

Configuration Manual

MSc Research Project MSc AI Top-up

Gabriel Amariei Student ID: 13130510

School of Computing National College of Ireland

Supervisor: Faithful Onwuegbuche

National College of Ireland

MSc Project Submission Sheet



School of Computing

Student Name:	Gabriel Amariei		
Student ID:	13130510		
Programme:	MSc Al Top-up	Year:	2024
Module:	MSc Research Project		
Lecturer:	Faithful Onwuegbuche		
Due Date:	12/08/2024		
Project Title:	Configuration Manual		

Word Count: xyz Page Count: 26

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

<u>ALL</u> internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

Signature:	Gabriel Amariei
Date:	12/08/2024

PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST

Attach a completed copy of this sheet to each project (including multiple copies)	
Attach a Moodle submission receipt of the online project submission, to each project (including multiple copies).	
You must ensure that you retain a HARD COPY of the project, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer.	

Assignments that are submitted to the Programme Coordinator Office must be placed into the assignment box located outside the office.

Office Use Only	
Signature:	
Date:	
Penalty Applied (if applicable):	

Configuration Manual

Gabriel Amariei Student ID: 13130510

Table of Contents

1	Introduction	2
2	Hardware Requirements	2
3	Software Reguirement	2
4	Package Reguirements	2
5	Data Download code	4
6	Initial exploratory data analysis	6
7	Convert doc and docx files to pdf	8
8	Extract text and metadata from pdf	8
9	Clean and vectorize the text	13
10) Check the right number of components and find the optimum n of clusters	18
11	Clustering implementation, evaluation and visualization	20

1 Introduction

This Configuration manual represents all the necessary requirements and details for the Project "*Document Clustering of Irish Government Circulars using Machine Learning Techniques*". It is structured as follows: the next section shows the hardware requirements followed by the software and library package requirements needed to run the project. The next section provides a description of the data and its location. The rest of the sections are divided into parts of workflow of the code, their execution, result and their evaluation.

2 Hardware Requirements

The configuration of the machine used for this project can be found below.

Hardware Model	Lenovo ThinkPad X1 Carbon 4th
Memory	16.0 GiB
Processor	Intel® Core™ i7-6600U CPU @ 2.60GHz × 4
Graphics	Mesa Intel® HD Graphics 520 (SKL GT2)
Disk Capacity	256.1 GB

3 Software Requirement

The operating system was Ubuntu.

OS Name	Ubuntu 22.04.1 LTS
ОЅ Туре	64-bit
GNOME Version	42.4
Windowing System	Wayland

Python 3.9.19 with the Jupyter lab development environment .

```
1 print (sys.version)
3.9.19 (main, May 6 2024, 19:43:03)
[GCC 11.2.0]
```

4 Package Requirements

The following python packages were imported:

import requests from bs4 import BeautifulSoup import os import re import pandas as pd from spire.doc import * from spire.doc import Document from spire.doc.common import * import os import numpy as np import pdfplumber import matplotlib import matplotlib.pyplot as plt from pylab import rcParams from sklearn.cluster import KMeans from sklearn.feature extraction.text import TfidfVectorizer from sklearn.decomposition import PCA from sklearn.manifold import TSNE from nltk.stem import WordNetLemmatizer from nltk.tokenize import RegexpTokenizer from nltk.corpus import stopwords from nltk import word_tokenize import nltk import re import string import warnings import random from sklearn.decomposition import TruncatedSVD from sklearn.pipeline import make_pipeline from sklearn.preprocessing import Normalizer from yellowbrick.cluster import KElbowVisualizer, SilhouetteVisualizer from sklearn.metrics import davies_bouldin_score from sklearn.metrics import silhouette samples import skfuzzy as fuzz

from sklearn.preprocessing import StandardScaler from tensorflow.keras.models import Sequential, Model from tensorflow.keras.layers import LSTM, Dense, Input, Dropout, BatchNormalization from tensorflow.keras.optimizers import Adam import tensorflow as tf from tensorflow.keras.layers import LayerNormalization

from sklearn.metrics import silhouette_score, davies_bouldin_score, calinski_harabasz_score from sklearn.metrics.cluster import normalized_mutual_info_score from sklearn.preprocessing import LabelEncoder

from sklearn.pipeline import make_pipeline from sklearn.preprocessing import Normalizer import plotly.express as px import plotly.graph_objects as go

import cufflinks as cf
from plotly.offline import iplot, init_notebook_mode
import gensim
from gensim.models import Word2Vec

from transformers import BertTokenizer, BertModel import torch import pickle

5 Data Download code

Step one: Download the individual links to the circulars File: **1.Download documents with dat info df.ipynb**

1. Download all individual page links

Step 2: Download pdf, doc and docx type documents

```
1 #create function to scrape data info
 3 def find_data(soup_tag, text):
         selected_tags = soup2.find_all(soup_tag)
abc_span = None
 4
 5
 6
         for tag in selected_tags:
           if text in tag.get_text(strip=True):
    # Remove the text from the tag's content
 7
 8
 9
                    abc_span = tag.get_text(strip=True).replace(text, '').strip()
10
                   break
         return abc_span
11
 1
    circulars info=pd.DataFrame(columns=('Link n','Circular n', 'page title', 'From', 'Published on',
                                                     'Effective_from',
'Effective_to', 'File_name', 'Sector', 'Author', 'Followed_by',
'Succeeds', 'Succeeded_by', 'Succeeds_L','Succeeded_by_L', 'url'))
 4
```

Iterate trough each link and download the documents and scrape further info for each circular from the website.

1	<pre>base_link='nttps://www.gov.le'</pre>
2	# Maximum allowed length for file name as some titles are too long to use as document name
3	MAX FILE NAME LENGTH = 70
4	
5	link n=0
6	
7	for item in circulars links:
8	url=base link+item
9	print(url)
10	# Requests URL and get response object
11	<pre>response2 = requests.get(url)</pre>
12	# Parse text obtained
13	<pre>soup2 = BeautifulSoup(response2.text, 'html.parser')</pre>
14	# Find all hyperlinks present on webpage
15	links = soup2.find all('a')
16	$\mathbf{i} = 0$
17	<pre>#page title = soup2.find('title')</pre>
18	<pre>page title = soup2.find('h1')</pre>
19	
20	#n of link included to avoid replacing documents with same truncated name
21	link_n=link_n+1
22	
23	Circular_n=find_data('span', 'Circular')
24	<pre>From=find_data('p', 'From')</pre>
25	Published_on=find_data('p', 'Published on')
26	Effective_from=find_data('p', 'Effective from')
27	Effective_to=find_data('p', 'Effective to')
28	Sector=find_data('p', 'Sector')
29	Author=find_data('p', 'Author')
30	Followed_by=find_data('p', 'Followed by')
31	Succeeds=find_data('p', 'Succeeds')
32	Succeeded_by=find_data('p', 'Succeeded by')
33	Succeeds_L=find_data('p', 'succeeds')
34	Succeeded_by_L=find_data('p', 'succeeded by')
35	

```
30
37
          for link in links:
38
                href=link.get('href')
39
                if href and href.endswith(('.pdf','.doc', '.docx')): # Check if href is not None before calling endswith
40
41
                      i += 1
42
                      print("Downloading file: ", i)
43
44
                      # Get response object for link
45
                      response = requests.get(link.get('href'))
46
47
                      cleaned_title= re.sub(r'[,\\/*?:"<>|]', '', page_title.string)
48
                      # Truncate the title to fit within the maximum file name length
truncated_title = cleaned_title[:MAX_FILE_NAME_LENGTH - len(str(i)) - 5] # Subtracting length of "_{
49
50
51
52
53
54
55
56
57
58
                      file_extension = os.path.splitext(href)[1]
                      pdf name = f"{link n} {truncated title} {i}{file extension}"
                      ### Write content in files
                      with open(pdf_name, 'wb') as pdf:
59
60
                            pdf.write(response.content)
61
                      pdf.close()
62
63
                      ## populate info circular info data frame
64
65
                      df_new_row=pd.DataFrame({'Link_n':[link_n],'Circular_n':[Circular_n],
                                   'page_title':[cleaned_title], 'From':[From],
                          'Published_on':[Published_on], 'Effective_from':[Effective_from],
                         'Effective_to':[Effective_to], 'File_name':[pdf_name],
                     'Sector':[Sector], 'Author':[Author], 'Followed_by':[Followed_by],
                    'Succeeds], 'Succeeded_by':[Succeeded_by],
                    'Succeeded_by_i:[Succeeded_by_i], 'Succeeded_by_i]
66
67
68
69
70
71
72
73
74
75
76
77
                                                           'Succeeds_L': [Succeeds_L], 'Succeeded_by_L': [Succeeded_by_L], 'url': [url]
                                                         })
                      circulars info=pd.concat([circulars info, df new row], ignore index=True)
                      print("File ", pdf_name, " downloaded")
78
79 ## write info into files (csv and xlsx)
80
81 circulars info.to csv('/home/ga/Desktop/project MSc AI/data/circulars info.csv')
82 circulars_info.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_info.xlsx',
83
                                       sheet_name='circulars_links', index=False)
84
85
```

