

# Unsupervised\_final\_project

December 10, 2024

## 1 Final Project for Unsupervised Learning

This project is based on a Kaggle competition of Disease Symptom Prediction and can be found here (<https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset>). The original competition was designed to see if students could come up with a disease prediction score associated with symptoms. Here, I am using the data a bit differently as to not come up with the same exact project. My project here is to just do an unsupervised clustering of symptoms to see if there are general trends in disease and symptom correlation.

I will use a hierarchical clustering and some other dimensionality reduction techniques to put the patients with similar symptoms together and to see if the clustering of symptoms leads to unique disease groups. Overall, the theory is that patients with the exact same set of symptoms will separate from each other. This in reality isn't practical, but we'll see how well the symptoms separate.

### 1.1 Data Acquisition and EDA

```
[153]: import kagglehub

# Download latest version
path = kagglehub.dataset_download("itachi9604/
↳disease-symptom-description-dataset")

print("Path to dataset files:", path)
```

Path to dataset files:

C:\Users\dvanbooven\.cache\kagglehub\datasets\itachi9604\disease-symptom-description-dataset\versions\2

```
[154]: # Import relevant libraries and then load in the dataset into data

import pandas as pd
import os
from sklearn.feature_extraction.text import TfidfVectorizer
from scipy.cluster.hierarchy import linkage, dendrogram
import matplotlib.pyplot as plt

data = pd.read_csv(os.path.join(path, "dataset.csv"))
```

```
print(data.head())

print(data.shape)

print(data['Disease'].nunique())
```

	Disease	Symptom_1	Symptom_2	Symptom_3	\
0	Fungal infection	itching	skin_rash	nodal_skin_eruptions	
1	Fungal infection	skin_rash	nodal_skin_eruptions	dischromic_patches	
2	Fungal infection	itching	nodal_skin_eruptions	dischromic_patches	
3	Fungal infection	itching	skin_rash	dischromic_patches	
4	Fungal infection	itching	skin_rash	nodal_skin_eruptions	

	Symptom_4	Symptom_5	Symptom_6	Symptom_7	Symptom_8	Symptom_9	\
0	dischromic_patches	NaN	NaN	NaN	NaN	NaN	
1	NaN	NaN	NaN	NaN	NaN	NaN	
2	NaN	NaN	NaN	NaN	NaN	NaN	
3	NaN	NaN	NaN	NaN	NaN	NaN	
4	NaN	NaN	NaN	NaN	NaN	NaN	

	Symptom_10	Symptom_11	Symptom_12	Symptom_13	Symptom_14	Symptom_15	\
0	NaN	NaN	NaN	NaN	NaN	NaN	
1	NaN	NaN	NaN	NaN	NaN	NaN	
2	NaN	NaN	NaN	NaN	NaN	NaN	
3	NaN	NaN	NaN	NaN	NaN	NaN	
4	NaN	NaN	NaN	NaN	NaN	NaN	

	Symptom_16	Symptom_17
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN

(4920, 18)

41

Here we have 18 columns in total. There is 1 disease column, and 17 symptom columns. There are 4920 patients.

[155]: # Notice that the symptoms are spread across different columns. Now take these columns and concatenate to get a single descriptor field

```
data['combined_symptoms'] = data['Symptom_1'] + " " + data['Symptom_2'] + " " +
data['Symptom_3']
```

```
data['combined_symptoms'] = data.loc[:, 'Symptom_1':'Symptom_17'].apply(
    lambda row: ' '.join(row.dropna()), axis=1
)
```

```
print(data)
```

		Disease	Symptom_1 \
0		Fungal infection	itching
1		Fungal infection	skin_rash
2		Fungal infection	itching
3		Fungal infection	itching
4		Fungal infection	itching
...		...	...
4915	(vertigo) Paroymsal	Positional Vertigo	vomiting
4916		Acne	skin_rash
4917	Urinary tract	infection	burning_micturition
4918		Psoriasis	skin_rash
4919		Impetigo	skin_rash

	Symptom_2	Symptom_3	Symptom_4 \
0	skin_rash	nodal_skin_eruptions	dischromic_patches
1	nodal_skin_eruptions	dischromic_patches	NaN
2	nodal_skin_eruptions	dischromic_patches	NaN
3	skin_rash	dischromic_patches	NaN
4	skin_rash	nodal_skin_eruptions	NaN
...	...	...	...
4915	headache	nausea	spinning_movements
4916	pus_filled_pimples	blackheads	scurring
4917	bladder_discomfort	foul_smell_of urine	continuous_feel_of_urine
4918	joint_pain	skin_peeling	silver_like_dusting
4919	high_fever	blister	red_sore_around_nose

	Symptom_5	Symptom_6	Symptom_7	Symptom_8 \
0	NaN	NaN	NaN	NaN
1	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN
...	...	...	...	...
4915	loss_of_balance	unsteadiness	NaN	NaN
4916	NaN	NaN	NaN	NaN
4917	NaN	NaN	NaN	NaN
4918	small_dents_in_nails	inflammatory_nails	NaN	NaN
4919	yellow_crust_ooze	NaN	NaN	NaN

