

# Introduction

```
In [1]: ▶ import numpy as np           #for large and multi-dime
import pandas as pd           #for data manipulation an
import os
import nltk                     #Natural Language process
```

```
In [2]: ▶ import warnings
warnings.filterwarnings("ignore")           #Ignoring unnecessary w
from gensim.models import Word2Vec
```

```
In [5]: ▶ data_path = "C:/Users/Sagnik_laptop/Documents/ML/ML1010/ML1010-A1/data/Review
data = pd.read_csv(data_path)
data_sel = data.head(10000)                #Considering only top
```

```
In [6]: ▶ data_sel.columns                #dataset column names
```

```
Out[6]: Index(['Id', 'ProductId', 'UserId', 'ProfileName', 'HelpfulnessNumerator',
              'HelpfulnessDenominator', 'Score', 'Time', 'Summary', 'Text'],
              dtype='object')
```

```
In [7]: ▶ data_score_removed = data_sel[data_sel['Score']!=3]           #removing neutral
```

## Converting Score values into class label either Positive or Negative.

```
In [8]: ▶ def partition(x):
            if x < 3:
                return 'positive'
            return 'negative'

score_upd = data_score_removed['Score']
t = score_upd.map(partition)
data_score_removed['Score']=t
```

## 1.Basic Cleaning-removing duplicates

```
In [31]: final_data = data_score_removed.drop_duplicates(subset={"UserId", "ProfileName", "Score", "HelpfulnessNumerator", "HelpfulnessDenominator", "Text"})
final = final_data[final_data['HelpfulnessNumerator'] <= final_data['HelpfulnessDenominator']]
summary_text = final['Text']
# text data
corpus = np.array(summary_text)
final_y = final['Score']
# print current data set
corpus
```

```
Out[31]: array(['I have bought several of the Vitality canned dog food products and have found them all to be of good quality. The product looks more like a stew than a processed meat and it smells better. My Labrador is finicky and she appreciates this product better than most.',
'Product arrived labeled as Jumbo Salted Peanuts...the peanuts were actually small sized unsalted. Not sure if this was an error or if the vendor intended to represent the product as "Jumbo".',
'This is a confection that has been around a few centuries. It is a light, pillowy citrus gelatin with nuts - in this case Filberts. And it is cut into tiny squares and then liberally coated with powdered sugar. And it is a tiny mouthful of heaven. Not too chewy, and very flavorful. I highly recommend this yummy treat. If you are familiar with the story of C.S. Lewis\' "The Lion, The Witch, and The Wardrobe" - this is the treat that seduces Edmund into selling out his Brother and Sisters to the Witch.',
...,
'I wanted to solely breastfeed but was unable to keep up and had to supplement formula. I chose Similac because they are a very reputable company and I had a ton of it from the hospital. We have used both the powder and ready made and he likes both. He got a little constipated at the beginning but has been on it now for 5 months and no problems. I read some other reviews about sucrose...well, sucrose is just sugar and ALL formula has sugar in it. Other companies just label it differently and since this one is organic it has to contain a pure sugar. As far as the hexane goes, I did my research, ALL formulas except one (Babys Only) use the hexane method. Babys Only contains brown rice syrup that has been shown to have high levels of arsenic..uh, no thanks. Also Similac has been tested and found to have no hexane present in the product. Overall my baby is healthy and smart so I am very happy with this formula.',
'i love the fact that i can get this delivered to my house with no delivery charge.it is so hard to find organic formula',
"We have a 7week old... He had gas and constipation problems for the first 5 weeks. We tried two different kinds of similac including for fussiness and gas and neither seemed to work. We switched to the organic a few weeks ago and saw quick improvement. I wish I could breast feed but I'm unable to, so for now this seems the best option especially since it was recommended we stick with a ready made formula for the gas problems.<br />I've read a lot of the reviews and took into consideration the information about sucrose. I plan on talking to the pediatrician and my midwife for additional information beyond the article written about it, especially since that is from 2008. I realize the concern and I am doing research on making my own formula so I know exactly what's in it and that it's organic, but in the meantime baby L eats great with this, is healthy, and has fewer stomach problems. It's middle of the road when it comes to $ - although Amazon is one of the more expensive places!!! Target has the best price. So for now it works and I recommend it!!"],
dtype=object)
```

