491396 |  | Why do so many people in the U.S. hate the sou... | My boyfriend doesn't feel guilty when he hurts ... | 0            | 59     | 146    | 12           | 30           | 0           | 32         |
| 151235 | 151235 | 237843 | 50930  |  | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag...  | 0            | 35     | 50     | 5            | 9            | 3           | 13         |

```
In [50]: 1 new_df['word_share'] = round(new_df['word_common']/new_df['word_total'],2)
2 new_df.head()
```

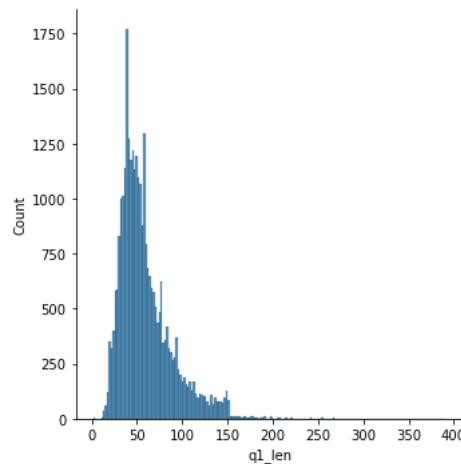
Out[50]:

|        |        |        |        |  | question1   | question2  | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_shar |
|--------|--------|--------|--------|--|---|--|--------------|--------|--------|--------------|--------------|-------------|------------|-----------|
| 398782 | 398782 | 496695 | 532029 |  | What is the best marketing automation tool for... | What is the best marketing automation tool for...  | 1            | 76     | 77     | 12           | 12           | 11          | 24         | 0.4       |
| 115086 | 115086 | 187729 | 187730 |  | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh...  | 0            | 49     | 57     | 12           | 15           | 7           | 23         | 0.3       |
| 327711 | 327711 | 454161 | 454162 |  | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers...  | 0            | 105    | 120    | 25           | 17           | 2           | 34         | 0.0       |
| 367788 | 367788 | 498109 | 491396 |  | Why do so many people in the U.S. hate the sou... | My boyfriend doesn't feel guilty when he hurts ... | 0            | 59     | 146    | 12           | 30           | 0           | 32         | 0.0       |
| 151235 | 151235 | 237843 | 50930  |  | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag...  | 0            | 35     | 50     | 5            | 9            | 3           | 13         | 0.2       |



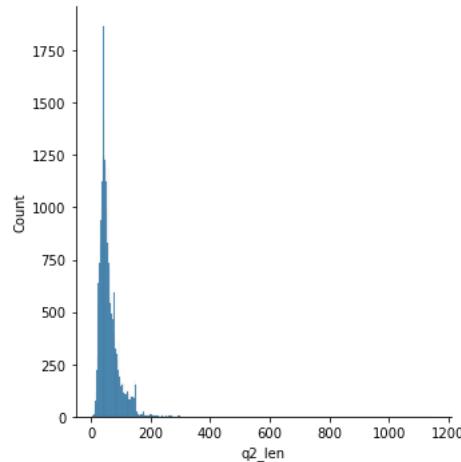
```
In [51]: 1 # Analysis of features
2 sns.displot(new_df['q1_len'])
3 print('minimum characters',new_df['q1_len'].min())
4 print('maximum characters',new_df['q1_len'].max())
5 print('average num of characters',int(new_df['q1_len'].mean()))
```

```
minimum characters 2
maximum characters 391
average num of characters 59
```



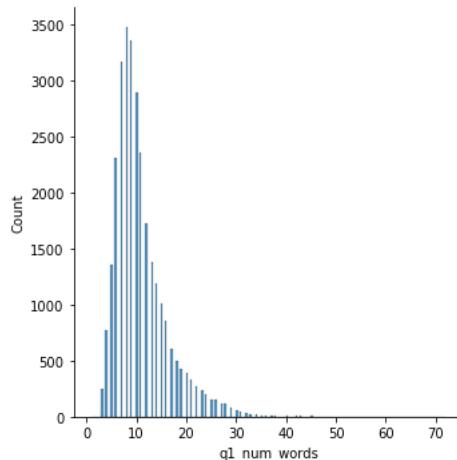
```
In [52]: 1 sns.displot(new_df['q2_len'])
2 print('minimum characters',new_df['q2_len'].min())
3 print('maximum characters',new_df['q2_len'].max())
4 print('average num of characters',int(new_df['q2_len'].mean()))
```

```
minimum characters 6
maximum characters 1151
average num of characters 60
```



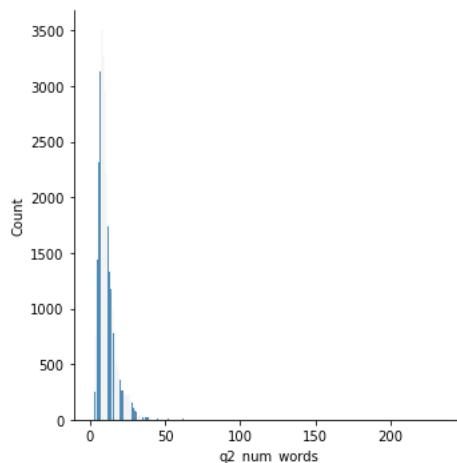
```
In [53]: 1 sns.distplot(new_df['q1_num_words'])
2 print('minimum words',new_df['q1_num_words'].min())
3 print('maximum words',new_df['q1_num_words'].max())
4 print('average num of words',int(new_df['q1_num_words'].mean()))
```

minimum words 1  
maximum words 72  
average num of words 10

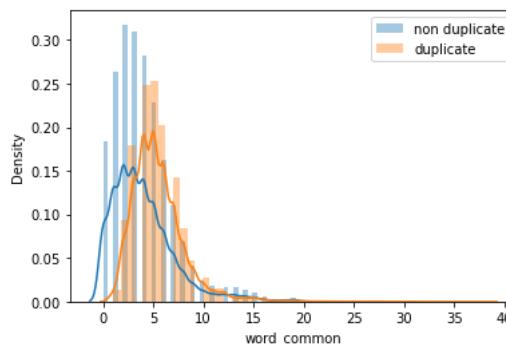


```
In [54]: 1 sns.distplot(new_df['q2_num_words'])
2 print('minimum words',new_df['q2_num_words'].min())
3 print('maximum words',new_df['q2_num_words'].max())
4 print('average num of words',int(new_df['q2_num_words'].mean()))
```

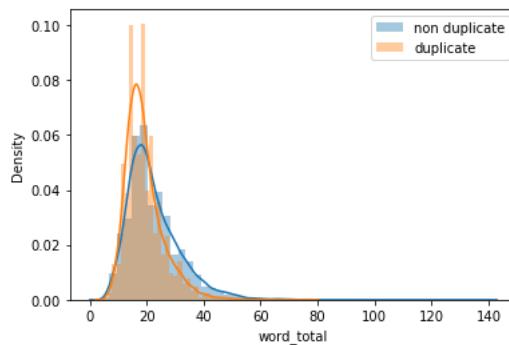
minimum words 1  
maximum words 237  
average num of words 11



```
In [55]: 1 # common words
2 sns.distplot(new_df[new_df['is_duplicate'] == 0]['word_common'],label='non duplicate')
3 sns.distplot(new_df[new_df['is_duplicate'] == 1]['word_common'],label='duplicate')
4 plt.legend()
5 plt.show()
```



```
In [56]: 1 # total words
2 sns.distplot(new_df[new_df['is_duplicate'] == 0]['word_total'],label='non duplicate')
3 sns.distplot(new_df[new_df['is_duplicate'] == 1]['word_total'],label='duplicate')
4 plt.legend()
5 plt.show()
```



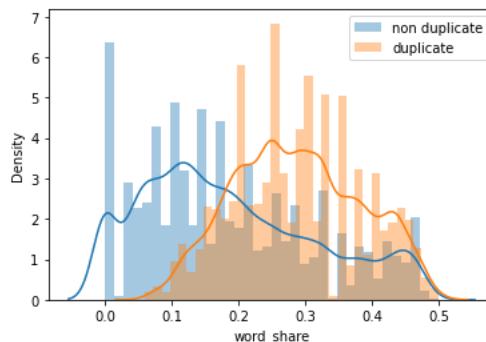
```
In [108]: 1 new_df[new_df['is_duplicate'] == 0]['word_total'][:5]
```

```
Out[108]: 115086    23
327711     34
367788     32
151235     13
244531     25
Name: word_total, dtype: int64
```

```
In [109]: 1 new_df[new_df['is_duplicate'] == 0]['word_common'][:5]
```

```
Out[109]: 115086    7
327711     2
367788     0
151235     3
244531     1
Name: word_common, dtype: int64
```

```
In [57]: 1 # word share
2 sns.distplot(new_df[new_df['is_duplicate'] == 0]['word_share'],label='non duplicate')
3 sns.distplot(new_df[new_df['is_duplicate'] == 1]['word_share'],label='duplicate')
4 plt.legend()
5 plt.show()
```



```
In [58]: 1 ques_df = new_df[['question1','question2']]
2 ques_df.head()
```

|               | question1   | question2   |
|---------------|---|---|
| <b>398782</b> | What is the best marketing automation tool for... | What is the best marketing automation tool for... |
| <b>115086</b> | I am poor but I want to invest. What should I do? | I am quite poor and I want to be very rich. Wh... |
| <b>327711</b> | I am from India and live abroad. I met a guy f... | T.I.E.T to Thapar University to Thapar Univers... |
| <b>367788</b> | Why do so many people in the U.S. hate the sou... | My boyfriend doesnt feel guilty when he hurts ... |
| <b>151235</b> | Consequences of Bhopal gas tragedy?               | What was the reason behind the Bhopal gas trag... |

```
In [59]: 1 final_df = new_df.drop(columns=['id','qid1','qid2','question1','question2'])
2 print(final_df.shape)
3 final_df.head()
```

(30000, 8)

Out[59]:

|        | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_share |
|--------|--------------|--------|--------|--------------|--------------|-------------|------------|------------|
| 398782 | 1            | 76     | 77     | 12           | 12           | 11          | 24         | 0.46       |
| 115086 | 0            | 49     | 57     | 12           | 15           | 7           | 23         | 0.30       |
| 327711 | 0            | 105    | 120    | 25           | 17           | 2           | 34         | 0.06       |
| 367788 | 0            | 59     | 146    | 12           | 30           | 0           | 32         | 0.00       |
| 151235 | 0            | 35     | 50     | 5            | 9            | 3           | 13         | 0.23       |

In [60]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 # merge texts
3 questions = list(ques_df['question1']) + list(ques_df['question2'])
4
5 cv = CountVectorizer(max_features=3000)
6 q1_arr, q2_arr = np.vsplit(cv.fit_transform(questions).toarray(), 2)
```

In [61]:

```
1 temp_df1 = pd.DataFrame(q1_arr, index=ques_df.index)
2 temp_df2 = pd.DataFrame(q2_arr, index=ques_df.index)
3 temp_df = pd.concat([temp_df1, temp_df2], axis=1)
4 temp_df.shape
```

Out[61]: (30000, 6008)

In [62]:

```
1 final_df = pd.concat([final_df, temp_df], axis=1)
2 print(final_df.shape)
3 final_df.head()
```

(30000, 6008)

Out[62]:

|        | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_share | 0 | 1 | ... | 2990 | 2991 | 2992 | 2993 | 2994 | 2995 | 2996 |
|--------|--------------|--------|--------|--------------|--------------|-------------|------------|------------|---|---|-----|------|------|------|------|------|------|------|
| 398782 | 1            | 76     | 77     | 12           | 12           | 11          | 24         | 0.46       | 0 | 0 | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 115086 | 0            | 49     | 57     | 12           | 15           | 7           | 23         | 0.30       | 0 | 0 | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 327711 | 0            | 105    | 120    | 25           | 17           | 2           | 34         | 0.06       | 0 | 0 | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    |
| 367788 | 0            | 59     | 146    | 12           | 30           | 0           | 32         | 0.00       | 0 | 0 | ... | 0    | 0    | 0    | 1    | 0    | 0    | 0    |
| 151235 | 0            | 35     | 50     | 5            | 9            | 3           | 13         | 0.23       | 0 | 0 | ... | 0    | 0    | 0    | 0    | 0    | 0    | 0    |

5 rows × 6008 columns

In [63]:

```
1 from sklearn.model_selection import train_test_split
2 X_train,X_test,y_train,y_test = train_test_split(final_df.iloc[:,1:],final_df.iloc[:,0].values,test_size=0.2,random_s
```

In [64]:

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score
3 rf = RandomForestClassifier()
4 rf.fit(X_train,y_train)
5 y_pred = rf.predict(X_test)
6 accuracy_score(y_test,y_pred)
```

Out[64]: 0.77

In [65]:

```
1 from xgboost import XGBClassifier
2 xgb = XGBClassifier()
3 xgb.fit(X_train,y_train)
4 y_pred = xgb.predict(X_test)
5 accuracy_score(y_test,y_pred)
```

Out[65]: 0.7645

In [66]:

```
1 from lightgbm import LGBMClassifier
2 lgb = LGBMClassifier()
3 lgb.fit(X_train,y_train)
4 y_pred = lgb.predict(X_test)
5 accuracy_score(y_test,y_pred)
```

Out[66]: 0.7691666666666666

## advance feature engineering!

```
In [110]: 1 # <!--
2 # 1. Token Features
3 # cwc_min: This is the ratio of the number of common words to the length of the smaller question
4 # cwc_max: This is the ratio of the number of common words to the length of the larger question
5 # csc_min: This is the ratio of the number of common stop words to the smaller stop word count among the two questions
6 # csc_max: This is the ratio of the number of common stop words to the larger stop word count among the two questions
7 # ctc_min: This is the ratio of the number of common tokens to the smaller token count among the two questions
8 # ctc_max: This is the ratio of the number of common tokens to the larger token count among the two questions
9 # last_word_eq: 1 if the last word in the two questions is same, 0 otherwise
10 # first_word_eq: 1 if the first word in the two questions is same, 0 otherwise
11
12 # 2. Length Based Features
13 # mean_len: Mean of the length of the two questions (number of words)
14 # abs_len_diff: Absolute difference between the length of the two questions (number of words)
15 # longest_substr_ratio: Ratio of the length of the longest substring among the two questions to the length of the smaller qu
16
17 # 3. Fuzzy Features
18 # fuzz_ratio: fuzz_ratio score from fuzzywuzzy
19 # fuzz_partial_ratio: fuzz_partial_ratio from fuzzywuzzy
20 # token_sort_ratio: token_sort_ratio from fuzzywuzzy
21 # token_set_ratio: token_set_ratio from fuzzywuzzy --> -->
```

```
In [116]: 1 import re
2 from bs4 import BeautifulSoup
3
4 import warnings
5 warnings.filterwarnings('ignore')
```

```
In [115]: 1 new_df = df.sample(30000, random_state=1)
2 new_df.head()
```

Out[115]:

|        | id     | qid1   | qid2   | question1   | question2   | is_duplicate |
|--------|--------|--------|--------|---|---|--------------|
| 237030 | 237030 | 33086  | 348102 | How can I stop playing video games?               | Should I stop playing video games with my child?  | 0            |
| 247341 | 247341 | 73272  | 8624   | Who is better Donald Trump or Hillary Clinton?    | Why is Hillary Clinton a better choice than Do... | 1            |
| 246425 | 246425 | 359482 | 359483 | What do you think is the chance that sometime ... | Do you think there will be another world war/n... | 1            |
| 306985 | 306985 | 1357   | 47020  | Why are so many questions posted to Quora that... | Why do people write questions on Quora that co... | 1            |
| 225863 | 225863 | 334315 | 334316 | Can there even be a movie ever rated 10/10 on ... | What are your 10/10 movies?                       | 0            |

```
In [146]: 1 new_df.nunique()
```

Out[146]:

|                      |       |
|----------------------|-------|
| id                   | 30000 |
| qid1                 | 28408 |
| qid2                 | 28335 |
| question1            | 28389 |
| question2            | 28311 |
| is_duplicate         | 2     |
| q1_len               | 248   |
| q2_len               | 281   |
| q1_num_words         | 64    |
| q2_num_words         | 79    |
| word_common          | 30    |
| word_total           | 75    |
| word_share           | 50    |
| cwc_min              | 125   |
| cwc_max              | 197   |
| csc_min              | 118   |
| csc_max              | 207   |
| ctc_min              | 410   |
| ctc_max              | 652   |
| last_word_eq         | 2     |
| first_word_eq        | 2     |
| abs_len_diff         | 52    |
| mean_len             | 98    |
| longest_substr_ratio | 2645  |
| fuzz_ratio           | 101   |
| fuzz_partial_ratio   | 92    |
| token_sort_ratio     | 101   |
| token_set_ratio      | 98    |
| dtype:               | int64 |



```
In [117]: 1 def preprocess(q):
2
3     q = str(q).lower().strip()
4
5     # Replace certain special characters with their string equivalents
6     q = q.replace('%', ' percent')
7     q = q.replace('$', ' dollar ')
8     q = q.replace('₹', ' rupee ')
9     q = q.replace('€', ' euro ')
10    q = q.replace('@', ' at ')
11
12    # The pattern '[math]' appears around 900 times in the whole dataset.
13    q = q.replace('[math]', '')
14
15    # Replacing some numbers with string equivalents (not perfect, can be done better to account for more cases)
16    q = q.replace(',000,000,000 ', 'b ')
17    q = q.replace(',000,000 ', 'm ')
18    q = q.replace(',000 ', 'k ')
19    q = re.sub(r'([0-9]+)00000000', r'\1b', q)
20    q = re.sub(r'([0-9]+)000000', r'\1m', q)
21    q = re.sub(r'([0-9]+)000', r'\1k', q)
22
23    # Decontracting words
24    # https://en.wikipedia.org/wiki/Wikipedia%3aList_of_English_contractions
25    # https://stackoverflow.com/a/19794953
26    contractions = {
27        "ain't": "am not",
28        "aren't": "are not",
29        "can't": "can not",
30        "can't've": "can not have",
31        "'cause": "because",
32        "could've": "could have",
33        "couldn't": "could not",
34        "couldn't've": "could not have",
35        "didn't": "did not",
36        "doesn't": "does not",
37        "don't": "do not",
38        "hadn't": "had not",
39        "hadn't've": "had not have",
40        "hasn't": "has not",
41        "haven't": "have not",
42        "he'd": "he would",
43        "he'd've": "he would have",
44        "he'll": "he will",
45        "he'll've": "he will have",
46        "he's": "he is",
47        "how'd": "how did",
48        "how'd'y": "how do you",
49        "how'll": "how will",
50        "how's": "how is",
51        "i'd": "i would",
52        "i'd've": "i would have",
53        "i'll": "i will",
54        "i'll've": "i will have",
55        "i'm": "i am",
56        "i've": "i have",
57        "isn't": "is not",
58        "it'd": "it would",
59        "it'd've": "it would have",
60        "it'll": "it will",
61        "it'll've": "it will have",
62        "it's": "it is",
63        "let's": "let us",
64        "ma'am": "madam",
65        "mayn't": "may not",
66        "might've": "might have",
67        "mightn't": "might not",
68        "mightn't've": "might not have",
69        "must've": "must have",
70        "mustn't": "must not",
71        "mustn't've": "must not have",
72        "needn't": "need not",
73        "needn't've": "need not have",
74        "o'clock": "of the clock",
75        "oughtn't": "ought not",
76        "oughtn't've": "ought not have",
77        "shan't": "shall not",
78        "sha'n't": "shall not",
79        "shan't've": "shall not have",
80        "she'd": "she would",
81        "she'd've": "she would have",
82        "she'll": "she will",
83        "she'll've": "she will have",
84        "she's": "she is",
85        "should've": "should have",
86        "shouldn't": "should not",
```

```

87 "shouldn't've": "should not have",
88 "so've": "so have",
89 "so's": "so as",
90 "that'd": "that would",
91 "that'd've": "that would have",
92 "that's": "that is",
93 "there'd": "there would",
94 "there'd've": "there would have",
95 "there's": "there is",
96 "they'd": "they would",
97 "they'd've": "they would have",
98 "they'll": "they will",
99 "they'll've": "they will have",
100 "they're": "they are",
101 "they've": "they have",
102 "to've": "to have",
103 "wasn't": "was not",
104 "we'd": "we would",
105 "we'd've": "we would have",
106 "we'll": "we will",
107 "we'll've": "we will have",
108 "we're": "we are",
109 "we've": "we have",
110 "weren't": "were not",
111 "what'll": "what will",
112 "what'll've": "what will have",
113 "what're": "what are",
114 "what's": "what is",
115 "what've": "what have",
116 "when's": "when is",
117 "when've": "when have",
118 "where'd": "where did",
119 "where's": "where is",
120 "where've": "where have",
121 "who'll": "who will",
122 "who'll've": "who will have",
123 "who's": "who is",
124 "who've": "who have",
125 "why's": "why is",
126 "why've": "why have",
127 "will've": "will have",
128 "won't": "will not",
129 "won't've": "will not have",
130 "would've": "would have",
131 "wouldn't": "would not",
132 "wouldn't've": "would not have",
133 "y'all": "you all",
134 "y'all'd": "you all would",
135 "y'all'd've": "you all would have",
136 "y'all're": "you all are",
137 "y'all've": "you all have",
138 "you'd": "you would",
139 "you'd've": "you would have",
140 "you'll": "you will",
141 "you'll've": "you will have",
142 "you're": "you are",
143 "you've": "you have"
144 }
145
146 q_decontracted = []
147
148 for word in q.split():
149     if word in contractions:
150         word = contractions[word]
151
152     q_decontracted.append(word)
153
154 q = ' '.join(q_decontracted)
155 q = q.replace("'ve", " have")
156 q = q.replace("n't", " not")
157 q = q.replace("'re", " are")
158 q = q.replace("'ll", " will")
159
160 # Removing HTML tags
161 q = BeautifulSoup(q)
162 q = q.get_text()
163
164 # Remove punctuations
165 pattern = re.compile('\W')
166 q = re.sub(pattern, ' ', q).strip()
167
168
169 return q