	Symptom_9	Symptom_10	Symptom_11	Symptom_12	Symptom_13	Symptom_14 \
0	NaN	NaN	NaN	NaN	NaN	NaN
1	NaN	NaN	NaN	NaN	NaN	NaN
2	NaN	NaN	NaN	NaN	NaN	NaN
3	NaN	NaN	NaN	NaN	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN	NaN

```

...      ...      ...      ...      ...      ...
4915      NaN      NaN      NaN      NaN      NaN      NaN
4916      NaN      NaN      NaN      NaN      NaN      NaN
4917      NaN      NaN      NaN      NaN      NaN      NaN
4918      NaN      NaN      NaN      NaN      NaN      NaN
4919      NaN      NaN      NaN      NaN      NaN      NaN

```

```

      Symptom_15 Symptom_16 Symptom_17 \
0      NaN      NaN      NaN
1      NaN      NaN      NaN
2      NaN      NaN      NaN
3      NaN      NaN      NaN
4      NaN      NaN      NaN

```

```

...      ...      ...      ...
4915      NaN      NaN      NaN
4916      NaN      NaN      NaN
4917      NaN      NaN      NaN
4918      NaN      NaN      NaN
4919      NaN      NaN      NaN

```

```

                                combined_symptoms
0      itching skin_rash nodal_skin_eruptions disc...
1      skin_rash nodal_skin_eruptions dischromic _...
2      itching nodal_skin_eruptions dischromic _pat...
3      itching skin_rash dischromic_patches
4      itching skin_rash nodal_skin_eruptions
...
4915      vomiting headache nausea spinning_movement...
4916      skin_rash pus_filled_pimples blackheads sc...
4917      burning_micturition bladder_discomfort foul...
4918      skin_rash joint_pain skin_peeling silver_l...
4919      skin_rash high_fever blister red_sore_arou...

```

[4920 rows x 19 columns]

Now the combined\_symptoms column contains a concatenated version of all of the symptoms in symptom columns 1 through 17

```

[156]: # The 19th column should now have some sort of data in it. For safety, let's
      ↪ check if it is still na, and if so remove it

      # Now inspect any rows with an Na

      rows_with_na = data[data['combined_symptoms'].isna()]

      print(rows_with_na)

```

Empty DataFrame

Columns: [Disease, Symptom\_1, Symptom\_2, Symptom\_3, Symptom\_4, Symptom\_5, Symptom\_6, Symptom\_7, Symptom\_8, Symptom\_9, Symptom\_10, Symptom\_11, Symptom\_12, Symptom\_13, Symptom\_14, Symptom\_15, Symptom\_16, Symptom\_17, combined\_symptoms]  
Index: []

This is an empty dataframe which means we have at least 1 symptom within each of the rows of the combined\_symptoms column.

```
[157]: #Let's look at the most and least frequent symptoms in the dataset.

from collections import Counter

# Step 1: Create the 'combined_symptoms' column if it doesn't exist
data['combined_symptoms'] = data.loc[:, 'Symptom_1':'Symptom_17'].apply(
    lambda row: ' '.join(row.dropna()), axis=1
)
all_symptoms = ' '.join(data['combined_symptoms'])

symptoms_list = all_symptoms.split()
symptom_counts = Counter(symptoms_list)
top_10_symptoms = symptom_counts.most_common(10)
print(top_10_symptoms)
lowest_10_symptoms = symptom_counts.most_common()[-10:] # Slice to get the
    ↳ last 10 (least common)
print(lowest_10_symptoms)
number_of_unique_symptoms = len(symptom_counts.keys())

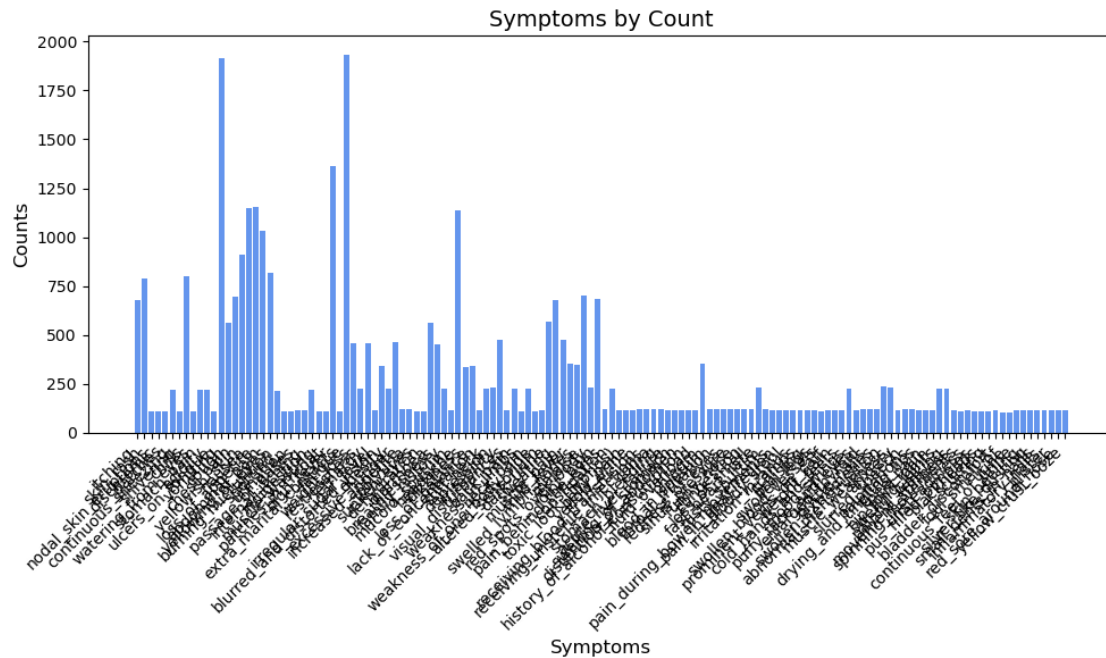
print(f"Number of unique symptoms: {number_of_unique_symptoms}")
```

```
[('fatigue', 1932), ('vomiting', 1914), ('high_fever', 1362),
('loss_of_appetite', 1152), ('nausea', 1146), ('headache', 1134),
('abdominal_pain', 1032), ('yellowish_skin', 912), ('yellowing_of_eyes', 816),
('chills', 798)]
[('dehydration', 108), ('weakness_in_limbs', 108), ('weakness_of_one_body_side',
108), ('swollen_blood_vessels', 108), ('spinning_movements', 108),
('pus_filled_pimples', 108), ('blackheads', 108), ('scurring', 108),
('foul_smell_of', 102), ('urine', 102)]
Number of unique symptoms: 134
```

```
[175]: # Plot the above distribution
symptoms, counts = zip(*symptom_counts.items())

# Create a bar chart
plt.figure(figsize=(10, 6))
plt.bar(symptoms, counts, color='cornflowerblue')
plt.xlabel('Symptoms', fontsize=12)
plt.ylabel('Counts', fontsize=12)
plt.title('Symptoms by Count', fontsize=14)
```

```
plt.xticks(rotation=45, ha='right') # Rotate x-axis labels for readability
plt.tight_layout()
plt.show()
```



