

MACHINE LEARNING

DISEASE PREDICTION FROM SYMPTOMS



TEAM

Name	Roll Number
D.KRISHNA MURTHY	CB.EN.U4CSE21016
MEENAKSHI S	CB.EN.U4CSE21035
PICHERI LIKITHA	CB.EN.U4CSE21044
RAGALA TEJDEEP	CB.EN.U4CSE21046

Introduction

- People are currently suffering from a variety of diseases.
- Many people are unsure if the symptoms they are experiencing are indicative of a certain disease, and hence they are unable to take the required safeguards.
- Disease prediction of a human means predicting the probability of a patient's disease after examining the combinations of the patient's symptoms.



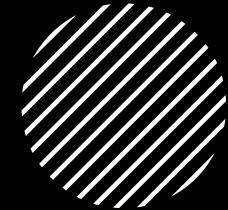
Motivation/significance

- The major goal of this project is to greatly aid physicians in predicting and diagnosing diseases at an early stage, as well as to design a model that will assist people in identifying diseases at an early stage so that risk may be decreased
- This analysis in the medical industry would lead to a streamlined and expedited treatment of patients

Dataset Considered



DATASET DESCRIPTION



- We used the below dataset from Kaggle.
- There are columns containing diseases, their symptoms and their weights.
- <https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset/data?select=dataset.csv>
<https://www.kaggle.com/datasets/itachi9604/disease-symptom-description-dataset/data?select=Symptom-severity.csv>

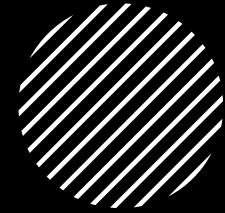
dataset.csv (632.2 kB) Download >

Detail Compact Column 10 of 18 columns

About this file

Disease with its symptoms.

Disease	Symptom_1	Symptom_2	Symptom_3	Symptom_4	Symptom_5	Symptom_6	Symptom_7	Symptom_8
Diseases that may be present	the symptoms experienced during the disease							
41 unique values	vomiting fatigue Other (3408)	17% 14% 69%	vomiting fatigue Other (3648)	18% 8% 74%	fatigue high_fever Other (3870)	15% 7% 79%	high_fever [null] Other (4194)	8% 7% 85%
Fungal infection	itching	skin_rash	nodal_skin_eruptions	dischromic_patches				
Fungal infection	skin_rash	nodal_skin_eruptions	dischromic_patches					



D a t a s e t d e s c r i p t i o n

- We used two csv files for predicting the disease. One csv file contains the symptoms and the disease columns and the next csv file contains the weights of the respective symptom.
- Disease csv file contains 18 features , in which 17 columns are the symptoms and the output is the disease of the respective symptoms
- The weight csv file contains the weight of the respective symptoms which is used for filling NAN symptoms values during encoding