6 Initial exploratory data analysis

File: 2. circulars_info_clean Exploratory data analysis.ipynb

```
1 circulars_info=pd.read_csv('/home/ga/Desktop/project MSc AI/data/circulars_info.csv', index_col = [0])
1 #set default parameters for chart
2
3 import matplotlib
4 import matplotlib.pyplot as plt
5 plt.style.use('fivethirtyeight')
6 from pylab import rcParams
7 matplotlib.rcParams['axtick.labelsize'] = 8
8 matplotlib.rcParams['ytick.labelsize'] = 7
9 matplotlib.rcParams['ytick.labelsize'] = 7
10 matplotlib.rcParams['text.color'] = 'k'
11
12
13 rcParams['figure.figsize'] = 10, 7
```

1 circulars_info.head()

	Link_n	Circular_n	page_title	From	Published_on	Effective_from	Effective_to	File_name	Sector	Author	Followed_by	Succeed
0	1	12/2024	Circular 122024 Arrangements for Occupational	Department of Public Expenditure, NDP Delivery	24 July 2024	NaN	NaN	1_Circular 122024 Arrangements for Occupationa	NaN	NaN	NaN	Na
1	2	LG	LG 08-2024 Guidelines for	Department of Housing, Local	17 July 2024	NaN	NaN	2_LG 08-2024 Guidelines for	NaN	LG 08-2024 Guidelines	NaN	Na

1 circulars_info.shape

(4017, 16)

1 circulars_info.info()

```
<class 'pandas.core.frame.DataFrame'>
Index: 4017 entries, 0 to 4016
Data columns (total 16 columns):
# Column Non-Null Count Dtype
```

plot NAs

```
def plot_nas(df: pd.DataFrame):
    if df.isnull().sum() != 0:
        na_df = (df.isnull().sum() / len(df)) * 100
        na_df = na_df.drop(na_df[=a_df == 0].index).sort_values(ascending=False)
        missing_data = pd.DataFrame({'Missing Ratio %' :na_df})
        missing_data.plot(kind = "barh")
        plt.show()
    else:
        print('No NAs found')
        plot_nas(circulars_info)
```

```
1 #drom empty columns
2 drop_col=['Succeeds', 'Succeeds_L', 'Succeeded_by', 'Succeeded_by_L']
3 circulars_info=circulars_info.drop(drop_col, axis=1)
```

1 circulars_info1=circulars_info
1 # Added column 'Year_published' for further analysis
2 circulars_info1['Year_published']= circulars_info1['Published_on'].str[-4:]

1 #transform the Published_on column to date time 2 circulars_infol['Published_on'] = pd.to_datetime(circulars_infol['Published_on']) 1 # Pivoted circulars_infol to get the number of circulars per department 2 circular_counts = circulars_infol.groupby(['From'])['Circular_n'].nunique().reset_index(name='Count') 4 # Print the result 5 print(circular_counts) 1 circular_counts['From'] = circular_counts['From'].str.replace('Department', 'Dep') 2 #circular_counts.plot(kind='bar', x='From', y='Count') 3 # Plot the results 5 plt.figure(figsize=(10, 6)) 6 ax = circular_counts.plot(kind='bar', x='From', y='Count') 7 plt.xticks(rotation=45) 8 # Add the number on top of each column 10 for p in ax.patches: 11 ax.annotate(str(p.get_height()), (p.get_x() + p.get_width() / 2., p.get_height()), ha='center', va='center',

```
13 plt.tight_layout()
14 plt.show()
15
```

: 1 #print the number of circulars per year 2 pd.set_option('display.max_rows', None) 4 df_c3=circulars_info1.pivot_table(values=['Circular_n'], 6 index=['Published_on'], 7 aggfunc=pd.Series.nunique 8) 9 df_c3 10 df_c3.resample('YE',level=0).sum().plot() 12 #print(df_c3.resample('Y',level=0).sum())

1 # Convert entire data frame as markdown and print 2 print(df_c3.resample('YE',level=0).sum().to_markdown())

```
2 # Add column 'Document_type'
3
4 # Use np.where to conditionally assign values based on a condition
a
10 # Plot the results
11 plt.figure(figsize=(10, 6))
2 ax = file type counts.plot(kind='bar')
13 plt.xlabel('File Type')
14 plt.ylabel('Number of Documents')
15 plt.title('Number of Documents by File Type')
16 plt.xticks(rotation=45)
18 # Add the number on top of each column
19 for p in ax.patches:
20
      ax.annotate(str(p.get height()), (p.get x() + p.get width() / 2., p.get height()), ha='center', va='center',
21
22 plt.tight_layout()
23 plt.show()
```

1 # Replace illegal characters in the entire DataFrame 2 circulars_infol_cleaned = circulars_infol.replace(r'\x02', '', regex=True) 3 # circulars_infol.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_info.xlsx', sheet_name='circulars_info' 5 6 7 # Replace illegal characters in the entire DataFrame 8 #circulars_infol_cleaned = circulars_infol.replace(r'\x02', '', regex=True) 9 # Save the cleaned DataFrame to Excel 1 circulars_infol_cleaned.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_infol.xlsx', sheet_name='circula 1 circulars_infol_cleaned.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_infol.xlsx', sheet_name='circula 1 circulars_infol_cleaned.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_infol.xlsx', sheet_name='circula 1 circulars_infol_cleaned.to_excel('/home/ga/Desktop/project MSc AI/data/circulars_infol.xlsx', sheet_name='circular 1 circulars_infol_to_csv('/home/ga/Desktop/project MSc AI/data/circulars_infol.xlsx', sheet_name='circular 1 circulars_infol.to_csv('/home/ga/Desktop/project MSc AI/data/circulars_infol.csv')