[176]: # Let's take this a step further and do a word cloud

```
from wordcloud import WordCloud

wordcloud = WordCloud(width=800, height=400, background_color='white').
    generate(all_symptoms)

plt.figure(figsize=(10, 5))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis('off') # Hide the axes
plt.title("Word Cloud of Disease Symptoms")
plt.show()
```

fatigue high fever cough high fever malaise phlegm  
dark urine loss of appetite itching fatigue  
loss\_of\_appetite abdominal\_pain  
breathlessness sweating abdominal\_pain yellowing\_of\_eyes continuous sneezing chills  
dizziness loss\_of\_balance depression irritability yellowing\_of\_eyes muscle pain  
nausea loss\_of\_appetite weight loss restlessness  
dark urine nausea weight loss restlessness  
loss\_of\_appetite mild fever headache nausea  
fatigue cough blurred\_and\_distorted\_vision excessive\_hunger yellowish\_skin nausea  
yellowing\_of\_eyes malaise joint\_pain dark\_urine  
fatigue lethargy abdominal\_pain diarrhea abnormal menstruation fatigue  
throat\_irritation redness\_of\_eyes swelled lymph nodes malaise  
vomiting high\_fever high\_fever breathlessness  
malaise red spots over body vomiting yellowish\_skin  
mild\_fever yellowing\_of\_eyes muscle pain itching  
increased appetite polyuria chills muscle pain chills  
phlegm chest\_pain high\_fever headache chills vomiting  
pain behind the eyes back pain fatigue yellowish\_skin

## 1.2 Convert text data to numerical form using TF-IDF

```
[177]: # Text data being unstructured cannot be directly processed by these clustering_
        ↪ algorithms unless it is transformed into numerical format
```

```
vectorizer = TfidfVectorizer(stop_words='english')
tfidf_matrix = vectorizer.fit_transform(data['combined_symptoms'])
```

### 1.3 Model Generation

```
[178]: from sklearn.cluster import AgglomerativeClustering
        from scipy.cluster.hierarchy import dendrogram, linkage

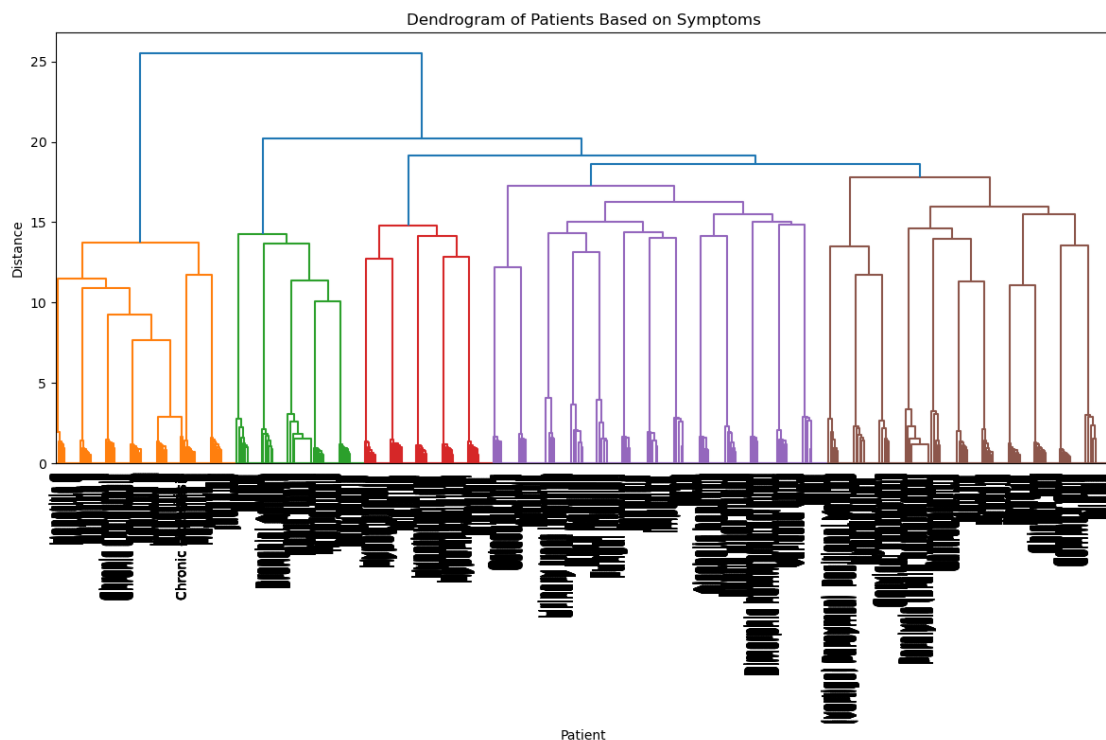
model = AgglomerativeClustering(
    n_clusters=None, # None means do not predefine the number of clusters
    distance_threshold=1.5, # Set a threshold for merging clusters
    linkage='ward'
)
model.fit(tfidf_matrix.toarray())
```

```
[178]: AgglomerativeClustering(distance_threshold=1.5, n_clusters=None)
```

