

Project_2_Quara Question Pairs Similarity

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

Dataset Description

The goal of this competition is to predict which of the provided pairs of questions contain two questions with the same meaning. The ground truth is the set of labels that have been supplied by human experts. The ground truth labels are inherently subjective, as the true meaning of sentences can never be known with certainty. Human labeling is also a 'noisy' process, and reasonable people will disagree. As a result, the ground truth labels on this dataset should be taken to be 'informed' but not 100% accurate, and may include incorrect labeling. We believe the labels, on the whole, to represent a reasonable consensus, but this may often not be true on a case by case basis for individual items in the dataset.

Please note:

as an anti-cheating measure, Kaggle has supplemented the test set with computer-generated question pairs. Those rows do not come from Quora, and are not counted in the scoring. All of the questions in the training set are genuine examples from Quora.

Data fields

id - the id of a training set question pair qid1, qid2 - unique ids of each question (only available in train.csv) question1, question2 - the full text of each question is_duplicate - the target variable, set to 1 if question1 and question2 have essentially the same meaning, and 0 otherwise.

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import nltk
4 from nltk.sentiment.vader import SentimentIntensityAnalyzer
5 import re
6 from textblob import TextBlob
7 from wordcloud import WordCloud
8 import seaborn as sns
9 import matplotlib.pyplot as plt
10 import cufflinks as cf
11 %matplotlib inline
12 from plotly.offline import init_notebook_mode, iplot
13 init_notebook_mode(connected = True)
14 cf.go_offline();
15 import plotly.graph_objs as go
16 from plotly.subplots import make_subplots
17
18
19 pd.set_option('display.max_columns', None)
20
21
22 from nltk.corpus import stopwords
23 from nltk.stem import SnowballStemmer
24 from sklearn.feature_extraction.text import CountVectorizer
25
26 from collections import Counter
27 from numpy import where
28
29
30 from sklearn.decomposition import PCA
31
32 from sklearn.preprocessing import OneHotEncoder
33 from sklearn.preprocessing import StandardScaler
34
35 from sklearn.model_selection import train_test_split
36 from sklearn.linear_model import LogisticRegression
37 from sklearn.metrics import accuracy_score
38 from sklearn.metrics import confusion_matrix, ConfusionMatrixDisplay
39 from sklearn import metrics
40
41 from scipy.sparse import hstack, vstack
42
43 from prettytable import PrettyTable
44 from scipy.stats import loguniform # Log-uniform is useful for searching
45 from sklearn.model_selection import RepeatedStratifiedKFold, RandomizedSe
46
47 import warnings
48 warnings.filterwarnings('ignore')
```

C:\ProgramData\anaconda3\Lib\site-packages\paramiko\transport.py:219: CryptographyDeprecationWarning:

Blowfish has been deprecated

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

Out[2]: (404290, 6)

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

Out[3]:

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

In [4]: 1 df.sample(10)

Out[4]:

	id	qid1	qid2	question1	question2	is_duplicate
368038	368038	498371	498372	How realistic are Tyler Durden's plans in "Fig...	[SPOILER ALERT] At the end of Fight Club, why ...	0
16732	16732	31869	31870	If goods are being transported from one state ...	If I want to sell goods from one state to anot...	0
82545	82545	139936	98268	What are some things new employees should know...	What are some things new employees should know...	0
253486	253486	368039	368040	Where does GoodRec.com get its data from?	Where does www.magnetic.com get its data from?	0
225230	225230	333516	333517	What was your latest discovery?	What is the latest discovery in finance?	0
357058	357058	9642	154412	How can you make a friend?	How do I make more friends? How do I?	1
11808	11808	22786	22787	What makes you proud to be Canadian?	How do I make perfect CV for job banks for Can...	0
300133	300133	361686	422891	Has India been succesful to divert world atten...	Why is India not conducting a surgical strike ...	0
88130	88130	148290	148291	How did Aditi Saini, once a careless girl (as ...	What are the perks of being an IES officer?	0
233509	233509	343774	343775	What is the use of shell scripting?	Where can I get Shell Script programs?	0

In [5]: 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 [6]: 1 df = df.sample(30000,random_state=2)

```
In [7]: 1 # missing values
        2 df.isnull().sum()
```

```
Out[7]: id          0
        qid1         0
        qid2         0
        question1    0
        question2    0
        is_duplicate  0
        dtype: int64
```

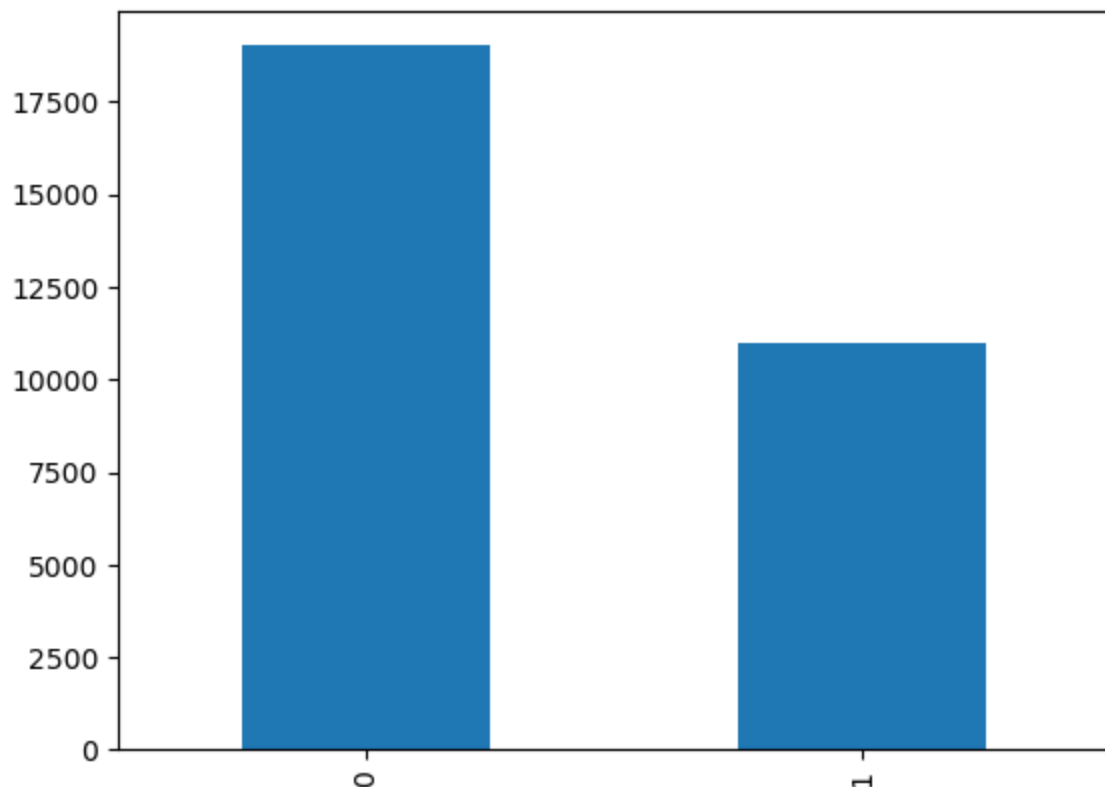
```
In [8]: 1 # duplicate rows
        2 df.duplicated().sum()
```

```
Out[8]: 0
```

```
In [9]: 1 # Diistribution of duplicate and non-duplicate questions
        2 print(df['is_duplicate'].value_counts())
        3 print(df['is_duplicate'].value_counts()/df['is_duplicate'].count()*100)
        4 df['is_duplicate'].value_counts().plot(kind = 'bar')
```

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

```
Out[9]: <Axes: >
```

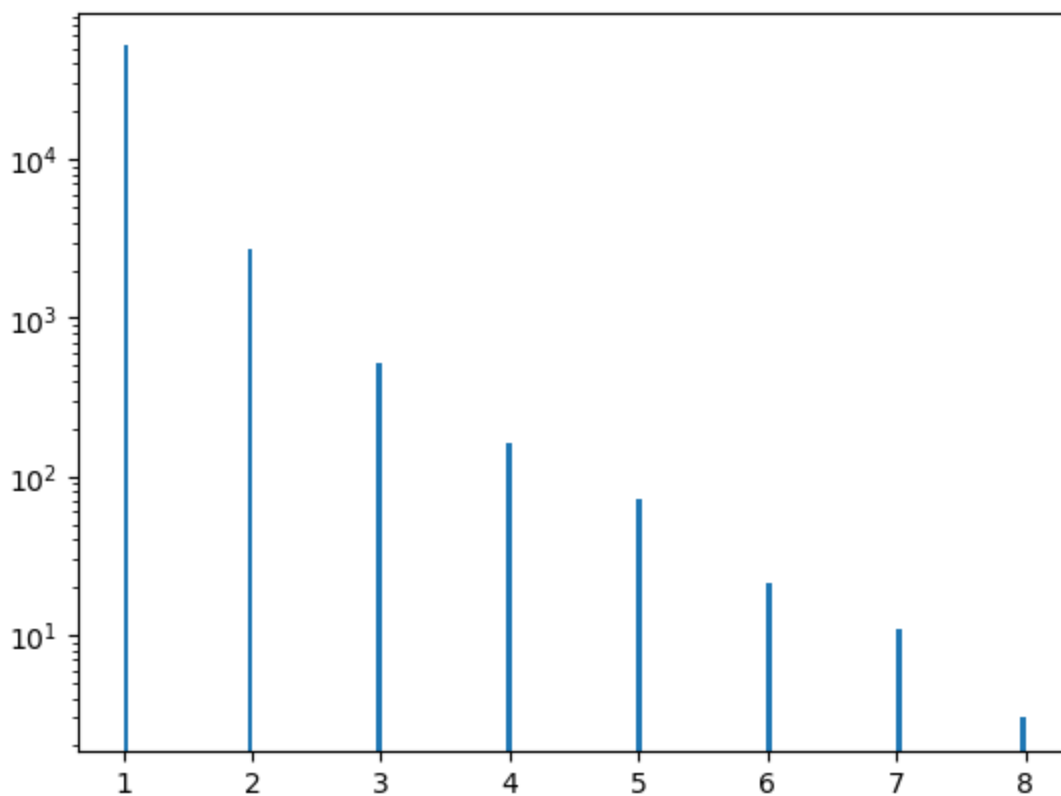


```
In [10]: 1 # Repeated questions
2
3 qid = pd.Series(df['qid1'].tolist() + df['qid2'].tolist())
4 print('Number of unique questions',np.unique(qid).shape[0])
5 x = qid.value_counts()>1
6 print('Number of questions getting repeated',x[x].shape[0])
```

Number of unique questions 55299

Number of questions getting repeated 3480

