# ISTANBUL TECHNICAL UNIVERSITY

# COMPUTER ENGINEERING DEPARTMENT

## BLG 202E

## NUMERICAL METHODS
## PROJECT 2 REPORT

**CRN**   : 21333

**LECTURER**   : Enver Özdemir

**PROJECT NO** : 2

**STUDENT**   : Madina Alzhanova 150220939

## SPRING 2024

# Contents

# 1   INTRODUCTION

This project aims to implement Latent Semantic Analysis (LSI) by using Singular Value Decomposition (SVD). LSI is a method for indexing and retrieving information from a big data set. SVD represents the matrix in form of USVt. U holds information about terms, V about documents, and S is matrix for singular values. We reduce the dimension of original matrix by SVD, then perform search of queries using the most optimal truncation with minimum errors.

# 2   MATERIALS AND METHODS

## 2.1   Dataset

The dataset is a large .csv document containing 4 columns: author, posted on date, rating and text. Since information posted prior to 2009 is not taken into account, they were deleted. Number of rows decreased from about 6189 to 5144 in VSCode editor using pandas code as provided below. It deletes empty rows as well as dates prior to 2009 as shown in code below. Document is a row of "text" column in the given .csv file. Since dataset is very large, running time of code is very long.

```
13    df = pd.read_csv('comcast_consumeraffairs_complaints.csv')  #read the file, create df(table)
14    df = df.dropna(subset=['text']) #delete null rows
15    df['posted_on'] = pd.to_datetime(df['posted_on'],errors='coerce',format='mixed')
16    df = df[df['posted_on'].dt.year >= 2009] #delete less than 2009
17    df.to_csv('comcast_consumeraffairs_complaints.csv', index=False) #update csv
```

Figure 1: csv file editing

## 2.2   Term-by-Document Matrix

I used numpy, pandas, nltk libraries in order to create term-by-document. It is shown in figure below:

```
1     import numpy as np
2     import pandas as pd
3     from nltk.tokenize import word_tokenize
4     from nltk.stem import PorterStemmer
5     from nltk.corpus import stopwords
6     import nltk
```

Figure 2: libraries

During preprocessing, I performed tokenization, removed stopwords, stemming, and lemmatization. so now tokens contain only one term for resembling words and non-

meaningful words(stopwords) are also not counted. This significantly reduces the size of term index matrix, thus reducing complexity of all project.

```
19   #preprocessing text function
20   def preprocess_text(text):
21       tokens = word_tokenize(text.lower())    #lowercasing
22       tokens = [word for word in tokens if word.isalpha()]    #check alphabetic
23       stop_words = set(stopwords.words('english'))    #create stopwords
24       tokens = [word for word in tokens if word not in stop_words]    #remove stopwords
25       stemmer = PorterStemmer()
26       tokens = [stemmer.stem(word) for word in tokens]
27       return tokens
28
29   df['tokens'] = df['text'].apply(preprocess_text)
```

Figure 3: preprocessing

After preprocessing text, and creating tokens, I create and fill term-index and term document matrices as in following code. Term index matrix maps each term to one id(index). I create an empty matrix, then fill it by each cycle of loop, with elements in tokens, so that they are unique. Term document matrix shows if document has a term in it.

```
31   #index of each term into matrix
32   term_index = {}
33   for tokens in df['tokens']:
34       for token in tokens:
35           if token not in term_index:
36               term_index[token] = len(term_index)
37
38   # # no of unique terms
39   # print("number of unique terms:", len(term_index))
40
41   # create td matrix and fill
42   td_matrix = np.zeros((len(term_index), len(df)))
43   for doc_id, tokens in enumerate(df['tokens']):
44       for term in tokens:
45           term_id = term_index.get(term, -1)
46           if term_id != -1:
47               td_matrix[term_id, doc_id] += 1
```

Figure 4: term-index and term-doc

## 2.3 SVD implementation

SVD is used to decompose matrix into form of USVt, so it can be reduced in size. U is an m×m orthogonal matrix whose columns are the left singular vectors of A. S is an m×n diagonal matrix whose non-zero elements are the singular values of A. Vt is an n×n orthogonal matrix whose rows are the right singular vectors of A. Since we are restricted from using SVD libraries, I implement SVD using Power method. Generally, Power method is a method to find a dominant(biggest) eigenvector and largest eigenvalue. In my code, it starts with a random vector, and by iterating b=Ab/norm many times, it eventually provides approximate eigenvector.

