

A decorative graphic in the top-left corner consisting of a grid of squares in shades of purple, blue, and green, arranged in a stepped pattern.

Topic Modeling

Neba Nfonsang
University of Denver



Objectives

- Find topics in documents
- Perform explicit semantics analysis
- Document clustering
- Latent semantic analysis
- Non-negative matrix factorization



Topic Modeling

- NLP is commonly used to search and extract information.
- Topic modeling involves automatically defining the subject of a document based on the important words in the document.
- We have already seen how to represent the relationship between words and documents using TF-IDF weights.
- A matrix (a collection of vectors) can be used to capture the statistical relationship between words and documents.



Term-Document Matrix

- A term-document matrix represents the relationship between words and documents.
- The rows are words while the columns are documents.
- A text or corpus of documents can be represented using a matrix.
- When we use a matrix to represent a corpus of documents, each column is a document and each row corresponds to a unique word in the text.

Term-Document Matrix (TDM)

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1

- The values in the term document matrix could be Boolean values indicating whether the word is present in the document or not
- The values in the TDM could be TF-IDF weights which indicate relative importance of the word in the document.
- Frequency counts could be used as well.

Explicit Semantic Analysis (ESA)

	d1	d2	d3
w1	w1d1	0	w1d3
w2	w2d2	w2d2	0
w3	0	w3d2	0
w4	0	0	0
w5	w5d2	w5d2	0
w6	w6d2	w6d2	w6d2

- ESA involves vectoral representation of text (words or documents) where the vector entries are typically tf-idf weightings.
- So, ESA starts with creating a Term Document Matrix of tf-idf weights. **w_id_j** are tf-idf values of word *i* on document *j* and zero values mean the word was absent in that document.

ESA: Word and Document Similarity

- Explicit Semantic Analysis is an approach used to find similarity (relatedness) between words and documents
- In ESA, TDM is used and typically, the entries of the vectors are the tf-idf weights (relative importance) of words on documents.
- The documents or words in the TDM can be expressed as vectors and the similarity between words or documents can be computed using:
 - $$\text{sim}(\text{vec1}, \text{vec2}) = \frac{\text{vec1} \bullet \text{vec2}}{|\text{vec1}| * |\text{vec2}|}$$

ESA: Documents as Vectors

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1

- The documents can be represented as column vectors:
 - $d1 = (1, 1, 0, 0, 1, 1)$
 - $d2 = (0, 1, 1, 0, 1, 1)$
 - $d3 = (1, 0, 0, 0, 0, 1)$
- Ones and zeros are used for simplicity. Practically, the actual tf-idf values should replace the ones.

ESA: Documents as Vectors

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1

- **Document similarity:** as already noted, document similarity can be analyzed using cosine similarity formula.
- Larger cosine values (lesser angles) indicates more similarity, which implies most words in one document appear in the other document, hence similarity in context. $\text{sim}(\text{vec1}, \text{vec2}) = \cos 90 = 0$



ESA and Document Similarity

- $d1 = (1, 1, 0, 0, 1, 1)$
- $d2 = (0, 1, 1, 0, 1, 1)$
- $d3 = (1, 0, 0, 0, 0, 1)$
- Given that the documents are represented as vectors, explicit semantic analysis finds the relationship between the document vectors.
- The idea is that, documents with similar corresponding vector entries are similar. So the cosine similarity is used to find the angle between the vector. Similar vectors have a smaller angle or “distance” between them.

ESA and Document Similarity

- $d1 = (1, 1, 0, 0, 1, 1)$
- $d2 = (0, 1, 1, 0, 1, 1)$
- $d3 = (1, 0, 0, 0, 0, 1)$
- Given the vectors representing documents in a corpus, the similarity between documents can be found using **cosine similarity**.
- The similarity between two vectors or documents could be found using the formula:
$$\text{sim}(\text{vec1}, \text{vec2}) = \frac{\text{vec1} \bullet \text{vec2}}{|\text{vec1}| |\text{vec2}|}$$
- Note that 1-cosine similarity is called **cosine distance**.