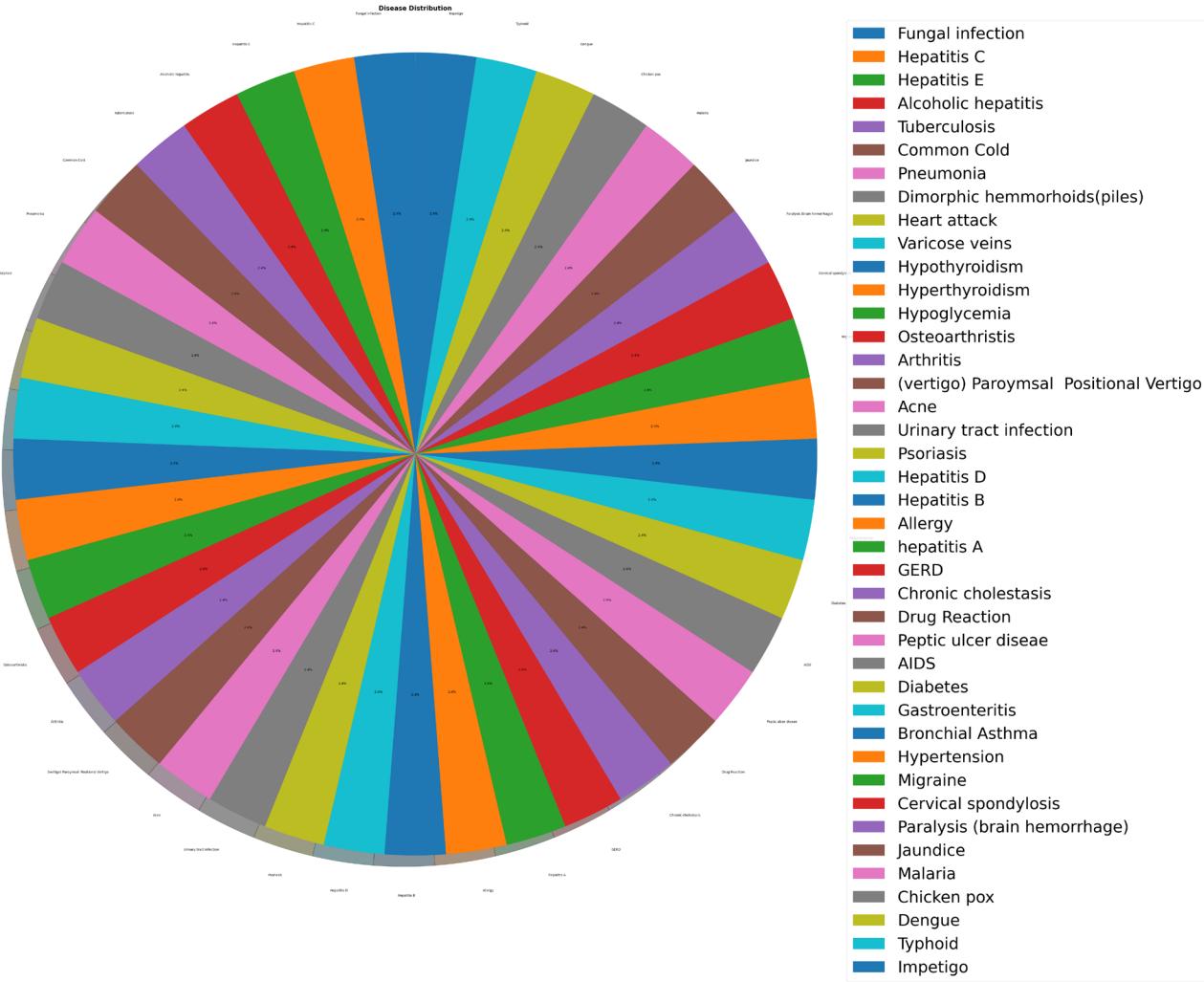
DATA PREPROCESSING

- After collecting that data, as that data is raw data we have to make it suitable for training our machine learning model. By using some python libraries like NumPy, and pandas, we have made that data suitable for machine learning models
- The collected data are preprocessed for the availability of missing values in most of the structured data. Hence, it is essential to fill out the missed data or remove or modify them to enhance the quality of the data set. The preprocessing step also eliminates the commas, punctuations, and white spaces. Once the preprocessing of data has been completed, it is then subjected to feature extraction followed by disease prediction
- The numerical features are scaled using StandardScaler and MinMaxScaler to ensure that they have the same scale and to improve the performance of the model.