```
In [11]: 1 # Repeated questions histogram
2
3 plt.hist(qid.value_counts().values,bins=160)
4 plt.yscale('log')
5 plt.show()
```



```
In [12]: 1 # Feature Engineering
2
3 df['q1_len'] = df['question1'].str.len()
4 df['q2_len'] = df['question2'].str.len()
```

In [13]:

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

Out[13]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len
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 doesnt 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 [14]: 1 import pandas as pd
2
3 # Assuming 'df' is your DataFrame
4 df['q1_num_words'] = df['question1'].apply(lambda row: len(str(row).split())
5 df['q2_num_words'] = df['question2'].apply(lambda row: len(str(row).split())
6 df.head()
7
```

Out[14]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_num
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 doesnt 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 [15]: 1 def common_words(row):
2         w1 = set(map(lambda word: word.lower().strip(), str(row['question1'])))
3         w2 = set(map(lambda word: word.lower().strip(), str(row['question2'])))
4         return len(w1 & w2)
5
```



```
In [16]: 1 df['word_common'] = df.apply(common_words, axis=1)
          2 df.head()
```

Out[16]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nu
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 doesnt 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 [17]: 1 def total_words(row):
          2     w1 = set(map(lambda word: word.lower().strip(), str(row['question1'])))
          3     w2 = set(map(lambda word: word.lower().strip(), str(row['question2'])))
          4     return (len(w1) + len(w2))
```

```
In [18]: 1 df['word_total'] = df.apply(total_words, axis=1)
         2 df.head()
```

Out[18]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nu
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 doesnt 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	

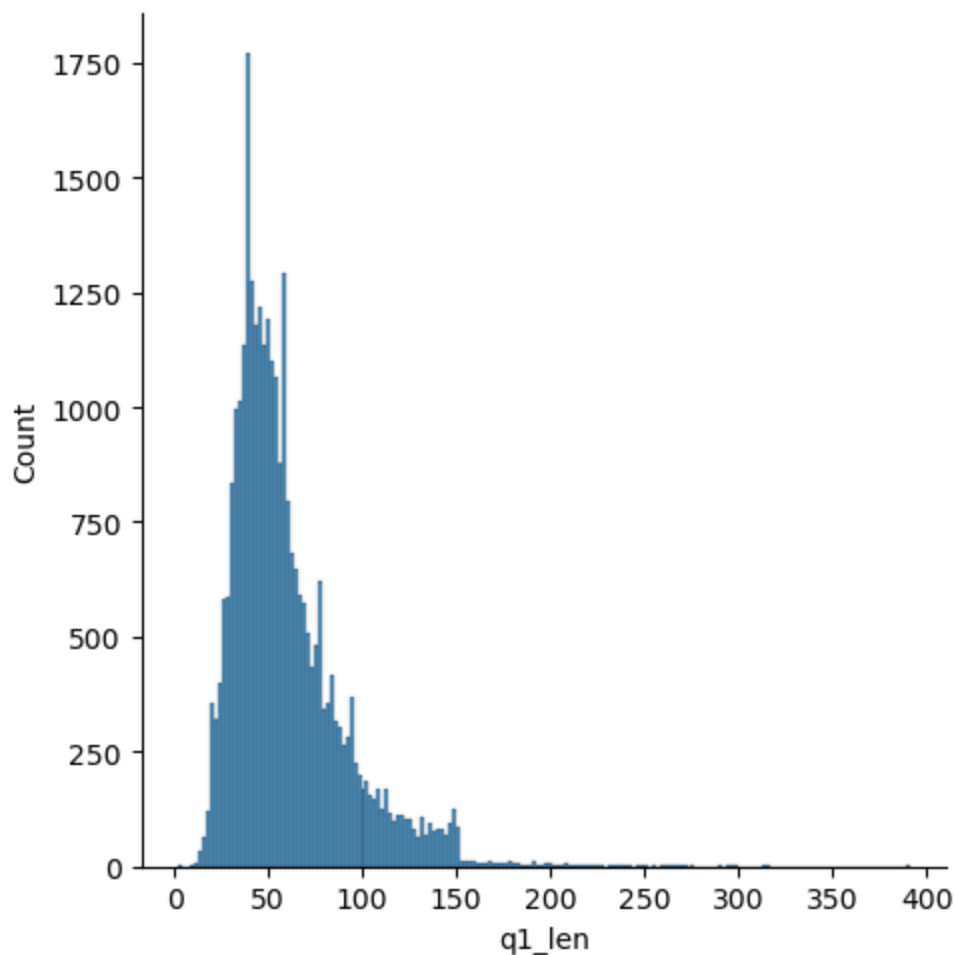
```
In [19]: 1 df['word_share'] = round(df['word_common']/df['word_total'],2)
        2 df.head()
```

Out[19]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nu
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 doesnt 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	

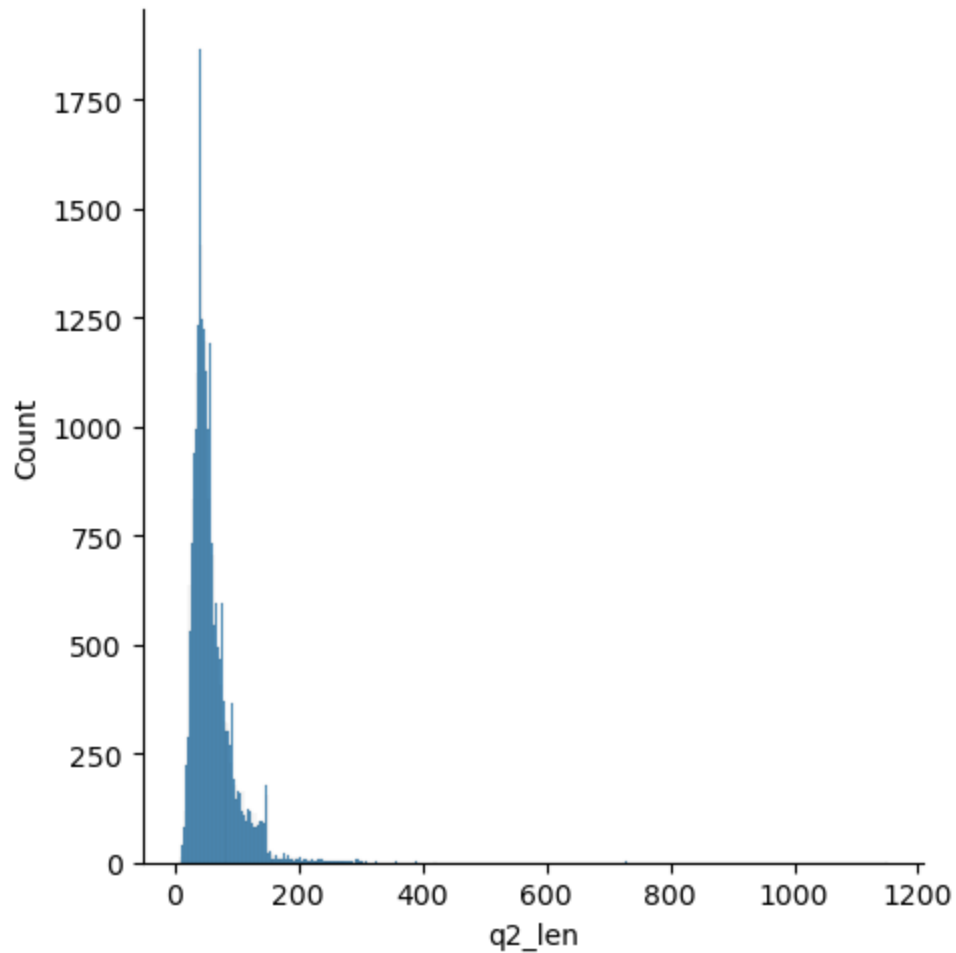
```
In [20]: 1 # Analysis of features
2 sns.displot(df['q1_len'])
3 print('minimum characters',df['q1_len'].min())
4 print('maximum characters',df['q1_len'].max())
5 print('average num of characters',int(df['q1_len'].mean()))
```