ESA: Words as Vectors

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1

- A word vector is a row vector with entries corresponding to the tf-idf weights of words across the documents.
- We can draw conclusions on the **relatedness** of two words based on their cosine similarity value, which captures similarity in the meaning of the two words

ESA: Word Vectors and Similarity

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1

- **Word similarity:** note that the meaning of a word can be inferred through the document or concept it contributes to.
- Terms or words with close meaning tend to be located together (appear across the same documents)



Using ESA for Document Classification

- So, ESA can be used for text or document classification. Mostly, content-based classification.
- If we have documents that are labelled or indexed in the database, we can compare a new document to the database documents, then classify the new document.
- Sometimes the documents are classified according to their attributes including genre, document type, author, year, etc.



Using ESA for Information Retrieval

- ESA is used for information retrieval.
- The Term Document Matrix of the corpus is used as the knowledge base.
- A new document is then used to query the knowledge base to find related documents.
- Basically, the new document is compared to all the documents in the knowledge base using the cosine similarity.
- Cosine similarity values range from 0 to 1 given that the vectors had positive entries.



Using ESA for Information Retrieval

- When a new document has been tokenized and is now represented as a list of words, we can convert this list of words into a vector by summing or averaging all the word vectors in the database for each of the words in the new document.
- The entire English Wikipedia is the knowledge base corpus , so the assumption here is that the words in the new documents are also in the database or TDM matrix.

Using ESA for Information Retrieval

- If new document was
- $\text{new_doc} = (w1, w2, w3)$
- Then the word vector for this new document will be: $\text{new_vector} = w1 + w2 + w3$
- This can be done in code by initializing the new vector with zero entries and dimensions = length of a row. Augmented statements can be used to add the vectors from the TDM

	d1	d2	d3
w1	1	0	1
w2	1	1	0
w3	0	1	0
w4	0	0	0
w5	1	1	0
w6	1	1	1



Using ESA for Information Retrieval

- Since the dimensions of the word vector correspond to documents, we can use the dimension with the highest tf-idf score to retrieve the document corresponding to this score as the related document
- On the other hand, we can find similar words by iterating through the rows of TDM and computing the cosine similarity of each row or word vector and the new word vector. The row in the database with highest cosine similarity represent the most similar words

Implementing Term Document Matrix

use the TFIDF matrix and the vocabulary size to get the Term Document Matrix

```
1 def term_document_matrix(TFIDF, word_list, word_dict):
2     # how many unique words are in the text
3     vocabulary_size = len(word_dict)
4     number_documents = len(TFIDF)
5
6     # create a numpy array with shape being
7     # vocab size for row num and number of doc for cols
8     TDM = np.zeros((vocabulary_size, number_documents))
9
10    # select each document in TFIDF
11    for doc_num in range(number_documents):
12        document = TFIDF[doc_num]
13
14        # for each word in the doc, get the
15        # int label from word_dict
16        for word in document.keys():
17            int_label = word_dict[word]
18
19            # get the value of TFIDF and
20            # assign it to the positions
21            # determined by int_label and doc_num
22            TDM[int_label, doc_num] = document[word]
23
24    return TDM
```

Note that the term document matrix can have words appearing in one document and not in the other. Wherever the word did not appear, there would be a numpy zero value

```
1 TDM = term_document_matrix(TFIDF, word_list, word_dict)
2 print((f"the dataset has {TDM.shape[0]} unique words and " +
3       f"{TDM.shape[1]} documents"))
```

the dataset has 927 unique words and 200 documents

```
1 TDM.shape
```

(927, 200)

```
1 # in the term document matrix, the rows are the words
2 len(word_list)
```

927

Term Document Matrix

```
1 TDM[:,10,:5]
```

```
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., 0., 5.29831737, 0., 0.],
       [0., 0., 0., 0., 0.],
       [0., 0., 0., 0., 0.]])
```