```

| In [ ]:   | 1                             |   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|-----------|-------------------------------|---|--------|---|---|--------------|-----------|-----------|--------------|--------------|-------------|--------------|--------------|------------------------------------|---|--------|--------|------------------------------------|---|---|---|---|----|--------|--------|--------|--------|---|---|---|---|---|------|-------|---|---|--------|--------|--------|---|---|---|---|---|----|-----|--------|--------|------|-------|---|---|------|-------|---|---|----|--------|--------|--------|--------|---|----------------------------|--------|--------|--------|---|----------------------------|---|----|----|----|---|---|
| In [118]: | 1                             | preprocess("I've already! wasn't <b>done</b>?")   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| Out[118]: | 'i have already was not done' |   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [119]: | 1                             | new_df['question1'] = new_df['question1'].apply(preprocess)   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | new_df['question2'] = new_df['question2'].apply(preprocess)   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [120]: | 1                             | new_df.head()   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| Out[120]: |                               | <table border="1"> <thead> <tr> <th></th><th>id</th><th>qid1</th><th>qid2</th><th>question1</th><th>question2</th><th>is_duplicate</th></tr> </thead> <tbody> <tr> <td>237030</td><td>237030</td><td>33086</td><td>348102</td><td>how can i stop playing video games</td><td>should i stop playing video games with my child</td><td>0</td></tr> <tr> <td>247341</td><td>247341</td><td>73272</td><td>8624</td><td>who is better donald trump or hillary clinton</td><td>why is hillary clinton a better choice than do...</td><td>1</td></tr> <tr> <td>246425</td><td>246425</td><td>359482</td><td>359483</td><td>what do you think is the chance that sometime ...</td><td>do you think there will be another world war n...</td><td>1</td></tr> <tr> <td>306985</td><td>306985</td><td>1357</td><td>47020</td><td>why are so many questions posted to quora that...</td><td>why do people write questions on quora that co...</td><td>1</td></tr> <tr> <td>225863</td><td>225863</td><td>334315</td><td>334316</td><td>can there even be a movie ever rated 10 10 on ...</td><td>what are your 10 10 movies</td><td>0</td></tr> </tbody> </table>   |        | id  | qid1  | qid2         | question1 | question2 | is_duplicate | 237030       | 237030      | 33086        | 348102       | how can i stop playing video games | should i stop playing video games with my child | 0      | 247341 | 247341                             | 73272   | 8624  | who is better donald trump or hillary clinton | why is hillary clinton a better choice than do... | 1  | 246425 | 246425 | 359482 | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1   | 306985  | 306985  | 1357 | 47020 | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1      | 225863 | 225863 | 334315  | 334316  | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | id                            | qid1  | qid2   | question1   | question2   | is_duplicate |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 237030    | 237030                        | 33086   | 348102 | how can i stop playing video games                | should i stop playing video games with my child   | 0            |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 247341    | 247341                        | 73272   | 8624   | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1            |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 246425    | 246425                        | 359482  | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1            |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 306985    | 306985                        | 1357  | 47020  | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1            |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 225863    | 225863                        | 334315  | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0            |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [121]: | 1                             | new_df['q1_len'] = new_df['question1'].str.len()  |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | new_df['q2_len'] = new_df['question2'].str.len()  |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [122]: | 1                             | new_df['q1_num_words'] = new_df['question1'].apply(lambda row: len(row.split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | new_df['q2_num_words'] = new_df['question2'].apply(lambda row: len(row.split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 3                             | new_df.head()   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| Out[122]: |                               | <table border="1"> <thead> <tr> <th></th><th>id</th><th>qid1</th><th>qid2</th><th>question1</th><th>question2</th><th>is_duplicate</th><th>q1_len</th><th>q2_len</th><th>q1_num_words</th><th>q2_num_words</th></tr> </thead> <tbody> <tr> <td>237030</td><td>237030</td><td>33086</td><td>348102</td><td>how can i stop playing video games</td><td>should i stop playing video games with my child</td><td>0</td><td>34</td><td>47</td><td>7</td><td>9</td></tr> <tr> <td>247341</td><td>247341</td><td>73272</td><td>8624</td><td>who is better donald trump or hillary clinton</td><td>why is hillary clinton a better choice than do...</td><td>1</td><td>45</td><td>56</td><td>8</td><td>10</td></tr> <tr> <td>246425</td><td>246425</td><td>359482</td><td>359483</td><td>what do you think is the chance that sometime ...</td><td>do you think there will be another world war n...</td><td>1</td><td>137</td><td>76</td><td>29</td><td>15</td></tr> <tr> <td>306985</td><td>306985</td><td>1357</td><td>47020</td><td>why are so many questions posted to quora that...</td><td>why do people write questions on quora that co...</td><td>1</td><td>85</td><td>85</td><td>16</td><td>16</td></tr> <tr> <td>225863</td><td>225863</td><td>334315</td><td>334316</td><td>can there even be a movie ever rated 10 10 on ...</td><td>what are your 10 10 movies</td><td>0</td><td>50</td><td>26</td><td>12</td><td>6</td></tr> </tbody> </table>  |        | id  | qid1  | qid2         | question1 | question2 | is_duplicate | q1_len       | q2_len      | q1_num_words | q2_num_words | 237030                             | 237030  | 33086  | 348102 | how can i stop playing video games | should i stop playing video games with my child | 0   | 34  | 47  | 7  | 9      | 247341 | 247341 | 73272  | 8624  | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1   | 45  | 56   | 8     | 10  | 246425  | 246425 | 359482 | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1   | 137   | 76  | 29 | 15  | 306985 | 306985 | 1357 | 47020 | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1    | 85    | 85  | 16  | 16 | 225863 | 225863 | 334315 | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies | 0      | 50     | 26     | 12  | 6                          |   |    |    |    |   |   |
|           | id                            | qid1  | qid2   | question1   | question2   | is_duplicate | q1_len    | q2_len    | q1_num_words | q2_num_words |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 237030    | 237030                        | 33086   | 348102 | how can i stop playing video games                | should i stop playing video games with my child   | 0            | 34        | 47        | 7            | 9            |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 247341    | 247341                        | 73272   | 8624   | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1            | 45        | 56        | 8            | 10           |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 246425    | 246425                        | 359482  | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1            | 137       | 76        | 29           | 15           |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 306985    | 306985                        | 1357  | 47020  | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1            | 85        | 85        | 16           | 16           |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 225863    | 225863                        | 334315  | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0            | 50        | 26        | 12           | 6            |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [123]: | 1                             | def common_words(row):  |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 3                             | w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 4                             | return len(w1 & w2)   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [124]: | 1                             | new_df['word_common'] = new_df.apply(common_words, axis=1)  |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | new_df.head()   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| Out[124]: |                               | <table border="1"> <thead> <tr> <th></th><th>id</th><th>qid1</th><th>qid2</th><th>question1</th><th>question2</th><th>is_duplicate</th><th>q1_len</th><th>q2_len</th><th>q1_num_words</th><th>q2_num_words</th><th>word_common</th></tr> </thead> <tbody> <tr> <td>237030</td><td>237030</td><td>33086</td><td>348102</td><td>how can i stop playing video games</td><td>should i stop playing video games with my child</td><td>0</td><td>34</td><td>47</td><td>7</td><td>9</td><td>5</td></tr> <tr> <td>247341</td><td>247341</td><td>73272</td><td>8624</td><td>who is better donald trump or hillary clinton</td><td>why is hillary clinton a better choice than do...</td><td>1</td><td>45</td><td>56</td><td>8</td><td>10</td><td>6</td></tr> <tr> <td>246425</td><td>246425</td><td>359482</td><td>359483</td><td>what do you think is the chance that sometime ...</td><td>do you think there will be another world war n...</td><td>1</td><td>137</td><td>76</td><td>29</td><td>15</td><td>13</td></tr> <tr> <td>306985</td><td>306985</td><td>1357</td><td>47020</td><td>why are so many questions posted to quora that...</td><td>why do people write questions on quora that co...</td><td>1</td><td>85</td><td>85</td><td>16</td><td>16</td><td>5</td></tr> <tr> <td>225863</td><td>225863</td><td>334315</td><td>334316</td><td>can there even be a movie ever rated 10 10 on ...</td><td>what are your 10 10 movies</td><td>0</td><td>50</td><td>26</td><td>12</td><td>6</td><td>1</td></tr> </tbody> </table> |        | id  | qid1  | qid2         | question1 | question2 | is_duplicate | q1_len       | q2_len      | q1_num_words | q2_num_words | word_common                        | 237030  | 237030 | 33086  | 348102                             | how can i stop playing video games              | should i stop playing video games with my child | 0   | 34  | 47 | 7      | 9      | 5      | 247341 | 247341  | 73272   | 8624  | who is better donald trump or hillary clinton | why is hillary clinton a better choice than do... | 1    | 45    | 56  | 8   | 10     | 6      | 246425 | 246425  | 359482  | 359483  | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1  | 137 | 76     | 29     | 15   | 13    | 306985  | 306985  | 1357 | 47020 | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1  | 85     | 85     | 16     | 16     | 5   | 225863                     | 225863 | 334315 | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies | 0 | 50 | 26 | 12 | 6 | 1 |
|           | id                            | qid1  | qid2   | question1   | question2   | is_duplicate | q1_len    | q2_len    | q1_num_words | q2_num_words | word_common |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 237030    | 237030                        | 33086   | 348102 | how can i stop playing video games                | should i stop playing video games with my child   | 0            | 34        | 47        | 7            | 9            | 5           |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 247341    | 247341                        | 73272   | 8624   | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1            | 45        | 56        | 8            | 10           | 6           |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 246425    | 246425                        | 359482  | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1            | 137       | 76        | 29           | 15           | 13          |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 306985    | 306985                        | 1357  | 47020  | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1            | 85        | 85        | 16           | 16           | 5           |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| 225863    | 225863                        | 334315  | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0            | 50        | 26        | 12           | 6            | 1           |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
| In [125]: | 1                             | def total_words(row):   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 2                             | w1 = set(map(lambda word: word.lower().strip(), row['question1'].split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 3                             | w2 = set(map(lambda word: word.lower().strip(), row['question2'].split(" ")))   |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |
|           | 4                             | return (len(w1) + len(w2))  |        |   |   |              |           |           |              |              |             |              |              |                                    |   |        |        |                                    |   |   |   |   |    |        |        |        |        |   |   |   |   |   |      |       |   |   |        |        |        |   |   |   |   |   |    |     |        |        |      |       |   |   |      |       |   |   |    |        |        |        |        |   |                            |        |        |        |   |                            |   |    |    |    |   |   |

```
In [126]: 1 new_df['word_total'] = new_df.apply(total_words, axis=1)
2 new_df.head()
```