minimum characters 2
maximum characters 391
average num of characters 59



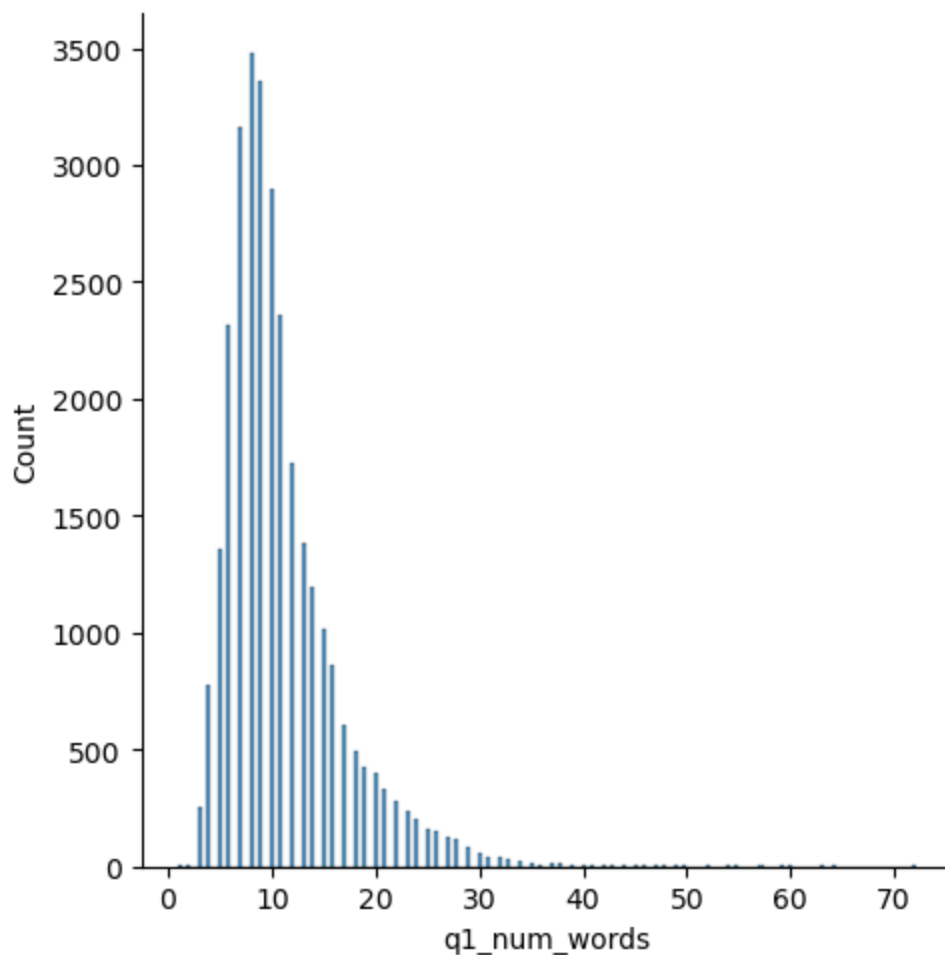
```
In [21]: 1 sns.displot(df['q2_len'])
2 print('minimum characters',df['q2_len'].min())
3 print('maximum characters',df['q2_len'].max())
4 print('average num of characters',int(df['q2_len'].mean()))
```

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



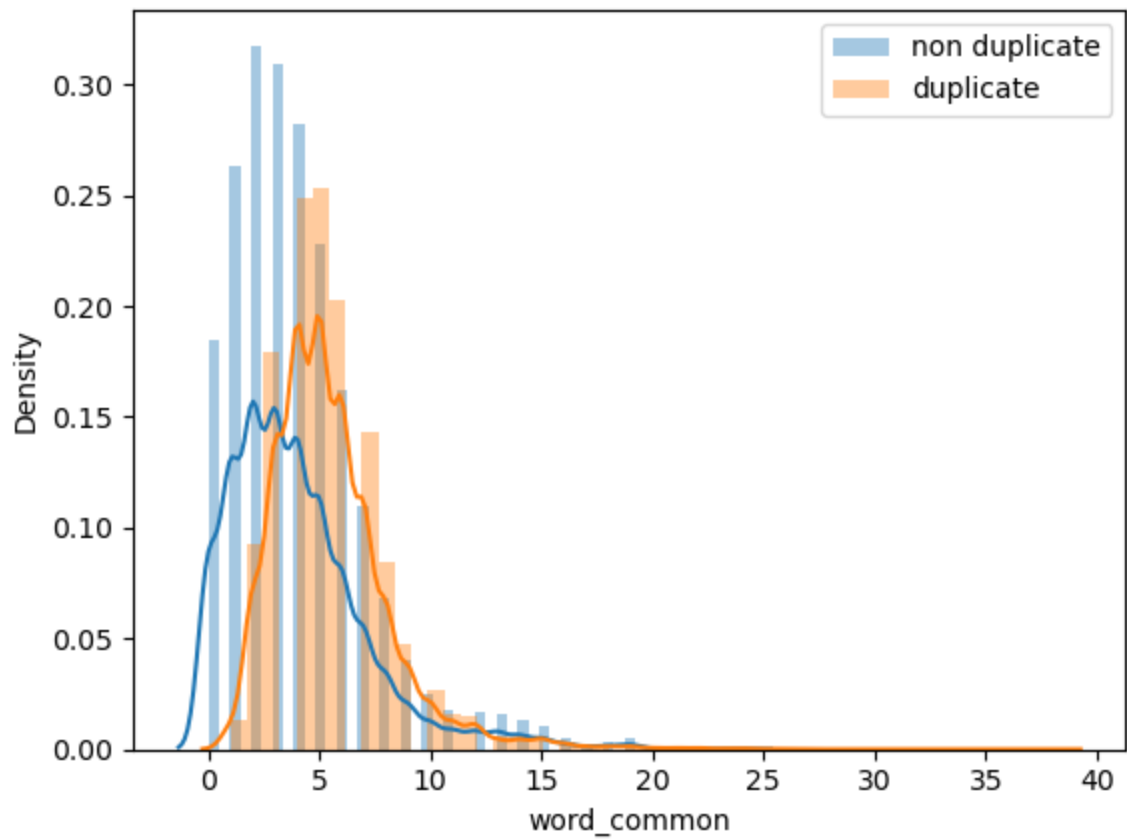
```
In [22]: 1 sns.displot(df['q1_num_words'])  
2 print('minimum words',df['q1_num_words'].min())  
3 print('maximum words',df['q1_num_words'].max())  
4 print('average num of words',int(df['q1_num_words'].mean()))
```

```
minimum words 1  
maximum words 72  
average num of words 10
```

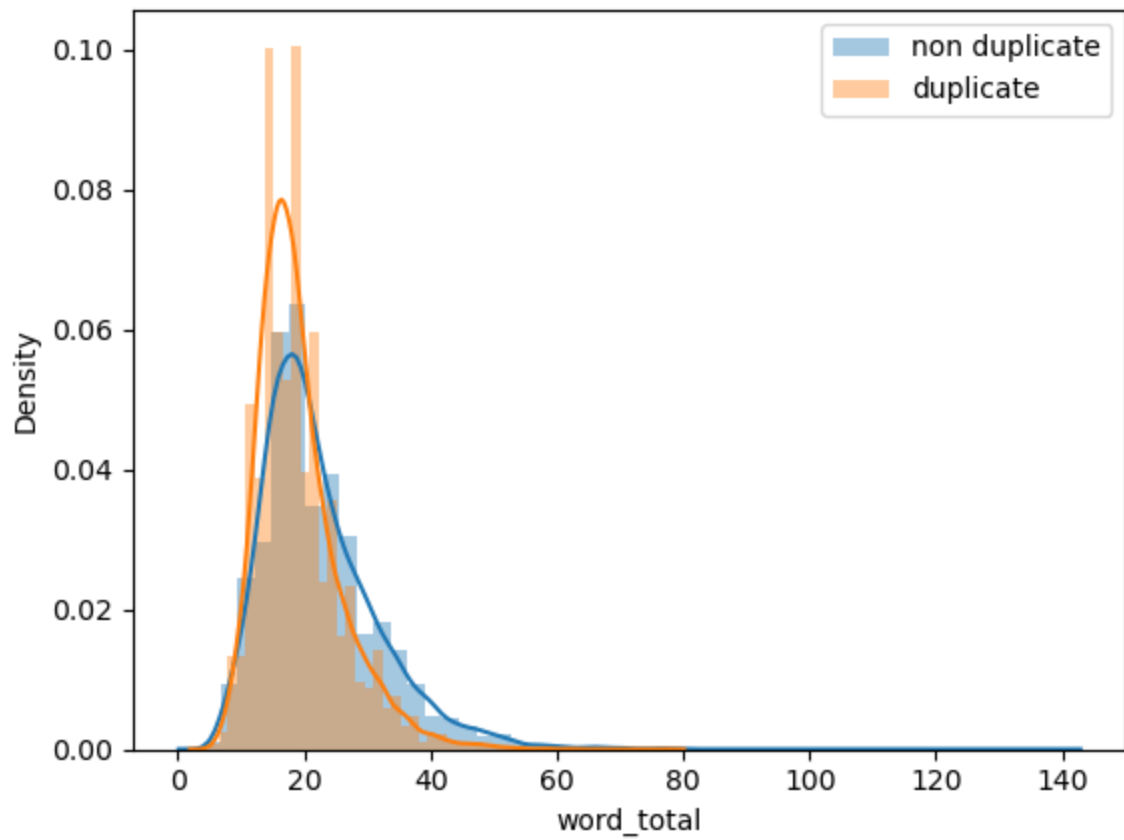


In [23]:

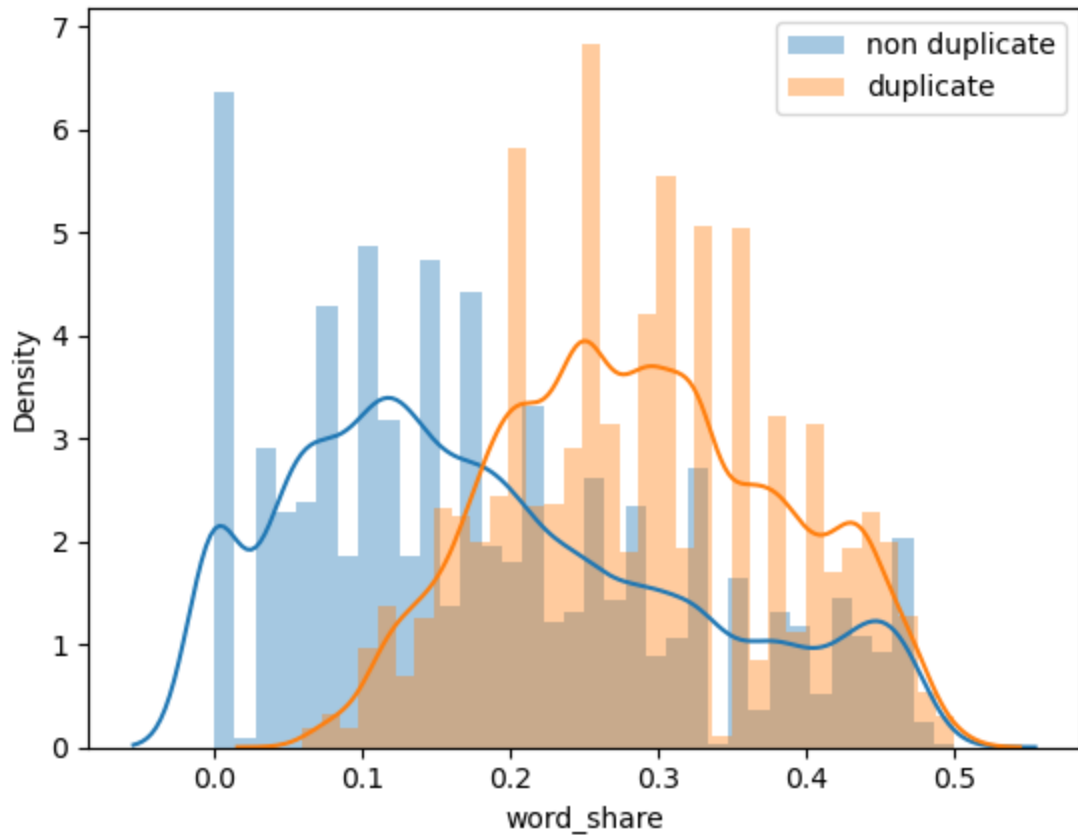
```
1 # common words
2 sns.distplot(df[df['is_duplicate'] == 0]['word_common'],label='non duplic
3 sns.distplot(df[df['is_duplicate'] == 1]['word_common'],label='duplicate'
4 plt.legend()
5 plt.show()
```



```
In [24]: 1 # total words
2 sns.distplot(df[df['is_duplicate'] == 0]['word_total'],label='non duplica
3 sns.distplot(df[df['is_duplicate'] == 1]['word_total'],label='duplicate')
4 plt.legend()
5 plt.show()
```




```
In [25]: 1 # word share
2 sns.distplot(df[df['is_duplicate'] == 0]['word_share'],label='non duplica
3 sns.distplot(df[df['is_duplicate'] == 1]['word_share'],label='duplicate')
4 plt.legend()
5 plt.show()
```



```
In [26]: 1 ques_df = df[['question1','question2']]
2 ques_df.head()
```

Out[26]:

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

```
In [27]: 1 final_df = df.drop(columns=['id', 'qid1', 'qid2', 'question1', 'question2'])
2 print(final_df.shape)
3 final_df.head()
```

(30000, 8)

Out[27]:

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

```
In [28]: 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 [29]: 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[29]: (30000, 6000)

```
In [30]: 1 final_df = pd.concat([final_df, temp_df], axis=1)
2 print(final_df.shape)
3 final_df.head()
```

(30000, 6008)

Out[30]:

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

```
In [31]: 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:].valu
3
```

RandomForestClassifier

```
In [33]: 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[33]: 0.7716666666666666

XGBClassifier

```
In [34]: 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[34]: 0.7661666666666667

Advanced Features

```
In [35]: 1 import re
2 from bs4 import BeautifulSoup
```


In [36]:

```

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
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]+)000000000', 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_contracti
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 [37]: 1 preprocess("I've already! wasn't <b>done</b>?")
```

```
Out[37]: 'i have already was not done'
```

```
In [38]: 1 df['question1'] = df['question1'].apply(preprocess)
2 df['question2'] = df['question2'].apply(preprocess)
```

```
In [39]: 1 df.head()
```

```
Out[39]:
```

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
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 doesnt 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 tragedy	0	35	50	

```
In [40]: 1 df['q1_len'] = df['question1'].str.len()
2 df['q2_len'] = df['question2'].str.len()
```



```
In [41]: 1 df['q1_num_words'] = df['question1'].apply(lambda row: len(row.split(" "))
2 df['q2_num_words'] = df['question2'].apply(lambda row: len(row.split(" "))
3 df.head()
```

Out[41]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

```
In [42]: 1 def common_words(row):
2 w1 = set(map(lambda word: word.lower().strip(), row['question1'].spli
3 w2 = set(map(lambda word: word.lower().strip(), row['question2'].spli
4 return len(w1 & w2)
```

```
In [43]: 1 df['word_common'] = df.apply(common_words, axis=1)
          2 df.head()
```

Out[43]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	



```
In [44]: 1 def total_words(row):
          2     w1 = set(map(lambda word: word.lower().strip(), row['question1'].spli
          3     w2 = set(map(lambda word: word.lower().strip(), row['question2'].spli
          4     return (len(w1) + len(w2))
```

```
In [45]: 1 df['word_total'] = df.apply(total_words, axis=1)
         2 df.head()
```

Out[45]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

```
In [46]: 1 df['word_share'] = round(df['word_common']/df['word_total'],2)
        2 df.head()
```

Out[46]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

```
In [47]: 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)))
41     token_features[1] = common_word_count / (max(len(q1_words), len(q2_words)))
42     token_features[2] = common_stop_count / (min(len(q1_stops), len(q2_stops)))
43     token_features[3] = common_stop_count / (max(len(q1_stops), len(q2_stops)))
44     token_features[4] = common_token_count / (min(len(q1_tokens), len(q2_tokens)))
45     token_features[5] = common_token_count / (max(len(q1_tokens), len(q2_tokens)))
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 [48]: 1 token_features = df.apply(fetch_token_features, axis=1)
2
3 df["cwc_min"] = list(map(lambda x: x[0], token_features))
4 df["cwc_max"] = list(map(lambda x: x[1], token_features))
5 df["csc_min"] = list(map(lambda x: x[2], token_features))
6 df["csc_max"] = list(map(lambda x: x[3], token_features))
7 df["ctc_min"] = list(map(lambda x: x[4], token_features))
8 df["ctc_max"] = list(map(lambda x: x[5], token_features))
9 df["last_word_eq"] = list(map(lambda x: x[6], token_features))
10 df["first_word_eq"] = list(map(lambda x: x[7], token_features))
```

```
In [49]: 1 df.head()
```

```
Out[49]:
```

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

```
In [50]: 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.lcs substrings(q1, q2))
24    length_features[2] = len(strs[0]) / (min(len(q1), len(q2)) + 1)
25
26    return length_features
```

```
In [51]: 1 # Apply the fetch_length_features function to each row
2 length_features = df.apply(fetch_length_features, axis=1)
3
4 # Extract features from the length_features list and assign them to new c
5 df['abs_len_diff'] = list(map(lambda x: x[0], length_features))
6 df['mean_len'] = list(map(lambda x: x[1], length_features))
7 df['longest_substr_ratio'] = list(map(lambda x: x[2], length_features))
8
```

In [52]:

1df.head()

Out[52]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

In [53]:

```
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 [54]:

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

In [55]:

```
1 print(df.shape)
2 df.head()
```

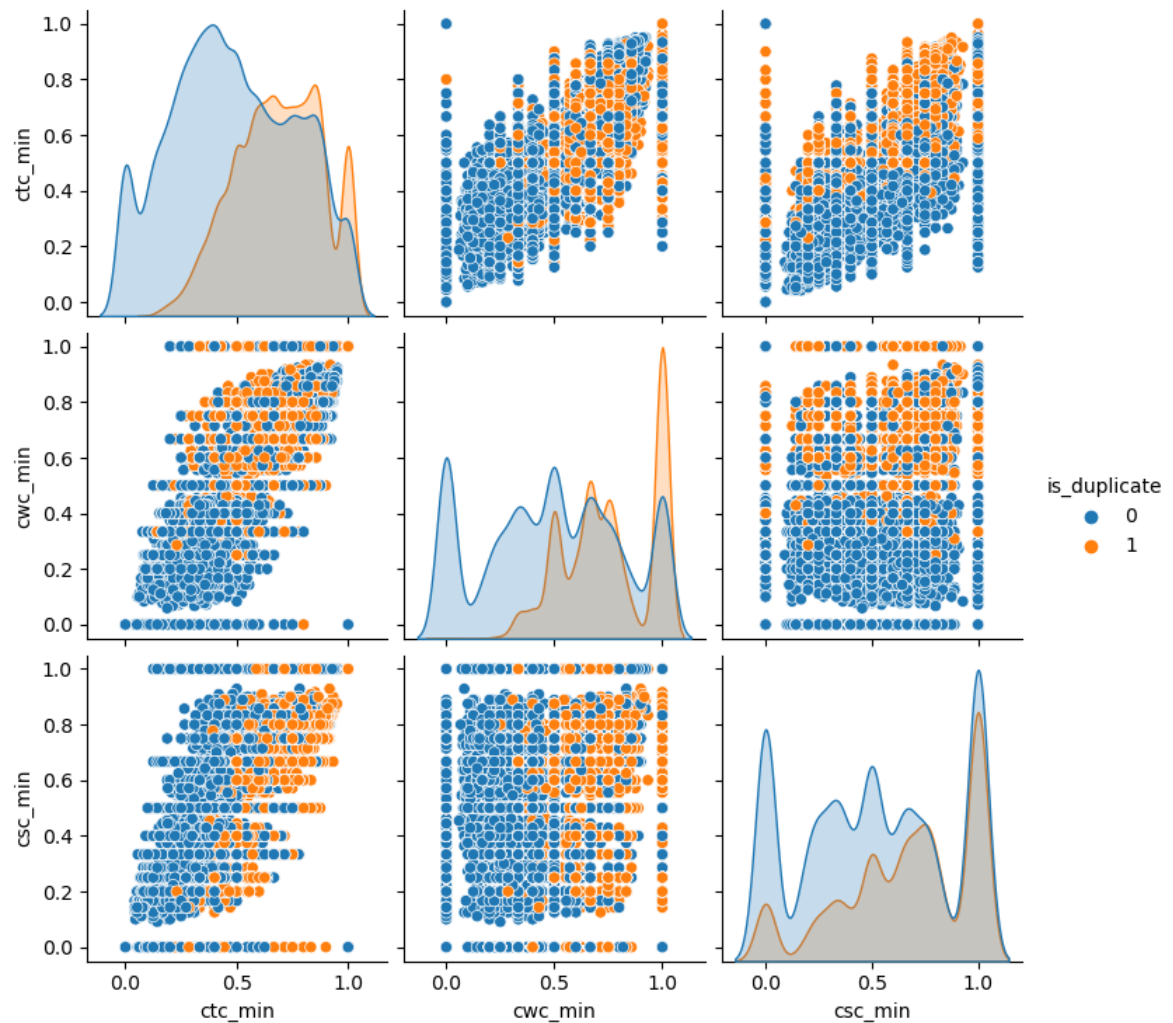
(30000, 28)

Out[55]:

	id	qid1	qid2	question1	question2	is_duplicate	q1_len	q2_len	q1_nur
398782	398782	496695	532029	what is the best marketing automation tool for...	what is the best marketing automation tool for...	1	75	76	
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	48	56	
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	104	119	
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	58	145	
151235	151235	237843	50930	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy	0	34	49	

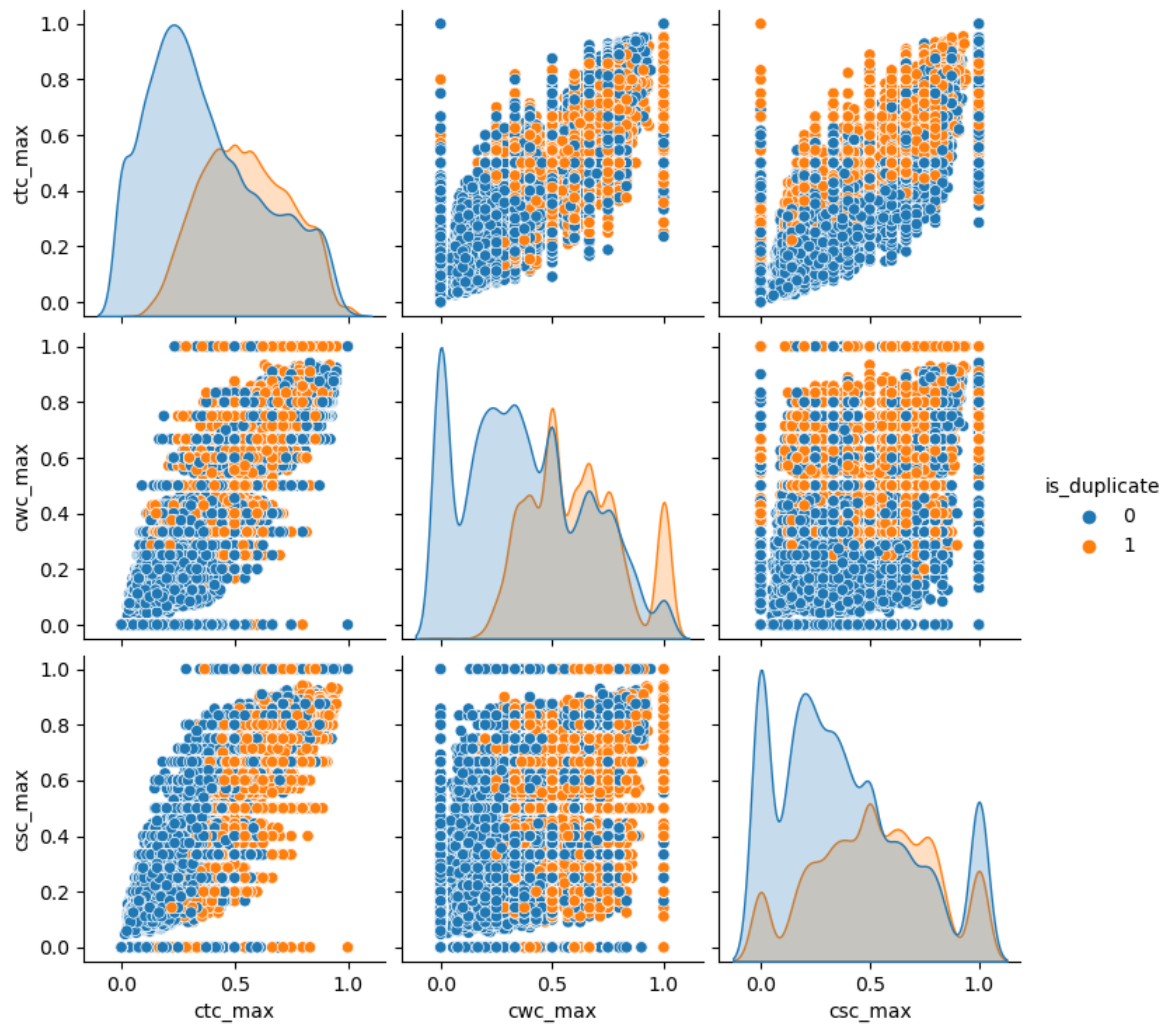
```
In [56]: 1 sns.pairplot(df[['ctc_min', 'cwc_min', 'csc_min', 'is_duplicate']],hue='i
```

```
Out[56]: <seaborn.axisgrid.PairGrid at 0x2fb06ddbb10>
```



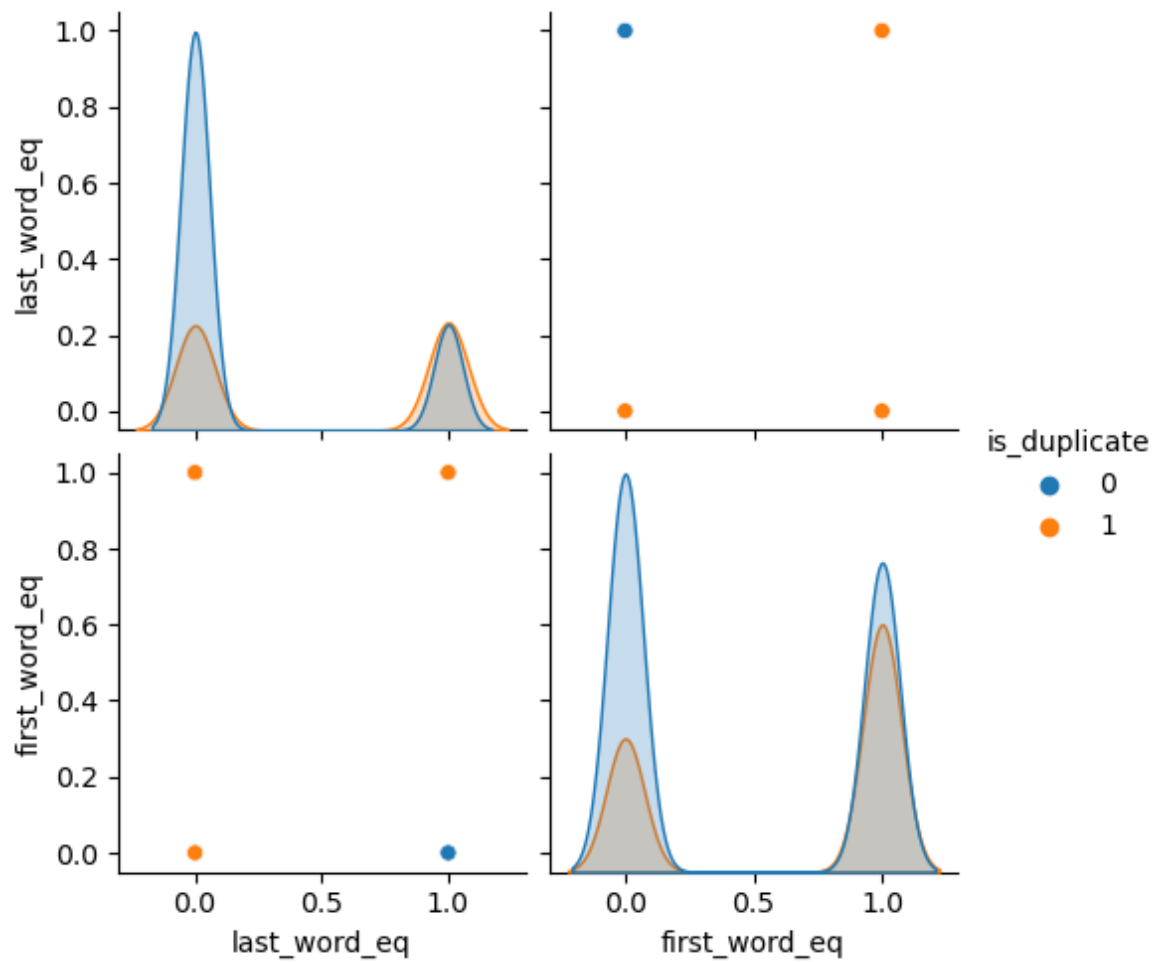
```
In [57]: 1 sns.pairplot(df[['ctc_max', 'cwc_max', 'csc_max', 'is_duplicate']],hue='i
```

```
Out[57]: <seaborn.axisgrid.PairGrid at 0x2fc584de490>
```