7 Convert doc and docx files to pdf

File: 3. Convert doc and docx files to pdf.ipynb

```
# Directory containing DOC and DOCX files
1
2 directory = '/home/ga/Desktop/project MSc AI/data/downloaded files'
   #function to convert file to pdf
def convert_to_pdf(file_path, pdf_file):
 1
 2
         try:
              # Load the document
5
              document = Document()
6
              document.LoadFromFile(file_path)
7
8
              # Save the document as PDF
9
              document.SaveToFile(pdf_file, FileFormat.PDF)
10
              return True
11
         except Exception as e:
              print(f"Error occurred while converting {file path} to PDF: {e}")
              return False
13
14
15 #function to convert batch to pdf
16 def batch_convert_to_pdf(directory):
17
         success count = 0
18
         error_count = 0
19
         for filename in os.listdir(directory):
    if filename.endswith(".docx") or filename.endswith(".doc"):
20
21
                   file_path = os.path.join(directory, filename)
pdf_file = os.path.splitext(file_path)[0] + '.pdf'
if convert_to_pdf(file_path, pdf_file):
23
24
25
                        success_count += 1
26
                   else:
27
                        error count += 1
28
29
         print(f"Successfully converted {success_count} files to PDF.")
30
         print(f"Failed to convert {error_count} files.")
31
```

```
1 # run the code
2 batch_convert_to_pdf(directory)
```

8 Extract text and metadata from pdf

File: 4. extract data from pdf.ipynb

```
1 # create Function to count words in a text
   #https://www.geeksforgeeks.org/count-words-in-a-given-string/
 4
   def countWords(s):
5
        # Check if the string is null
6
       # cneck if the string is nu
# or empty then return zero
if s.strip() == "":
8
 9
             return 0
10
        # Splitting the string
        words = s.split()
13
14
        return len(words)
15
16 #def count_words(text):
17
      #return len(text.split())
```

1 # create Function to count words in a text 2 #https://www.geeksforgeeks.org/count-words-in-a-given-string/ def countWords(s): 4 5 # Check if the string is null 6 # or empty then return zero
if s.strip() == "": 7 8 9 return 0 # Splitting the string 10 words = s.split() return len(words)

```
1 #create function to extract pdf info including text
 3 def extract_pdf_info(pdf_file):
 4
           try:
                # Ensure pdf_file is a string representing the file path
if not isinstance(pdf_file, str):
    raise ValueError("pdf_file must be a string representing the file path.")
 5
 6
 7
 8
 9
                 # Open the PDF file
                with pdfplumber.open(pdf_file) as pdf:
    # Extract document information
10
11
12
                       doc info = {
13
                             'Document Name': os.path.basename(pdf_file),
                             'Number of Pages': len(pdf.pages),
'File Size (bytes)': os.path.getsize(pdf_file),
14
15
                            'Author_pdf': pdf.metadata.get('Author', None),
'Creator_pdf': pdf.metadata.get('Creator', None),
'Producer_pdf': pdf.metadata.get('Producer', None),
'Subject_pdf': pdf.metadata.get('Subject', None),
16
17
18
19
20
                            'Title_pdf': pdf.metadata.get('Title', None)
21
22
23
                      }
                      # Extract text from each page
24
25
26
                      results = []
for page in pdf.pages:
                             text = page.extract text()
27
                             if text:
28
29
                                  results.append(text)
                      end_text = ''.join(results)
30
                      doc_info['word_count'] = countWords(end_text)
doc_info['Text'] = end_text
32
34
           except Exception as e:
35
                 print(f"Error occurred while reading {pdf_file}: {e}")
                 doc_info = None
36
38
          return doc_info
```

```
1 # Directory containing PDF files
2 pdf_directory = '/home/ga/Desktop/project MSc AI/data/downloaded files'
3
4 # List to store extracted information
5 pdf_info_list = []
6 failed_count = 0
7
```

12 print("Number of files failed to read:", failed_count)

8

```
2 # Filter out None values from the list if needed
3 #pdf_info_list1 = [info for info in pdf_info_list if info is not None]
4 5 # Create DataFrame from the list of extracted information
6 pdf_df = pd.DataFrame(pdf_info_list)
7
```

1 # Save DataFrame to CSV file
2 csv_file = '/home/ga/Desktop/project MSc AI/data/pdf_info_and_text.csv'
3 xlsx_file = '/home/ga/Desktop/project MSc AI/data/pdf_info_and_text.xlsx'

1 circulars_info.to_csv(csv_file)
2 circulars_info.to_excel(xlsx_file, sheet_name='data', index=False)

```
1 # Check for duplicate rows
2 duplicate_rows = pdf_df[pdf_df.duplicated()]
3 print("Number of duplicate rows:", duplicate_rows.shape[0])
4
5 # Display duplicate rows if any
6 if not duplicate_rows.empty:
7 print(duplicate_rows)
```

Number of duplicate rows: 0

1 print(len(pdf_df))

3841

```
4 # Replace Evaluation Warning: The document was created with Spire.Doc for Python. with in pdf_df
5 pdf_df2 = pdf_df.astype(str).replace(
7 (?i)Evaluation Warning: The document was created with Spire.Doc for Python.", "",
7 regex=True).astype(pdf_df.dtypes.to_dict())
1 # Apply the count_words function to each row of the Text column and create a new column with word counts
2 pdf_df2['word_count_2'] = pdf_df2['Text'].apply(lambda x: countWords(x))
1 # Group the DataFrame by word counts and summarize the data
2 summary_df = pdf_df2.groupby('word_count_2').agg({'Text': 'count'}).reset_index()
3 summary_df.columns = ['Word Count', 'Number of Texts']
```

```
5 #check to see the number of words
6 print(summary_df.to_string())
```

```
1 #Histogram for word count
2 # Modify the 'word_count' column to group documents with over 10000 words
3 pdf_df2['word_count_grouped'] = np.where(pdf_df2['word_count'] > 10000, pdf_df2['word_count'])
4 # Plot the histogram with the modified data, removing the string conversion
6 pdf_df2['word_count_grouped'].plot(kind='hist',bins=100, title='word_count') # Removed .astype(str)
7 plt.gca().spines[['top', 'right']].set_visible(False)
8 plt.show()
```

```
/]: 1 less_than_100 = pdf_df2[pdf_df2['word_count'] < 100]['word_count'].count()
    print("Number of documents with less than 100 words:", less_than_100)
Number of documents with less than 100 words: 69
3]: 1 over_10000 = pdf_df2[pdf_df2['word_count'] > 10000]['word_count'].count()
    print("Number of documents more than than 10000 words:", over_10000)
Number of documents more than than 10000 words: 79
5]: 1 #as thee are 69 documents with less than 100 words and 79 with over 10000 words
    2 # we select only documents with word count of more than 100 words and which contain the first document in the
    3 pdf_df3 = pdf_df2[(pdf_df2['word_count_2'] > 100) & (pdf_df2['Document Name'].str.contains('_1.pdf'))]
7]: 1 len(pdf_df3)
```