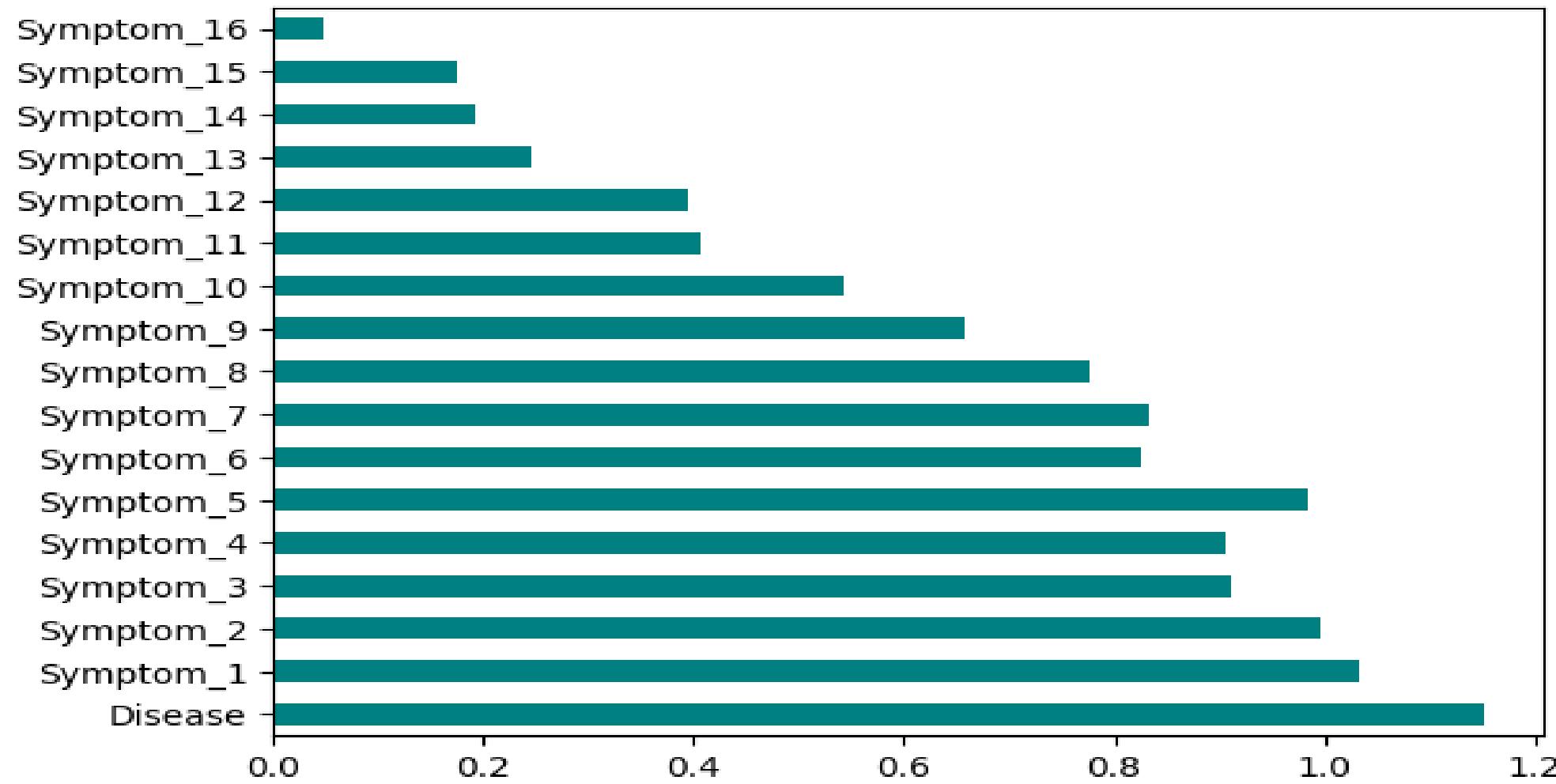


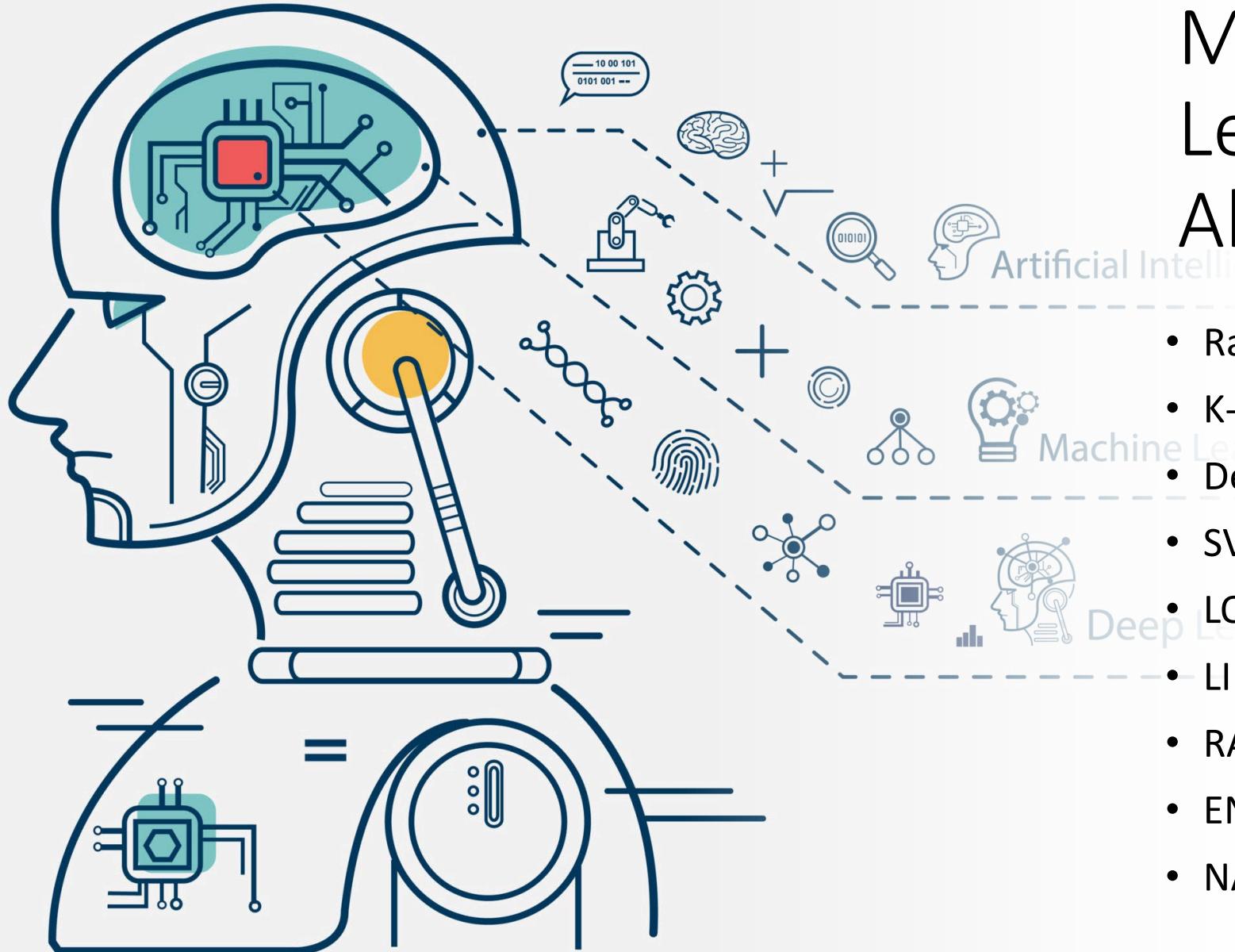
Some Visualizations

Visualization of Diseases



FEATURE SELECTION



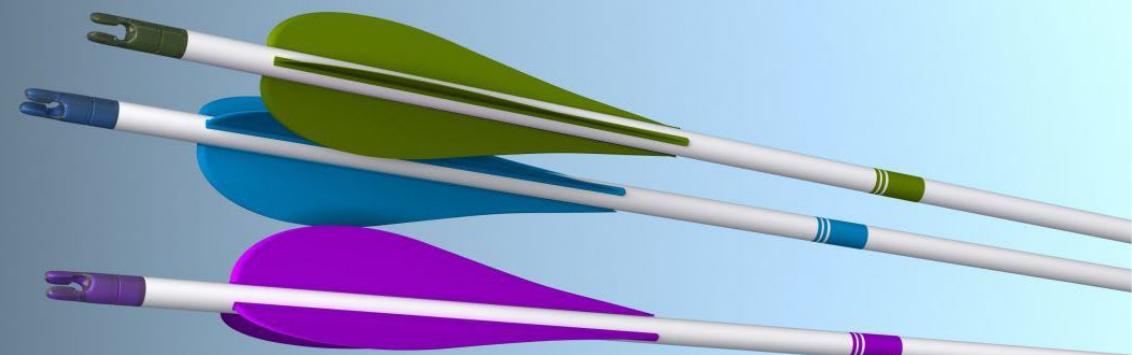


Machine Learning Algorithms Used

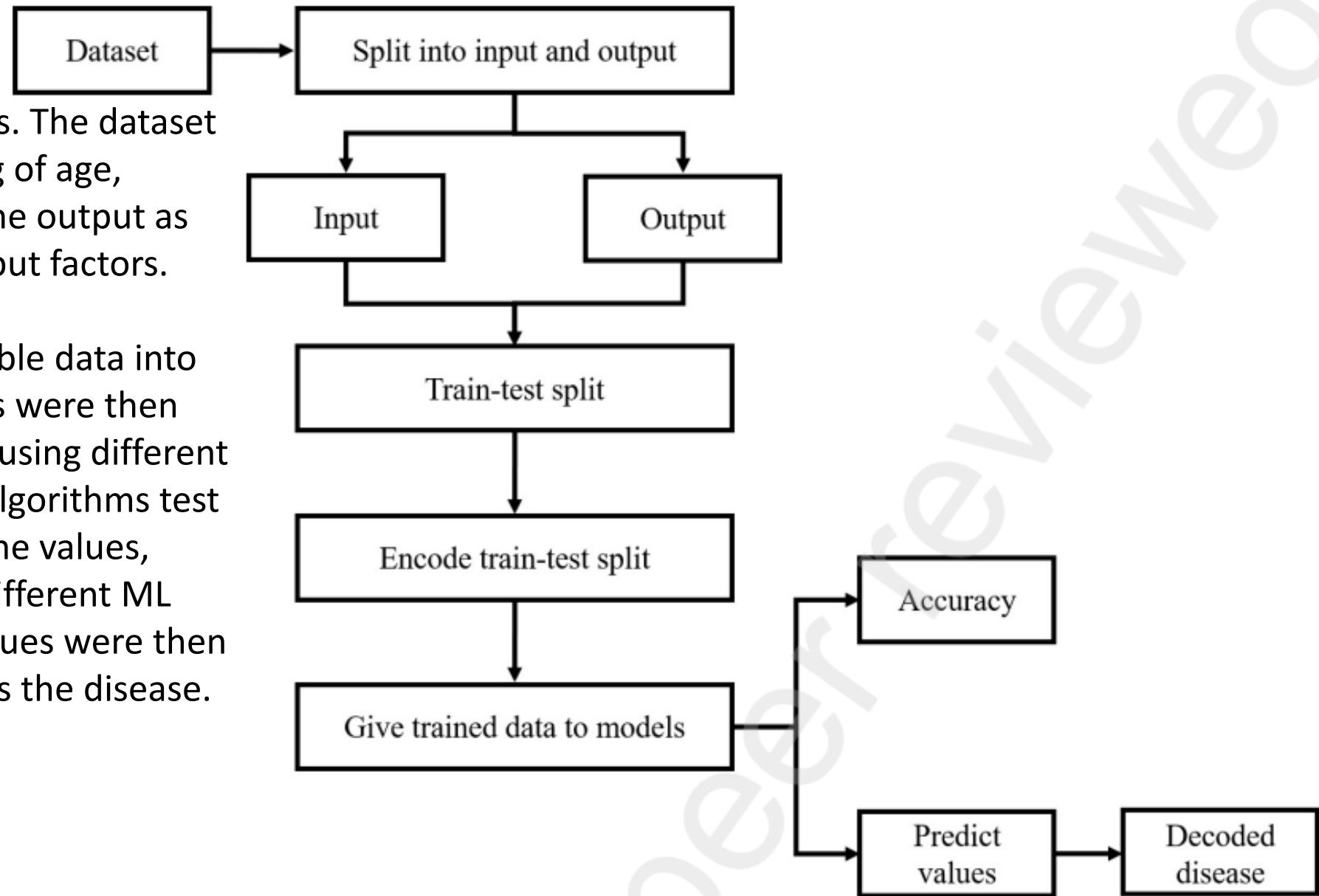
- Random Forest Classifier
- K-Nearest Neighbour Classifier
- Decision Tree Classifier
- SVM
- LOGISTIC REGRESSION
- LINEAR REGRESSION
- RANDOM FOREST
- ENSEMBLE LEARNING
- NAÏVE BAYES

Performance Metrics Considered

- Accuracy
- Precision
- Recall
- ROC-AUC Score
- ROC Graphs