Now we will perform a simple hierarchical clustering to see how the patients separate from each other.

```
[179]: linkage_matrix = linkage(tfidf_matrix.toarray(), method='ward')

plt.figure(figsize=(12, 8))
dendrogram(linkage_matrix, labels=data['Disease'].values, leaf_rotation=90,
            leaf_font_size=10)
plt.title("Dendrogram of Patients Based on Symptoms")
plt.xlabel("Patient")
plt.ylabel("Distance")
plt.tight_layout()
plt.show()
```



```
[180]: # Sanity check and to put the clusters within the main table
data['Cluster'] = model.labels_
print(data[['Disease', 'Cluster']])
```

	Disease	Cluster
0	Fungal infection	18
1	Fungal infection	114
2	Fungal infection	18
3	Fungal infection	61
4	Fungal infection	86
...	...	...
4915	(vertigo) Paroymsal Positional Vertigo	5
4916	Acne	50



4917	Urinary tract infection	30
4918	Psoriasis	27
4919	Impetigo	19

[4920 rows x 2 columns]

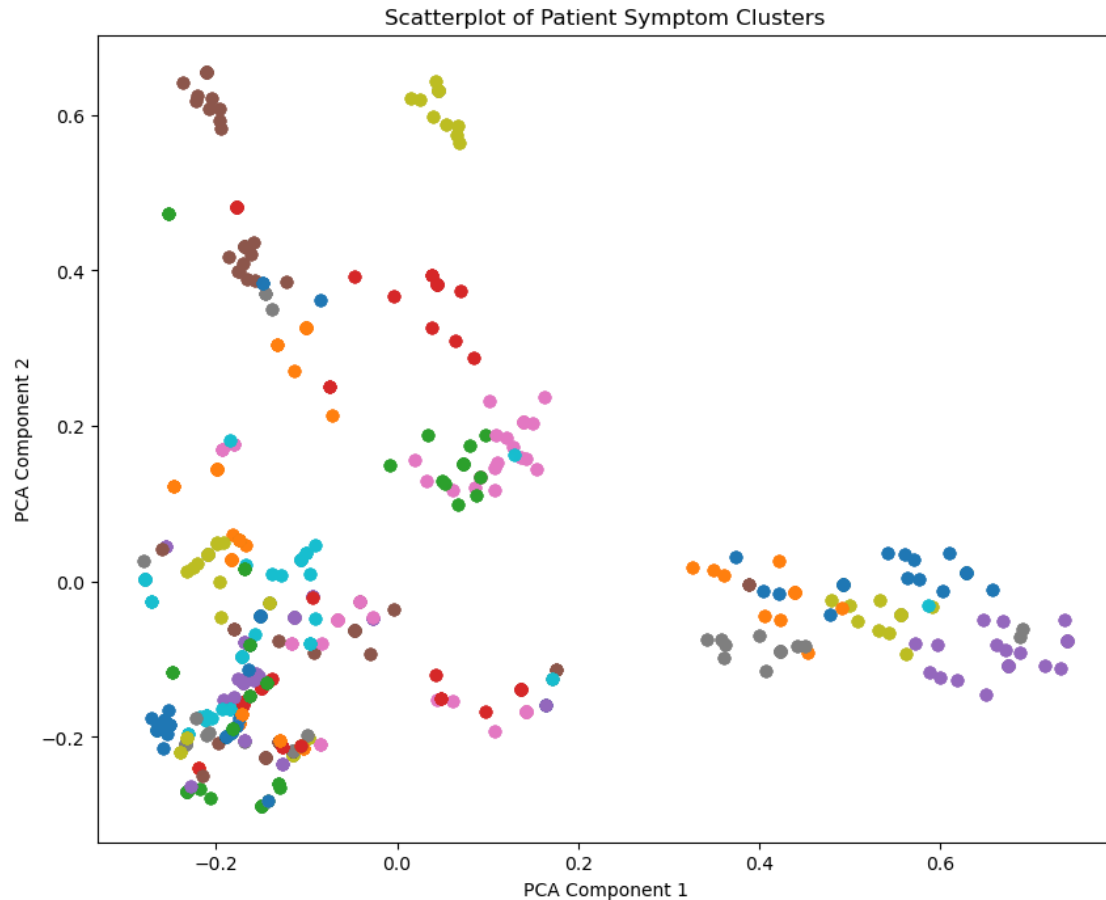
```
[181]: # The dendrogram is nice, but let's do a PCA that will show us a scatterplot of
      ↪ the disease clusters

from sklearn.decomposition import PCA

pca = PCA(n_components=2)
reduced_data = pca.fit_transform(tfidf_matrix.toarray())

plt.figure(figsize=(10, 8))
for cluster in sorted(data['Cluster'].unique()):
    cluster_data = reduced_data[data['Cluster'] == cluster]
    plt.scatter(
        cluster_data[:, 0], cluster_data[:, 1],
        label=f'Cluster {cluster}', alpha=0.7
    )

plt.title("Scatterplot of Patient Symptom Clusters")
plt.xlabel("PCA Component 1")
plt.ylabel("PCA Component 2")
plt.show()
```



