

In [ ]:

```
import pandas as pd
import matplotlib.pyplot as plt
import re
import time
import warnings
import numpy as np
from nltk.corpus import stopwords
from sklearn.preprocessing import normalize
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
warnings.filterwarnings("ignore")
import sys
import os
import pandas as pd
import numpy as np
from tqdm import tqdm

import spacy
```

In [ ]:

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

df['question1'] = df['question1'].apply(lambda x: str(x))
df['question2'] = df['question2'].apply(lambda x: str(x))
```

In [ ]:

```
df.head()
```

Out[ ]:

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

In [ ]:

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

tfidf = TfidfVectorizer(lowercase=False, )
tfidf.fit_transform(questions)

word2tfidf = dict(zip(tfidf.get_feature_names(), tfidf.idf_))
```

In [ ]:

```
nlp = spacy.load('en_core_web_sm')

vecs1 = []
for qu1 in tqdm(list(df['question1'])):
    doc1 = nlp(qu1)
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 404290/404290 [2:13:51<00:00, 50.34it/s]
```

```
100%|██████████████████████████████████████████████████████████████████████████████| 404290/404290 [1:47:52<00:00, 62.46it/s]
```

```
if os.path.isfile('nlp_features_train.csv'):
    dfnlp = pd.read_csv("nlp_features_train.csv",encoding='latin-1')
else:
    print("download nlp_features_train.csv from drive or run previous notebook")

if os.path.isfile('df_fe_without_preprocessing_train.csv'):
    dfppro = pd.read_csv("df_fe_without_preprocessing_train.csv",encoding='latin-1')
else:
    print("download df_fe_without_preprocessing_train.csv from drive or run previous notebook")
```

```
df1 = dfnlp.drop(['qid1', 'qid2', 'question1', 'question2'], axis=1)
df2 = dfppro.drop(['qid1', 'qid2', 'question1', 'question2', 'is_duplicate'], axis=1)
df3 = df.drop(['qid1', 'qid2', 'question1', 'question2', 'is_duplicate'], axis=1)
df3_q1 = pd.DataFrame(df3.q1_feats_m.values.tolist(), index= df3.index)
df3_q2 = pd.DataFrame(df3.q2_feats_m.values.tolist(), index= df3.index)
```

Out[ ]:

id	is_duplicate	cwc_min	cwc_max	csc_min	csc_max	ctc_min	ctc_max	last_word_eq	first_word_eq	abs_len_diff	measures
1	0	1	1	1	1	1	1	0	0	0	0.000000
2	0	1	1	1	1	1	1	0	0	0	0.000000
3	0	1	1	1	1	1	1	0	0	0	0.000000
4	0	1	1	1	1	1	1	0	0	0	0.000000
5	0	1	1	1	1	1	1	0	0	0	0.000000
6	0	1	1	1	1	1	1	0	0	0	0.000000
7	0	1	1	1	1	1	1	0	0	0	0.000000
8	0	1	1	1	1	1	1	0	0	0	0.000000
9	0	1	1	1	1	1	1	0	0	0	0.000000
10	0	1	1	1	1	1	1	0	0	0	0.000000
11	0	1	1	1	1	1	1	0	0	0	0.000000
12	0	1	1	1	1	1	1	0	0	0	0.000000
13	0	1	1	1	1	1	1	0	0	0	0.000000
14	0	1	1	1	1	1	1	0	0	0	0.000000
15	0	1	1	1	1	1	1	0	0	0	0.000000
16	0	1	1	1	1	1	1	0	0	0	0.000000
17	0	1	1	1	1	1	1	0	0	0	0.000000
18	0	1	1	1	1	1	1	0	0	0	0.000000
19	0	1	1	1	1	1	1	0	0	0	0.000000
20	0	1	1	1	1	1	1	0	0	0	0.000000
21	0	1	1	1	1	1	1	0	0	0	0.000000
22	0	1	1	1	1	1	1	0	0	0	0.000000
23	0	1	1	1	1	1	1	0	0	0	0.000000
24	0	1	1	1	1	1	1	0	0	0	0.000000
25	0	1	1	1	1	1	1	0	0	0	0.000000
26	0	1	1	1	1	1	1	0	0	0	0.000000
27	0	1	1	1	1	1	1	0	0	0	0.000000
28	0	1	1	1	1	1	1	0	0	0	0.000000
29	0	1	1	1	1	1	1	0	0	0	0.000000
30	0	1	1	1	1	1	1	0	0	0	0.000000
31	0	1	1	1	1	1	1	0	0	0	0.000000
32	0	1	1	1	1	1	1	0	0	0	0.000000
33	0	1	1	1	1	1	1	0	0	0	0.000000
34	0	1	1	1	1	1	1	0	0	0	0.000000
35	0	1	1	1	1	1	1	0	0	0	0.000000
36	0	1	1	1	1	1	1	0	0	0	0.000000
37	0	1	1	1	1	1	1	0	0	0	0.000000
38	0	1	1	1	1	1	1	0	0	0	0.000000
39	0	1	1	1	1	1	1	0	0	0	0.000000
40	0	1	1	1	1	1	1	0	0	0	0.000000
41	0	1</									

