

✓ Worksheet 11

Name: Jeong Yong Yang, Junho Son UID: U95912941, U64222022

Topics

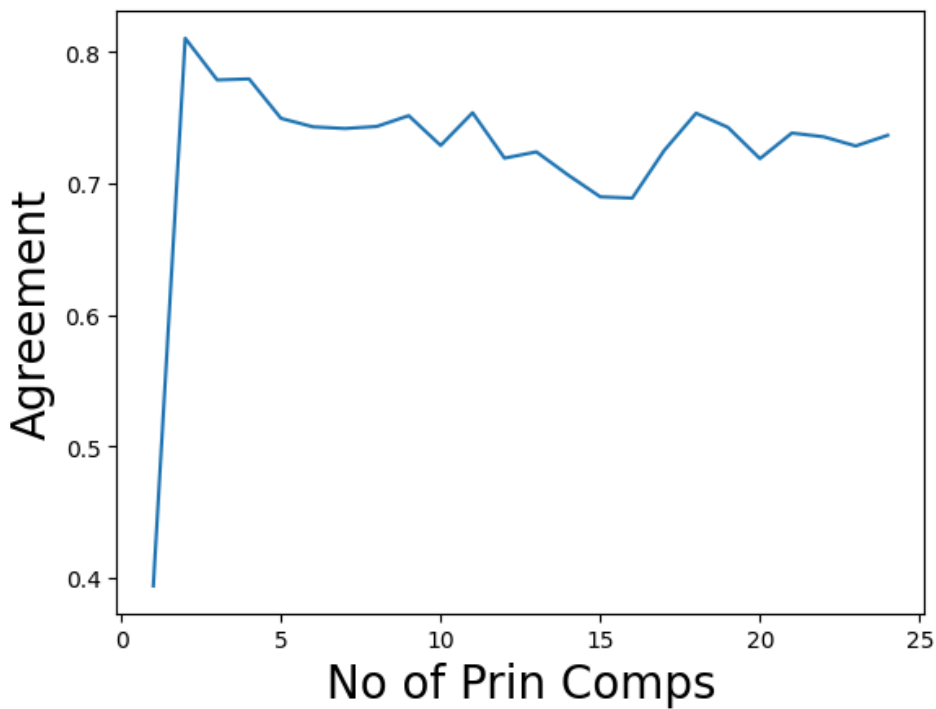
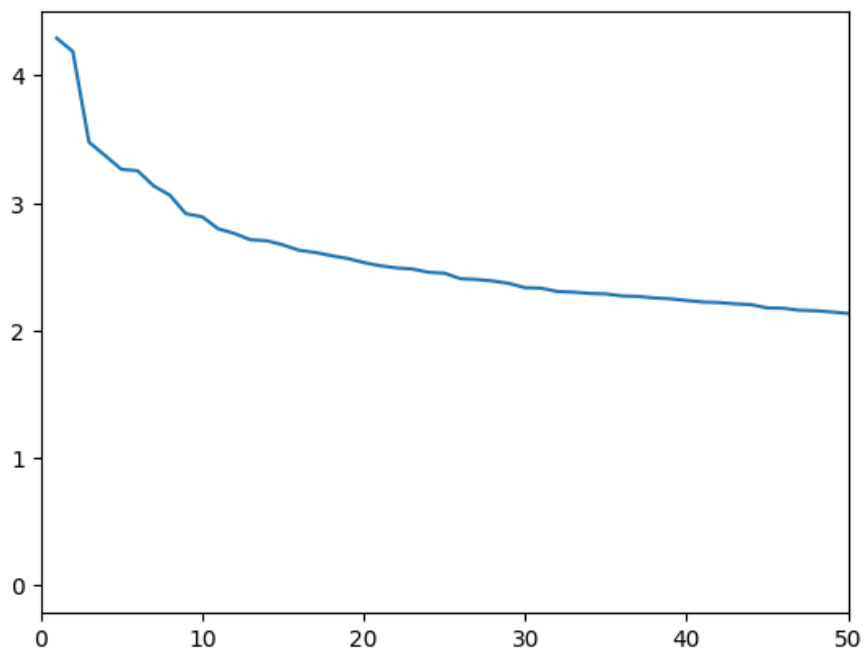
- Latent Semantic Analysis