]: 3304

1 # Modify the 'word_count' column to group documents with over 10000 words 2 pdf_df3['word_count_grouped'] = np.where(pdf_df3['word_count_2'] > 10000, 10000, pdf_df3['word_count_2']) 4 # Plot the histogram with the modified data, removing the string conversion
5 pdf_df3['word_count_grouped'].plot(kind='hist',bins=100, title='word_count') #
6 plot of the string for the string for the string for the string conversion
6 plot of the string for the string for the string for the string conversion
6 plot of the string for the string for the string for the string conversion
6 plot of the string for the string for the string conversion
6 plot of the string for the 6 plt.gca().spines[['top', 'right']].set_visible(False) 7 plt.show() 1 #most pdf documents did not have an author with the highest stated being the Dep of Finance 2 # Count the occurrences of each Author_pdf 3 author_counts = pdf_df3['Author_pdf'].value_counts() 5 # Select the top 10 authors
6 top_10_authors = author_counts.head(10) 8 # Plot the bar chart 9 plt.figure(figsize=(10, 6)) b pt://guterigs/ze=(16, 6)/ 0 ax=top 10_authors.plot(kind='bar') 1 plt.title('Top 10 Authors') 2 plt.xlabel('Author') 1 plt.ylabel('Number of Documents') 14 plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability 15 16 # Add the number on top of each column 18 for p in ax.patches: 19 ax.annotate(str(p.get_height()), (p.get_x() + p.get_width() / 2., p.get_height()), ha='center', va='center', 21 plt.tight layout()# Adjust layout to prevent labels from overlapping plt.show()

merge the circulars info and the pdf_df3

```
1 file_to_merge = '/home/ga/Desktop/project MSc AI/data/circulars_infol.csv'
2 #xlsx_file = '/home/ga/Desktop/readings master prep/project/data/circulars_info.xlsx'

1
2
3 circulars_info=pd.read_csv(file_to_merge)
4 circulars_info = circulars_info.drop(columns=['Unnamed: 0'])

1 # the file name works as unique identifier of the documents
2 #in order to merge the circulars_info and the data from the Pdf_df3 we will replace the extension of the doc and
3
4 def change_extension(filename):
5 if filename.endswith('.doc'):
6 return filename[:-4] + '.pdf' # Replace .doc with .pdf
9 else:
10 return filename[:-5] + '.pdf' # Replace .docx with .pdf
9 else:
11 circulars_info['Document Name'] = circulars_info['File_name'].apply(change_extension)

1 # Perform the merge
2 merged_df = pdf_df3.merge(circulars_info, on='Document Name', how='left')
```

1 merged_df[['Sector','Author']].head(5)

1 # Remove ':' and leading spaces from 'Sector' and 'Author' columns 2 merged_df['Sector'] = merged_df['Sector'].str.lstrip(': ').str.strip() 3 merged_df['Author'] = merged_df['Author'].str.lstrip(': ').str.strip() 1 # Replace illegal characters with a suitable replacement (e.g., a space) 2 merged_df = merged_df.replace('\x02', ' ', regex=True) 1 #transform the Published on column to date time 2 merged_df['Published_on'] = pd.to_datetime(merged_df['Published_on']) 1 # Save DataFrame to CSV file 2 csv_file = '/home/ga/Desktop/project MSc AI/data/pdf_df3_1.csv' 3 xlsx_file = '/home/ga/Desktop/project MSc AI/data/pdf_df3_1.xlsx' 1 merged_df.to_excel(xlsx_file, sheet_name='files_detail', index=False) 2 merged_df.to_csv(csv_file, index=False)

4 print("PDF information saved to:", csv_file)

```
1 #Create chart for the departments that published the Circulars as per cleaned data
 2 merged_df['From']=merged_df['From'].str.replace('Department', 'Dep')
 3
4 # Count the occurrences of each Department
5 department counts = merged_df['From'].value_counts()
6
 7 # Select the top 10 departments
 8 top 10 departments = department counts.head(10)
10 # Plot the bar chart
11 plt.figure(figsize=(10, 6))
12 bars = top_10_departments.plot(kind='bar')
13 plt.title('Top 10 Departments')
14 plt.xlabel('Department')
15 plt.ylabel('Number of Documents')
16 plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for better readability
17
18 # Add count labels on top of each bar
19 for bar in bars.patches:
       20
21
22
                     ha='center', va='center'
                     size=10, xytext=(0, 8),
24
                     textcoords='offset points')
25
26 plt.tight_layout() # Adjust layout to prevent labels from overlapping
1 #plot the number of documents in the corpus per year published
 2
 3 pd.set_option('display.max_rows', None)
 4
 5 df=merged_df.pivot_table(values=['Circular_n'],
 6
                             index=[ 'Published on'],
                           aggfunc=pd.Series.nunique
 8
                        )
 9
11 df.resample('YE',level=0).sum().plot()
```

1 # Convert entire data frame as markdown and print year 2 print(df.resample('YE',level=0).sum().to_markdown())

9 Clean and vectorize the text

File: 5. Clean and vectorize the text.ipynb

```
1 #import the data
2
3
4 csv_file3 = '/home/ga/Desktop/project MSc AI/data/pdf_df3_1.csv'
5
6 pdf df4=pd.read csv(csv file3, index col='Document Name')
```

2. Text pre-processing

```
1

2 ## Download NLTK resources (stopwords and WordNet corpus) if they were not downloaded

3 #nltk.download('stopwords')

4 #nltk.download('wordnet')

5 #nltk.download('punkt')
```

```
1 #https://stackoverflow.com/questions/54396405/how-can-i-preprocess-nlp-text-lowercase-remove-sp
 2
 3
 4
 5
 6 lemmatizer = WordNetLemmatizer()
 8 # set the stopwords dictionary
9 custom_stopwords = set(stopwords.words('english') + ['circular', 'education'])
10
11 #create function to clean the text with the arguments being the text to tokenize,
12 #the tokenizer and the stopwords dictionary
13 #function returns the cleaned ext
14
15 def clean text(text, tokenizer, stopwords):
16
17
          text = str(text).lower() # Lowercase words
         text = str(text).tower() # Lowercase words
text = re.sub(r"\[[(.*?)\]", "", text) # Remove [+XYZ chars] in content
text = re.sub(r'\truck) + ", ", text) # remove web links
text = re.sub(r"\\s+", " ", text) # Remove multiple spaces tabs, and line breaks in content
text = re.sub(r"\\s+", " ", text) # Remove ellipsis (and last word)
text = re.sub(r"(?<=\w)-(?=\w)", " ", text) # Replace dash between words
text = re.sub('[0-9]+', '', text) #Remove numbers
18
19
20
21
22
23
24
          text = re.sub(f"[{re.escape(string.punctuation)}]", "", text) # Remove punctuation
25
26
          tokens = tokenizer(text) # Get tokens from text
          tokens = [t for t in tokens if not t in stopwords] # Remove stopwords
tokens = ["" if t.isdigit() else t for t in tokens] # Remove digits
27
28
          #Lemmatization - to return the base or dictionary form of a word, which is known as the lem
29
          tokens=[lemmatizer.lemmatize(w) for w in tokens]
30
31
32
          tokens = [t for t in tokens if len(t) > 1] # Remove short tokens
          return " ".join(tokens)
33
```

```
1 #run the function
2 pdf_df4['cleanText'] = pdf_df4['Text'].map(lambda x: clean_text(x, word_tokenize, custom_stopwords))
1 # Tokenize the text data
2 pdf_df4['tokenizedText'] = pdf_df4['cleanText'].apply(word_tokenize)
3
```

```
1 #function to display HTML
2
3 from IPython.core.display import HTML, display
4 def visualisel(title, content):
5   text = ' '
6   display(HTML(f'<hl>{title}</hl>'))
7   text += ' ' + content.replace('\n', '<br>')
8   display(HTML(f''{text}'''))
```