Implementing ESA

```
1 # define a function that finds related documents
2 # to the new document
3
4 def find_related_docs(tweet, TDM):
5     # create a vector representing the new document
6     vec_dim = TDM.shape[1] # dimension of vector
7     # fill the vectors with zeroes
8     # this would be a row vector representing the word
9     new_vector = np.zeros(vec_dim)
10
11     for word in tweet:
12         # get the integer label or position of word
13         # the position matches with the row position of word on TDM
14         position = word_dict[word]
15         # extract rows with the words in tweets and add them
16         # the entries with the highest tf-idf
17         # corresponds to the most related document
18         new_vector = new_vector + TDM[position, :]
19
20     # label the tf-idfs of the new vector with integers
21     # then sort the tf-idf entries of the new vector
22     doc_list = sorted(zip(range(vec_dim), new_vector),
23                       key=lambda x:x[1], reverse=True)
24     return doc_list
```

```
1 # the most related document has the highest total
2 # for tf-idf. This is if all or most words in the
3 # new document appear in a particular document, that document
4 # will have values of tf-idf leading to higher tf-idf totals
5
6 related_docs = find_related_docs(new_tweet, TDM)
7 related_docs[:10]
```

```
[(17, 6.725433722188182),
 (34, 6.725433722188182),
 (2, 5.298317366548036),
 (166, 3.506557897319982),
 (167, 3.506557897319982),
 (190, 3.506557897319982),
 (191, 3.506557897319982),
 (154, 3.2188758248682006),
 (164, 3.2188758248682006),
 (165, 3.2188758248682006)]
```

Documents Related to Search Query

```
1 # print the top 5 documents related to the
2 # new document used for search query
3
4 print("Relatiated Documents:")
5
6 for int_label, total_tfidf in related_docs[0:5]:
7     print("-----")
8     print(int_label, tweets[int_label])
```

Relatiated Documents:

17 ['jh', 'hines', 'staff', 'newly', 'issued', '@apple', '#connected', 'macbook', 'ipad', 'mini', '#txed']

34 ['rt', '@tra_hall', 'jh', 'hines', 'staff', 'newly', 'issued', '@apple', '#connected', 'macbook', 'ipad', 'mini', '#txed',
'#txed']

2 ['#aapl:5', 'rocket', 'stocks', 'buy', 'december', 'gains', 'apple']

166 ['ipad', 'mini', 'unboxing', 'via', '@youtube', '@apple', '#ipadmini', '#ipad', '#macbook', '#macbookpro', '#startup', '#h
ipster', '#unboxing']

167 ['ipad', 'mini', 'first', 'time', 'startup', 'via', '@youtube', '@apple', '#ipadmini', '#ipad', '#macbook', '#macbookpro',
'#startup', '#hipster', '#unbox']

Alternatively, Find Related Words to New Vector: Cosine Similarities

```
1 # Alternatively, use similarity cosine to find related documents
2
3 def similarity(vec1, vec2):
4     """
5     compute cosine similarity between vectors
6     """
7     sim = np.dot(vec1, vec2)
8     norm1 = np.sqrt(np.dot(vec1, vec1))
9     norm2 = np.sqrt(np.dot(vec2, vec2))
10    cosine_sim = sim/(norm1*norm2)
11
12    return cosine_sim
13
14 def find_similar_words(tweet, TDM):
15     vec_dim = TDM.shape[1] # dimension of vector
16     # initialize new vector for new tweet
17     new_vector = np.zeros(vec_dim)
18
19     # create new vector (word vector) for new tweet
20     for word in tweet:
21         position = word_dict[word]
22         new_vector = new_vector + TDM[position, :]
23
24     # similarity between new vector and each word vector
25     sim_words = [similarity(new_vector, TDM[i, :]) for i in range(TDM.shape[0])]
26     sim_words_sorted = sorted(zip(range(TDM.shape[0]), sim_words),
27                               key=lambda x:x[1], reverse=True)
28
29     return sim_words_sorted
30
```

```
1 similar_words = find_similar_words(new_tweet, TDM)
2 similar_words[:10]
```

```
[(593, 0.7640071602554539),
 (346, 0.7398422630773179),
 (518, 0.6849608378181927),
 (863, 0.66219576837537),
 (526, 0.6273056189756162),
 (581, 0.6273056189756162),
 (844, 0.6273056189756162),
 (381, 0.6273056189756161),
 (460, 0.6273056189756161),
 (537, 0.6273056189756161)]
```