Out[126]:

|        |        |        |        |  | question1   | question2   | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total |
|--------|--------|--------|--------|--|---|---|--------------|--------|--------|--------------|--------------|-------------|------------|
| 237030 | 237030 | 33086  | 348102 |  | how can i stop playing video games                | should i stop playing video games with my child   | 0            | 34     | 47     | 7            | 9            | 5           | 16         |
| 247341 | 247341 | 73272  | 8624   |  | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1            | 45     | 56     | 8            | 10           | 6           | 18         |
| 246425 | 246425 | 359482 | 359483 |  | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1            | 137    | 76     | 29           | 15           | 13          | 40         |
| 306985 | 306985 | 1357   | 47020  |  | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1            | 85     | 85     | 16           | 16           | 5           | 30         |
| 225863 | 225863 | 334315 | 334316 |  | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0            | 50     | 26     | 12           | 6            | 1           | 16         |

```
In [127]: 1 new_df['word_share'] = round(new_df['word_common']/new_df['word_total'],2)
2 new_df.head()
3
```

Out[127]:

|        |        |        |        |  | question1   | question2   | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_share |
|--------|--------|--------|--------|--|---|---|--------------|--------|--------|--------------|--------------|-------------|------------|------------|
| 237030 | 237030 | 33086  | 348102 |  | how can i stop playing video games                | should i stop playing video games with my child   | 0            | 34     | 47     | 7            | 9            | 5           | 16         | 0.31       |
| 247341 | 247341 | 73272  | 8624   |  | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1            | 45     | 56     | 8            | 10           | 6           | 18         | 0.33       |
| 246425 | 246425 | 359482 | 359483 |  | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1            | 137    | 76     | 29           | 15           | 13          | 40         | 0.32       |
| 306985 | 306985 | 1357   | 47020  |  | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1            | 85     | 85     | 16           | 16           | 5           | 30         | 0.17       |
| 225863 | 225863 | 334315 | 334316 |  | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0            | 50     | 26     | 12           | 6            | 1           | 16         | 0.06       |

```
In [128]: 1 # Advanced Features
2 from nltk.corpus import stopwords
3
4 def fetch_token_features(row):
5
6     q1 = row['question1']
7     q2 = row['question2']
8
9     SAFE_DIV = 0.0001
10
11    STOP_WORDS = stopwords.words("english")
12
13    token_features = [0.0]*8
14
15    # Converting the Sentence into Tokens:
16    q1_tokens = q1.split()
17    q2_tokens = q2.split()
18
19    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
20        return token_features
21
22    # Get the non-stopwords in Questions
23    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
24    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])
25
26    #Get the stopwords in Questions
27    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
28    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])
29
30    # Get the common non-stopwords from Question pair
31    common_word_count = len(q1_words.intersection(q2_words))
32
33    # Get the common stopwords from Question pair
34    common_stop_count = len(q1_stops.intersection(q2_stops))
35
36    # Get the common Tokens from Question pair
37    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))
38
39
40    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)) + SAFE_DIV)
41    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)) + SAFE_DIV)
42    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)) + SAFE_DIV)
43    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)) + SAFE_DIV)
44    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
45    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
46
47    # Last word of both question is same or not
48    token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])
49
50    # First word of both question is same or not
51    token_features[7] = int(q1_tokens[0] == q2_tokens[0])
52
53    return token_features
```

```
In [129]: 1 token_features = new_df.apply(fetch_token_features, axis=1)
2
3 new_df["cwc_min"]      = list(map(lambda x: x[0], token_features))
4 new_df["cwc_max"]      = list(map(lambda x: x[1], token_features))
5 new_df["csc_min"]      = list(map(lambda x: x[2], token_features))
6 new_df["csc_max"]      = list(map(lambda x: x[3], token_features))
7 new_df["ctc_min"]      = list(map(lambda x: x[4], token_features))
8 new_df["ctc_max"]      = list(map(lambda x: x[5], token_features))
9 new_df["last_word_eq"] = list(map(lambda x: x[6], token_features))
10 new_df["first_word_eq"] = list(map(lambda x: x[7], token_features))
```

In [130]: 1 new\_df

Out[130]:

|  |  |  |  |  | id     | qid1   | qid2   | question1 | question2   | is_duplicate                                      | q1_len | q2_len | q1_num_words | q2_num_words | ... | word_total | word_share | cwc_min | c...     |   |
|--|--|--|--|--|--------|--------|--------|-----------|---|---|--------|--------|--------------|--------------|-----|------------|------------|---------|----------|---|
|  |  |  |  |  | 237030 | 237030 | 33086  | 348102    | how can i stop playing video games                | should i stop playing video games with my child   | 0      | 34     | 47           | 7            | 9   | ...        | 16         | 0.31    | 0.999975 | 0 |
|  |  |  |  |  | 247341 | 247341 | 73272  | 8624      | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... | 1      | 45     | 56           | 8            | 10  | ...        | 18         | 0.33    | 0.999980 | 0 |
|  |  |  |  |  | 246425 | 246425 | 359482 | 359483    | what do you think is the chance that sometime ... | do you think there will be another world war n... | 1      | 137    | 76           | 29           | 15  | ...        | 40         | 0.32    | 0.857131 | 0 |
|  |  |  |  |  | 306985 | 306985 | 1357   | 47020     | why are so many questions posted to quora that... | why do people write questions on quora that co... | 1      | 85     | 85           | 16           | 16  | ...        | 30         | 0.17    | 0.374995 | 0 |
|  |  |  |  |  | 225863 | 225863 | 334315 | 334316    | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        | 0      | 50     | 26           | 12           | 6   | ...        | 16         | 0.06    | 0.499975 | 0 |
|  |  |  |  |  | 298455 | 298455 | 88357  | 148625    | why is donald trump still ducking his income...   | why does not trump begin speaking the truth ...   | 1      | 63     | 74           | 13           | 15  | ...        | 26         | 0.19    | 0.285710 | 0 |
|  |  |  |  |  | 16366  | 16366  | 31205  | 31206     | how often do women get wet when they see a rea... | how often do women initiate a conversation wit... | 0      | 57     | 90           | 13           | 17  | ...        | 29         | 0.28    | 0.499994 | 0 |
|  |  |  |  |  | 379818 | 379818 | 11786  | 511366    | how do i disable voice data for text to speec...  | how do i convert call center recordings into t... | 0      | 72     | 148          | 15           | 28  | ...        | 42         | 0.17    | 0.428565 | 0 |
|  |  |  |  |  | 54795  | 54795  | 96691  | 96692     | how does mildew form on carpets                   | how can you prevent mildew from forming on car... | 0      | 31     | 50           | 6            | 9   | ...        | 15         | 0.27    | 0.666644 | 0 |
|  |  |  |  |  | 335036 | 335036 | 462289 | 462290    | how do i leave a legacy                           | how do you leave a legacy to your children        | 0      | 23     | 42           | 6            | 9   | ...        | 15         | 0.33    | 0.999950 | 0 |

30000 rows × 21 columns

In [133]: 1 # !pip install distance

```
In [134]: 1 import distance
2
3 def fetch_length_features(row):
4
5     q1 = row['question1']
6     q2 = row['question2']
7
8     length_features = [0.0]*3
9
10    # Converting the Sentence into Tokens:
11    q1_tokens = q1.split()
12    q2_tokens = q2.split()
13
14    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
15        return length_features
16
17    # Absolute Length features
18    length_features[0] = abs(len(q1_tokens) - len(q2_tokens))
19
20    #Average Token Length of both Questions
21    length_features[1] = (len(q1_tokens) + len(q2_tokens))/2
22
23    strs = list(distance.lcsubstrings(q1, q2))
24    length_features[2] = len(strs[0]) / (min(len(q1), len(q2)) + 1)
25
26
27    return length_features
```

```
In [135]: 1 length_features = new_df.apply(fetch_length_features, axis=1)
2
3 new_df['abs_len_diff'] = list(map(lambda x: x[0], length_features))
4 new_df['mean_len'] = list(map(lambda x: x[1], length_features))
5 new_df['longest_substr_ratio'] = list(map(lambda x: x[2], length_features))
```

```
In [136]: 1 new_df.head()
```

Out[136]:

|        |        |        |        |  |  | how can i<br>stop<br>playing<br>video<br>games                      | should i<br>stop<br>playing<br>video<br>games<br>with my<br>child |  | 0 | 34  | 47 |  | 7  |  | 9  | ... | 0.799984 | 0.333322 | 0.249994 | 0.714271 |  |  |  |  |
|--------|--------|--------|--------|--|--|---|---|--|---|-----|----|--|----|--|----|-----|----------|----------|----------|----------|--|--|--|--|
| 237030 | 237030 | 33086  | 348102 |  |  |   |   |  |   |     |    |  |    |  |    |     |          |          |          |          |  |  |  |  |
| 247341 | 247341 | 73272  | 8624   |  |  | who is<br>better<br>donald<br>trump or<br>hillary<br>clinton        | why is<br>hillary<br>clinton a<br>better<br>choice<br>than do...  |  | 1 | 45  | 56 |  | 8  |  | 10 | ... | 0.833319 | 0.333322 | 0.249994 | 0.74999  |  |  |  |  |
| 246425 | 246425 | 359482 | 359483 |  |  | what do<br>you think<br>is the<br>chance<br>that<br>sometime<br>... | do you<br>think there<br>will be<br>another<br>world war<br>n...  |  | 1 | 137 | 76 |  | 29 |  | 15 | ... | 0.499996 | 0.999986 | 0.538457 | 0.86666  |  |  |  |  |
| 306985 | 306985 | 1357   | 47020  |  |  | why are<br>so many<br>questions<br>posted to<br>quora<br>that...    | why do<br>people<br>write<br>questions<br>on quora<br>that co...  |  | 1 | 85  | 85 |  | 16 |  | 16 | ... | 0.333330 | 0.333328 | 0.285710 | 0.312491 |  |  |  |  |
| 225863 | 225863 | 334315 | 334316 |  |  | can there<br>even be a<br>movie<br>ever rated<br>10 10 on<br>...    | what are<br>your 10<br>10 movies                                  |  | 0 | 50  | 26 |  | 12 |  | 6  | ... | 0.166664 | 0.000000 | 0.000000 | 0.16666  |  |  |  |  |

5 rows × 24 columns

```
In [139]: 1 # !pip install fuzzywuzzy
```

```
In [140]: 1 # Fuzzy Features
2 from fuzzywuzzy import fuzz
3
4 def fetch_fuzzy_features(row):
5
6     q1 = row['question1']
7     q2 = row['question2']
8
9     fuzzy_features = [0.0]*4
10
11     # fuzz_ratio
12     fuzzy_features[0] = fuzz.QRatio(q1, q2)
13
14     # fuzz_partial_ratio
15     fuzzy_features[1] = fuzz.partial_ratio(q1, q2)
16
17     # token_sort_ratio
18     fuzzy_features[2] = fuzz.token_sort_ratio(q1, q2)
19
20     # token_set_ratio
21     fuzzy_features[3] = fuzz.token_set_ratio(q1, q2)
22
23     return fuzzy_features
```

```
In [141]: 1 fuzzy_features = new_df.apply(fetch_fuzzy_features, axis=1)
2
3 # Creating new feature columns for fuzzy features
4 new_df['fuzz_ratio'] = list(map(lambda x: x[0], fuzzy_features))
5 new_df['fuzz_partial_ratio'] = list(map(lambda x: x[1], fuzzy_features))
6 new_df['token_sort_ratio'] = list(map(lambda x: x[2], fuzzy_features))
7 new_df['token_set_ratio'] = list(map(lambda x: x[3], fuzzy_features))
```

```
In [142]: 1 print(new_df.shape)
2 new_df.head()
```