1 pdf_df4['Text'][2]

```
1 #see the text
2 visualise1(pdf_df4['File_name'][5], pdf_df4['Text'][5])
```

```
1 # see the cleaned text
2 visualise1(pdf_df4['File_name'][5], pdf_df4['cleanText'][5])
3
2 #plot the document length destribution
4
  pdf df4['word count clean'].iplot(
5
       kind='hist',
6
       bins=100.
7
       xTitle='Clean text word count',
       yTitle='Count'
8
9
       title='Cleaned Text Word Count Distribution',
LO
       colors='green'
11 )
```

```
1 #print the top 20 words in the original text|
2
3 common_words = get_top_n_words(pdf_df4['Text'], 20)
4
5 #for word, freq in common_words:
6 #print(word, freq)
7
8 df1 = pd.DataFrame(common_words, columns = ['Words', 'count'])
9
10 custom iplot(df1, 'Words', title='Top 20 words in review before cleaning the text')
```

```
1 #print most frequent bigrams in te original text
2 common_words = get_top_n_bigram(pdf_df4['Text'], 20,2)
3
4 #for word, freq in common_words:
5 #print(word, freq)
6 df3 = pd.DataFrame(common_words, columns = ['Words', 'count'])
7 custom iplot(df3, 'Words', title='Top 20 bigrams before cleaning the text')
```

```
1 #print most frequent bigrams in the cleaned text
2
3 common words = get top n bigram(pdf_df4['cleanText'], 20,2)
4
5 #for word, freq in common words:
6
      #print(word, freq)
7 df4 = pd.DataFrame(common_words, columns = ['Words', 'count'])
8 custom_iplot(df4, 'Words', title='Top 20 bigrams after cleaning the text')
1 ##get ngrams before and after removing stop words
2 # plot Top 20 trigrams before cleaning the text
3
4 common words = get top n bigram(pdf df4['Text'], 20,3)
5
6 #for word, freq in common words:
7
      #print(word, freq)
8 df5 = pd.DataFrame(common_words, columns = ['Words', 'count'])
9 custom iplot(df5, 'Words', title='Top 20 trigrams before cleaning the text')
```

```
1 #print most frequent trigrams in the cleaned text
2 common_words = get_top_n_bigram(pdf_df4['cleanText'], 20,3)
3
4 #for word, freq in common_words:
5     #print(word, freq)
6 df6 = pd.DataFrame(common_words, columns = ['Words', 'count'])
7 custom_iplot(df6, 'Words', title='Top 20 trigrams in the cleaned the text')
```

```
1 #save the file for further reference
2 csv_file_pdf_df4 = '/home/ga/Desktop/project MSc AI/data/pdf_df4.csv'
3 xlsx_file_pdf_df4 = '/home/ga/Desktop/project MSc AI/data/pdf_df4.xlsx'
```

3. Text vectorization

3.1 Text vectorization using TF-IDF

```
1 # TF-IDF document frequency function
2 vectorizer = TfidfVectorizer(
3 min_df = 0.01, #to eliminate words that are not too common - ignore terms that app
4 max_df = 0.97, #to eliminate words that are too common - ignore terms that appear
5 #max_features = 5000,
6 #stop_words = 'english', #remove the english stopwords - done already
7 #token_pattern=u'(?ui)\\b\\w*[a-z]+\\w*\\b' # eliminate words that contain number;
8 )
```

1 # Fit and transform the text data using TF-IDF vectorizer 2 tfidf_matrix = vectorizer.fit_transform(pdf_df4['cleanText'])

```
1 tfidf_matrix
```

```
<3304x4125 sparse matrix of type '<class 'numpy.float64'>'
with 966753 stored elements in Compressed Sparse Row format>
```

```
1 # Print how many rows and columns of the TF-IDF matrix consists
2 print("n_documents: %d, n_features: %d" % tfidf_matrix.shape)
```

n_documents: 3304, n_features: 4125

```
1 # Select the first ten documents from the data set
2 #see https://okan.cloud/posts/2022-01-16-text-vectorization-using-python-tf-idf/
3 
4 tf_idf_df = pd.DataFrame(tfidf_matrix.todense())
5 tf_idf_df.columns = vectorizer.get_feature_names_out()
6 #tf_idf_df.index=pdf_df4['Document_name']
```

```
1 tf_idf_df.tail(3)
```



3.2 Text vectorization word2vec

https://www.tensorflow.org/text/tutorials/word2vec/https://www.geeksforgeeks.org/python-word-embedding-using-word2vec/https://dylancastillo.co/nlp-snippets-cluster-documents-using-word2vec/

9]: 1 import gensim 2 from gensim.models import Word2Vec

by using min_count=50 we reduced the Vocabulary Size: 63817 to 4846

There are several ways to vectorize a Word2Vec model for document representation which include Averaging Word Vectors, TF-IDF Weighted Averaging where each word vector is weighted by its TF-IDF score, Concatenation of Min, Max, and Average Vectors for the words in a document which captures different aspects of the word distributions and Using Pre-trained Sentence Embeddings.

```
# vectorize using the min max and average
1
2
   def vectorize word2vec(line):
3
       words = [1]
       for word in line: # line - iterable list of tokens
4
5
           try:
6
               w2v_idx = word2vec_indices[word]
           except KeyError: # if we do not have a vector for this word in your w2v model, continue
8
               continue
           words.append(word2vec vectors[w2v idx])
9
10
       if words:
           words = np.asarray(words)
           min_vec = words.min(axis=0)
           max_vec = words.max(axis=0)
14
           avg vec = words.mean(axis=0)
           return np.concatenate((min_vec, max_vec, avg_vec))
16
       else:
18
           return None
```

4]	:	<pre>1 word2vec_matrix_vectors = pd.DataFrame(vectorized_docs) 2</pre>
:	1	<pre>vectorized_docs = [vectorize_word2vec(doc) for doc in tokenized_text]</pre>

```
1 # save matix for further reference
```

2 word2vec_matrix_vectors.to_csv('/home/ga/Desktop/project MSc AI/data/word2vec_matrix_vectors.csv', index=False)

3.3 Text vectorization BERT

4 # Load pre-trained BERT model and tokenizer

- 5 tokenizer_BERT = BertTokenizer.from_pretrained('bert-base-uncased')
- 6 model_BERT = BertModel.from_pretrained('bert-base-uncased')

```
1 #basic Pre-cleaning for BERT
 3 def text_clean_BERT(text):
 4
               #text = text.lower() # lowercase everything
x = text.encode('ascii', 'ignore').decode() # remove unicode characters
x = re.sub(r'https*\S+', ' ', x) # remove links
x = re.sub(r'http*\S+', ' ', x)
x = re.sub(r'\s{2,}', ' ', x) #replaces multiple consecutive spaces with a single space
 5
      1
 6
 8
 9
10
               # cleaning up text
#x = re.sub(r'\'\w+', '', x) # removes words that start with an apostrophe (contractions or possessives).
#x = re.sub(r'\w*\d+\w*', '', x) #removes words containing digits
#x = re.sub(r'\s{2,}', ' ', x) #replaces multiple consecutive spaces with a single space
#x = re.sub(r'_, '', x) #remove the underscore character (.)
#x = re.sub(r'(.)', '', x) #remove the underscore character
14
16
17
               # Consider keeping punctuation that might be relevant to context
# x = re.sub(r'\s[^\w\s]\s', '', x) # Commenting out to keep punctuation surrounded by spaces
18
                # Handle specific cases where spacing might be important
               x = re.sub(r'\.([a-zA-Z])', r'. \1', x)
x = re.sub(r'\,([a-zA-Z])', r', \1', x)
24
25
                return x.strip()
```