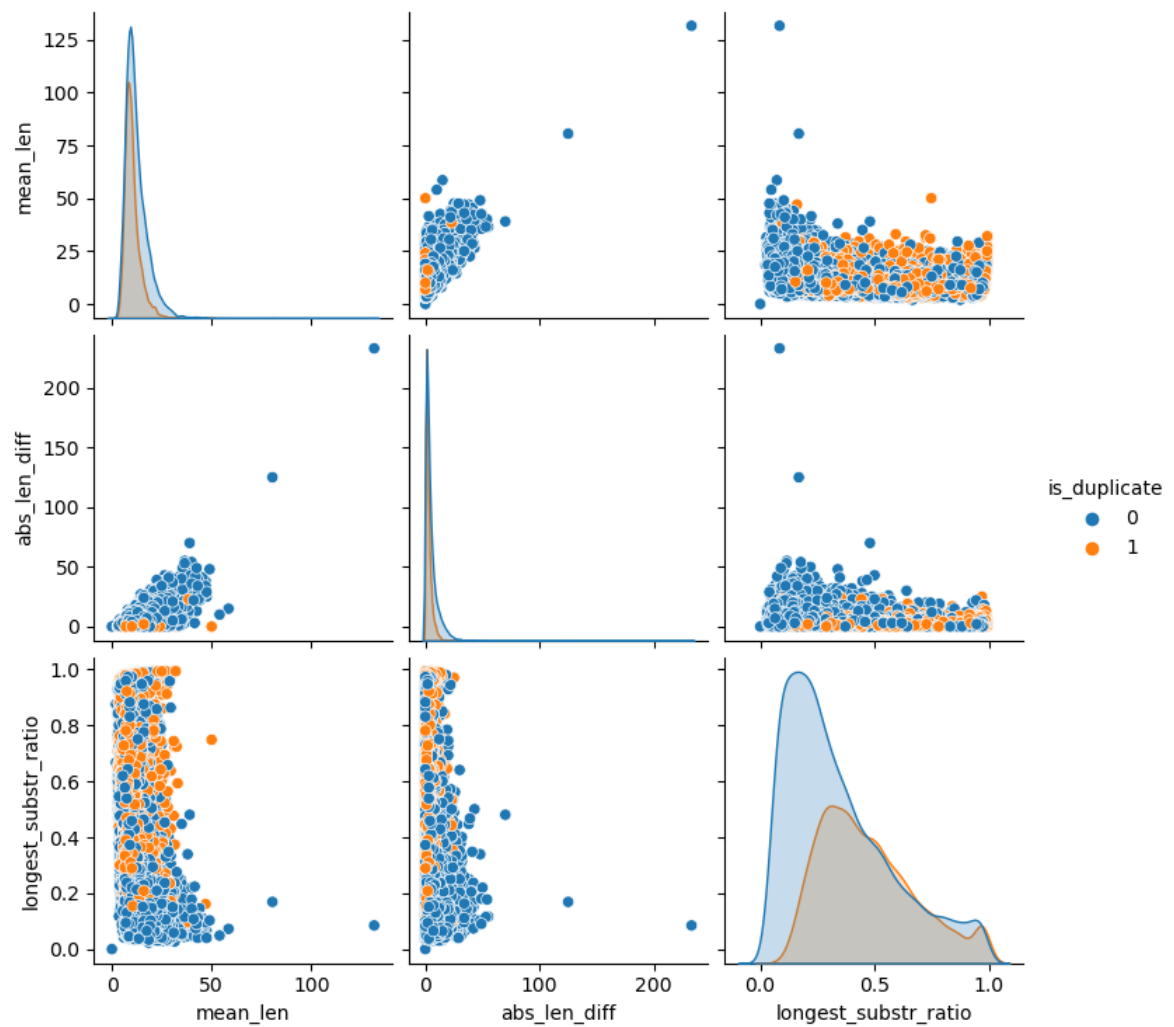
```
In [58]: 1 sns.pairplot(df[['last_word_eq', 'first_word_eq', 'is_duplicate']],hue='i
```

```
Out[58]: <seaborn.axisgrid.PairGrid at 0x2fb06e1a0d0>
```



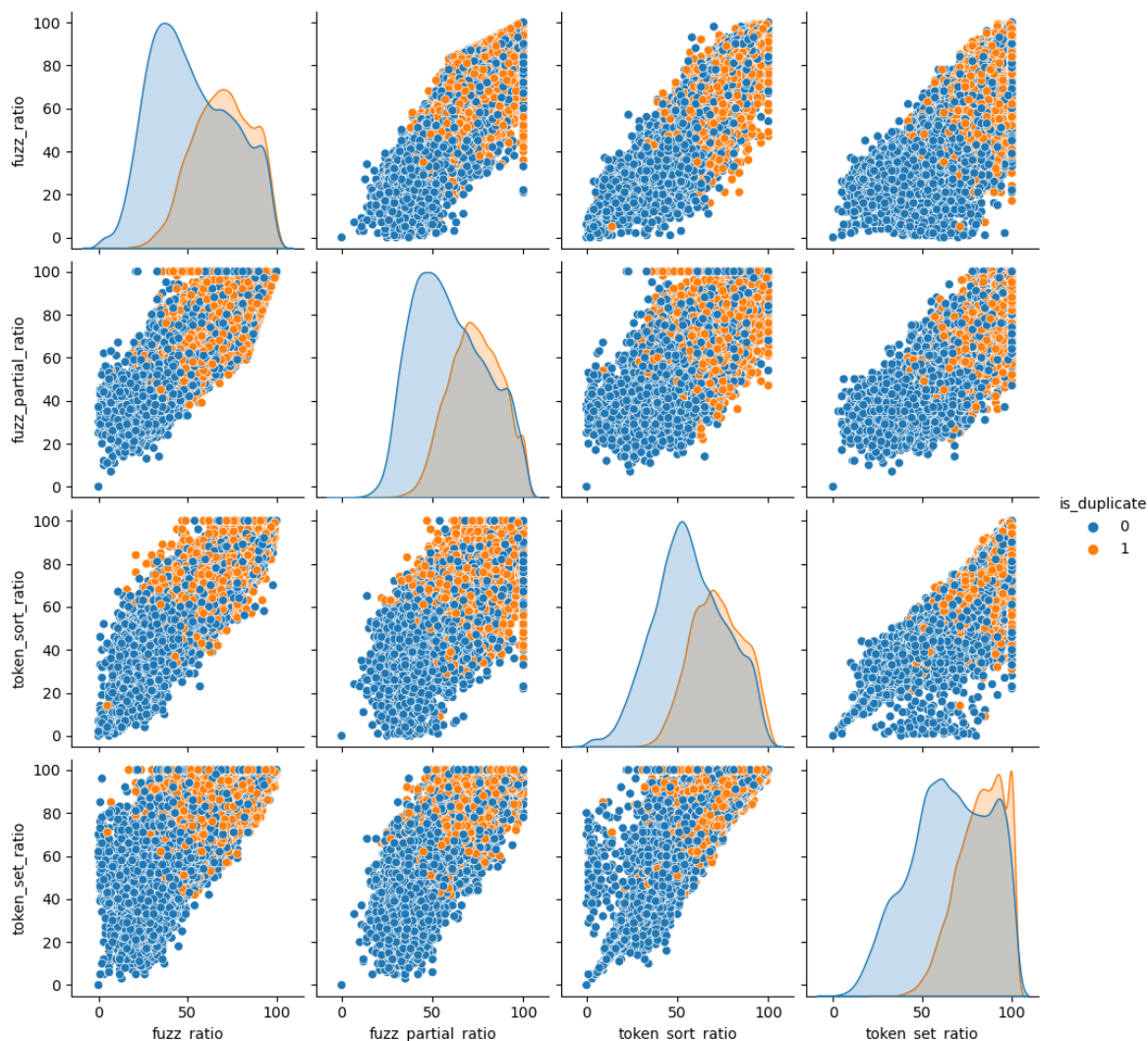
```
In [59]: 1 sns.pairplot(df[['mean_len', 'abs_len_diff', 'longest_substr_ratio', 'is_d
```

```
Out[59]: <seaborn.axisgrid.PairGrid at 0x2fb06606490>
```