- Confusion Matrix



Functioning of the ML models. The dataset was split into input consisting of age, gender, and symptoms and the output as the diseases based on the input factors.

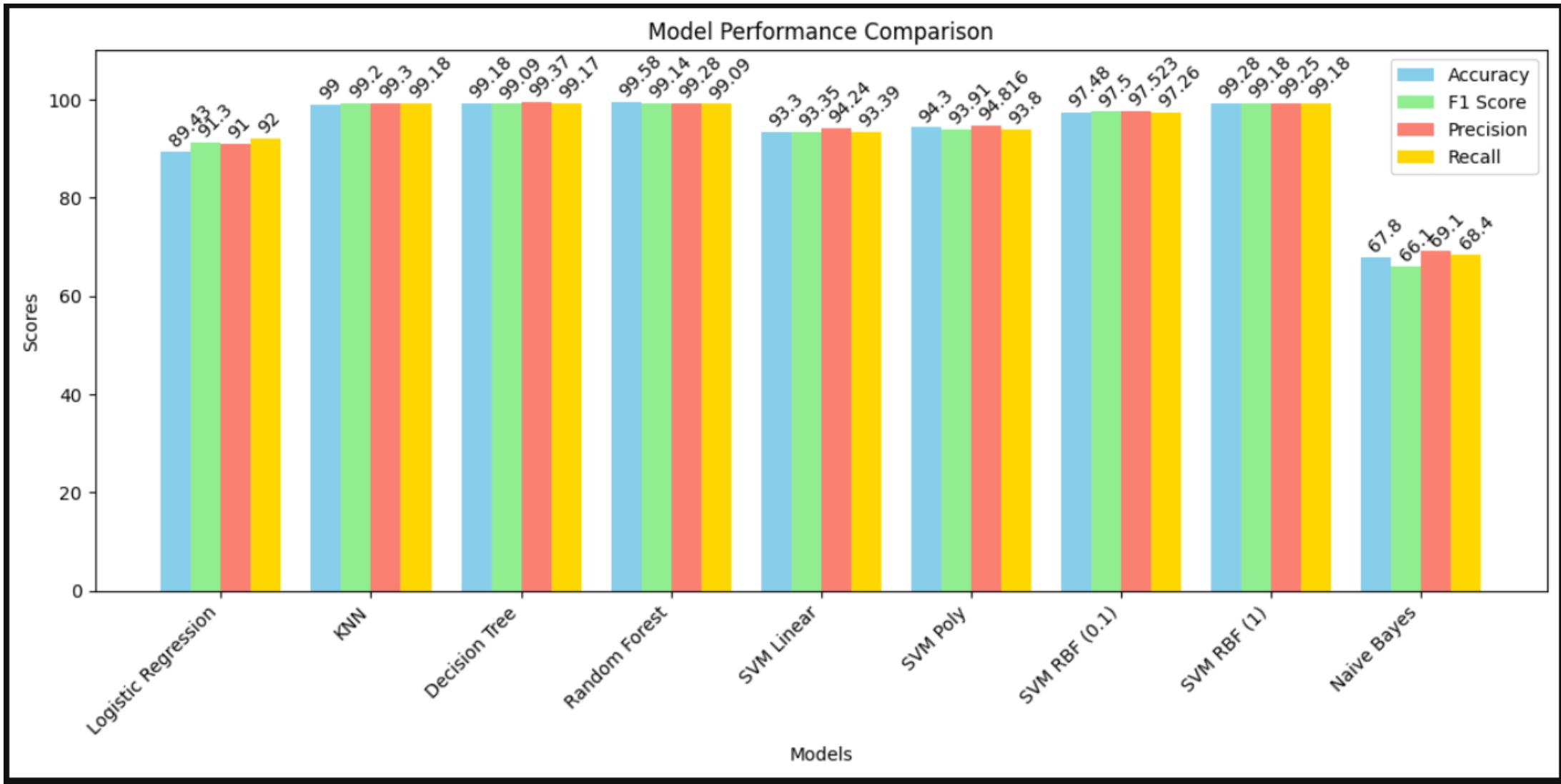


We randomly split the available data into train and test sets. These sets were then encoded and further trained using different algorithms. After which the algorithms test the training set and predict the values, resulting in the accuracy of different ML algorithms. The predicted values were then decoded to give the output as the disease.