## 1.4 Model Refinement

This is a lot of clusters. You can kind of see the patterns if you look closely, but this is very noisy. So let's see if we can reduce the number of clusters.

[182]: *# Let's apply the elbow method to find the optimal number of clusters.*

```
from sklearn.cluster import KMeans

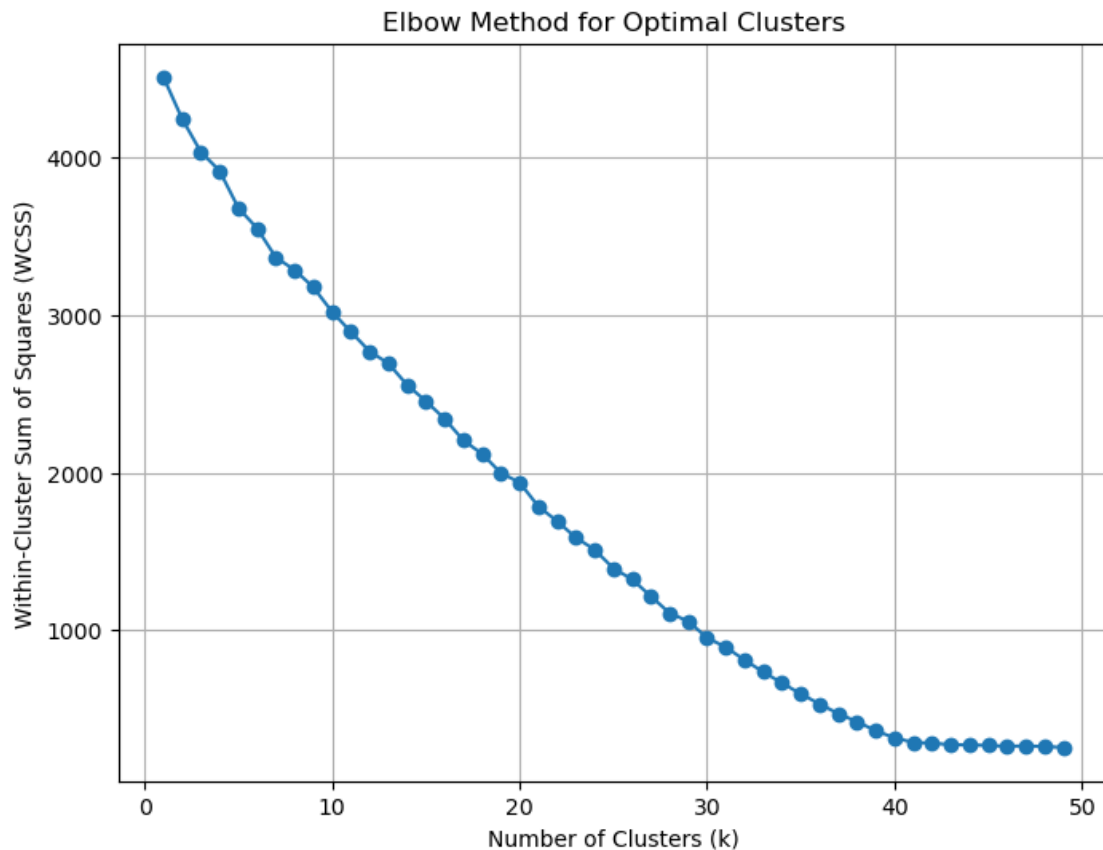
wcss = [] # List to store WCSS for each k
k_values = range(1, 50) # Number of clusters to try (1 to 50)

for k in k_values:
    kmeans = KMeans(n_clusters=k, random_state=42)
    kmeans.fit(tfidf_matrix)
    wcss.append(kmeans.inertia_) # Inertia is the WCSS

plt.figure(figsize=(8, 6))
plt.plot(k_values, wcss, marker='o')
```

```
plt.title("Elbow Method for Optimal Clusters")
plt.xlabel("Number of Clusters (k)")
plt.ylabel("Within-Cluster Sum of Squares (WCSS)")
plt.grid()
plt.show()
```