In [60]: 1 sns.pairplot(df[['fuzz_ratio', 'fuzz_partial_ratio', 'token_sort_ratio', 't

Out[60]: <seaborn.axisgrid.PairGrid at 0x2fc5871bd10>



In [61]: 1 # Using TSNE for Dimentionality reduction for 15 Features(Generated after
2
3 from sklearn.preprocessing import MinMaxScaler
4
5 X = MinMaxScaler().fit_transform(df[['cwc_min', 'cwc_max', 'csc_min', 'cs
6 y = df['is_duplicate'].values

```
In [62]: 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)
```

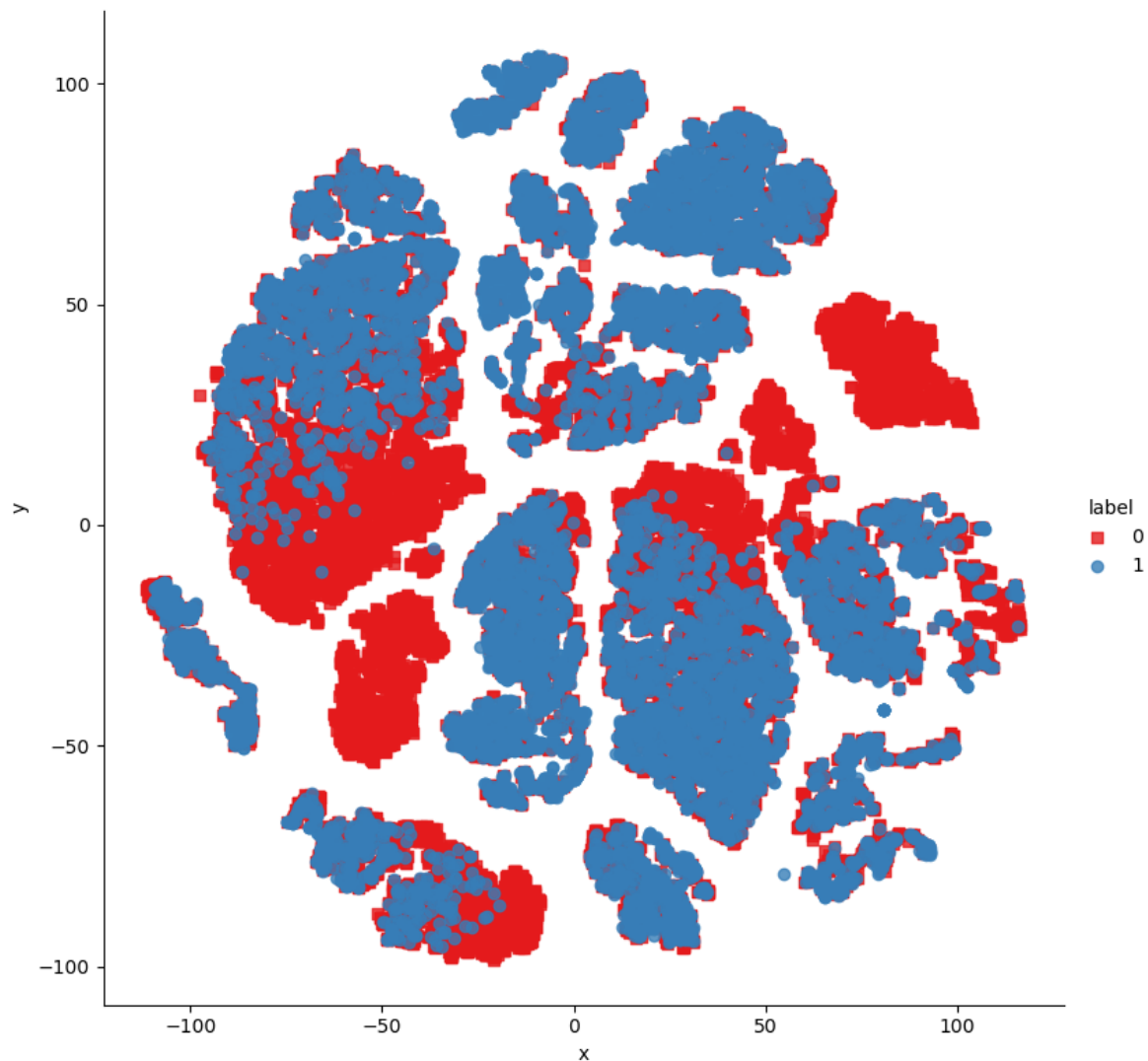


```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 30000 samples in 0.068s...
[t-SNE] Computed neighbors for 30000 samples in 5.413s...
[t-SNE] Computed conditional probabilities for sample 1000 / 30000
[t-SNE] Computed conditional probabilities for sample 2000 / 30000
[t-SNE] Computed conditional 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: 0.080413
[t-SNE] Computed conditional probabilities in 0.941s
[t-SNE] Iteration 50: error = 110.0354767, gradient norm = 0.0222168 (50 iterations in 15.025s)
[t-SNE] Iteration 100: error = 89.8034515, gradient norm = 0.0077262 (50 iterations in 11.897s)
[t-SNE] Iteration 150: error = 85.8413696, gradient norm = 0.0047961 (50 iterations in 12.044s)
[t-SNE] Iteration 200: error = 83.9956284, gradient norm = 0.0039413 (50 iterations in 10.530s)
[t-SNE] Iteration 250: error = 82.8466949, gradient norm = 0.0029376 (50 iterations in 13.284s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 82.846695
[t-SNE] Iteration 300: error = 3.3983612, gradient norm = 0.0073742 (50 iterations in 9.378s)
[t-SNE] Iteration 350: error = 2.7810087, gradient norm = 0.0073581 (50 iterations in 9.022s)
[t-SNE] Iteration 400: error = 2.4468734, gradient norm = 0.0068702 (50 iterations in 7.910s)
[t-SNE] Iteration 450: error = 2.2388027, gradient norm = 0.0064638 (50 iterations in 7.764s)
[t-SNE] Iteration 500: error = 2.0952334, gradient norm = 0.0061122 (50 iterations in 7.882s)
```

```
[t-SNE] Iteration 550: error = 1.9906163, gradient norm = 0.0057656 (50 iterations in 9.103s)
[t-SNE] Iteration 600: error = 1.9118668, gradient norm = 0.0053742 (50 iterations in 8.223s)
[t-SNE] Iteration 650: error = 1.8510072, gradient norm = 0.0050041 (50 iterations in 7.879s)
[t-SNE] Iteration 700: error = 1.8027258, gradient norm = 0.0046457 (50 iterations in 7.946s)
[t-SNE] Iteration 750: error = 1.7636205, gradient norm = 0.0043452 (50 iterations in 8.119s)
[t-SNE] Iteration 800: error = 1.7313575, gradient norm = 0.0040350 (50 iterations in 8.592s)
[t-SNE] Iteration 850: error = 1.7042049, gradient norm = 0.0037814 (50 iterations in 8.236s)
[t-SNE] Iteration 900: error = 1.6811382, gradient norm = 0.0035771 (50 iterations in 8.122s)
[t-SNE] Iteration 950: error = 1.6614714, gradient norm = 0.0033116 (50 iterations in 8.398s)
[t-SNE] Iteration 1000: error = 1.6444725, gradient norm = 0.0031130 (50 iterations in 10.454s)
[t-SNE] KL divergence after 1000 iterations: 1.644472
```

```
In [63]: 1 x_df = pd.DataFrame({'x': tsne2d[:, 0], 'y': tsne2d[:, 1], 'label': y})  
2  
3 # draw the plot in the appropriate place in the grid  
4 sns.lmplot(data=x_df, x='x', y='y', hue='label', fit_reg=False, height=8,  
5
```