Words that are related
occur together across
documents

Related Words

```
1 similar_words = find_similar_words(new_tweet, TDM)
2 similar_words[:10]
```

```
[(593, 0.7640071602554539),
 (346, 0.7398422630773179),
 (518, 0.6849608378181927),
 (863, 0.66219576837537),
 (526, 0.6273056189756162),
 (581, 0.6273056189756162),
 (844, 0.6273056189756162),
 (381, 0.6273056189756161),
 (460, 0.6273056189756161),
 (537, 0.6273056189756161)]
```

```
1 # print similar words
2 for int_label, sim_score in similar_words[0:10]:
3     print(word_list[int_label], sim_score)
4
```

```
macbook 0.7640071602554539
mini 0.7398422630773179
ipad 0.6849608378181927
#macbook 0.66219576837537
#startup 0.6273056189756162
@youtube 0.6273056189756162
#hipster 0.6273056189756162
newly 0.6273056189756161
jh 0.6273056189756161
issued 0.6273056189756161
```


Document Clustering

- We can group documents based on their content. If the documents have similar contents, we can group them together.
- Before we group documents we need to specify how we want to measure the distance between them. – use cosine similarity

- Using the cosine similarity:

$$\text{sim}(\vec{u}, \vec{v}) = \frac{\vec{u} \cdot \vec{v}}{||\vec{u}|| ||\vec{v}||}$$

a simple approach is to impose a minimum similarity cut off and any value above the cut off implies the documents in question are similar



Document Clustering

- This is how the algorithm works:
- Every document is assigned to its own cluster
- The matrix similarity are calculated.
- The two most similar vectors are merged.
- Similarity is recalculated for the reduced set of clusters until a desired number of clusters is reached.

Latent Semantic Analysis (LSA)

- LSA, sometimes called latent semantic indexing, is a data reduction technique that uses singular value decomposition to extract latent dimensions from the data, which are compact representation of the

Symmetric S

$$S = Q\Lambda Q^T = \lambda_1 \mathbf{q}_1 \mathbf{q}_1^T + \lambda_2 \mathbf{q}_2 \mathbf{q}_2^T + \cdots + \lambda_r \mathbf{q}_r \mathbf{q}_r^T$$

Any matrix A

$$A = U\Sigma V^T = \sigma_1 \mathbf{u}_1 \mathbf{v}_1^T + \sigma_2 \mathbf{u}_2 \mathbf{v}_2^T + \cdots + \sigma_r \mathbf{u}_r \mathbf{v}_r^T$$

$$\begin{matrix} M \\ m \times n \end{matrix} = \begin{matrix} U \\ m \times m \end{matrix} \begin{matrix} \Sigma \\ m \times n \end{matrix} \begin{matrix} V^\dagger \\ n \times n \end{matrix}$$

where

U is a unitary matrix ($UU^\dagger = 1$)

Σ is diagonal with non-negative elements

V^\dagger is a unitary matrix



LSA and Latent Dimensions

- With singular value decomposition, the term-document matrix is decomposed into components and the first few components (k dimensions) with the maximum variance are extracted such that the variance in the data is preserved.
- The k dimensions selected represent k latent topics or dimensions of the underlying dataset or TDM
- Only the most significant dimensions are kept to get rid of the noise. We want to represent the most information in the dimensions we select.

LSA and SVD

- In LSA, singular value decomposition (SDV) is applied to the term-document matrix to obtain a SVD term matrix U , the diagonal matrix Σ and SDV document matrix V .