Latent Semantic Analysis

In this section we will fetch news articles from 3 different categories. We will perform Tf-idf vectorization on the corpus of documents and use SVD to represent our corpus in the feature space of topics that we've uncovered from SVD. We will attempt to cluster the documents into 3 clusters as we vary the number of singular vectors we use to represent the corpus, and compare the output to the clustering created by the news article categories. Do we end up with a better clustering the more singular vectors we use?

```
1  import nltk
2  nltk.download('punkt')
3  nltk.download('stopwords')
4  import numpy as np
5  from sklearn import metrics
6  import matplotlib.pyplot as plt
7  from sklearn.cluster import KMeans
8  from sklearn.datasets import fetch_20newsgroups
9  from sklearn.feature_extraction.text import TfidfVectorizer
10 from nltk.stem.snowball import SnowballStemmer
11 from nltk.tokenize import word_tokenize, sent_tokenize
12
13
14 categories = ['comp.os.ms-windows.misc', 'sci.space', 'rec.sport.baseball']
15 news_data = fetch_20newsgroups(subset='train', categories=categories)
16 vectorizer = TfidfVectorizer(stop_words='english', min_df=4, max_df=0.8)
17
18 stemmed_data = [" ".join(SnowballStemmer("english", ignore_stopwords=True).stem(word)
19                          for sent in sent_tokenize(message)
20                          for word in word_tokenize(sent))
21                  for message in news_data.data]
22
23 dtm = vectorizer.fit_transform(stemmed_data)
24 terms = vectorizer.get_feature_names_out()
25 centered_dtm = dtm - np.mean(dtm, axis=0)
26
27 u, s, vt = np.linalg.svd(centered_dtm)
28 plt.xlim([0, 50])
29 plt.plot(range(1, len(s)+1), s)
30 plt.show()
31
32 ag = []
33 max = len(u)
34 for singular_vectors in range(1, 25):
35     vectorsk = u.dot(np.diag(s))[:, :singular_vectors]
36     kmeans = KMeans(n_clusters=3, init='k-means++', max_iter=100, n_init=10, random_state=0)
37     kmeans.fit_predict(np.asarray(vectorsk))
38     labelsk = kmeans.labels_
39     ag.append(metrics.v_measure_score(labelsk, news_data.target)) # closer to 1 means closer to news categories
40
41 plt.plot(range(1, 25), ag)
42 plt.ylabel('Agreement', size=20)
43 plt.xlabel('No of Prin Comps', size=20)
44 plt.show()
45
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.
```



No we do not end up with a better clustering the more singular vectors we use despite the agreement being pretty high. The best number of singular vectors to use seems to be 3 from the graph above.

✎ Embeddings

The data comes from the [Yelp Dataset](#). Each line is a review that consists of a label (0 for negative reviews and 1 for positive reviews) and a set of words.

```
1 i will never forget this single breakfast experience in mad...
0 the search for decent chinese takeout in madison continues ...
0 sorry but me julio fell way below the standard even for med...
1 so this is the kind of food that will kill you so there s t...
```

In order to transform the set of words into vectors, we will rely on a method of feature engineering called word embeddings (Tfidf is one way to get these embeddings). Rather than simply indicating which words are present, word embeddings represent each word by "embedding" it in a low-dimensional vector space which may carry more information about the semantic meaning of the word. (for example in this space, the words "King" and "Queen" would be close).

word2vec.txt contains the word2vec embeddings for about 15 thousand words. Not every word in each review is present in the provided word2vec.txt file. We can treat these words as being "out of vocabulary" and ignore them.

Example

Let x_i denote the sentence "a hot dog is not a sandwich because it is not square" and let a toy word2vec dictionary be as follows:

hot	0.1	0.2	0.3
not	-0.1	0.2	-0.3
sandwich	0.0	-0.2	0.4
square	0.2	-0.1	0.5

we would first trim the sentence to only contain words in our vocabulary: "hot not sandwich not square" then embed x_i into the feature space:

$$\phi_2(x_i) = \frac{1}{5}(\text{word2vec}(\text{hot}) + 2 \cdot \text{word2vec}(\text{not}) + \text{word2vec}(\text{sandwich}) + \text{word2vec}(\text{square})) = [0.02 \ 0.06 \ 0.12]^T$$

a) Implement a function to trim out-of-vocabulary words from the reviews. Your function should return an nd array of the same dimension and dtype as the original loaded dataset.

```

1  import csv
2  import numpy as np
3
4  VECTOR_LEN = 300    # Length of word2vec vector
5  MAX_WORD_LEN = 64   # Max word length in dict.txt and word2vec.txt
6
7  def load_tsv_dataset(file):
8      """
9      Loads raw data and returns a tuple containing the reviews and their ratings.
10
11      Parameters:
12          file (str): File path to the dataset tsv file.
13
14      Returns:
15          An np.ndarray of shape N. N is the number of data points in the tsv file.
16          Each element dataset[i] is a tuple (label, review), where the label is
17          an integer (0 or 1) and the review is a string.
18      """
19      dataset = np.loadtxt(file, delimiter='\t', comments=None, encoding='utf-8',
20                          dtype='l,O')
21      return dataset
22
23  def load_feature_dictionary(file):
24      """
25      Creates a map of words to vectors using the file that has the word2vec
26      embeddings.
27
28      Parameters:
29          file (str): File path to the word2vec embedding file.
30
31      Returns:
32          A dictionary indexed by words, returning the corresponding word2vec
33          embedding np.ndarray.
34      """
35      word2vec_map = dict()
36      with open(file) as f:
37          read_file = csv.reader(f, delimiter='\t')
```

```

38         for row in read_file:
39             word, embedding = row[0], row[1:]
40             word2vec_map[word] = np.array(embedding, dtype=float)
41     return word2vec_map
42
43
44 def trim_reviews(path_to_dataset):
45     sentences = load_tsv_dataset(path_to_dataset)
46     categories = []
47     all_sentences = []
48     for x in sentences:
49         categories.append(x[0])
50         all_sentences.append(x[1].split())
51
52     word2vec_map = load_feature_dictionary("./word2vec.txt")
53     final_value = []
54     index = 0
55     for x in all_sentences:
56
57         value = [0] * VECTOR_LEN
58         for y in x:
59             if y in word2vec_map:
60                 factor = [element * 0.2 for element in word2vec_map[y]]
61                 value = [a + b for a, b in zip(value, factor)]
62             value = np.insert(value, 0, categories[index])
63             value = [round(float(element), 6) for element in value]
64
65         final_value.append(value)
66         index += 1
67
68     return final_value
69
70 trim_train = trim_reviews("./train_small.tsv")
71 trim_test = trim_reviews("./test_small.tsv")
72
73

```