

In [1]:

```
1 import os
2 import pandas as pd
3 import numpy as np
4 import bs4 as bs
5 import urllib.request
6 import re
7 import spacy
8 import re, string, unicodedata
9 import nltk
10 from nltk.corpus import stopwords
11 from nltk.stem import WordNetLemmatizer
```

2023-07-21 02:12:53.199692: I tensorflow/core/platform/cpu_feature_guard.cc:193] This TensorFlow binary is optimized with oneAPI Deep Neural Network Library (oneDNN) to use the following CPU instructions in performance-critical operations: SSE4.1 SSE4.2 AVX AVX2 FMA
To enable them in other operations, rebuild TensorFlow with the appropriate compiler flags.

In [2]:

```
1 lem=WordNetLemmatizer()
```

In [3]:

```
1 os.chdir(r'/Users/myyntiimac/Desktop/Combine xml file')
```

In [4]:

```
1 from glob import glob
2
3 path=r'/Users/myyntiimac/Desktop/Combine xml file'
4 all_files = glob(os.path.join(path, "*.xml"))
```

In [5]:

```
1 import xml.etree.ElementTree as ET
```

In [7]:

```
1 dfs = []
2 for filename in all_files:
3     tree = ET.parse(filename)
4     root = tree.getroot()
5     root=ET.tostring(root, encoding='utf8').decode('utf8')
6     dfs.append(root)
7
8
9 parsed_article = bs.BeautifulSoup(dfs[0], 'xml')
10
11 paragraphs = parsed_article.find_all('p')
12
13
14
15
16 article_text_full = ""
17
18 for p in paragraphs:
19     article_text_full += p.text
20     print(p.text)
21
22
23
24
25
26 def data_preprocessing(each_file):
27
28     parsed_article = bs.BeautifulSoup(each_file, 'xml')
29
30     paragraphs = parsed_article.find_all('para')
31
32
33
34     article_text_full = ""
35
36     for p in paragraphs:
37         article_text_full += p.text
38         print(p.text)
39
40     return article_text_full
```

In [8]:

```
1 data=[data_preprocessing(each_file) for each_file in dfs]
```

Difficulties in social situations (in suspected social anxiety disorder)

Irrational and out-of-proportion fear or avoidance of particular objects or situations (in suspected specific phobia)

Intense anxiety reactions with exposure to specific situations (in suspected agoraphobia)

Anxiety is the most common feature in phobic disorders. Manifestations include the following:

Elevated heart rate

Elevated blood pressure

Tremor

Palpitations

Diarrhea

Sweating

Dyspnea

Paresthesias

Dizziness

Because anxiety manifests with a number of physical symptoms, any patient who presents with a de novo complaint of physical symptoms suggestive of an anxiety disorder should undergo a physical examination to he

In [9]:

```

1  def remove_stop_word(file):
2      # Load the spaCy English language model
3      nlp = spacy.load("en_core_web_sm")
4
5      # Define punctuation characters
6      punctuations = string.punctuation
7
8      # Load stopwords from NLTK library
9      from nltk.corpus import stopwords
10     stopwords = stopwords.words('english')
11
12     # Additional symbols to remove
13     SYMBOLS = " ".join(string.punctuation).split(" ") + ["-", "...", "'", "\""]
14     stopwords = nltk.corpus.stopwords.words('english') + SYMBOLS
15
16     # Apply spaCy NLP pipeline to the file text, disabling parser and named entities
17     doc = nlp(file, disable=['parser', 'ner'])
18
19     # Lemmatize tokens and remove pronouns ("-PRON-")
20     tokens = [tok.lemma_.lower().strip() for tok in doc if tok.lemma_ != '-PRON-']
21
22     # Remove stopwords and punctuation from the tokens
23     tokens = [tok for tok in tokens if tok not in stopwords and tok not in punctuations]
24
25     # Lemmatize tokens again using NLTK WordNetLemmatizer
26     s = [lem.lemmatize(word) for word in tokens]
27     tokens = ' '.join(s)
28
29     # Remove square brackets and extra whitespaces
30     article_text = re.sub(r'\[[0-9]*\]', ' ', tokens)
31     article_text = re.sub(r'\s+', ' ', article_text)
32
33     # Remove non-alphabetic characters
34     formatted_article_text = re.sub('[^a-zA-Z]', ' ', article_text)
35     formatted_article_text = re.sub(r'\s+', ' ', formatted_article_text)
36
37     # Remove words with 3 characters or less
38     formatted_article_text = re.sub(r'\W*\b\w{1,3}\b', '', formatted_article_text)
39
40     return formatted_article_text
41