Comparision Of Performance metrics of different models

S.NO	MODEL	ACCURACY	F1 SCORE	PRECISION	RECALL
1.	LOGISTIC REGRESSION	89.43%	91.3	91	92
2.	KNN(K=4)	99%	99.2	99.3	99.18
3.	DECISION TREE	99.18%	99.09	99.37	99.17
4.	RANDOM FOREST	99.58%	99.14	99.28	99.09
5.	SVM-LINEAR	93.3%	93.35	94.24	93.39
6.	SVM-POLYNOMIAL (DEGREE=3)	94.3%	93.91	94.816	93.8
7.	SVM-RBF KERNEL (GAMMA=0.1)	97.48%	97.5	97.523	97.26
8.	SVM-RBF KERNEL (GAMMA=1)	99.28%	99.18	99.25	99.18
9.	NAIVE BAYES	67.8%	66.1	69.1	68.4

• COMPARATIVE ANALYSIS



INFERENCE

Different machine learning models were used to examine the prediction of disease for available input dataset.

The table shows the comparative analysis of the accuracy of the training models. From the earlier necessities, Random Forest Algorithm was best with a model accuracy of 99.58%.

Following the Random Forest is the Decision Tree with an accuracy of 99.18%.

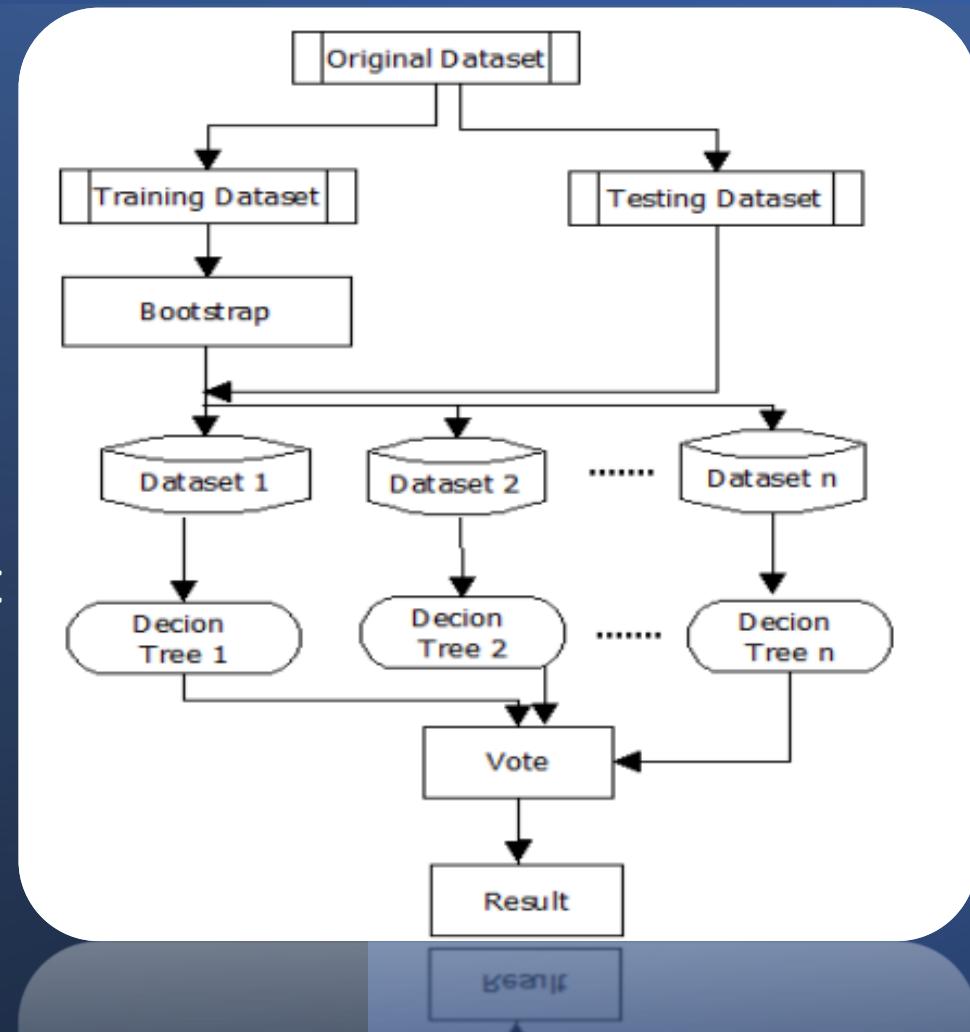
The SVM and KNN models were also very close.

However, the suggested model, that is Random forest model, yields the most accurate result, 99.5% as compared to earlier methods.

Among all the models, we gained the highest accuracy for the Random Forest model of 99.58 %.