2

Next function does decomposition using previous power method function. It finds most significant eigenvector and eigenvalue many times iterating num-singular-values times and fills U, S and Vt matrices. Elements of U are calculated by Av/s. The function returns U, S and Vt matrices, so decomposes original matrix.

```python
58    def power_method(matrix, num_iterations=100):
59        b = np.random.rand(matrix.shape[1]) #random 1D vector
60        for _ in range(num_iterations): #converges as iterates
61            b = np.dot(matrix, b)   #A @ b
62            b /= np.linalg.norm(b)  #normalize
63        return b    #eigvector for
64
65    def svd_power_method(matrix, num_singular_values):
66        m, n = matrix.shape
67        U = np.zeros((m, num_singular_values))
68        S = np.zeros(num_singular_values)
69        Vt = np.zeros((num_singular_values, n))
70        matrix_copy = matrix.copy()  #to preserve original matrix
71        for i in range(num_singular_values):
72            v = power_method(matrix_copy.T @ matrix_copy) #input to power method func is
73            sigma = np.linalg.norm(matrix_copy @ v)
74            u = (matrix_copy @ v) / sigma    #Av/s
75            matrix_copy -= np.outer(u, v) * sigma
76            U[:, i] = u #shoten
77            S[i] = sigma
78            Vt[i, :] = v
79        return U, S, Vt
```

Figure 5: SVD using power method

## 2.4   Question 1

Question 1 aims to create a lower-rank approximation of original matrix with minimum loss. Lower-rank approximation is done by selecting most significant k columns of V matrix(most significant k eigenvectors), and creating new matrix called A-hat.
k is found using two methods: mean square error method and Frobenius Norm method. k is said to be optimal when MSE and FN are minimal, thus A-hat is closest to original matrix A. During iteration for both methods, k is between 10 and min(t,d)/10 and increments by 20. Formulas for mean square error and Frobenius Norm are given below.

$$MSE = \tfrac{1}{td}\Sigma_{i,j}(A_{ij} - \hat{A}_{ij})^2$$

$$||A - \hat{A}||_F = \sqrt{\Sigma_{i,j}(A_{ij} - \hat{A}_{ij})^2}$$

Figure 6: MSE and FN formula

According to that formula, the function of calculating MSE and FN is created as shown in code. It calculates MSE and FN for different value of k, and appends result into an array. Returns arrays containing all values of MSE and FN.

3

```
81  def calculate_mse_and_fn(matrix, U, S, VT, k_values):
82      mse_results = {}
83      fn_results = {}
84      for k in k_values:
85          U_k = U[:, :k]   #shorten
86          S_k = S[:k]   #shorten
87          VT_k = VT[:k, :]   #shorten
88          A_hat = U_k @ np.diag(S_k) @ VT_k   #low-approx Ahat
89          mse = np.mean((matrix - A_hat) ** 2)     #formula
90          fn = np.linalg.norm(matrix - A_hat) #formula
91          mse_results[k] = mse     #fill array
92          fn_results[k] = fn   #fill array
93      return mse_results, fn_results
94
95  def find_optimal_k(mse_results, fn_results):
96      optimal_k_mse = min(mse_results, key=mse_results.get)     #minimal error
97      optimal_k_fn = min(fn_results, key=fn_results.get)
98      return optimal_k_mse, optimal_k_fn
```

Figure 7: MSE and FN functions, optimal k function

And then using the given arrays, by finding minimal MSE and minimal FN, and its index(corresponding k), optimal k is calculated and returned.

## 2.5   Question 2

Question 2 requires implementing cosine similarity function, and finding most relevant document to given queries using it. When q is query vector, and d is document vector, similarity function is defined as below:

$$\textbf{Equation.10: } sim(q,d) = \frac{\vec{q} \cdot \vec{d}}{||q|| \cdot ||d||}$$