Example of Term Document Matrix

		d_1	d_2	d_3	d_4	d_5	d_6	
	ship	1	0	1	0	0	0	
	boat	0	1	0	0	0	0	
	ocean	1	1	0	0	0	0	
	voyage	1	0	0	1	1	0	
	trip	0	0	0	1	0	1	

<https://nlp.stanford.edu/IR-book/html/htmledition/latent-semantic-indexing-1.html>

LSA and SVD

SVD term matrix (U)

	1	2	3	4	5
ship	-0.44	-0.30	0.57	0.58	0.25
boat	-0.13	-0.33	-0.59	0.00	0.73
ocean	-0.48	-0.51	-0.37	0.00	-0.61
voyage	-0.70	0.35	0.15	-0.58	0.16
trip	-0.26	0.65	-0.41	0.58	-0.09

Diagonal matrix

2.16	0.00	0.00	0.00	0.00
0.00	1.59	0.00	0.00	0.00
0.00	0.00	1.28	0.00	0.00
0.00	0.00	0.00	1.00	0.00
0.00	0.00	0.00	0.00	0.39

SVD document matrix

	d_1	d_2	d_3	d_4	d_5	d_6
1	-0.75	-0.28	-0.20	-0.45	-0.33	-0.12
2	-0.29	-0.53	-0.19	0.63	0.22	0.41
3	0.28	-0.75	0.45	-0.20	0.12	-0.33
4	0.00	0.00	0.58	0.00	-0.58	0.58
5	-0.53	0.29	0.63	0.19	0.41	-0.22

Now, we can multiply the matrices and consider only a few, say two singular values. That means we are selecting only the first two components. It is like “zeroing” out the rest of the singular values in the diagonal matrix after the second one, then multiply the matrices. So the term document matrix is approximated using k-dimensional singular values

LSA: projecting a new document to a singular space

- Documents and words can be easily classified in the singular space.
- A new document (or query) is mapped or projected into the singular space using the transformation:
$$\hat{v} = \Sigma_k^{-1} U_k^\dagger v$$
- v is the original document
- \bar{v} is the transformed version of the original document
- U_k and Σ_k are the truncated versions of U and Σ . When we select $k=2$ components, all other rows after the 2 rows are zeroed out to get the truncated versions.

Examples of Truncated Σ_k

2.16	0.00	0.00	0.00	0.00
0.00	1.59	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00
0.00	0.00	0.00	0.00	0.00

- If $k=2$, only the first two rows of the matrices are maintained. The truncated version therefore has the rest of the rows filled with zeros.

Implementing LSA

- Most packages such as sklearn have the svd feature for decomposing a matrix into its three components: In LSA, SDV needs to be first applied to the TDM:

```
u, sigma, vt = svd(TDM)
```

Shapes of Matrices

- `m, n = TDM.shape`
- `print(sigma.shape)`
- `print(vt.shape)`
- Note that: $TDM = U\Sigma V^T$
- `sdv` can be imported as follows:
- `from scipy.linalg import svd`

Implementing SVD

- When we decompose the term document matrix, we get

$$U (m \times m) \cdot \text{Sigma} (n \times n) \cdot V^T (n \times n)$$

- But we actually require:

$$U (m \times m) \cdot \text{Sigma} (m \times n) \cdot V^T (n \times n)$$

- So we can reconstruct the diagonal matrix into an $m \times n$ matrix
- We could decide to keep sigma as $n \times n$ and reconstruct u to $m \times n$. So we can have $m \times n$, $n \times n$, $n \times n$

LSA Implementation in Python

Latent Semantic Analysis

This is the extraction of latent dimensions or topics from a singular decomposition of the term-document matrix

```
1 # starts with the term-document matrix  
2 TDM
```

```
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.]])
```