Advantage of Random Forest:

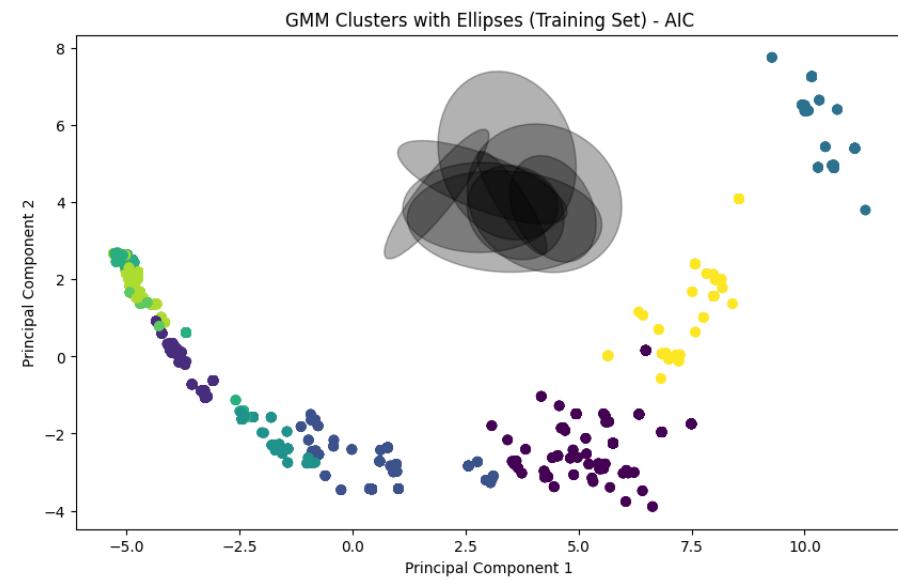
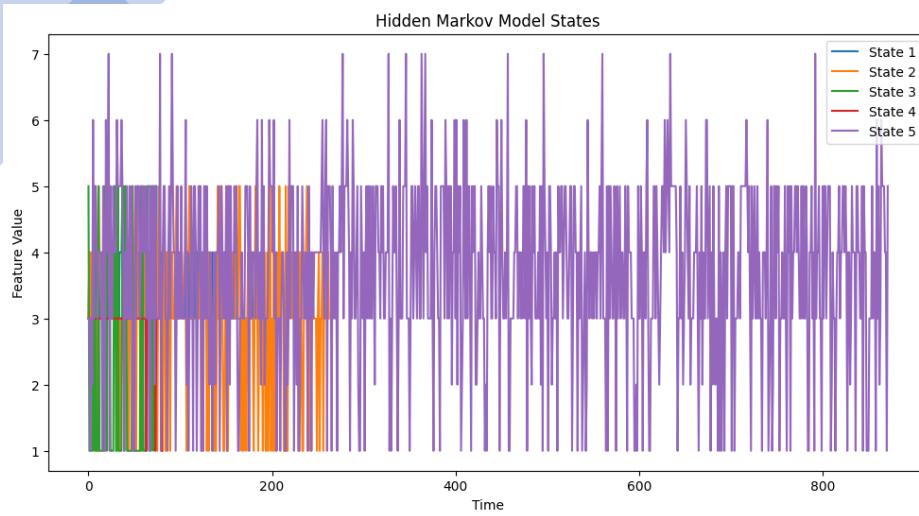
- Different machine learning models were used to examine the prediction of disease for available input dataset.
- While training the model, the decision forests that are formed while concluding are pruned as soon as they encounter a weak symptom or a symptom that does not occur in a location.
- Thus, Random Forest Algorithm minimizes the cost whilst predicting a more realistic model.