C:\Users\dvanbooven\AppData\Local\anaconda3\lib\site-packages\sklearn\cluster\\_kmeans.py:870: FutureWarning: The default value of `n\_init` will change from 10 to 'auto' in 1.4. Set the value of `n\_init` explicitly to suppress the warning  
 warnings.warn(



```
[183]: # Find the max value of cluster...

n_clusters = len(set(data['Cluster'])) # Count unique cluster labels

print(f"Number of clusters found: {n_clusters}")
```

Number of clusters found: 116

```
[184]: # Let's reduce the number of clusters from 116 to 40 and redo the
        ↪ agglomerative clustering to see what that does to the analysis.
```

```
model = AgglomerativeClustering(
    n_clusters=40, # None means do not predefine the number of clusters
    linkage='ward'
)
model.fit(tfidf_matrix.toarray())
data['Cluster'] = model.labels_
print(data[['Disease', 'Cluster']])
```

	Disease	Cluster
0	Fungal infection	3
1	Fungal infection	3
2	Fungal infection	3
3	Fungal infection	3
4	Fungal infection	3
...	...	...
4915	(vertigo) Paroymsal Positional Vertigo	12
4916	Acne	8
4917	Urinary tract infection	1
4918	Psoriasis	33
4919	Impetigo	21

[4920 rows x 2 columns]

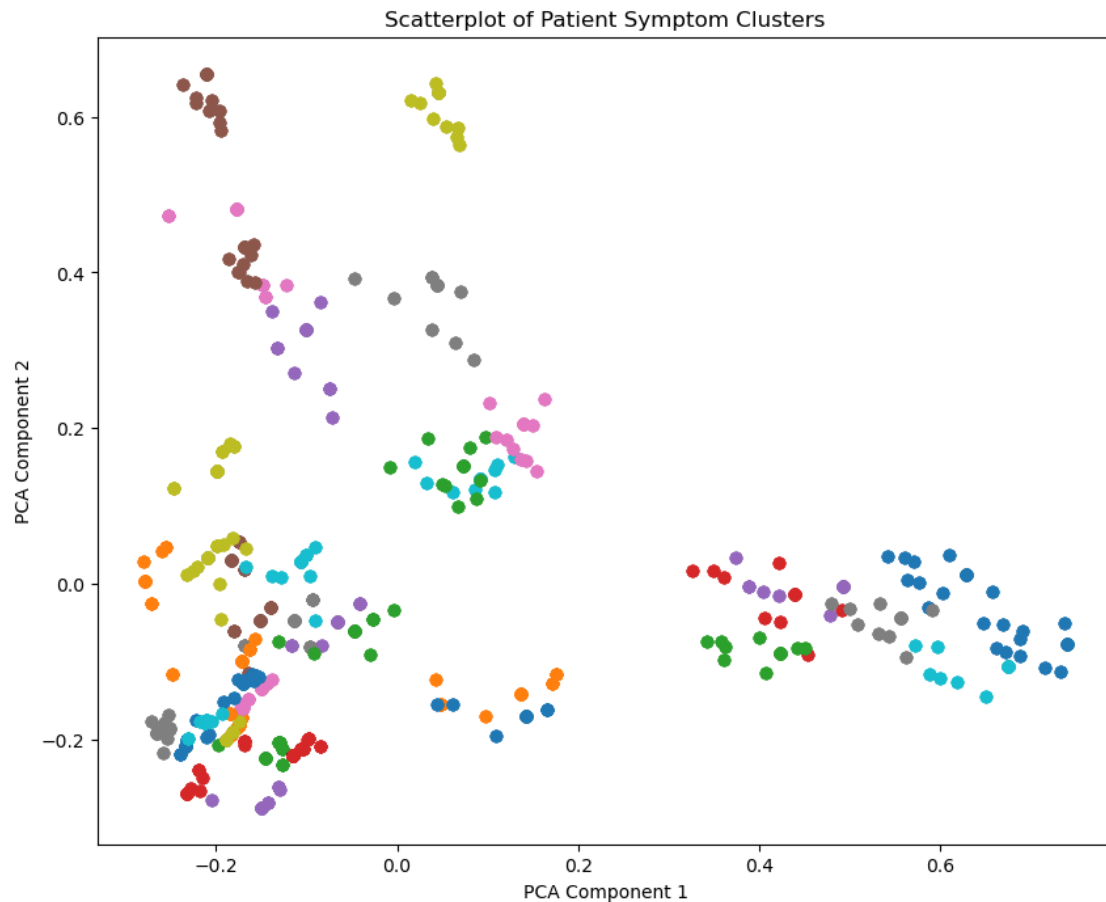
The top 5 values for “fungal infection” now all appear to be headed towards cluster 3 which is a good sign

```
[185]: # Let's redo the scatterplot PCA to show the 40 cluster model
```

```
pca = PCA(n_components=2)
reduced_data = pca.fit_transform(tfidf_matrix.toarray())

plt.figure(figsize=(10, 8))
for cluster in sorted(data['Cluster'].unique()):
    cluster_data = reduced_data[data['Cluster'] == cluster]
    plt.scatter(
        cluster_data[:, 0], cluster_data[:, 1],
        label=f'Cluster {cluster}', alpha=0.7
    )

plt.title("Scatterplot of Patient Symptom Clusters")
plt.xlabel("PCA Component 1")
plt.ylabel("PCA Component 2")
plt.show()
```