(30000, 28)

Out[142]:

| 237030 | 237030 | 33086  | 348102 | how can i stop playing video games                | should i stop playing video games with my child   |  | 0 | 34  | 47 |  | 7  |  | 9  | ... | 0.555549 |  | 0.0 |  | 0.0 |  |  |  |  |  |
|--------|--------|--------|--------|---|---|--|---|-----|----|--|----|--|----|-----|----------|--|-----|--|-----|--|--|--|--|--|
| 247341 | 247341 | 73272  | 8624   | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... |  | 1 | 45  | 56 |  | 8  |  | 10 | ... | 0.599994 |  | 0.0 |  | 0.0 |  |  |  |  |  |
| 246425 | 246425 | 359482 | 359483 | what do you think is the chance that sometime ... | do you think there will be another world war n... |  | 1 | 137 | 76 |  | 29 |  | 15 | ... | 0.464284 |  | 0.0 |  | 0.0 |  |  |  |  |  |
| 306985 | 306985 | 1357   | 47020  | why are so many questions posted to quora that... | why do people write questions on quora that co... |  | 1 | 85  | 85 |  | 16 |  | 16 | ... | 0.312498 |  | 0.0 |  | 1.0 |  |  |  |  |  |
| 225863 | 225863 | 334315 | 334316 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        |  | 0 | 50  | 26 |  | 12 |  | 6  | ... | 0.083333 |  | 0.0 |  | 0.0 |  |  |  |  |  |

5 rows × 28 columns

In [143]: 1 new\_df

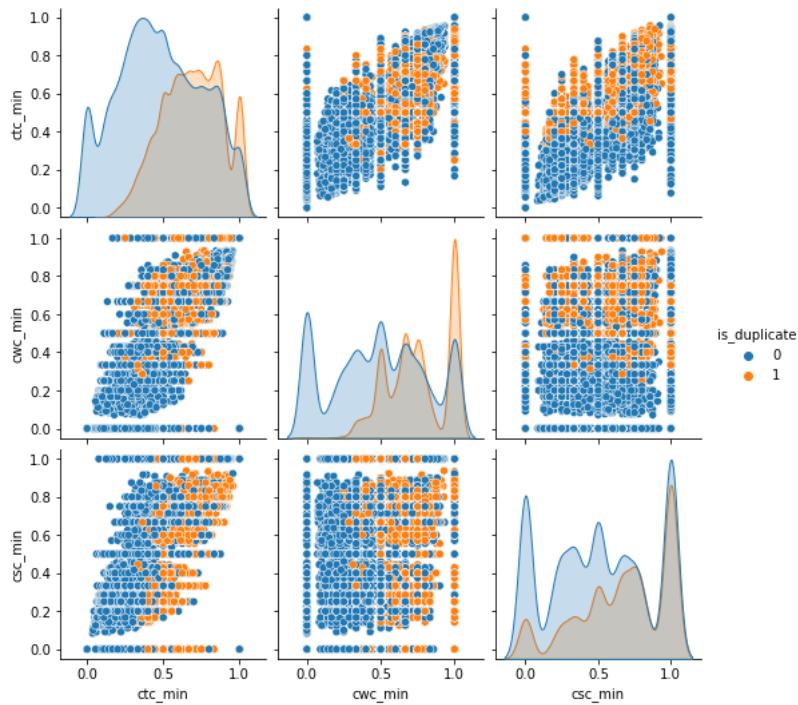
Out[143]:

|     | num_words | ...      | ctc_max | last_word_eq | first_word_eq | abs_len_diff | mean_len | longest_substr_ratio | fuzz_ratio | fuzz_partial_ratio | token_sort_ratio | token_set_ratio |
|-----|-----------|----------|---------|--------------|---------------|--------------|----------|----------------------|------------|--------------------|------------------|-----------------|
| 9   | ...       | 0.555549 | 0.0     | 0.0          | 2.0           | 8.0          | 0.771429 | 72                   | 85         | 69                 | 87               |                 |
| 10  | ...       | 0.599994 | 0.0     | 0.0          | 2.0           | 9.0          | 0.347826 | 42                   | 49         | 83                 | 92               |                 |
| 15  | ...       | 0.464284 | 0.0     | 0.0          | 13.0          | 21.5         | 0.298701 | 46                   | 55         | 70                 | 94               |                 |
| 16  | ...       | 0.312498 | 0.0     | 1.0          | 0.0           | 16.0         | 0.139535 | 53                   | 53         | 51                 | 59               |                 |
| 6   | ...       | 0.083333 | 0.0     | 0.0          | 6.0           | 9.0          | 0.259259 | 42                   | 46         | 50                 | 46               |                 |
| ... | ...       | ...      | ...     | ...          | ...           | ...          | ...      | ...                  | ...        | ...                | ...              |                 |
| 15  | ...       | 0.307690 | 0.0     | 1.0          | 2.0           | 12.0         | 0.171875 | 53                   | 46         | 50                 | 53               |                 |
| 17  | ...       | 0.470585 | 0.0     | 1.0          | 4.0           | 15.0         | 0.327586 | 59                   | 58         | 61                 | 77               |                 |
| 28  | ...       | 0.222221 | 0.0     | 1.0          | 14.0          | 20.0         | 0.123288 | 31                   | 40         | 41                 | 53               |                 |
| 9   | ...       | 0.444440 | 1.0     | 1.0          | 3.0           | 7.5          | 0.343750 | 72                   | 74         | 64                 | 81               |                 |
| 9   | ...       | 0.555549 | 0.0     | 1.0          | 3.0           | 7.5          | 0.625000 | 68                   | 87         | 65                 | 95               |                 |



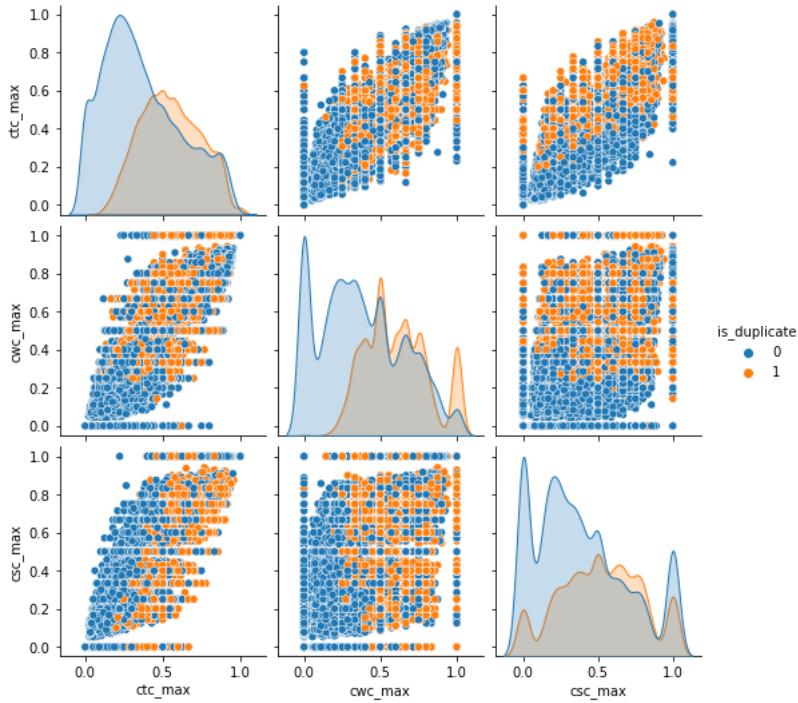
```
In [153]: 1 sns.pairplot(new_df[['ctc_min', 'cwc_min', 'csc_min', 'is_duplicate']],hue='is_duplicate')
2
```

Out[153]: <seaborn.axisgrid.PairGrid at 0x18704e6a850>



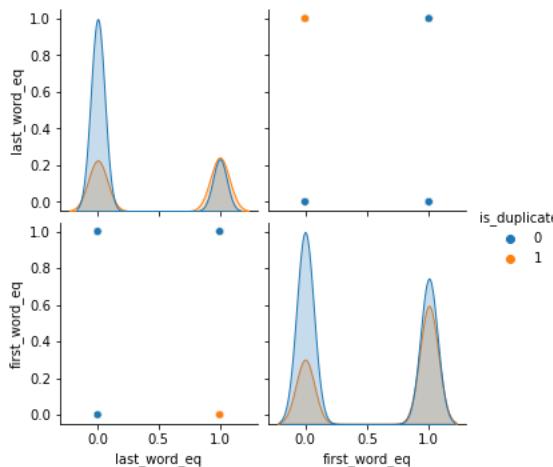
```
In [154]: 1 sns.pairplot(new_df[['ctc_max', 'cwc_max', 'csc_max', 'is_duplicate']],hue='is_duplicate')
2
```

Out[154]: <seaborn.axisgrid.PairGrid at 0x18705597220>



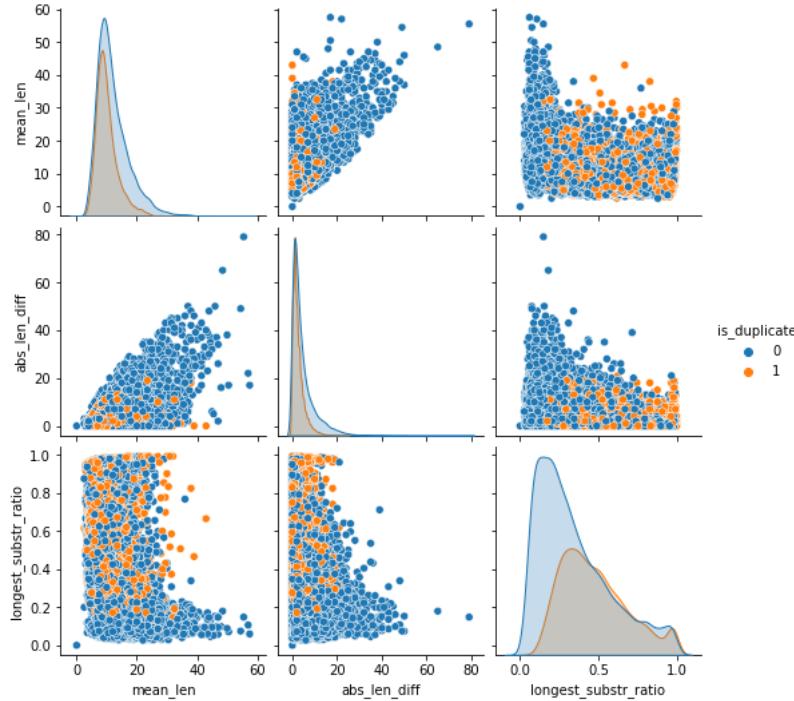
```
In [155]: 1 sns.pairplot(new_df[['last_word_eq', 'first_word_eq', 'is_duplicate']], hue='is_duplicate')
2
```

