## 3.6 Featurizing text data with tfidf weighted word-vectors

```python
In [2]: 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

        # exctract word2vec vectors
        # https://github.com/explosion/spaCy/issues/1721
        # http://landinghub.visualstudio.com/visual-cpp-build-tools
        import spacy
```

```
C:\Users\Himanshu Pc\Anaconda3\lib\site-packages\sklearn\feature_extraction\text.py:17: DeprecationWarning: Using or importing the ABCs from 'col
lections' instead of from 'collections.abc' is deprecated, and in 3.8 it will stop working
  from collections import Mapping, defaultdict
```

```python
In [3]: # avoid decoding problems
        df = pd.read_csv("train.csv")

        # encode questions to unicode
        # https://stackoverflow.com/a/6812069
        # ----------------- python 2 --------------------
        # df['question1'] = df['question1'].apply(lambda x: unicode(str(x),"utf-8"))
        # df['question2'] = df['question2'].apply(lambda x: unicode(str(x),"utf-8"))
        # ----------------- python 3 --------------------
        df['question1'] = df['question1'].apply(lambda x: str(x))
        df['question2'] = df['question2'].apply(lambda x: str(x))
```

```
In [4]:  df.head()
```

Out[4]:

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

```
In [5]:  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)

         # dict key:word and value:tf-idf score
         word2tfidf = dict(zip(tfidf.get_feature_names(), tfidf.idf_))
```

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

```
In [6]:  import en_core_web_sm
```

```
In [7]:   # en_vectors_web_lg, which includes over 1 million unique vectors.
          #nlp = spacy.load('en_core_web_sm')
          nlp = en_core_web_sm.load()

          vecs1 = []
          # https://github.com/noamraph/tqdm
          # tqdm is used to print the progress bar
          for qu1 in tqdm(list(df['question1'])):
              doc1 = nlp(qu1)
              # 384 is the number of dimensions of vectors
              mean_vec1 = np.zeros([len(doc1), len(doc1[0].vector)])
              for word1 in doc1:
                  # word2vec
                  vec1 = word1.vector
                  # fetch df score
                  try:
                      idf = word2tfidf[str(word1)]
                  except:
                      idf = 0
                  # compute final vec
                  mean_vec1 += vec1 * idf
              mean_vec1 = mean_vec1.mean(axis=0)
              vecs1.append(mean_vec1)
          df['q1_feats_m'] = list(vecs1)
```

100%|████████████████████████████████████████████████| 404290/404290 [47:08<00:00, 142.94it/s]

```
In [20]:  vecs2 = []
          for qu2 in tqdm(list(df['question2'])):
              doc2 = nlp(qu2)
              mean_vec1 = np.zeros([len(doc1), len(doc2[0].vector)])
              for word2 in doc2:
                  # word2vec
                  vec2 = word2.vector
                  # fetch df score
                  try:
                      idf = word2tfidf[str(word2)]
                  except:
                      #print word
                      idf = 0
                  # compute final vec
                  mean_vec2 += vec2 * idf
              mean_vec2 = mean_vec2.mean(axis=0)
              vecs2.append(mean_vec2)
          df['q2_feats_m'] = list(vecs2)
```

100%|████████████████████████████████████████████████| 404290/404290 [1:16:04<00:00, 88.57it/s]

```
In [21]: #prepro_features_train.csv (Simple Preprocessing Feartures)
         #nlp_features_train.csv (NLP Features)
         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")
```

```
In [22]: 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)
```

```
In [23]: # dataframe of nlp features
         df1.head()
```

Out[23]:

| | id | is_duplicate | cwc_min | cwc_max | csc_min | csc_max | ctc_min | ctc_max | last_word_eq | first_word_eq | abs_len_diff | mean_len | token_set_ratio | token_sort_ratio | fuzz_ratio |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0.999980 | 0.833319 | 0.999983 | 0.999983 | 0.916659 | 0.785709 | 0.0 | 1.0 | 2.0 | 13.0 | 100 | 93 | 93 |
| 1 | 1 | 0 | 0.799984 | 0.399996 | 0.749981 | 0.599988 | 0.699993 | 0.466664 | 0.0 | 1.0 | 5.0 | 12.5 | 86 | 63 | 66 |
| 2 | 2 | 0 | 0.399992 | 0.333328 | 0.399992 | 0.249997 | 0.399996 | 0.285712 | 0.0 | 1.0 | 4.0 | 12.0 | 66 | 66 | 54 |
| 3 | 3 | 0 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.000000 | 0.0 | 0.0 | 2.0 | 12.0 | 36 | 36 | 35 |
| 4 | 4 | 0 | 0.399992 | 0.199998 | 0.999950 | 0.666644 | 0.571420 | 0.307690 | 0.0 | 1.0 | 6.0 | 10.0 | 67 | 47 | 46 |

```
In [24]:  # data before preprocessing
          df2.head()
```

Out[24]:

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

```
In [25]:  # Questions 1 tfidf weighted word2vec
          df3_q1.head()
```

Out[25]:

|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... | 86 | 87 | 88 | 89 |
|---|---|---|---|---|---|---|---|---|---|---|-----|----|----|----|----|
| 0 | -6.179507 | 37.450731 | -67.929894 | 32.224274 | 143.348826 | 135.374574 | 17.865208 | 54.562352 | 81.618936 | 232.909839 | ... | -71.834689 | -60.222858 | -22.026407 | 103.336720 | -68.477 |
| 1 | 9.236668 | -80.371416 | -45.785907 | 78.291656 | 183.568221 | 100.894077 | 74.344804 | 48.360802 | 127.297421 | 112.987302 | ... | -32.130515 | -98.080325 | 19.113790 | -20.507508 | -76.981 |
| 2 | 97.546832 | 22.972194 | -39.558379 | 18.723413 | 56.928618 | 48.307643 | 8.719268 | 36.893738 | 106.899947 | 226.283077 | ... | -66.835018 | 87.592131 | 4.032431 | 56.851710 | -43.625 |
| 3 | 57.586978 | -22.017089 | -4.599294 | -88.939271 | -4.732171 | -54.209048 | 74.614947 | 106.533737 | 15.520611 | 39.009709 | ... | 28.362970 | 41.981222 | -11.204987 | 16.833428 | -36.372 |
| 4 | 83.185784 | -40.506985 | -83.403923 | -52.648658 | 79.074884 | -19.038248 | 53.728722 | 97.648612 | 160.555822 | 290.541356 | ... | -4.390959 | 109.604406 | -91.160167 | -25.739913 | 133.123 |

5 rows × 96 columns

```
In [26]:  # Questions 2 tfidf weighted word2vec
          df3_q2.head()
```

Out[26]:

|   | 0 |
|---|---|
| 0 | 1.187354 |
| 1 | 3.790469 |
| 2 | 4.870265 |
| 3 | 6.003683 |
| 4 | 7.616745 |

```
In [27]: 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 : 96
Number of features in question2 w2v  dataframe : 1
Number of features in final dataframe  : 126
```

```
In [28]: # 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')
```