Implement SVD

```
1 # decompose the term-document matrix
2 u, sigma, vt = svd(TDM)
```

```
1 # print the svd term matrix
2 u
```

```
array([[ 7.36806317e-05,  3.40770148e-04, -7.77823847e-03, ...,
         1.05289194e-02, -7.46998095e-05,  6.69916034e-04],
       [ 3.44239397e-05,  3.98821846e-05, -2.23763614e-05, ...,
         3.75624127e-04,  1.25668498e-04,  8.11037736e-05],
       [ 7.07206337e-03,  8.57493541e-04, -1.24079164e-02, ...,
         1.04724698e-03,  5.16283990e-04,  6.36249565e-04],
       ...,
       [ 3.18213574e-04,  1.54391576e-03, -1.40497485e-02, ...,
         8.24330670e-01, -6.80895387e-05,  1.20092051e-04],
       [ 5.63686605e-05,  4.49917970e-04, -1.57398765e-03, ...,
        -2.17853697e-03,  8.78420001e-01, -1.53654818e-04],
       [ 4.61032665e-05,  6.28857248e-05, -4.27434654e-05, ...,
         1.25491796e-04,  9.38039129e-03,  9.17042436e-01]])
```

```
1 # print sigma
2 sigma[:5]
```

```
array([34.12143118, 28.69781728, 25.40239066, 24.53357897, 22.68779674])
```

```
1 # print the svd document matrix
2 vt
```

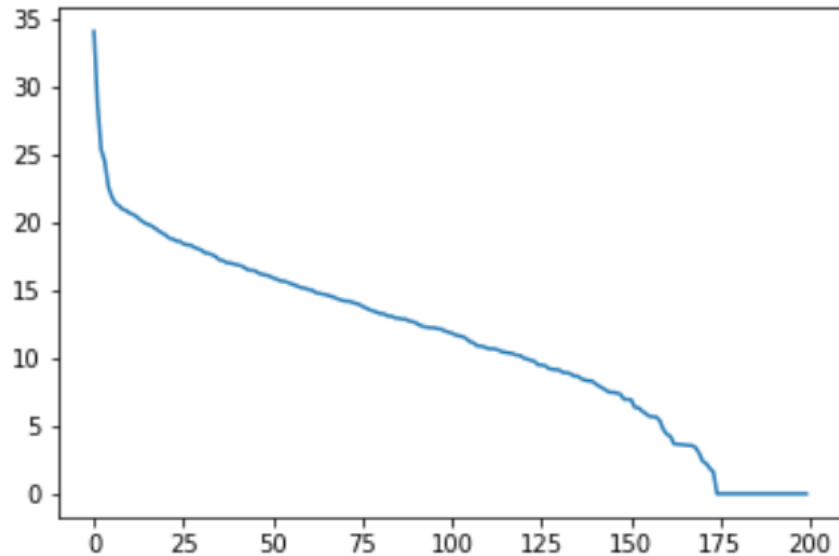
```
array([[ 2.33523602e-04,  2.10340123e-04,  1.88044005e-04, ...,
         1.57925095e-03,  2.19588538e-03,  7.71782278e-04],
       [ 2.33069252e-04,  1.31105677e-03,  2.16875708e-03, ...,
         1.70972476e-02,  1.20388328e-01,  5.43102303e-03],
       [-1.18582648e-04, -2.57411317e-02, -2.27348622e-02, ...,
        -1.45939720e-03, -3.54534831e-02, -1.02066596e-01],
       ...,
       [ 0.00000000e+00, -9.73800447e-17, -1.32165671e-16, ...,
         1.14491749e-16,  5.20417043e-17,  1.73472348e-17],
       [ 0.00000000e+00,  1.86049093e-16,  2.06215253e-16, ...,
         5.81132364e-17, -4.51028104e-17, -3.46944695e-18],
       [ 0.00000000e+00,  8.84708973e-17,  6.93889390e-17, ...,
        -9.71445147e-17, -5.37764278e-17, -1.38777878e-17]])
```

```
1 # print the shapes of the matrices
2 print(TDM.shape, u.shape, sigma.shape, vt.shape)
```