```

In [10]:

```
1 clean_data=[remove_stop_word(file) for file in data]
2 clean_data
```

erapy social anxiety disorder include following paroxetine sertraline approved venlafaxine approved escitalopram fluoxetine fluvoxaminephene lzinemoclobemide approve united states tricyclic antidepressant tcas b eta blocker propranolol selected anticonvulsant gabapentin pregabalin valproic acid topiramate tiagabine control study demonstrate efficacy psychopharmacologic intervention specific phobia need administration s hort benzodiazepine useful temporary anxiety relief specific situation agent consider agoraphobia include following ssri escitalopram citalop ram fluoxetine fluvoxamine paroxetine sertraline venlafaxine reboxetin esome clomipramine imipramine some benzodiazepine alprazolam lorazepam diazepam clonazepam mirtazapinemoclobemidepsychotherapeutic interventi on helpful treat phobic disorder include following social anxiety diso rder social phobia self exposure monotherapy computer base exposure tr aining clinician lead exposure combination therapy self exposure cogni tive behavioral therapy self help manual specific phobia base approach include gradual desensitization relaxation breathing control technique exposure therapy agoraphobia combination exposure therapy relaxation b reathe retrainingsee treatment medication detail phobia define irratio nal fear produce conscious avoidance fear subject activity situation a ffected person usually recognize reaction excessive collectively phobi

In []:

```
1 #vectorizing the clean data
```

In [11]:

```
1 from sklearn.feature_extraction.text import CountVectorizer
2 from sklearn.preprocessing import normalize
3 import pandas
```

In [12]:

```
1 vectorizer = CountVectorizer(stop_words=stopwords.words('english'), ngram_range=(
2 vectorizer
```

Out[12]:

```
CountVectorizer(ngram_range=(2, 2),
                stop_words=['i', 'me', 'my', 'myself', 'we', 'our', 'ours',
                            'ourselves', 'you', "you're", "you've", "you'll",
                            "you'd", 'your', 'yours', 'yourself', 'yourselves',
                            'he', 'him', 'his', 'himself', 'she', "she's",
                            'her', 'hers', 'herself', 'it', "it's", 'its',
                            'itself', ...])
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [17]:

```
1 X=vectorizer.transform(clean_data)
2 X
```

Out[17]:

```
<30x66449 sparse matrix of type '<class 'numpy.int64'>'
  with 72731 stored elements in Compressed Sparse Row format>
```

In [18]:

```
1 # Get the feature names using get_feature_names_out
2 feature_names = vectorizer.get_feature_names_out()
```

In [19]:

```

1 # Convert the sparse matrix to a DataFrame using the feature names
2 data_final = pd.DataFrame(X.toarray(), columns=feature_names)
3 data_final

```

Out[19]:

	aase smith	aasly aharon	abandon lead	abbass premaxillary	abbott barrett	abbott melhem	abbreviation head	abdelrahman abbass	abdoman chest
0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0
3	1	0	0	1	0	0	0	1	0
4	0	0	0	0	0	0	0	0	0
5	0	0	0	0	0	0	0	0	0
6	0	1	0	0	0	0	0	0	0
7	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	1	0	0
10	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0
12	0	0	0	0	0	0	0	0	1
13	0	0	0	0	0	0	0	0	0
14	0	0	0	0	1	0	0	0	0
15	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0
17	0	0	0	0	0	0	0	0	0
18	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0
20	0	0	1	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0
22	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0
24	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0
26	0	0	0	0	0	0	0	0	0
27	0	0	0	0	0	1	0	0	0
28	0	0	0	0	0	0	0	0	0
29	0	0	0	0	0	0	0	0	0

30 rows × 66449 columns

In [20]:

```
1 from sklearn.feature_extraction.text import TfidfTransformer
2
3 tran=TfidfTransformer().fit(data_final.values)
4
5 X=tran.transform(X).toarray()
6
7 X = normalize(X)
```

In [21]:

```
1 X
```

Out[21]:

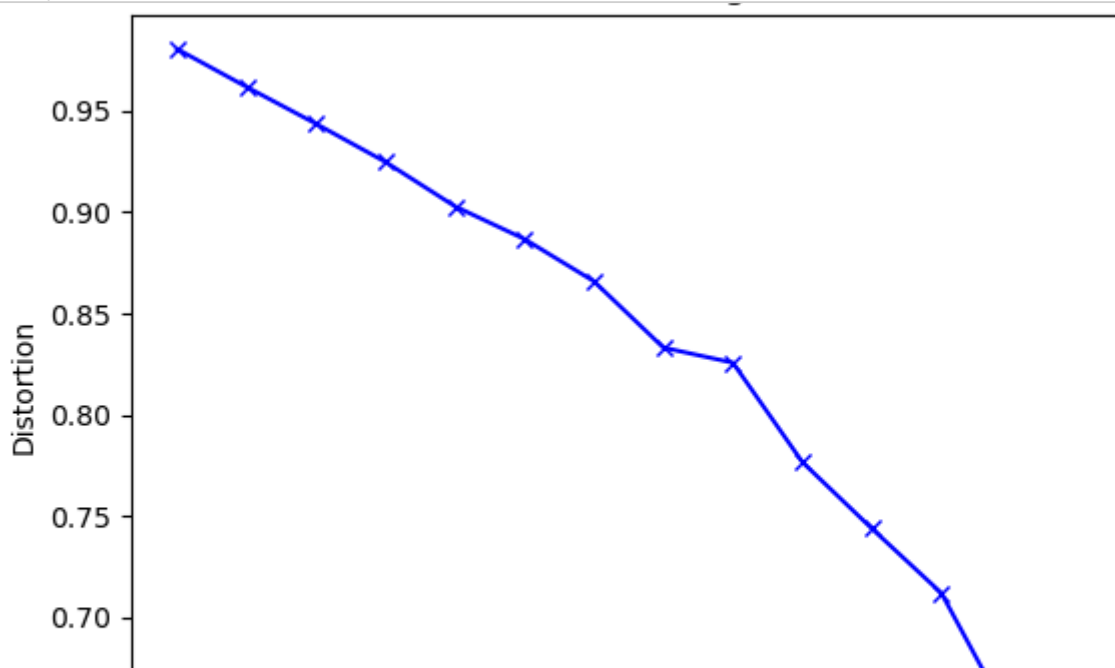
```
array([[0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       ...,
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.],
       [0., 0., 0., ..., 0., 0., 0.]])
```


In [22]:

```

1  from sklearn.cluster import KMeans
2  from sklearn import metrics
3  from scipy.spatial.distance import cdist
4  import numpy as np
5  import matplotlib.pyplot as plt
6
7
8  distortions = []
9  inertias = []
10 mapping1 = {}
11 mapping2 = {}
12 K = range(1,15)
13
14 for k in K:
15     #Building and fitting the model
16     kmeanModel = KMeans(n_clusters=k).fit(X)
17     kmeanModel.fit(X)
18
19     distortions.append(sum(np.min(cdist(X, kmeanModel.cluster_centers_,
20                                     'euclidean'),axis=1)) / X.shape[0])
21     inertias.append(kmeanModel.inertia_)
22
23     mapping1[k] = sum(np.min(cdist(X, kmeanModel.cluster_centers_,
24                                     'euclidean'),axis=1)) / X.shape[0]
25     mapping2[k] = kmeanModel.inertia_
26
27 for key,val in mapping1.items():
28     print(str(key)+' : '+str(val))
29
30 plt.plot(K, distortions, 'bx-')
31 plt.xlabel('Values of K')
32 plt.ylabel('Distortion')
33 plt.title('The Elbow Method using Distortion')
34 plt.show()
35

```



In [34]:

```

1  from sklearn.feature_extraction.text import CountVectorizer
2  from sklearn.cluster import KMeans
3  import pandas as pd
4
5  # Assuming you have a variable 'text_data' that contains the text data as a list
6
7  # Initialize the CountVectorizer
8  vectorizer = CountVectorizer()
9
10 # Fit and transform the text data 'text_data'
11 X_vectorized = vectorizer.fit_transform(clean_data)
12
13 # Define the number of clusters ('true_k') and create the KMeans model
14 true_k = 12
15 model = KMeans(n_clusters=true_k, init='k-means++', max_iter=100, n_init=1)
16
17 # Fit the KMeans model to the vectorized data
18 model.fit(X_vectorized)
19
20 # Create a DataFrame to store the vectorized data and cluster labels
21 cluster_result_data = pd.DataFrame(clean_data, columns=['text'])
22 cluster_result_data['group'] = model.predict(X_vectorized)
23
24 # Get the feature names using get_feature_names_out
25 terms = vectorizer.get_feature_names_out()
26
27 # Get the cluster centers and print the top 50 terms for each cluster
28 order_centroids = model.cluster_centers_.argsort()[:, :-1]
29
30 for i in range(true_k):
31     print('Cluster %d:' % i)
32     for ind in order_centroids[i, :50]:
33         print(' %s' % terms[ind])
34

```

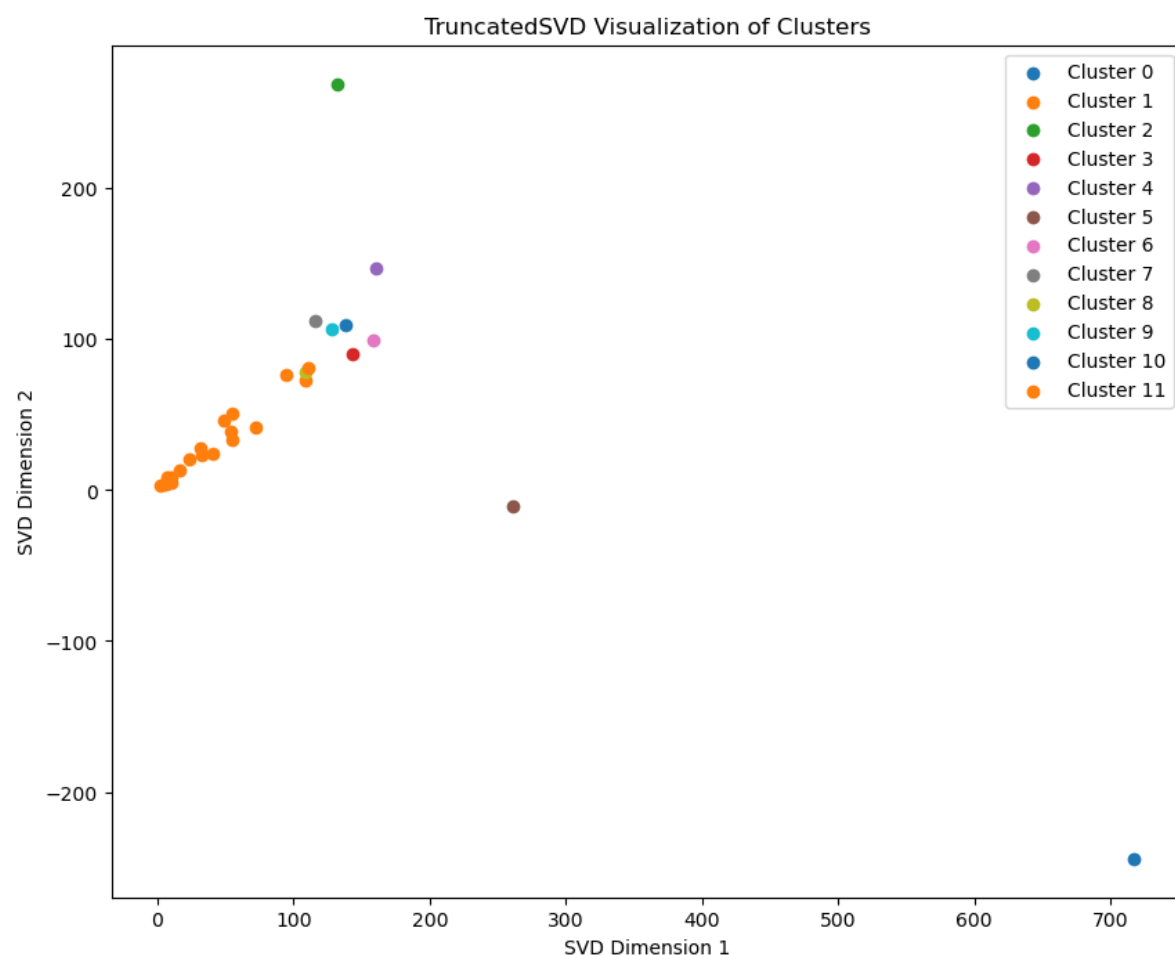
```

failure
wall
develop
surgical
early
infarct
anterior
severe
increase
usually
muscle
intervention
treat
cardiogenic
first
mitral
vfwr
Cluster 11:
heat
patient

```

In [38]:

```
1 import matplotlib.pyplot as plt
2 from sklearn.decomposition import TruncatedSVD
3
4 # Dimensionality reduction using TruncatedSVD
5 svd = TruncatedSVD(n_components=2, random_state=42)
6 X_svd = svd.fit_transform(X_vectorized)
7
8 # Get the cluster labels from the fitted KMeans model
9 cluster_labels = model.labels_
10
11 # Plot the clusters using a scatter plot
12 plt.figure(figsize=(10, 8))
13 for i in range(true_k):
14     plt.scatter(X_svd[cluster_labels == i, 0], X_svd[cluster_labels == i, 1], la
15
16 plt.title('TruncatedSVD Visualization of Clusters')
17 plt.xlabel('SVD Dimension 1')
18 plt.ylabel('SVD Dimension 2')
19 plt.legend()
20 plt.show()
21
```



In []:

1