1 circulars_text_BERT = pdf_df4['Text'].apply(text_clean_BERT).tolist()

```
1 # BERT can work with text of maximum length of 512 tokens. This is why if the document is too long
2 #we have to split and encode chunks, then average
3 #Function to encode documents using BERT with handling long documents
 4
     def encode documents(docs, max length=512):
           embeddings = []
 5
           for doc in docs:
 6
                 inputs = tokenizer_BERT(doc, return_tensors='pt', max_length=max_length, truncation=True, padding=True)
                 with torch.no_grad():
    outputs = model BERT(**inputs)
 8
10
                 doc embedding = outputs.last hidden state.mean(dim=1).squeeze().numpy()
                # If document is too long, split and encode chunks, then average
#it does that by converting the doc in a tnsor format, truncating the document if it exceeds the maximum
#or padding the document if it is shorter than the maximm length
if len(inputs['input_ids'][0]) > max_length:
    church corections.
14
15
                       chunk_embeddings = []
for i in range(0, len(doc), max_length):
17
18
                             chunk = doc[i:i+max_length]
19
                             inputs = tokenizer_BERT(chunk, return_tensors='pt', max_length=max_length,
                             truncation=True, padding=True)
# Passes the tokenized inputs through the BERT model without computing gradients (saves memory an
#Computes the average embedding of the tokens in the last hidden state for the entire document an
22
                             with torch.no_grad():
                                   outputs = model(**inputs)
24
                             chunk_embeddings.append(outputs.last_hidden_state.mean(dim=1).squeeze().numpy())
                       doc_embedding = np.mean(chunk_embeddings, axis=0)
                 embeddings.append(doc_embedding)
           return np.array(embeddings)
29
```

```
1 # Encode the documents using BERT
2 document_vectors_BERT = encode_documents(circulars_text_BERT)
1 document_vectors_BERT.shape
(3304, 768)
1 #save for further reference
2 BERT_matrix_vectors = pd.DataFrame(document_vectors_BERT)
```

```
3 BERT_matrix_vectors.to_csv('/home/ga/Desktop/project MSc AI/data/BERT_matrix_vectors.csv', index=False)
```

10 Check the right number of components and find the optimum n of clusters

File: 6 Check the right number of components and find the optimum n of clusters.ipynb

```
1 #import the data
2 csv file_pdf_df4 = '/home/ga/Desktop/project MSc AI/data/pdf_df4.csv'
3 tf_idf_df=pd.read_csv('/home/ga/Desktop/project MSc AI/data/tf_idf_df.csv', index col='Unnamed: 0')
```

2. TF-IDF performing dimensionality reduction using TruncatedSVD

https://scikit-learn.org/stable/auto examples/text/plot document clustering.html

```
# tf idf df is the DataFrame with TF-IDF values
2 X = tf idf df.values
    #Try different number of components
  n components range = range(1, min(X.shape[0], X.shape[1]), 50) # Adjust step size as needed
  explained variance ratios = []
# Finding optimal number of features for clustering
4
5 for n_components in n_components_range:
       svd = TruncatedSVD(n_components=n_components)
6
       svd.fit(X)
7
8
       explained_variance = svd.explained_variance_ratio_.sum()
9
       explained_variance_ratios.append(explained_variance)
10
       print(f"Number of components = {n_components}, Explained Variance = {explained_variance}")
```

```
1 # Plot explained variance vs. number of components
2 plt.plot(n_components_range, explained_variance_ratios)
3 plt.xlabel('Number of Components')
4 plt.ylabel('Explained Variance Ratio')
5 plt.title('Explained Variance vs. Number of Components (Truncated SVD)')
6 plt.show()
```

It looks like the first few 1000 components explain 94% of the variance. After this adding more components does not signiffically improve the representation of the data.

```
1 svd = TruncatedSVD(n_components=1000)
2 X SVD=svd.fit transform(X)
```

3. Find the optimum n of clusters

```
1 # Use Elbow Method to find optimal number of clusters
2 # Try different ranges of clusters
3 visualizer = KElbowVisualizer(KMeans(max_iter=500, n_init=20, random_state=SEED), k=(2,40))
4 visualizer.fit(X_SVD)
5 visualizer.show()
1 #optimum n of clusters using the silhouette metric
2 visualizer = KElbowVisualizer(KMeans(max_iter=500, n_init=20, random_state=SEED), k=(2,40), metric='silhouette')
3 visualizer.fit(X_SVD)
```

4 visualizer.show()

- 1 #optimum n of clusters using the calinsky harabasz
 2 visualizer = KElbowVisualizer(KMeans(max_iter=500, n_init=20, random_state=SEED), k=(2,40), metric='calinski_hara visualizer.fit(X SVD)
- 4 visualizer.show()

1	#Calculate Davies-Bouldin scores for different numbers of clusters
2	from sklearn.metrics import davies bouldin score
3	davies bouldin scores = []
4	for n in range(2, 41):
5	kmeans = KMeans(n clusters=n, max iter=500, n init=20, random state=SEED)
6	labels = kmeans.fit predict(X SVD)
7	score = davies bouldin score(\overline{X} SVD, labels)
8	davies bouldin scores.append(score)
9	<pre>print(f"Number of clusters = {n}, Davies-Bouldin score = {score}")</pre>

1 #plot Davies-Bouldin scores vs. number of clusters

- plt.plot(range(2, 41), davies_bouldin_scores)
 plt.xlabel('Number of Clusters')
 plt.ylabel('Davies-Bouldin Score')
- 3
- 5 plt.title('Davies-Bouldin Score')
 5 plt.title('Davies-Bouldin Score vs. Number of Clusters')
 6 plt.show() 4

```
1 #plot the silhouette plot for the selected 21 n of clusters
   #https://stackabuse.com/k-means-elbow-method-and-silhouette-analysis-with-yellowbrick-and-scikit-learn/
 3 from sklearn.metrics import silhouette samples
5 model_3clust = KMeans(n_clusters = 21, max_iter=500, n_init=20, random_state=SEED)
6 model 3clust.fit(X SVD)
 7 labels = model_3clust.labels
8 silhouette_vals = silhouette_samples(X_SVD, labels)
10 visualizer = SilhouetteVisualizer(model 3clust)
11 visualizer.fit(X_SVD)
                                # Fit the data to the visualizer
                           # Finalize and render the figure
12 visualizer.show()
```

4. TruncatedSVD for word2vec

```
1 #load data
2 word2vec df=pd.read csv('/home/ga/Desktop/project MSc AI/data/word2vec matrix vectors.csv')
```

```
1 #Try different number of components
   n_components_range = range(1, min(X.shape[0], X.shape[1]), 50) # Adjust step size as needed
 3 explained_variance_ratios = []
4 # Finding optimal number of features for clustering
   for n components in n components range:
         svd = TruncatedSVD(n_components=n_components)
         svd.fit(X)
         explained_variance = svd.explained_variance_ratio_.sum()
explained_variance_ratios.append(explained_variance)
 8
10
         print(f"Number of components = {n_components}, Explained Variance = {explained_variance}")
 1 # Plot explained variance vs. number of components
2 plt.plot(n_components_range, explained_variance_ratios)
3 plt.xlabel('Number of Components')
 4 plt.ylabel('Explained Variance Ratio')
5 plt.title('Explained Variance vs. Number of Components (Truncated SVD)')
6 plt.show()
```

5. TruncatedSVD for BERT

#load data BERT dfl=pd.read csv('/home/ga/Desktop/project MSc AI/data/BERT matrix vectors.csv') 2

1	<pre>X = BERT_dfl.values</pre>
1 2 3 4 5 6 7 8 9 10	<pre>#Try different number of components n_components_range = range(1, min(X.shape[0], X.shape[1]), 50) # Adjust step size as needed explained_variance_ratios = [] # Finding optimal number of features for clustering for n_components in n_components_range: svd = TruncatedSVD(n_components=n_components) svd.fit(X) explained_variance = svd.explained_variance_ratiosum() explained_variance_ratios.append(explained_variance) print(f"Number of components = {n_components}, Explained Variance = {explained_variance}")</pre>
1 2 3 4 5 6	<pre># Plot explained variance vs. number of components plt.plot(n_components_range, explained_variance_ratios) plt.xlabel('Number of Components') plt.ylabel('Explained Variance Ratio') plt.title('Explained Variance vs. Number of Components (Truncated SVD)') plt.show()</pre>