Out[155]: <seaborn.axisgrid.PairGrid at 0x187073538e0>



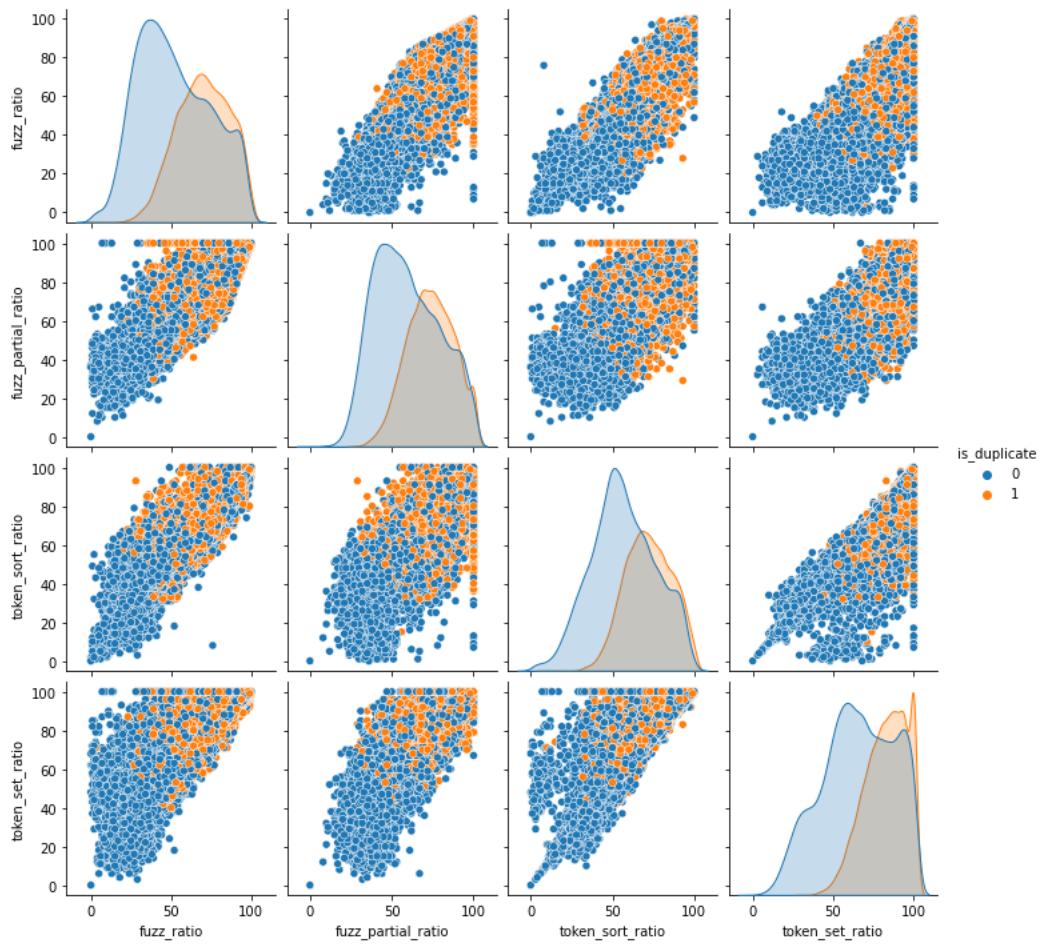
```
In [156]: 1 sns.pairplot(new_df[['mean_len', 'abs_len_diff', 'longest_substr_ratio', 'is_duplicate']], hue='is_duplicate')
2
```

Out[156]: <seaborn.axisgrid.PairGrid at 0x18708709d30>



```
In [157]: 1 sns.pairplot(new_df[['fuzz_ratio', 'fuzz_partial_ratio','token_sort_ratio','token_set_ratio', 'is_duplicate']],hue='is_duplicate')
```

Out[157]: <seaborn.axisgrid.PairGrid at 0x1870a1d2b80>



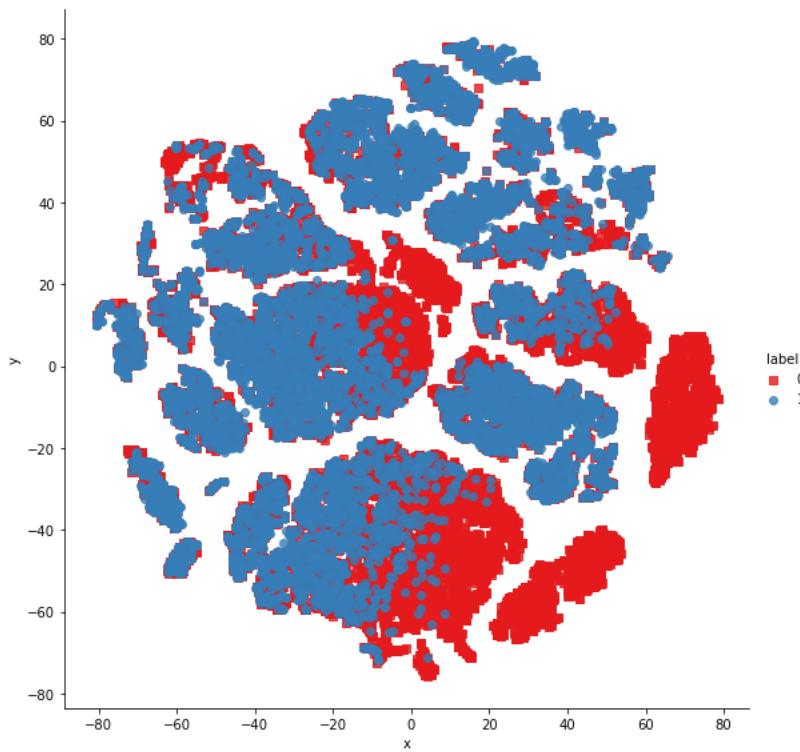
```
In [158]: 1 # Using TSNE for Dimensionality reduction for 15 Features(Generated after cleaning the data) to 3 dimension
2
3 from sklearn.preprocessing import MinMaxScaler
4
5 X = MinMaxScaler().fit_transform(new_df[['cwc_min', 'cwc_max', 'csc_min', 'csc_max', 'ctc_min', 'ctc_max', 'last_word_eq']])
6 y = new_df['is_duplicate'].values
```

```
In [159]: 1 from sklearn.manifold import TSNE
2
3 tsne2d = TSNE(
4     n_components=2,
5     init='random', # pca
6     random_state=101,
7     method='barnes_hut',
8     n_iter=1000,
9     verbose=2,
10    angle=0.5
11 ).fit_transform(X)

[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 30000 samples in 0.159s...
[t-SNE] Computed neighbors for 30000 samples in 6.390s...
[t-SNE] Computed conditional probabilities for sample 1000 / 30000
[t-SNE] Computed conditional probabilities for sample 2000 / 30000
[t-SNE] Computed conditional probabilities for sample 3000 / 30000
[t-SNE] Computed conditional probabilities for sample 4000 / 30000
[t-SNE] Computed conditional probabilities for sample 5000 / 30000
[t-SNE] Computed conditional probabilities for sample 6000 / 30000
[t-SNE] Computed conditional probabilities for sample 7000 / 30000
[t-SNE] Computed conditional probabilities for sample 8000 / 30000
[t-SNE] Computed conditional probabilities for sample 9000 / 30000
[t-SNE] Computed conditional probabilities for sample 10000 / 30000
[t-SNE] Computed conditional probabilities for sample 11000 / 30000
[t-SNE] Computed conditional probabilities for sample 12000 / 30000
[t-SNE] Computed conditional probabilities for sample 13000 / 30000
[t-SNE] Computed conditional probabilities for sample 14000 / 30000
[t-SNE] Computed conditional probabilities for sample 15000 / 30000
[t-SNE] Computed conditional probabilities for sample 16000 / 30000
[t-SNE] Computed conditional probabilities for sample 17000 / 30000
[t-SNE] Computed conditional probabilities for sample 18000 / 30000
[t-SNE] Computed conditional probabilities for sample 19000 / 30000
[t-SNE] Computed conditional probabilities for sample 20000 / 30000
[t-SNE] Computed conditional probabilities for sample 21000 / 30000
[t-SNE] Computed conditional probabilities for sample 22000 / 30000
[t-SNE] Computed conditional probabilities for sample 23000 / 30000
[t-SNE] Computed conditional probabilities for sample 24000 / 30000
[t-SNE] Computed conditional probabilities for sample 25000 / 30000
[t-SNE] Computed conditional probabilities for sample 26000 / 30000
[t-SNE] Computed conditional probabilities for sample 27000 / 30000
[t-SNE] Computed conditional probabilities for sample 28000 / 30000
[t-SNE] Computed conditional probabilities for sample 29000 / 30000
[t-SNE] Computed conditional probabilities for sample 30000 / 30000
[t-SNE] Mean sigma: 0.087420
[t-SNE] Computed conditional probabilities in 1.120s
[t-SNE] Iteration 50: error = 111.1139450, gradient norm = 0.0000528 (50 iterations in 9.618s)
[t-SNE] Iteration 100: error = 92.4731140, gradient norm = 0.0029080 (50 iterations in 10.295s)
[t-SNE] Iteration 150: error = 86.6217270, gradient norm = 0.0014014 (50 iterations in 8.987s)
[t-SNE] Iteration 200: error = 84.3923798, gradient norm = 0.0009428 (50 iterations in 9.675s)
[t-SNE] Iteration 250: error = 83.2190552, gradient norm = 0.0007072 (50 iterations in 9.664s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 83.219055
[t-SNE] Iteration 300: error = 3.8108232, gradient norm = 0.0011424 (50 iterations in 9.858s)
[t-SNE] Iteration 350: error = 3.2906470, gradient norm = 0.0006687 (50 iterations in 10.132s)
[t-SNE] Iteration 400: error = 2.9285789, gradient norm = 0.0004417 (50 iterations in 9.777s)
[t-SNE] Iteration 450: error = 2.6810791, gradient norm = 0.0003204 (50 iterations in 9.914s)
[t-SNE] Iteration 500: error = 2.5010517, gradient norm = 0.0002474 (50 iterations in 9.856s)
[t-SNE] Iteration 550: error = 2.3637815, gradient norm = 0.0001989 (50 iterations in 9.691s)
[t-SNE] Iteration 600: error = 2.2556167, gradient norm = 0.0001645 (50 iterations in 9.516s)
[t-SNE] Iteration 650: error = 2.1680160, gradient norm = 0.0001390 (50 iterations in 9.498s)
[t-SNE] Iteration 700: error = 2.0954065, gradient norm = 0.0001198 (50 iterations in 9.805s)
[t-SNE] Iteration 750: error = 2.0341637, gradient norm = 0.0001048 (50 iterations in 10.480s)
[t-SNE] Iteration 800: error = 1.9818006, gradient norm = 0.0000927 (50 iterations in 10.568s)
[t-SNE] Iteration 850: error = 1.9364457, gradient norm = 0.0000829 (50 iterations in 10.624s)
[t-SNE] Iteration 900: error = 1.8966687, gradient norm = 0.0000748 (50 iterations in 10.716s)
[t-SNE] Iteration 950: error = 1.8615599, gradient norm = 0.0000680 (50 iterations in 10.481s)
[t-SNE] Iteration 1000: error = 1.8303556, gradient norm = 0.0000622 (50 iterations in 10.986s)
[t-SNE] KL divergence after 1000 iterations: 1.830356
```

```
In [160]: 1 x_df = pd.DataFrame({'x':tsne2d[:,0], 'y':tsne2d[:,1] , 'label':y})  
2  
3 # draw the plot in appropriate place in the grid  
4 sns.lmplot(data=x_df, x='x', y='y', hue='label', fit_reg=False, size=8, palette="Set1", markers=['s', 'o'])
```