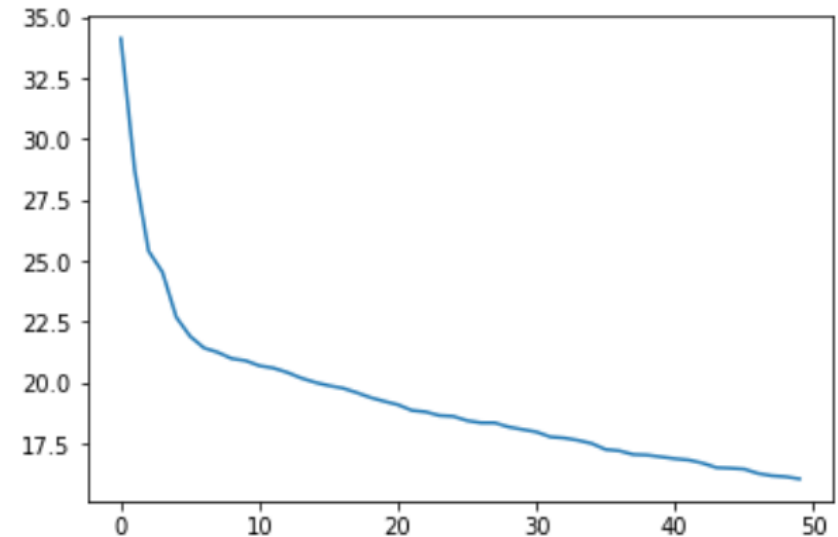
```
(927, 200) (927, 927) (200,) (200, 200)
```

How Many Dimensions Should be Extracted?

```
1 # it is difficult to see the precise "elbow" value
2 # when there are several values of sigma
3 plt.plot(sigma)
4 plt.show()
```



```
1 plt.plot(sigma[:50])
2 plt.show()
```



The graph shows that $k=5$ would be appropriate, from the "elbow" value

Creating Truncated Matrices with K Dimensions

```
1  # from the shape of sigma, it can be noticed that sigma is a vector,  
2  # sigma is supposed to be a matrix mxn, we can reconstruct it if we wanted  
3  # Let's go straight to converting sigma into a truncated diagonal matrix  
4  
5  k = 5  
6  sigma_k = sigma[:k]  
7  sk = np.diag(sigma_k)  
8  sk  
9  # this is nxn matrix
```

```
array([[34.12143118,  0.,          0.,          0.,          0.],  
       [ 0.,          28.69781728,  0.,          0.,          0.],  
       [ 0.,          0.,          25.40239066,  0.,          0.],  
       [ 0.,          0.,          0.,          24.53357897,  0.],  
       [ 0.,          0.,          0.,          0.,          22.68779674]])
```

```
1  # we need a 927x5 for uk to keep only k-dimensions  
2  uk = u[:, :k]
```

```
1  # we need 5x200 for vtk to keep only k-dimensions  
2  vtk = vt[:, :5, :200]
```

```
1  print(uk.shape, sk.shape, vtk.shape)
```

```
(927, 5) (5, 5) (5, 200)
```

View a few Unique Words in the Dataset

```
1 # print the word list, just a few words out of 927
2 print(word_list[:20])
```

```
['best-designed', 'well', 'broadway', '#iwantjobsback', 'th
ree', 'increases', 'thankful', 'stocks', 'wilson', 'competi
tion', 'charging', 'im', 'bloomberg', '@baba', 'analyst',
'music', 'iphones', 'shareholders', 'reveal', 'self']
```

```
1 len(word_list)
```

```
927
```