0	id	is_duplicate	0.799984	0.333328	0.399992	0.249997	0.399996	0.285712	last_word_eq	first_word_eq	abs_len_diff	me
1	1	0	0.799984	0.399996	0.749981	0.599988	0.699993	0.466664	0.0	1.0	5.0	
2	2	0	0.399992	0.333328	0.399992	0.249997	0.399996	0.285712	0.0	1.0	4.0	
3	3	0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.0	0.0	2.0	
4	4	0	0.399992	0.199998	0.999950	0.666644	0.571420	0.307690	0.0	1.0	6.0	

In [ ]:

```
# data before preprocessing
df2.head()
```

Out[ ]:

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

In [ ]:

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

Out[ ]:

	0	1	2	3	4	5	6	7	8	9	...
0	121.929927	100.083900	72.497894	115.641800	48.370870	34.619058	172.057787	-92.502617	113.223315	50.562441	...
1	-78.070939	54.843781	82.738482	98.191872	51.234859	55.013510	-39.140730	-82.692352	45.161489	-9.556289	...
2	-5.355015	73.671810	14.376365	104.130241	1.433537	35.229116	148.519385	-97.124595	41.972195	50.948731	...
3	5.778359	-34.712038	48.999631	59.699204	40.661263	41.658731	-36.808594	24.170655	0.235600	29.407290	...
4	51.138220	38.587312	123.639488	53.333041	47.062739	37.356212	298.722753	106.421119	106.248914	65.880707	...

**5 rows x 384 columns**

In [ ]:

```
# Questions 2 tfidf weighted word2vec
df3 %>% q2.head()
```

Out[ ]:

	0	1	2	3	4	5	6	7	8	9	...
0	125.983301	95.636485	42.114702	95.449980	-37.386295	39.400078	-148.116070	-87.851475	110.371966	62.272814	... 16
1	-106.871904	80.290331	79.066297	59.302092	-42.175328	117.616655	-144.364237	-127.131513	22.962533	25.397575	... -4
2	7.072875	15.513378	1.846914	85.937583	-33.808811	94.702337	-122.256856	-114.009530	53.922293	60.131814	... 8

3	0	1	2	3	4	5	6	7	8	9	...	3
39.421531	44.136989	24.010929	85.265863	-0.339022	-9.323137	-60.499651	-37.044763	49.407848	23.350150			
4					-	-	-					-2
31.950101	62.854106	1.778164	36.218768	45.130875	66.674880	106.342341	-22.901008	59.835938	62.663961	...		

5 rows x 384 columns



In [ ]:

```
print("Number of features in nlp dataframe :", df1.shape[1])
print("Number of features in preprocessed dataframe :", df2.shape[1])
print("Number of features in question1 w2v dataframe :", df3_q1.shape[1])
print("Number of features in question2 w2v dataframe :", df3_q2.shape[1])
print("Number of features in final dataframe :", df1.shape[1]+df2.shape[1]+df3_q1.shape[1]+df3_q2.shape[1])
```

```
Number of features in nlp dataframe : 17
Number of features in preprocessed dataframe : 12
Number of features in question1 w2v dataframe : 384
Number of features in question2 w2v dataframe : 384
Number of features in final dataframe : 794
```

In [ ]:

```
# storing the final features to csv file
if not os.path.isfile('final_features.csv'):
    df3_q1['id']=df1['id']
    df3_q2['id']=df1['id']
    df1 = df1.merge(df2, on='id',how='left')
    df2 = df3_q1.merge(df3_q2, on='id',how='left')
    result = df1.merge(df2, on='id',how='left')
    result.to_csv('final_features.csv')
```