Out[160]: <seaborn.axisgrid.FacetGrid at 0x1870bdb3dc0>

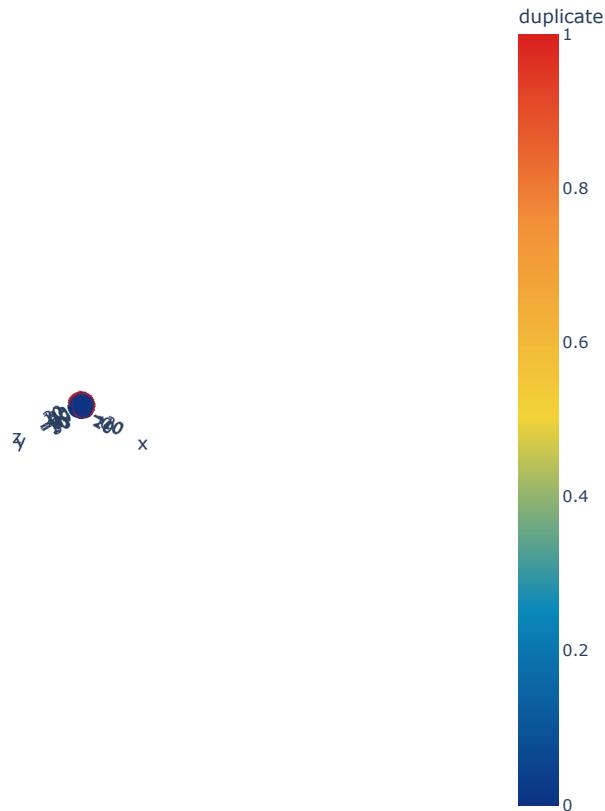


```
In [161]: 1 tsne3d = TSNE(
2     n_components=3,
3     init='random', # pca
4     random_state=101,
5     method='barnes_hut',
6     n_iter=1000,
7     verbose=2,
8     angle=0.5
9 ).fit_transform(X)

[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 30000 samples in 0.170s...
[t-SNE] Computed neighbors for 30000 samples in 8.881s...
[t-SNE] Computed conditional probabilities for sample 1000 / 30000
[t-SNE] Computed conditional probabilities for sample 2000 / 30000
[t-SNE] Computed conditional probabilities for sample 3000 / 30000
[t-SNE] Computed conditional probabilities for sample 4000 / 30000
[t-SNE] Computed conditional probabilities for sample 5000 / 30000
[t-SNE] Computed conditional probabilities for sample 6000 / 30000
[t-SNE] Computed conditional probabilities for sample 7000 / 30000
[t-SNE] Computed conditional probabilities for sample 8000 / 30000
[t-SNE] Computed conditional probabilities for sample 9000 / 30000
[t-SNE] Computed conditional probabilities for sample 10000 / 30000
[t-SNE] Computed conditional probabilities for sample 11000 / 30000
[t-SNE] Computed conditional probabilities for sample 12000 / 30000
[t-SNE] Computed conditional probabilities for sample 13000 / 30000
[t-SNE] Computed conditional probabilities for sample 14000 / 30000
[t-SNE] Computed conditional probabilities for sample 15000 / 30000
[t-SNE] Computed conditional probabilities for sample 16000 / 30000
[t-SNE] Computed conditional probabilities for sample 17000 / 30000
[t-SNE] Computed conditional probabilities for sample 18000 / 30000
[t-SNE] Computed conditional probabilities for sample 19000 / 30000
[t-SNE] Computed conditional probabilities for sample 20000 / 30000
[t-SNE] Computed conditional probabilities for sample 21000 / 30000
[t-SNE] Computed conditional probabilities for sample 22000 / 30000
[t-SNE] Computed conditional probabilities for sample 23000 / 30000
[t-SNE] Computed conditional probabilities for sample 24000 / 30000
[t-SNE] Computed conditional probabilities for sample 25000 / 30000
[t-SNE] Computed conditional probabilities for sample 26000 / 30000
[t-SNE] Computed conditional probabilities for sample 27000 / 30000
[t-SNE] Computed conditional probabilities for sample 28000 / 30000
[t-SNE] Computed conditional probabilities for sample 29000 / 30000
[t-SNE] Computed conditional probabilities for sample 30000 / 30000
[t-SNE] Mean sigma: 0.087420
[t-SNE] Computed conditional probabilities in 1.157s
[t-SNE] Iteration 50: error = 111.1141052, gradient norm = 0.0000095 (50 iterations in 20.843s)
[t-SNE] Iteration 100: error = 92.1124878, gradient norm = 0.0019646 (50 iterations in 24.411s)
[t-SNE] Iteration 150: error = 85.0801544, gradient norm = 0.0007227 (50 iterations in 21.225s)
[t-SNE] Iteration 200: error = 83.1425705, gradient norm = 0.0004393 (50 iterations in 19.145s)
[t-SNE] Iteration 250: error = 82.1798553, gradient norm = 0.0002817 (50 iterations in 19.086s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 82.179855
[t-SNE] Iteration 300: error = 3.4739718, gradient norm = 0.0008218 (50 iterations in 21.519s)
[t-SNE] Iteration 350: error = 2.8796194, gradient norm = 0.0003712 (50 iterations in 25.331s)
[t-SNE] Iteration 400: error = 2.5150728, gradient norm = 0.0002101 (50 iterations in 24.924s)
[t-SNE] Iteration 450: error = 2.2753220, gradient norm = 0.0001347 (50 iterations in 26.005s)
[t-SNE] Iteration 500: error = 2.1078000, gradient norm = 0.0000941 (50 iterations in 24.315s)
[t-SNE] Iteration 550: error = 1.9846343, gradient norm = 0.0000698 (50 iterations in 25.387s)
[t-SNE] Iteration 600: error = 1.8901992, gradient norm = 0.0000538 (50 iterations in 26.536s)
[t-SNE] Iteration 650: error = 1.8153672, gradient norm = 0.0000431 (50 iterations in 25.098s)
[t-SNE] Iteration 700: error = 1.7544904, gradient norm = 0.0000353 (50 iterations in 26.003s)
[t-SNE] Iteration 750: error = 1.7041844, gradient norm = 0.0000294 (50 iterations in 26.094s)
[t-SNE] Iteration 800: error = 1.6622106, gradient norm = 0.0000257 (50 iterations in 24.552s)
[t-SNE] Iteration 850: error = 1.62666277, gradient norm = 0.0000225 (50 iterations in 25.561s)
[t-SNE] Iteration 900: error = 1.5962178, gradient norm = 0.0000198 (50 iterations in 24.661s)
[t-SNE] Iteration 950: error = 1.5701591, gradient norm = 0.0000179 (50 iterations in 24.392s)
[t-SNE] Iteration 1000: error = 1.5478457, gradient norm = 0.0000165 (50 iterations in 27.312s)
[t-SNE] KL divergence after 1000 iterations: 1.547846
```

```
In [162]: 1 import plotly.graph_objs as go
2 import plotly.tools as tls
3 import plotly.offline as py
4 py.init_notebook_mode(connected=True)
5
6 trace1 = go.Scatter3d(
7     x=tsne3d[:,0],
8     y=tsne3d[:,1],
9     z=tsne3d[:,2],
10    mode='markers',
11    marker=dict(
12        sizemode='diameter',
13        color = y,
14        colorscale = 'Portland',
15        colorbar = dict(title = 'duplicate'),
16        line=dict(color='rgb(255, 255, 255)'),
17        opacity=0.75
18    )
19 )
20
21 data=[trace1]
22 layout=dict(height=800, width=800, title='3d embedding with engineered features')
23 fig=dict(data=data, layout=layout)
24 py.iplot(fig, filename='3DBubble')
```

3d embedding with engineered features



```
In [163]: 1 ques_df = new_df[['question1','question2']]
2 ques_df.head()
```

Out[163]:

|        | question1   | question2   |
|--------|---|---|
| 237030 | how can i stop playing video games                | should i stop playing video games with my child   |
| 247341 | who is better donald trump or hillary clinton     | why is hillary clinton a better choice than do... |
| 246425 | what do you think is the chance that sometime ... | do you think there will be another world war n... |
| 306985 | why are so many questions posted to quora that... | why do people write questions on quora that co... |
| 225863 | can there even be a movie ever rated 10 10 on ... | what are your 10 10 movies                        |

```
In [164]: 1 final_df = new_df.drop(columns=['id','qid1','qid2','question1','question2'])
2 print(final_df.shape)
3 final_df.head()
```

Out[164]: (30000, 23)

Out[164]:

|        | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_share | cwc_min  | cwc_max  | ... ctc_max  | last_word_eq |
|--------|--------------|--------|--------|--------------|--------------|-------------|------------|------------|----------|----------|--------------|--------------|
| 237030 | 0            | 34     | 47     | 7            | 9            | 5           | 16         | 0.31       | 0.999975 | 0.799984 | ... 0.555549 | 0.0          |
| 247341 | 1            | 45     | 56     | 8            | 10           | 6           | 18         | 0.33       | 0.999980 | 0.833319 | ... 0.599994 | 0.0          |
| 246425 | 1            | 137    | 76     | 29           | 15           | 13          | 40         | 0.32       | 0.857131 | 0.499996 | ... 0.464284 | 0.0          |
| 306985 | 1            | 85     | 85     | 16           | 16           | 5           | 30         | 0.17       | 0.374995 | 0.333330 | ... 0.312498 | 0.0          |
| 225863 | 0            | 50     | 26     | 12           | 6            | 1           | 16         | 0.06       | 0.499975 | 0.166664 | ... 0.083333 | 0.0          |

5 rows × 23 columns

```
In [165]: 1 from sklearn.feature_extraction.text import CountVectorizer
2 # merge texts
3 questions = list(ques_df['question1']) + list(ques_df['question2'])
4
5 cv = CountVectorizer(max_features=3000)
6 q1_arr, q2_arr = np.vsplit(cv.fit_transform(questions).toarray(),2)
```

```
In [166]: 1 temp_df1 = pd.DataFrame(q1_arr, index= ques_df.index)
2 temp_df2 = pd.DataFrame(q2_arr, index= ques_df.index)
3 temp_df = pd.concat([temp_df1, temp_df2], axis=1)
4 temp_df.shape
```

Out[166]: (30000, 6000)

```
In [167]: 1 final_df = pd.concat([final_df, temp_df], axis=1)
2 print(final_df.shape)
3 final_df.head()
```

(30000, 6023)