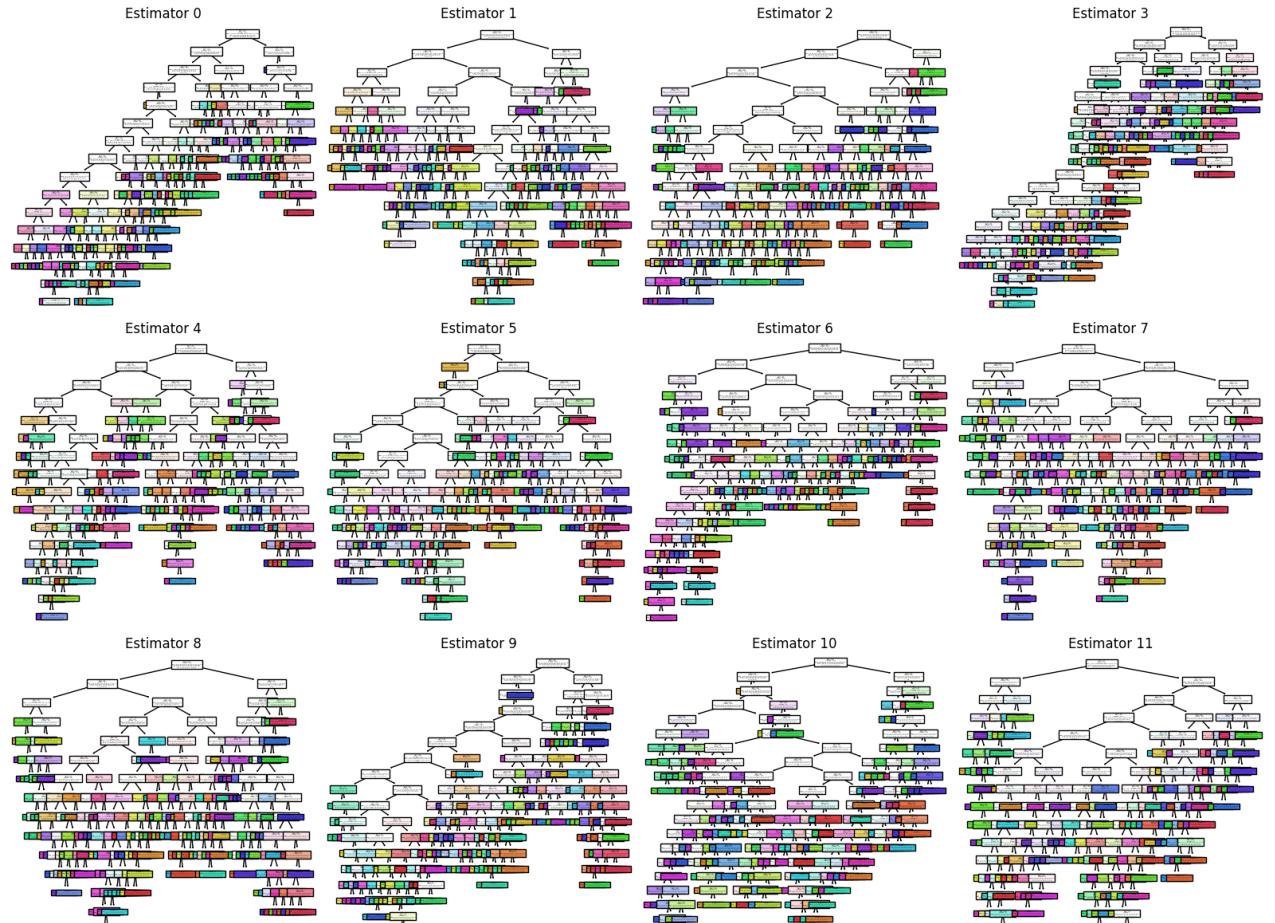
REVIEW 2

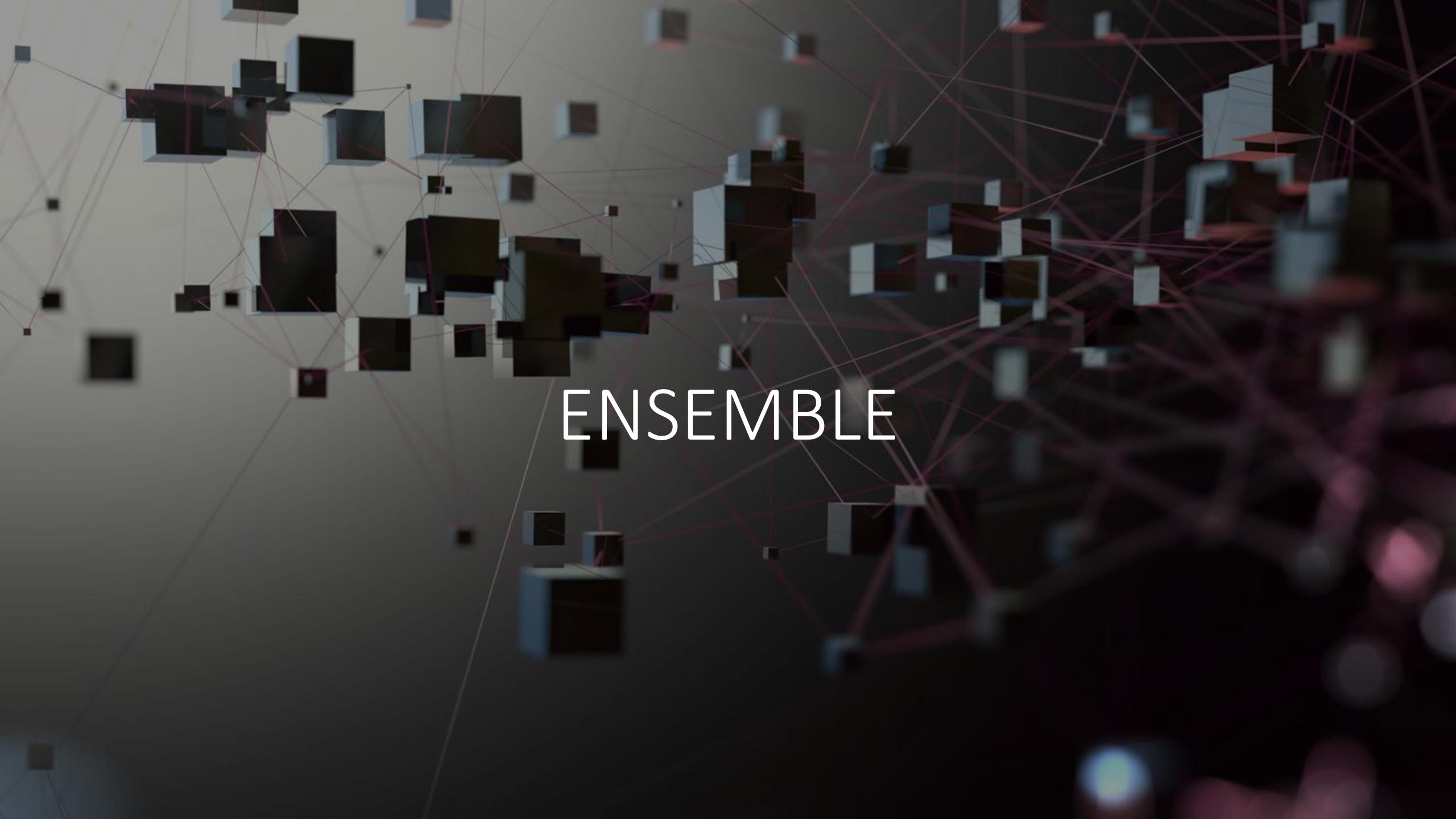


PROBABILISTIC MODEL

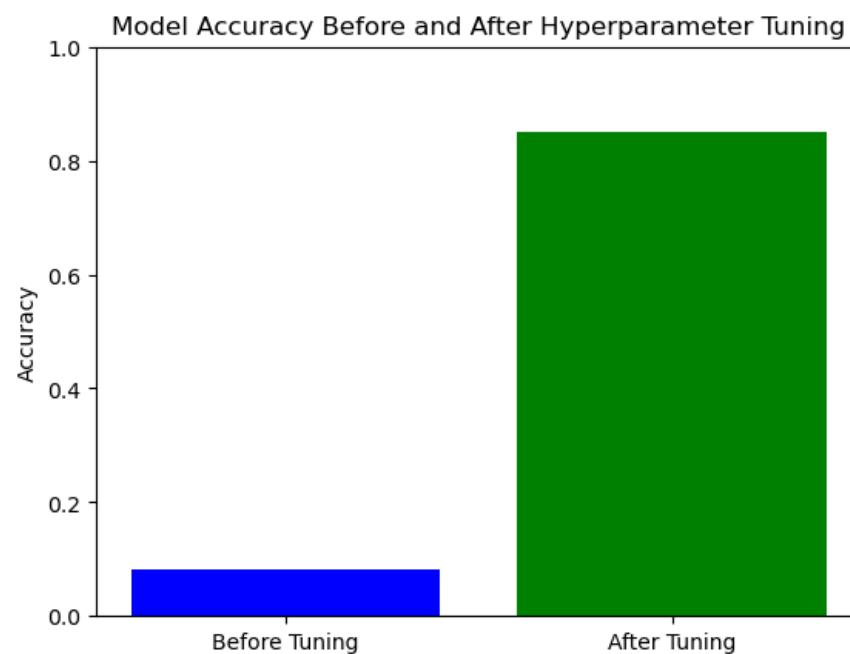