## ## Text pre-processing

```
In [27]: > from nltk.corpus import stopwords
stop_words = nltk.corpus.stopwords.words('english')
```

```
In [33]: > # from nltk.stem import PorterStemmer                # Stemmer
# import re
# temp = []
# snow = nltk.stem.SnowballStemmer('english')
# for sentence in final_X:
#     sentence = sentence.lower()                # Converting to Lowercase
#     cleanr = re.compile('<.*?>')
#     sentence = re.sub(cleanr, ' ', sentence)    #Removing HTML tags
#     sentence = re.sub(r'[?|!|\'|\"|#]',r'',sentence)
#     sentence = re.sub(r'[.,|)|(|\\|/]',r' ',sentence)    #Removing Punctuations
#     words = [snow.stem(word) for word in sentence.split() if word not in stop_words]
#     temp.append(words)

# final_X = temp
wpt = nltk.WordPunctTokenizer()
def normalize_document(doc):
    # lower case and remove special characters\whitespaces
    doc = re.sub(r'^a-zA-Z0-9\s', '', doc, re.I)
    # covert to lower case
    doc = doc.lower()
    doc = doc.strip()
    cleanr = re.compile('<.*?>')
    # remove html tags
    doc = re.sub(cleanr, ' ', doc)
    #remove punctuations
    doc = re.sub(r'[?|!|\'|\"|#]',r'', doc)
    doc = re.sub(r'[.,|)|(|\\|/]',r' ', doc)
    # tokenize document
    tokens = wpt.tokenize(doc)
    # filter stopwords out of document
    filtered_tokens = [token for token in tokens if token not in stop_words]
    # re-create document from filtered tokens
    doc = ' '.join(filtered_tokens)
    return doc
normalize_corpus = np.vectorize(normalize_document)
```

```
In [34]: norm_corpus = normalize_corpus(corpus)
norm_corpus
```

```
Out[34]: array(['bought several vitality canned dog food products found good quality
product looks like stew processed meat smells better labrador finicky appreciates
product better',
'product arrived labeled jumbo salted peanuts peanuts actually small sized
unsalted sure error vendor intended represent product jumbo',
'confection around centuries light pillowy citrus gelatin nuts - case filberts
cut tiny squares liberally coated powdered sugar tiny mouthful heavenly chewy
flavorful highly recommend yummy treat familiar story c lewis lion witch
wardrobe - treat seduces edmund selling brother sisters witch',
...,
'wanted solely breastfeed unable keep supplement formula chose similar reputable
company ton hospital used powder ready made likes got little constipated beginning
5 months problems read reviews sucrose well sucrose sugar formula sugar
companies label differently since this one organic contain pure sugar far hexane
goes research formulas except one babys use hexane method babys contains brown
rice syrup shown high levels arsenic uh thanks also similac tested found hexane
present product overall baby healthy smart happy formula',
'love fact get relieved house delivery charge it hard find organic formula',
'7 week old gas constipation problems first 5 weeks tried two different kinds
similac including fussiness gas neither seemed work switched organic weeks ago
saw quick improvement wish could breast feed im unable seems best option
especially since recommended stick ready made formula gas problems ive read lot
reviews took consideration information sucrose plan talking pediatrician midwife
additional information beyond article written especially since 2008 realize concern
research making formula know exactly whats organic mean time baby 1 eats great
healthy fewer stomach problems middle road comes $ - although amazon one
expensive places target best price works recommend'],
dtype='<U6094')
```

```
In [25]: # save the data
import dill
dill.dump_session('notebook_env.db')

# Load the data
# dill.load_session('notebook_env.db')
```

## Bag of words

```
In [52]: ▶ from sklearn.feature_extraction.text import CountVectorizer #For Bag
cv = CountVectorizer()
cv_matrix = cv.fit_transform(norm_corpus)
cv_matrix = cv_matrix.toarray()
cv_matrix
```

```
Out[52]: array([[0, 0, 0, ..., 0, 0, 0],
               [0, 0, 0, ..., 0, 0, 0],
               [0, 0, 0, ..., 0, 0, 0],
               ...,
               [0, 0, 0, ..., 0, 0, 0],
               [0, 0, 0, ..., 0, 0, 0],
               [0, 0, 0, ..., 0, 0, 0]], dtype=int64)
```

```
In [58]: vocab = cv.get_feature_names()
pd.DataFrame(cv_matrix, columns=vocab)
```

Out[58]:

	00	000	0003	000kwh	002	008	0100	0174	02	03	...	zoo	zoom	zotz	zs	zucch
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
5	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
6	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
7	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
9	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
10	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
12	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
13	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
14	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
15	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
16	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
17	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
18	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
19	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
20	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
21	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
22	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
23	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
24	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
25	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
26	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
27	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
28	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
29	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
8688	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8689	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	

	00	000	0003	000kwh	002	008	0100	0174	02	03	...	zoo	zoom	zotz	zs	zucch
8690	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8691	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8692	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8693	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8694	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8695	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8696	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8697	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8698	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8699	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8700	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8701	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8702	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8703	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8704	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8705	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8706	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8707	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8708	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8709	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8710	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8711	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8712	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8713	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8714	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8715	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8716	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	
8717	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	

8718 rows × 18608 columns



```
In [59]: # Bag of n-grams model
bv = CountVectorizer(ngram_range=(2,2))
bv_matrix = bv.fit_transform(norm_corpus)
bv_matrix = bv_matrix.toarray()
vocab = bv.get_feature_names()
pd.DataFrame(bv_matrix, columns=vocab)
```

Out[59]:

	00 10lb	00 12	00 15	00 22	00 24	00 30	00 47	00 49	00 50	00 51	...	zukes soft	zukes thing	zukes think	zukes top	zukes treats	zukes w
0	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
1	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
2	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
3	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
4	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
5	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
6	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
7	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
8	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
9	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
10	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
11	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
12	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
13	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
14	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
15	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
16	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
17	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
18	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
19	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
20	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
21	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
22	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
23	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
24	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
25	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
26	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
27	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
28	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	
29	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	



	00 10lb	00 12	00 15	00 22	00 24	00 30	00 47	00 49	00 50	00 51	...	zukes soft	zukes thing	zukes think	zukes top	zukes treats	zukes w/
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
8688	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8689	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8690	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8691	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8692	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8693	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8694	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8695	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8696	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8697	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8698	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8699	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8700	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8701	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8702	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8703	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8704	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8705	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8706	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8707	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8708	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8709	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8710	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8711	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8712	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8713	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8714	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8715	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8716	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0
8717	0	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0

8718 rows × 216937 columns



```
In [62]: #Term Frequency - Inverse Document Frequency(TF-IDF)
from sklearn.feature_extraction.text import TfidfVectorizer
tv = TfidfVectorizer(min_df=0., max_df=1., use_idf=True)
tv_matrix = tv.fit_transform(norm_corpus)
tv_matrix = tv_matrix.toarray()
vocab = tv.get_feature_names()
pd.DataFrame(np.round(tv_matrix, 2), columns=vocab)
```

Out[62]:

	00	000	0003	000kwh	002	008	0100	0174	02	03	...	zoo	zoom	zotz	zs	zuc
0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
2	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
10	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
11	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
12	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
13	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
14	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
15	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
16	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
17	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
18	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
19	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
20	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
21	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
22	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
23	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
24	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
25	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
26	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
27	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
28	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	
29	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	

	00	000	0003	000kwh	002	008	0100	0174	02	03	...	zoo	zoom	zotz	zs	zuc
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
8688	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8689	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8690	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8691	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8692	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8693	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8694	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8695	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8696	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8697	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8698	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8699	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8700	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8701	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8702	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8703	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8704	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8705	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8706	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8707	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8708	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8709	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8710	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8711	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8712	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8713	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8714	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8715	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8716	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0
8717	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0

8718 rows × 18608 columns



```
In [65]: ▶ import gensim
w2v_num_features = 500
w2v_model = gensim.models.Word2Vec(norm_corpus, size=w2v_num_features, window=
min_count=10, sample=1e-3)
```

```
In [66]: ▶ def averaged_word2vec_vectorizer(corpus, model, num_features):
    vocabulary = set(model.wv.index2word)

    def average_word_vectors(words, model, vocabulary, num_features):
        feature_vector = np.zeros((num_features,), dtype="float64")
        nwords = 0.

        for word in words:
            if word in vocabulary:
                nwords = nwords + 1.
                feature_vector = np.add(feature_vector, model[word])
        if nwords:
            feature_vector = np.divide(feature_vector, nwords)

        return feature_vector

    features = [average_word_vectors(tokenized_sentence, model, vocabulary, num_features)
                for tokenized_sentence in corpus]
    return np.array(features)
```

```
In [68]: ▶ features = averaged_word2vec_vectorizer(corpus=corpus, model=w2v_model,
num_features=500)
```