Words that Contribute Most, to a Topic

Find the Meaning of an Extrated Topic

When a topic or a vector (column vector) is obtained from the extracted SVD term matrix, we can define what that topic is by mapping the vector (whose entries are latent tf-idf scores) to the wordlist (unique words in the text in the order in which we have them on the term-document matrix) then sorting the list by the latent scores. We want to identify the words that contribute most to a particular topic

```
1 # map the entries of the latent scores from the
2 # topic vector to their corresponding words
3 def top_words(wordlist, topic_vector):
4     doc_list = sorted(zip(wordlist, topic_vector),
5                        key=lambda x:x[1], reverse=True)
6
7     return doc_list
8
```

```
1 # let's find the words that contribute most to topic 1
2 # this will help us define topic 1
3 topwords = top_words(word_list, uk[:,0])
4 for word, weight in topwords[:10]:
5     print(word, ":", weight)
```

```
iphones : 0.2987441726579015
federal : 0.2987441726579015
cites : 0.2987441726579015
department : 0.2987441726579015
century : 0.2987441726579015
unlock : 0.2987441726579015
law : 0.2987441726579015
18th : 0.2987441726579015
justice : 0.2942978100155947
@thehill : 0.26993335660102
```


Transform a Given Document into the Singular Space

```
1 # given a sentence or a word document,
2 # Let's define the vector in a regular word space
3 document = ["ipad", "iphone", "mini"]
4 doc_vector = np.zeros(len(word_list))
5
6 # view the first 10 entries
7 doc_vector[:10]
```

```
array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0.])
```

```
1 for word in document:
2     # extract the index of the word
3     integer_index = word_dict[word]
4     # update the entry corresponding to
5     # the integer index
6     doc_vector[integer_index] = 1
7
8     # view a portion of the vector
9 doc_vector[:200]
```

```
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., 0., 0., 0., 0., 0., 0., 0., 0., 0., 0.,  
       0., 0.]
```

Let's project the vector into a singular or latent space

$$\hat{v} = \Sigma_k^{-1} U_k^\dagger v$$

```
1 # sigma_k^-1 is a singular vector
2 sigma_k
```

```
array([34.12143118, 28.69781728, 25.40239066, 24.53357897, 22.687
79674])
```

```
1 1/sigma_k
```

```
array([0.02930709, 0.03484586, 0.03936637, 0.04076046, 0.04407656])
```

1 | uk.T

```
array([[ 7.36806317e-05,  3.44239397e-05,  7.07206337e-03, ...,
         3.18213574e-04,  5.63686605e-05,  4.61032665e-05],
       [ 3.40770148e-04,  3.98821846e-05,  8.57493541e-04, ...,
         1.54391576e-03,  4.49917970e-04,  6.28857248e-05],
       [-7.77823847e-03, -2.23763614e-05, -1.24079164e-02, ...,
        -1.40497485e-02, -1.57398765e-03, -4.27434654e-05]).
```

Transform a Given Document into the Singular Space

```
1  # the projected document would be  
2  singular_doc = 1/sigma_k*np.dot(uk.T, doc_vector)  
3  singular_doc
```

```
array([ 0.00017605,  0.01445142,  0.00050387, -0.0004712 , -0.00121434])
```

Get the Topic Most Related to the Document

Topic most relevant for the document

```
1  # the entries of the singular vector are on specific dimensions
2  # in the singular space. We want to find the dimension with the largest
3  # value, so we use np.argmax which returns the dimension with the largest value
4  topic = np.argmax(singular_doc)
5  topic
```

Retrieving Top Words

Retrieve the top words for the topic related to the document

```
1 topwords = top_words(word_list, uk[:,topic])  
2 for word, weights in topwords[:10]:  
3     print(word, ":", weight)
```

```
#macbook : 0.26993335660102  
#startup : 0.26993335660102  
@youtube : 0.26993335660102  
#hipster : 0.26993335660102  
#macbookpro : 0.26993335660102  
via : 0.26993335660102  
mini : 0.26993335660102  
ipad : 0.26993335660102  
#ipadmini : 0.26993335660102  
startup : 0.26993335660102
```



ESA and LSA

- ESA uses a corpus, then defines the meaning of a word based on the documents the word contributes the most.
- Original Tf-idf weights are used from the TDM to find words in the corpus related to a query document or find similarity between query doc and docs in corpus
- LSA finds the underlying representation or latent topics in a text. The latent topics can then be labelled based on the words that contribute most to the topics
- A query document can be mapped to the latent space to see which topic is relevant to the document.