Out[63]: <seaborn.axisgrid.FacetGrid at 0x2fc5cc3a810>



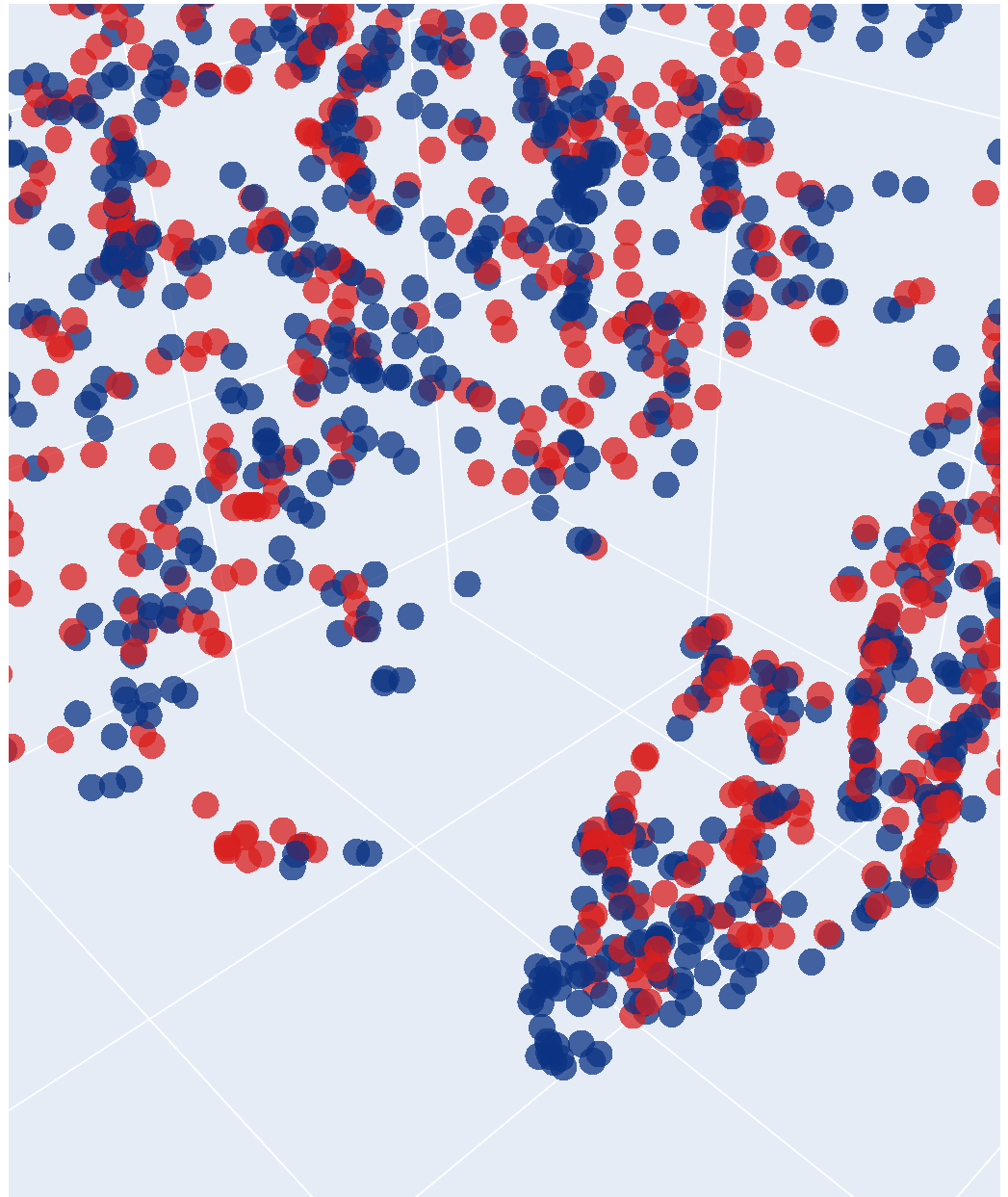
```
In [64]: 1 tsne3d = TSNE(  
2         n_components=3,  
3         init='random', # pca  
4         random_state=101,  
5         method='barnes_hut',  
6         n_iter=1000,  
7         verbose=2,  
8         angle=0.5  
9     ).fit_transform(X)
```

```
[t-SNE] Computing 91 nearest neighbors...
[t-SNE] Indexed 30000 samples in 0.071s...
[t-SNE] Computed neighbors for 30000 samples in 4.859s...
[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: 0.080413
[t-SNE] Computed conditional probabilities in 0.858s
[t-SNE] Iteration 50: error = 110.5925522, gradient norm = 0.0148754 (50 iterations in 33.847s)
[t-SNE] Iteration 100: error = 86.6675262, gradient norm = 0.0046031 (50 iterations in 24.531s)
[t-SNE] Iteration 150: error = 83.4906464, gradient norm = 0.0022670 (50 iterations in 24.854s)
[t-SNE] Iteration 200: error = 82.2365875, gradient norm = 0.0015197 (50 iterations in 25.397s)
[t-SNE] Iteration 250: error = 81.5207214, gradient norm = 0.0011329 (50 iterations in 20.254s)
[t-SNE] KL divergence after 250 iterations with early exaggeration: 81.520721
[t-SNE] Iteration 300: error = 2.9802334, gradient norm = 0.0038841 (50 iterations in 27.950s)
[t-SNE] Iteration 350: error = 2.3667295, gradient norm = 0.0030882 (50 iterations in 29.651s)
[t-SNE] Iteration 400: error = 2.0555992, gradient norm = 0.0024692 (50 iterations in 31.315s)
[t-SNE] Iteration 450: error = 1.8723609, gradient norm = 0.0020611 (50 iterations in 29.924s)
[t-SNE] Iteration 500: error = 1.7512760, gradient norm = 0.0017748 (50 iterations in 30.053s)
```

```
[t-SNE] Iteration 550: error = 1.6657951, gradient norm = 0.0015336 (50 iterations in 27.281s)
[t-SNE] Iteration 600: error = 1.6035148, gradient norm = 0.0013240 (50 iterations in 28.066s)
[t-SNE] Iteration 650: error = 1.5577862, gradient norm = 0.0011421 (50 iterations in 30.451s)
[t-SNE] Iteration 700: error = 1.5237987, gradient norm = 0.0009500 (50 iterations in 28.701s)
[t-SNE] Iteration 750: error = 1.4991173, gradient norm = 0.0007596 (50 iterations in 30.606s)
[t-SNE] Iteration 800: error = 1.4807075, gradient norm = 0.0006571 (50 iterations in 33.086s)
[t-SNE] Iteration 850: error = 1.4661463, gradient norm = 0.0005565 (50 iterations in 36.494s)
[t-SNE] Iteration 900: error = 1.4544582, gradient norm = 0.0004825 (50 iterations in 38.770s)
[t-SNE] Iteration 950: error = 1.4449804, gradient norm = 0.0004071 (50 iterations in 32.307s)
[t-SNE] Iteration 1000: error = 1.4379606, gradient norm = 0.0003029 (50 iterations in 31.926s)
[t-SNE] KL divergence after 1000 iterations: 1.437961
```

```
In [65]: 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 fe
23 fig=dict(data=data, layout=layout)
24 py.iplot(fig, filename='3DBubble')
```

3d embedding with engineered features




```
In [66]: 1 ques_df = df[['question1', 'question2']]
         2 ques_df.head()
```

Out[66]:

	question1	question2
398782	what is the best marketing automation tool for...	what is the best marketing automation tool for...
115086	i am poor but i want to invest what should i do	i am quite poor and i want to be very rich wh...
327711	i am from india and live abroad i met a guy f...	t i e t to thapar university to thapar univers...
367788	why do so many people in the u s hate the sou...	my boyfriend doesnt feel guilty when he hurts ...
151235	consequences of bhopal gas tragedy	what was the reason behind the bhopal gas tragedy

```
In [67]: 1 final_df = df.drop(columns=['id', 'qid1', 'qid2', 'question1', 'question2'])
         2 print(final_df.shape)
         3 final_df.head()
```

(30000, 23)

Out[67]:

	is_duplicate	q1_len	q2_len	q1_num_words	q2_num_words	word_common	word_total
398782	1	75	76	13	13	12	26
115086	0	48	56	13	16	8	24
327711	0	104	119	28	21	4	38
367788	0	58	145	14	32	1	34
151235	0	34	49	5	9	3	13

```
In [68]: 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 [69]: 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[69]: (30000, 6000)

```
In [70]: 1 final_df = pd.concat([final_df, temp_df], axis=1)
2 print(final_df.shape)
3 final_df.head()
```

(30000, 6023)

Out[70]:

	is_duplicate	q1_len	q2_len	q1_num_words	q2_num_words	word_common	word_total
398782	1	75	76	13	13	12	26
115086	0	48	56	13	16	8	24
327711	0	104	119	28	21	4	38
367788	0	58	145	14	32	1	34
151235	0	34	49	5	9	3	13

```
In [71]: 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:].valu
3
```

KNN

```
In [72]: 1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.metrics import accuracy_score
3 knn_classifier = KNeighborsClassifier(n_neighbors=3)
4 knn_classifier.fit(X_train, y_train)
5 KNeighborsClassifier(n_neighbors=3)
6 y_pred = knn_classifier.predict(X_test)
7 accuracy = accuracy_score(y_test, y_pred)
8 accuracy
```

Out[72]: 0.6956666666666667

Decision Tree

```
In [73]: 1 from sklearn.tree import DecisionTreeClassifier
2 clf = DecisionTreeClassifier()
3 clf.fit(X_train, y_train)
4 DecisionTreeClassifier()
5 y_pred = clf.predict(X_test)
6 accuracy = metrics.accuracy_score(y_test, y_pred)
7 accuracy
```

Out[73]: 0.7313333333333333

In []: 1

Random Forest

```
In [74]: 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[74]: 0.7838333333333334

XGBoost (Extreme Gradient Boosting)

```
In [75]: 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[75]: 0.7946666666666666

Boosting Algorithm

```
In [79]: 1 from sklearn.ensemble import AdaBoostClassifier
2 base_classifier = DecisionTreeClassifier(max_depth=3)
3 adaboost_classifier = AdaBoostClassifier(base_classifier,
4 n_estimators=50, random_state=42)
5 adaboost_classifier.fit(X_train, y_train)
6 AdaBoostClassifier(estimator=DecisionTreeClassifier(max_depth=1),
7 random_state=42)
8 y_pred = adaboost_classifier.predict(X_test)
9 accuracy = metrics.accuracy_score(y_test, y_pred)
10 accuracy
```

Out[79]: 0.7725

Logistic Regression

```
In [80]: 1 from sklearn import linear_model
          2 lrg = linear_model.LogisticRegression()
          3 lrg.fit(X_train, y_train)
          4 LogisticRegression()
          5 lrg.score(X_test, y_test)
```

Out[80]: 0.7151666666666666

```
In [84]: 1 from sklearn.metrics import confusion_matrix
```

```
In [85]: 1 # for random forest model
          2 confusion_matrix(y_test,y_pred)
```

Out[85]: array([[3191, 621],
 [744, 1444]], dtype=int64)

```
In [86]: 1 # for xgboost model
          2 confusion_matrix(y_test,y_pred1)
```

Out[86]: array([[3237, 575],
 [657, 1531]], dtype=int64)