## 1.5 Model Evaluation

```
[186]: # Let's look at silhouette scores

from sklearn.metrics import silhouette_score

# get model labels
labels = model.labels_

# Calculate silhouette score
score = silhouette_score(tfidf_matrix, labels)
print(f"Silhouette Score: {score}")

# Remember a silhouette score near 1 are well-separated and cohesive, near 0 is
# unclear, and -1 shows misassignment
```

Silhouette Score: 0.7698943712079074

```
[187]: import numpy as np

top_terms = {}

feature_names = np.array(vectorizer.get_feature_names_out())

for cluster in np.unique(labels):
    cluster_indices = np.where(labels == cluster)[0]
    cluster_data = tfidf_matrix[cluster_indices]
    avg_tfidf_per_term = cluster_data.mean(axis=0).A1
    top_terms[cluster] = feature_names[np.argsort(avg_tfidf_per_term)[::-1][:
↪10]] # Top 10 terms

# Display top terms for each cluster
for cluster, terms in top_terms.items():
    print(f"Cluster {cluster}: {' '.join(terms)}")
```

Cluster 0: dark\_urine, yellowing\_of\_eyes, joint\_pain, yellowish\_skin, abdominal\_pain, nausea, loss\_of\_appetite, vomiting, mild\_fever, muscle\_pain

Cluster 1: bladder\_discomfort, continuous\_feel\_of\_urine, urine, foul\_smell\_of, burning\_micturition, extra\_marital\_contacts, family\_history, fast\_heart\_rate, fatigue, increased\_appetite

Cluster 2: spotting\_, urination, burning\_micturition, stomach\_pain, itching, skin\_rash, family\_history, foul\_smell\_of, fluid\_overload, fatigue

Cluster 3: nodal\_skin\_eruptions, \_patches, dischromic, itching, skin\_rash, lethargy, high\_fever, distention\_of\_abdomen, dizziness, drying\_and\_tingling\_lips

Cluster 4: dehydration, sunken\_eyes, diarrhoea, vomiting, dizziness, drying\_and\_tingling\_lips, enlarged\_thyroid, excessive\_hunger, extra\_marital\_contacts, hip\_joint\_pain

Cluster 5: shivering, watering\_from\_eyes, continuous\_sneezing, chills, fast\_heart\_rate, high\_fever, headache, foul\_smell\_of, fluid\_overload, fatigue

Cluster 6: swollen\_legs, cramps, prominent\_veins\_on\_calf, bruising, swollen\_blood\_vessels, obesity, fatigue, excessive\_hunger, extra\_marital\_contacts, family\_history

Cluster 7: altered\_sensorium, weakness\_of\_one\_body\_side, headache, vomiting, fast\_heart\_rate, hip\_joint\_pain, high\_fever, foul\_smell\_of, fluid\_overload, fatigue

Cluster 8: blackheads, scurring, pus\_filled\_pimples, skin\_rash, yellowish\_skin, high\_fever, headache, foul\_smell\_of, fluid\_overload, fatigue

Cluster 9: belly\_pain, toxic\_look\_, typhos, constipation, diarrhoea, chills, abdominal\_pain, headache, nausea, high\_fever

Cluster 10: weakness\_in\_limbs, neck\_pain, back\_pain, dizziness, loss\_of\_balance, fast\_heart\_rate, high\_fever, headache, foul\_smell\_of, fluid\_overload

Cluster 11: passage\_of\_gases, internal\_itching, indigestion, abdominal\_pain, loss\_of\_appetite, vomiting, fatigue, high\_fever, headache, foul\_smell\_of

Cluster 12: unsteadiness, spinning\_movements, loss\_of\_balance, headache, nausea, vomiting, yellowish\_skin, family\_history, foul\_smell\_of, fluid\_overload

Cluster 13: weight\_loss, dark\_urine, itching, yellowish\_skin, abdominal\_pain,

high\_fever, vomiting, fatigue, headache, foul\_smell\_of

Cluster 14: knee\_pain, hip\_joint\_pain, swelling\_joints, painful\_walking, neck\_pain, joint\_pain, fast\_heart\_rate, headache, foul\_smell\_of, fluid\_overload

Cluster 15: patches\_in\_throat, muscle\_wasting, extra\_marital\_contacts, high\_fever, yellowish\_skin, fatigue, hip\_joint\_pain, headache, foul\_smell\_of, fluid\_overload

Cluster 16: breathlessness, chest\_pain, sweating, vomiting, fatigue, high\_fever, headache, foul\_smell\_of, fluid\_overload, yellowish\_skin

Cluster 17: muscle\_pain, diarrhoea, sweating, chills, headache, nausea, high\_fever, vomiting, excessive\_hunger, extra\_marital\_contacts

Cluster 18: ulcers\_on\_tongue, stomach\_pain, acidity, cough, chest\_pain, vomiting, yellowish\_skin, fluid\_overload, high\_fever, headache

Cluster 19: itching, yellowing\_of\_eyes, yellowish\_skin, abdominal\_pain, nausea, loss\_of\_appetite, vomiting, high\_fever, headache, foul\_smell\_of

Cluster 20: distention\_of\_abdomen, history\_of\_alcohol\_consumption, fluid\_overload, swelling\_of\_stomach, yellowish\_skin, abdominal\_pain, vomiting, loss\_of\_appetite, loss\_of\_smell, dizziness