11 Clustering implementation, evaluation and visualization

file: 7. Implementing the clustering.ipynb

2. Load the data and reduce dimensionality

: 1 # Load the data
2
3 base_address_data='/home/ga/Desktop/project MSc AI/data/'
4 tfidf_df = pd.read_csv(base_address_data+'tf_idf_df.csv', index_col='Unnamed: 0')
5 word2vec_df = pd.read_csv(base_address_data+'word2vec_matrix_vectors.csv')
6 bert_df = pd.read_csv(base_address_data+'BERT_matrix_vectors.csv')
1 # Reduce dimensionality to xyz components

```
2 svd = TruncatedSVD(n_components=750)
3 tfidf_df_r= svd.fit_transform(tfidf_df)
4 # Convert the NumPy array back to a Pandas DataFrame
5 tfidf_df_r = pd.DataFrame(tfidf_df_r, index=tfidf_df.index)
6 # We are adding back the index from the original word2vec_df DataFrame
```

```
1
2 # Reduce dimensionality to xyz components
3 svd = TruncatedSVD(n_components=135)
4 word2vec_df_r= svd.fit_transform(word2vec_df)
5 # Convert the NumPy array back to a Pandas DataFrame
6 word2vec_df_r = pd.DataFrame(word2vec_df_r, index=word2vec_df.index)
7 # We are adding back the index from the original word2vec_df DataFrame
```

```
1 # Reduce dimensionality to xyz components
2 svd = TruncatedSVD(n_components=200)
3 bert_df_r= svd.fit_transform(bert_df)
4 # Convert the NumPy array back to a Pandas DataFrame
```

3. Set clustering functions

```
1 #set k
2 k=21
1 def kmeans clustering(df, n clusters=k):
2
    Performs k-means clustering on the given dataframe.
3
4
5
6
      df: The dataframe containing the vectorized data.
7
       n clusters: The number of clusters to form.
8
9
    Returns:
10
      The KMeans model fitted to the data.
11
      kmeans = KMeans(n_clusters=n_clusters, max_iter=500, n_init=20, init='k-means++',
12
13
                       random state=SEED)#n init=20,
14
       kmeans.fit(df)
15
       return kmeans
```

```
1 def efcm_clustering(df, n_clusters=k, m=1.1):
2
3
       Performs Enhanced Fuzzy C-Means clustering on the given dataframe.
4
5
       Args:
6
           df: The dataframe containing the vectorized data.
7
           n clusters: The number of clusters to form.
8
           m: fuziness factor
9
10
       Returns:
       The FCM model fitted to the data.
11
       scaler = StandardScaler()
14
       scaled data = df.values
15
16
17
       # Use cmeans from skfuzzy
18
19
       cntr, u, u0, d, jm, p, fpc = fuzz.cluster.cmeans(
20
           scaled data.T, n clusters, m, error=0.005, maxiter=1000, init=None, seed=SEED
21
       cluster_membership = np.argmax(u, axis=0)
23
       return cluster membership
```

```
1 #lstm clustering
 3
    def build autoencoder(input shape, dense units, dropout rate=0.2):
 4
 5
          Builds a simple autoencoder for tabular data using Dense layers.
 6
 7
          Args:
 8
               input shape: The shape of the input data.
 9
               dense_units: A list containing the number of units in each Dense layer.
10
               dropout rate: Dropout rate for regularization.
11
          Returns:
12
          The autoencoder model, input layer, and the encoded (bottleneck) layer output.
13
14
15
          input layer = Input(shape=input shape)
132
133
          # Scale the data
134
          scaler = StandardScaler()
235
          df scaled = scaler.fit transform(df)
236
2<sup>37</sup>
2<sup>38</sup>
2<sup>39</sup>
          # 1. Pretrain Autoencoder:
input_shape = (df_scaled.shape[1],)
autoencoder, input_layer, encoded_output = build_autoencoder(input_shape, dense_units, dropout_rate)
autoencoder.compile(optimizer=Adam(learning_rate=learning_rate), loss='mse')
2:40
241
          autoencoder.fit(df_scaled, df_scaled, epochs=pretrain_epochs, batch_size=batch_size, verbose=0)
242
 43
          # 2. Extract Encoder:
          encoder = Model(inputs=input_layer, outputs=encoded_output)
 44
 45
          encoded_X = encoder.predict(df_scaled)
 46
 47
          # 3. Initialize Cluster Centers (using K-Means++ for better initialization):
          kmeans = KMeans(n_clusters=n_clusters, n_init=20, init='k-means++')
y_pred = kmeans.fit_predict(encoded_X)
 48
 49
 50
 51
          # 4. Build Clustering Layer:
 52
          clustering_layer = Dense(n_clusters, activation='softmax', name='clustering')(encoded_output)
          model = Model(inputs=input_layer, outputs=clustering_layer)
 53
 54
 55
          def target_distribution(q):
    weight = q ** 2 / q.sum(0)
    return (weight.T / weight.sum(1)).T
 56
 57
 58
 59
          # 5. Train Clustering Model:
 60
          model.compile(optimizer=Adam(learning_rate=learning_rate), loss='kld')
          y_pred_last = np.copy(y_pred)
y_pred = tf.keras.utils.to_categorical(y_pred, n_clusters)
 61
 63
 64
          for epoch in range(epochs):
              q = model.predict(df_scaled, verbose=0) # Reduced verbosity
p = target_distribution(q)
 65
 66
 67
               model.fit(df_scaled, p, epochs=1, batch_size=batch_size, verbose=0) # Reduced verbosity
 69
                _pred = q.argmax(1)
              delta_label = np.sum(y_pred != y_pred_last).astype(np.float32) / y_pred.shape[0]
y_pred_last = np.copy(y_pred)
if epoch > 0 and delta_label < 0.001:</pre>
 70
 71
72
  73
                   print(f'Converged at epoch {epoch}')
  74
                   break
  75
```

76

return y_pred

4. Run Clustering

```
1 def run clustering(df):
 3
      Runs all clustering methods on the given dataframe and evaluates their performance.
 4
 5
      Args:
 6
        df: The dataframe containing the vectorized data.
 7
 8
      Returns:
      A dictionary containing the evaluation metrics for each clustering method.
 9
10
11
      results = {}
12
      # K-means clustering
13
      kmeans_model = kmeans_clustering(df)
kmeans_predictions = kmeans_model.labels_
14
15
16
      results['K-means'] = kmeans_predictions #
17
    # Enhanced Fuzzy C-Means clustering
efcm_predictions = efcm_clustering(df)
#efcm_model = efcm_clustering(df)
#efcm_predictions = efcm_model.predict(df.values)
18
19
20
      results['EFCM'] = efcm predictions
23
24
25
      # LSTM clustering
      lstm_model = lstm_clustering(df)
      results['LSTM'] = lstm_model #lstm_predictions
26
27
28
29
     return results
```