Out[167]:

|        | is_duplicate | q1_len | q2_len | q1_num_words | q2_num_words | word_common | word_total | word_share | cwc_min  | cwc_max  | ... 2990 | 2991 | 2992 | 2993 |
|--------|--------------|--------|--------|--------------|--------------|-------------|------------|------------|----------|----------|----------|------|------|------|
| 237030 | 0            | 34     | 47     | 7            | 9            | 5           | 16         | 0.31       | 0.999975 | 0.799984 | ... 0    | 0    | 0    | 0    |
| 247341 | 1            | 45     | 56     | 8            | 10           | 6           | 18         | 0.33       | 0.999980 | 0.833319 | ... 0    | 0    | 0    | 0    |
| 246425 | 1            | 137    | 76     | 29           | 15           | 13          | 40         | 0.32       | 0.857131 | 0.499996 | ... 0    | 1    | 0    | 0    |
| 306985 | 1            | 85     | 85     | 16           | 16           | 5           | 30         | 0.17       | 0.374995 | 0.333330 | ... 0    | 0    | 0    | 0    |
| 225863 | 0            | 50     | 26     | 12           | 6            | 1           | 16         | 0.06       | 0.499975 | 0.166664 | ... 0    | 0    | 0    | 0    |

5 rows × 6023 columns

```
In [168]: from sklearn.model_selection import train_test_split
train,X_test,y_train,y_test = train_test_split(final_df.iloc[:,1:],final_df.iloc[:,0].values,test_size=0.2,random_state=1)
```

```
In [169]: 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.metrics import accuracy_score
3 rf = RandomForestClassifier()
4 rf.fit(X_train,y_train)
5 y_pred = rf.predict(X_test)
6 accuracy_score(y_test,y_pred)
```

Out[169]: 0.785

```
In [218]: 1 # rf = RandomForestClassifier()
2 # rf.fit(X_train,y_train)
3 y_pred = rf.predict(X_test)
4 accuracy_score(y_train[:6000],y_pred)
```

Out[218]: 0.5371666666666667

In [ ]:

1

```
In [212]: 1 print(classification_report(y_test, y_pred))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.82      | 0.85   | 0.83     | 3810    |
| 1            | 0.72      | 0.68   | 0.70     | 2190    |
| accuracy     |           |        | 0.79     | 6000    |
| macro avg    | 0.77      | 0.76   | 0.77     | 6000    |
| weighted avg | 0.78      | 0.79   | 0.78     | 6000    |

```
In [170]: 1 from xgboost import XGBClassifier
2 xgb = XGBClassifier()
3 xgb.fit(X_train,y_train)
4 y_pred1 = xgb.predict(X_test)
5 accuracy_score(y_test,y_pred1)
```

Out[170]: 0.785

```
In [211]: 1 print(classification_report(y_test, y_pred1))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.83      | 0.83   | 0.83     | 3810    |
| 1            | 0.70      | 0.71   | 0.71     | 2190    |
| accuracy     |           |        | 0.79     | 6000    |
| macro avg    | 0.77      | 0.77   | 0.77     | 6000    |
| weighted avg | 0.79      | 0.79   | 0.79     | 6000    |

```
In [171]: 1 from lightgbm import LGBMClassifier
2 lgb = LGBMClassifier()
3 lgb.fit(X_train,y_train)
4 y_pred2 = lgb.predict(X_test)
5 accuracy_score(y_test,y_pred2)
```

Out[171]: 0.7911666666666667

```
In [210]: 1 print(classification_report(y_test, y_pred2))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.84      | 0.84   | 0.84     | 3810    |
| 1            | 0.71      | 0.71   | 0.71     | 2190    |
| accuracy     |           |        | 0.79     | 6000    |
| macro avg    | 0.77      | 0.77   | 0.77     | 6000    |
| weighted avg | 0.79      | 0.79   | 0.79     | 6000    |

```
In [172]: 1 from sklearn.metrics import confusion_matrix
```

```
In [173]: 1 # for random forest model
2 confusion_matrix(y_test,y_pred)
```

Out[173]: array([[3224, 586],
 [ 704, 1486]], dtype=int64)

```
In [174]: 1 # for xgboost model
2 confusion_matrix(y_test,y_pred1)
```

Out[174]: array([[3162, 648],
 [ 642, 1548]], dtype=int64)

```
In [176]: 1 # for Lightgbm model
2 confusion_matrix(y_test,y_pred2)
```

Out[176]: array([[3182, 628],
 [ 625, 1565]], dtype=int64)

*We will choose the random forest model because in our case we should consider low False positive like it was actual not duplicate question but our model predicted as duplicate. so it will make a bad user experience !*

```
In [ ]: 1

In [181]: 1 def test_common_words(q1,q2):
2     w1 = set(map(lambda word: word.lower().strip(), q1.split(" ")))
3     w2 = set(map(lambda word: word.lower().strip(), q2.split(" ")))
4     return len(w1 & w2)

In [182]: 1 def test_total_words(q1,q2):
2     w1 = set(map(lambda word: word.lower().strip(), q1.split(" ")))
3     w2 = set(map(lambda word: word.lower().strip(), q2.split(" ")))
4     return (len(w1) + len(w2))

In [183]: 1 def test_fetch_token_features(q1,q2):
2
3     SAFE_DIV = 0.0001
4
5     STOP_WORDS = stopwords.words("english")
6
7     token_features = [0.0]*8
8
9     # Converting the Sentence into Tokens:
10    q1_tokens = q1.split()
11    q2_tokens = q2.split()
12
13    if len(q1_tokens) == 0 or len(q2_tokens) == 0:
14        return token_features
15
16    # Get the non-stopwords in Questions
17    q1_words = set([word for word in q1_tokens if word not in STOP_WORDS])
18    q2_words = set([word for word in q2_tokens if word not in STOP_WORDS])
19
20    #Get the stopwords in Questions
21    q1_stops = set([word for word in q1_tokens if word in STOP_WORDS])
22    q2_stops = set([word for word in q2_tokens if word in STOP_WORDS])
23
24    # Get the common non-stopwords from Question pair
25    common_word_count = len(q1_words.intersection(q2_words))
26
27    # Get the common stopwords from Question pair
28    common_stop_count = len(q1_stops.intersection(q2_stops))
29
30    # Get the common Tokens from Question pair
31    common_token_count = len(set(q1_tokens).intersection(set(q2_tokens)))
32
33
34    token_features[0] = common_word_count / (min(len(q1_words), len(q2_words)) + SAFE_DIV)
35    token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)) + SAFE_DIV)
36    token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)) + SAFE_DIV)
37    token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)) + SAFE_DIV)
38    token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
39    token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)) + SAFE_DIV)
40
41    # Last word of both question is same or not
42    token_features[6] = int(q1_tokens[-1] == q2_tokens[-1])
43
44    # First word of both question is same or not
45    token_features[7] = int(q1_tokens[0] == q2_tokens[0])
46
47    return token_features
```

```
In [184]: 1 def test_fetch_length_features(q1,q2):
2
3     length_features = [0.0]*3
4
5     # Converting the Sentence into Tokens:
6     q1_tokens = q1.split()
7     q2_tokens = q2.split()
8
9     if len(q1_tokens) == 0 or len(q2_tokens) == 0:
10        return length_features
11
12    # Absolute length features
13    length_features[0] = abs(len(q1_tokens) - len(q2_tokens))
14
15    #Average Token Length of both Questions
16    length_features[1] = (len(q1_tokens) + len(q2_tokens))/2
17
18    strs = list(distance.lcsubstrings(q1, q2))
19    length_features[2] = len(strs[0]) / (min(len(q1), len(q2)) + 1)
20
21    return length_features
```

```
In [185]: 1 def test_fetch_fuzzy_features(q1,q2):
2
3     fuzzy_features = [0.0]*4
4
5     # fuzz_ratio
6     fuzzy_features[0] = fuzz.QRatio(q1, q2)
7
8     # fuzz_partial_ratio
9     fuzzy_features[1] = fuzz.partial_ratio(q1, q2)
10
11    # token_sort_ratio
12    fuzzy_features[2] = fuzz.token_sort_ratio(q1, q2)
13
14    # token_set_ratio
15    fuzzy_features[3] = fuzz.token_set_ratio(q1, q2)
16
17    return fuzzy_features
```

```
In [186]: 1 def query_point_creator(q1,q2):
2
3     input_query = []
4
5     # preprocess
6     q1 = preprocess(q1)
7     q2 = preprocess(q2)
8
9     # fetch basic features
10    input_query.append(len(q1))
11    input_query.append(len(q2))
12
13    input_query.append(len(q1.split(" ")))
14    input_query.append(len(q2.split(" ")))
15
16    input_query.append(test_common_words(q1,q2))
17    input_query.append(test_total_words(q1,q2))
18    input_query.append(round(test_common_words(q1,q2)/test_total_words(q1,q2),2))
19
20    # fetch token features
21    token_features = test_fetch_token_features(q1,q2)
22    input_query.extend(token_features)
23
24    # fetch length based features
25    length_features = test_fetch_length_features(q1,q2)
26    input_query.extend(length_features)
27
28    # fetch fuzzy features
29    fuzzy_features = test_fetch_fuzzy_features(q1,q2)
30    input_query.extend(fuzzy_features)
31
32    # bow feature for q1
33    q1_bow = cv.transform([q1]).toarray()
34
35    # bow feature for q2
36    q2_bow = cv.transform([q2]).toarray()
37
38
39
40    return np.hstack((np.array(input_query).reshape(1,22),q1_bow,q2_bow))
```

```
In [187]: 1 q1 = 'Where is the capital of India?'
2 q2 = 'What is the current capital of Pakistan?'
3 q3 = 'Which city serves as the capital of India?'
4 q4 = 'What is the business capital of India?'
```

```
In [188]: 1 rf.predict(query_point_creator(q1,q4))
2
```

```
Out[188]: array([1], dtype=int64)
```

```
In [189]: 1 cv
```

```
Out[189]: CountVectorizer(max_features=3000)
```

```
In [191]: 1 import pickle
2
3 pickle.dump(rf,open('quora_duplicate_question_model.pkl','wb'))
4 pickle.dump(cv,open('cv.pkl','wb'))
```

## deployment

```
In [205]: 1 # !pip install streamlit
2 # !pip install helper
3 # !pip install grp
4 # !pip uninstall celery
5 # !pip install celery==5.0.5
```

```
In [207]: 1 # import streamlit as st
2 # import helper
3
4 # import pickle
5
6 # model = pickle.load(open('model.pkl','rb'))
7
8 # st.header('Duplicate Question Pairs')
9
10 # q1 = st.text_input('Enter question 1')
11 # q2 = st.text_input('Enter question 2')
12
13 # if st.button('Find'):
14 #     query = helper.query_point_creator(q1,q2)
15 #     result = model.predict(query)[0]
16
17 #     if result:
18 #         st.header('Duplicate')
19 #     else:
20 #         st.header('Not Duplicate')
```

```
In [ ]: 1 q1
```

```
In [3]: 1 # pip install nbconvert
```

```
In [7]: 1 # !pip install sudo apt-get install pandoc
```

```
In [ ]: 1
```