Figure 8: Cosine similarity formula

```
140  def cosine_similarity(v1, v2):
141      norm_v1 = np.linalg.norm(v1)     #normailze q
142      norm_v2 = np.linalg.norm(v2)     #normailze v
143      if norm_v1 == 0 or norm_v2 == 0:
144          return 0  # Return 0 similarity if either vector has zero length to avoid division by zero
145      return np.dot(v1, v2) / (norm_v1 * norm_v2) #formula
146
147  def query_vector(query, term_index, U, S, k):
148      q = np.zeros(len(term_index))
149      for word in query:
150          if word in term_index:
151              q[term_index[word]] = 1 #create a query vector so it has 1 where it has terms for term-index matrix
152      return q @ U[:, :k] @ np.linalg.inv(np.diag(S[:k]))
153
154  def find_most_relevant_document(query, term_index, U, S, VT, k):
155      q_vec = query_vector(query, term_index, U, S, k)
156      VT_reduced = VT[:k, :]   # Here VT should have rows up to k
157      document_similarities = [cosine_similarity(q_vec, VT_reduced[:, i]) for i in range(VT_reduced.shape[1])] #array with cosine values
158      return np.argmax(document_similarities),df.iloc[np.argmax(document_similarities)]['text'] #the more cos -> the more similar
```

Figure 9: cos-similarity, query-vector and finding document functions

In the code, i have 3 functions: cosine-similarity, query-vector and find-most-relevant-document. First, i normalize q and v vectors, then apply the formula given in Figure 8. Then in second function, create a query vector, same as term-document matrices'. So i fill 1-s where it has elements of term-index matrix.
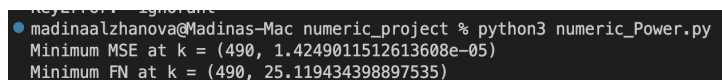
Third function calls second function, and creates query vector. Then second argument

4

given as d in equation is row of reduced Vt. Given q and d, we call first function, and get array or cosine values. We take maximum of it, because bigger cosine value corresponds to bigger similarity. So, third function returns index of maximum similarity and its content. After that, we call the function according to the given queries, and wait for the code to run.

# 3 OUTPUTS

## 3.1 Question 1

According to results of Question 1, optimal k is same for MSE and FN methods. It should be 490 after .csv file has been modified. MSE = 1.4249011512613608e-05 whereas FN = 25.119434398897535, which shows that SVD approximated the matrix with a little loss.



Figure 10: Results for Question 1

## 3.2 Question 2

Outputs of Question 2 say that query1's relevant document is at index 8, query2's is 4912, query3 and query4's is 1030 in modified dataset. The terminal screenshots are provided below.



Figure 11: Results for Question 2

In second question, the document has same words as in query, or other forms of that word or synonyms. There are some examples:

Document for query1 contains both overwhelming and ignorant.

Document for query2 contains "adapter" 4 times, also has "router" and etc.

Document for query3 contains lie(=liar), advertising(=promotion), crappy(=horrible) and etc.

Document for query4 contains intrenet, kindergartener(kindergarten), supervisor/managers(=clerk), crappy(=horrible) and etc.

# 4 DISCUSSION

In general, in this project I used latent semantic indexing method by using SVD. Before every process, it is crucial to reduce number of rows in dataset. Even reduced, it takes much time to run the code, so without getting rid of unnecessary texts, compile time will be too long. At first, I erased last several rows starting from 2008 year manually, but it was not enough, because file still contained empty rows and complains of older years randomly in file. So, using pandas library for this was efficient. ntlk library was used to perform preprocessing of text, and creation of term-index, term-document matrices.

In order to perform SVD I used Power method. Power method works by finding approximate eigenvector, and from that calculates eigenvalues. Having eigenvectors, we fill matrix V, with eigenvalues we calculate singualrvalues and fill diagonal matrix S. Using formula, we fill matrix U. Thus, I got SVD.

Question 1 is very useful for finding most optimal reduction size. It also can be used for compression of any text data, as well as images, so that it takes much less memory but has comparably same content. I used MSE and FN methods, and they gave same k, which means that original matrix of size about 5000 can be represented as size of 490. It reduces size more than 10 times, while its MSE is only 1.4249011512613608e-05(very negligible).

Question 2 works as a search engine. Provided queries, it finds most relevant document from large dataset. For example, google search works similarly. Searchbar gets queries, and it outputs webpages(documents) in order of relevance. In our case, we outputted only most relevant one. There are some terms in queries that are not found in a document, but it is counted as most relevant, because other terms fit the document very well. In our assignment, queries are given as terms, but they should be also preprocessed. If this code is used in larger scope, user's input must be checked for stopwords, then lemmatized, tokenized and etc.

# 5 CONCLUSION

In conclusion, one of efficient applications of SVD - LSI was implemented in this project. It can be used efficiently for reducing size of big dataset and to work as search engine. This project was difficult to implement because running the code takes much time, also I was not familiar with python and its libraries. But, the project was implemented successfully and results look correct.