2 # Run clustering and evaluate for each dataset 3 tfidf_results = run_clustering(tfidf_df_r) 4 word2vec_results = run_clustering(word2vec_df_r) 5 bert_results = run_clustering(bert_df_r)

```
1 # put the results together
2 tfidf_results_df = pd.DataFrame(tfidf_results)
3 word2vec_results_df = pd.DataFrame(word2vec_results)
4 bert_results_df = pd.DataFrame(word2vec_results)
5
6 # Concatenate the results horizontally to avoid doubling the index length
7 combined_results = pd.concat([
8 tfidf_results_df.add_suffix('_tfidf'),
9 word2vec_results_df.add_suffix('_word2vec'),
10 bert_results_df.add_suffix('_bert')
11 ], axis=1) # Change axis to 1 for horizontal concatenation
12
13 # Now you should be able to add the clustering results without errors
14 combined_results['K-means_tfidf'] = tfidf_results['K-means']
15 combined_results['EFCM_tfidf'] = tfidf_results['LSTM']
17
18 # Display the combined dataframe
19 print(combined_results.head())
```

5. Calculate metrics

```
1 # function to measure the quality of the clustering results
3 def calculate metrics(df, predicted labels):
4
       Calculates clustering metrics.
5
6
      Args:
df: The dataframe containing the vectorized data.
predicted_labels: The predicted cluster labels.
 7
8
 9
10
11
      Returns:
      A dictionary containing the Calinski-Harabasz Index, Silhouette Score, and Entropy.
12
13
       # Check if number of unique labels is greater than 1
if len(np.unique(predicted_labels)) > 1:
14
15
         chi = calinski harabasz_score(df, predicted labels)
silhouette = silhouette_score(df, predicted_labels)
16
17
18
         dbs=davies_bouldin_score(df, predicted_labels)
19
20
         # Calculate entropy - need laels for this
         #unique labels, counts = np.unique(predicted_labels, return_counts=True)
#probabilities = counts / counts.sum()
#entropy = -np.sum(probabilities * np.log2(probabilities))
21
22
24
25
         return {'Calinski-Harabasz Index': chi, 'Silhouette Score': silhouette, 'Davies Bouldin Index': dbs} #, 'Er
26
       else:
         return {'Calinski-Harabasz Index': None, 'Silhouette Score': None, 'Davies Bouldin Index': None} # Return /
27
28
1 # Calculate metrics for each clustering method and dataset
2 metrics = {}
 3 for dataset_name, results in [('TF-IDF', tfidf_results), ('Word2Vec', word2vec_results), ('BERT', bert_results
      for metrics[dataset_name] = {}
for method_name, predictions in results.items():
    metrics[dataset_name][method_name] = calculate_metrics(
        tfidf_df if dataset_name == 'TF-IDF' else (word2vec_df if dataset_name == 'Word2Vec' else bert_df),
 4
 5
 6
 7
 8
              predictions
 9
         1
 1 #save the metrics as a data frame
7 metrics df
 8
```

6 Plot the clusters using PCA

```
1 # plot the clusters for each of the dataset and for each of the clustering method used using PCA
2
3 #import matplotlib.pyplot as plt
 4 #from sklearn.decomposition import PCA
 5
6 ....
      Plots the clusters using PCA for dimensionality reduction.
 7
8
9
      Args:
        df: The dataframe containing the vectorized data.
labels: The predicted cluster labels.
10
11
         title: The title of the plot.
13 """
15 def plot clusters(df, labels, title):
         # Reduce dimensionality to 2D using PCA
         pca = PCA(n_components=2)
17
18
         reduced_data = pca.fit_transform(df)
19
20
       # Plot the clusters
        fig, ax = plt.subplots(figsize=(8, 6))
scatter = ax.scatter(reduced_data[:, 0], reduced_data[:, 1], c=labels, cmap='viridis', alpha=0.7)
23
         ax.set title(title)
        ax.set_xlabel("PCA Component 1")
ax.set_ylabel("PCA Component 2")
24
25
26
27
      # Create a legend with all clusters
28
29
        unique_labels = np.unique(labels)
        handles = [plt.Line2D([0], [0], marker='o', color='w', label=f"Cluster {label}",
markerfacecolor=plt.cm.viridis(label / unique_labels.max()), markersize=10)
30
31
                    for label in unique_labels]
32
33
      # Position the leaend outside to the right
34
        ax.legend(handles=handles, loc='center left', bbox_to_anchor=(1, 0.5))
35
36  # Plot clusters for each dataset and clustering method
37  for dataset_name, results in [('TF-IDF', tfidf_results), ('Word2Vec', word2vec_results),
38  ('BERT', bert_results)]:
39
         for method name, predictions in results.items():
40
           plot_clusters(
    tfidf_df if dataset_name == 'TF-IDF' else (word2vec_df_r if dataset_name == 'Word2Vec' else bert_df_r),
41
                predictions,
42
43
                f"{dataset_name} - {method_name} Clustering"
44
         )
```

7. Plot the clusters using t-SNE

```
1 #plot the clusters for each of the dataset and for each of the clustering method used using t-SNE
    #import matplotlib.pyplot as plt
 3 #from sklearn.manifold import TSNE
 4
 5 def plot clusters(df, labels, title):
     # Reduce dimensionality to 2D using TSNE
tsne = TSNE(n_components=2, random_state=SEED)
 6
 8
         reduced data = tsne.fit transform(df)
 9
10
      # Plot the clusters
        fig, ax = plt.subplots(figsize=(8, 6))
         scatter = ax.scatter(reduced_data[:, 0], reduced_data[:, 1], c=labels, cmap='viridis', alpha=0.7)
         ax.set title(title)
        ax.set_xlabel("t-SNE Component 1")
ax.set_ylabel("t-SNE Component 2")
14
15
16
17
      # Create a legend with all clusters
18
         unique_labels = np.unique(labels)
        handles = [plt.Line2D([0], [0], marker='o', color='w', label=f"Cluster {label}",
markerfacecolor=plt.cm.viridis(label / unique_labels.max()), markersize=10)
19
20
21
                     for label in unique labels]
22
      # Position the legend outside to the right
ax.legend(handles=handles, loc='center left', bbox_to_anchor=(1, 0.5))
24
25
26 # Plot clusters for each dataset and clustering method
27 for dataset_name, results in [('TF-IDF', tfidf_results), ('Word2Vec', word2vec_results), ('BERT', bert_results)]:
28 for method_name, predictions in results.items():
         plot_clusters(
29
              tfidf_df if dataset_name == 'TF-IDF' else (word2vec_df_r if dataset_name == 'Word2Vec' else bert_df_r),
30
              predictions.
31
32
              f"{dataset name} - {method name} Clustering"
33
         )
```

8. Plot interactive clusters for Word2Vec with k-means with both the PCA and the t-SNA

```
1 # For interactive plots
 2 import plotly.express as px
 3 import plotly.graph_objects as go
 4
 5 pdf_df4 = pd.read_csv(base_address_data+'pdf_df4.csv')
 6 word2vec_results_df=pd.DataFrame(word2vec_results)
 8 tsne = TSNE(n_components=2, random_state=SEED)
 9 word2vec_tsne = tsne.fit_transform(word2vec_df)
10
11
12 # Create a DataFrame for plotting
12 # Create a batarrame for plotting
13 plot_df = pd.DataFrame({
14 'x': word2vec_tsne[:, 0],
15 'y': word2vec_tsne[:, 1],
16 'Cluster': word2vec_results_df['K-means'],
17 'Document Name': pdf_df4['Document Name']
18 })
19
# Interactive Plot using Plotly for t-SNE results with hover data
fig_tsne = px.scatter(plot_df, x='x', y='y', color='Cluster', hover_data=['Document Name'],
labels={'x': 't-SNE Component 1', 'y': 't-SNE Component 2', 'color': 'Cluster'},
title='Cluster Assignments (t-SNE)')
24 fig_tsne.show()
```