Cluster 21: yellow\_crust\_ooze, blister, red\_sore\_around\_nose, skin\_rash, high\_fever, yellowish\_skin, fast\_heart\_rate, headache, foul\_smell\_of, fluid\_overload

Cluster 22: red\_spots\_over\_body, swelled\_lymph\_nodes, mild\_fever, lethargy, malaise, itching, skin\_rash, headache, loss\_of\_appetite, high\_fever

Cluster 23: movement\_stiffness, swelling\_joints, painful\_walking, stiff\_neck, muscle\_weakness, fast\_heart\_rate, headache, foul\_smell\_of, fluid\_overload, fatigue

Cluster 24: mucoid\_sputum, family\_history, breathlessness, cough, high\_fever, fatigue, yellowish\_skin, fluid\_overload, hip\_joint\_pain, headache

Cluster 25: loss\_of\_smell, throat\_irritation, redness\_of\_eyes, sinus\_pressure, congestion, runny\_nose, continuous\_sneezing, phlegm, swelled\_lymph\_nodes, muscle\_pain

Cluster 26: irritation\_in\_anus, bloody\_stool, pain\_in\_anal\_region, pain\_during\_bowel\_movements, constipation, yellowish\_skin, high\_fever, headache, foul\_smell\_of, fluid\_overload

Cluster 27: visual\_disturbances, indigestion, acidity, stiff\_neck, depression, blurred\_and\_distorted\_vision, excessive\_hunger, irritability, headache, loss\_of\_appetite

Cluster 28: muscle\_weakness, abnormal\_menstruation, mood\_swings, restlessness, fast\_heart\_rate, irritability, weight\_loss, excessive\_hunger, diarrhoea, sweating

Cluster 29: palpitations, slurred\_speech, drying\_and\_tingling\_lips, anxiety, blurred\_and\_distorted\_vision, excessive\_hunger, irritability, sweating, headache, nausea

Cluster 30: increased\_appetite, polyuria, irregular\_sugar\_level, restlessness, obesity, blurred\_and\_distorted\_vision, lethargy, weight\_loss, excessive\_hunger, fatigue

Cluster 31: lack\_of\_concentration, loss\_of\_balance, dizziness, chest\_pain, headache, yellowish\_skin, fatigue, high\_fever, foul\_smell\_of, fluid\_overload

Cluster 32: receiving\_blood\_transfusion, receiving\_unsterile\_injections,

yellow\_urine, lethargy, dark\_urine, malaise, itching, yellowing\_of\_eyes, yellowish\_skin, abdominal\_pain  
Cluster 33: skin\_peeling, silver\_like\_dusting, small\_dents\_in\_nails, inflammatory\_nails, joint\_pain, skin\_rash, extra\_marital\_contacts, fluid\_overload, fatigue, fast\_heart\_rate  
Cluster 34: family\_history, yellowish\_skin, yellowing\_of\_eyes, nausea, loss\_of\_appetite, fatigue, abnormal\_menstruation, fluid\_overload, history\_of\_alcohol\_consumption, hip\_joint\_pain  
Cluster 35: rusty\_sputum, fast\_heart\_rate, phlegm, breathlessness, cough, chest\_pain, sweating, malaise, chills, high\_fever  
Cluster 36: pain\_behind\_the\_eyes, back\_pain, red\_spots\_over\_body, muscle\_pain, joint\_pain, malaise, skin\_rash, chills, headache, nausea  
Cluster 37: coma, stomach\_bleeding, acute\_liver\_failure, dark\_urine, joint\_pain, yellowing\_of\_eyes, abdominal\_pain, yellowish\_skin, loss\_of\_appetite, nausea  
Cluster 38: blood\_in\_sputum, swelled\_lymph\_nodes, mild\_fever, phlegm, breathlessness, weight\_loss, cough, chest\_pain, malaise, sweating  
Cluster 39: brittle\_nails, enlarged\_thyroid, swollen\_extremities, puffy\_face\_and\_eyes, weight\_gain, cold\_hands\_and\_feets, depression, abnormal\_menstruation, mood\_swings, dizziness

```
[188]: # proof of concept, let's look at a cluster and see if we can determine what
      ↪ the disease was...
```

```
cluster_0_data = data[data['Cluster'] == 0]

# give me a unique set of diseases for this cluster

print(data['Disease'].unique())
```

```
['Fungal infection' 'Allergy' 'GERD' 'Chronic cholestasis' 'Drug Reaction'
 'Peptic ulcer disease' 'AIDS' 'Diabetes' 'Gastroenteritis'
 'Bronchial Asthma' 'Hypertension' 'Migraine' 'Cervical spondylosis'
 'Paralysis (brain hemorrhage)' 'Jaundice' 'Malaria' 'Chicken pox'
 'Dengue' 'Typhoid' 'hepatitis A' 'Hepatitis B' 'Hepatitis C'
 'Hepatitis D' 'Hepatitis E' 'Alcoholic hepatitis' 'Tuberculosis'
 'Common Cold' 'Pneumonia' 'Dimorphic hemmorhoids(piles)' 'Heart attack'
 'Varicose veins' 'Hypothyroidism' 'Hyperthyroidism' 'Hypoglycemia'
 'Osteoarthritis' 'Arthritis' '(vertigo) Parosmal Positional Vertigo'
 'Acne' 'Urinary tract infection' 'Psoriasis' 'Impetigo']
```

```
[ ]:
```

```
[ ]:
```