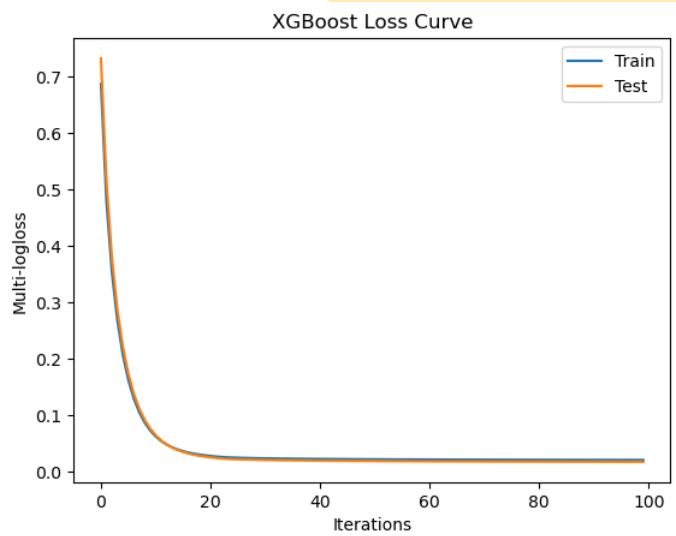
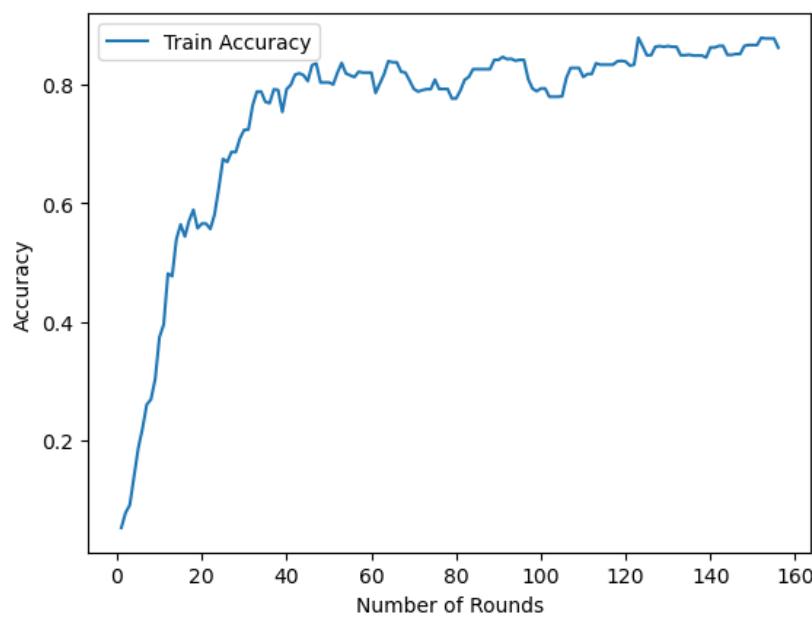


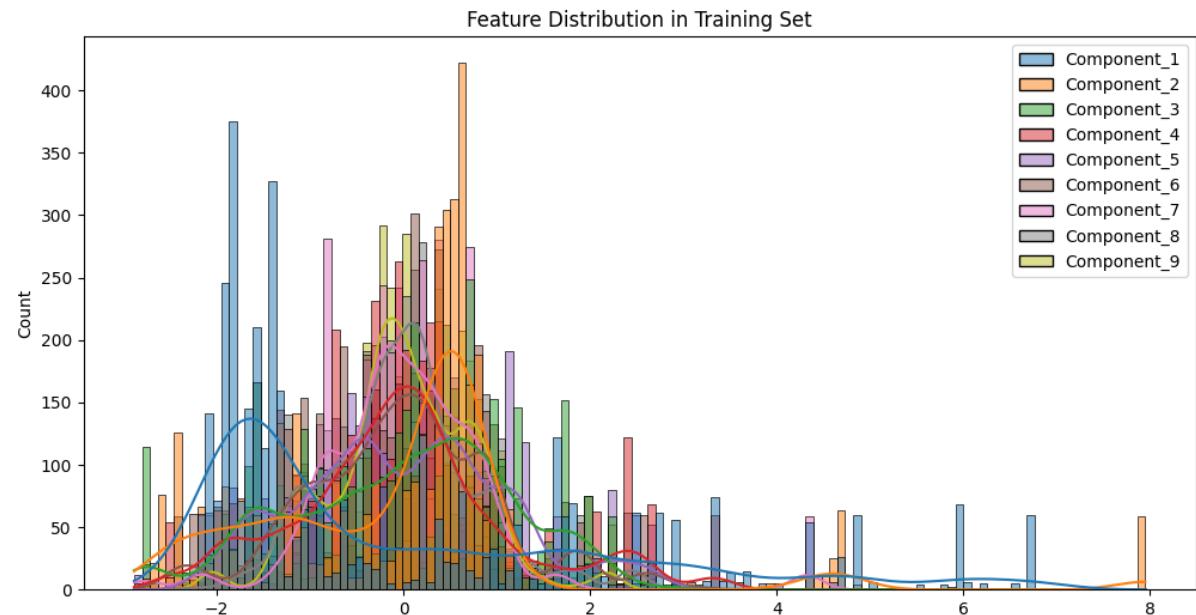
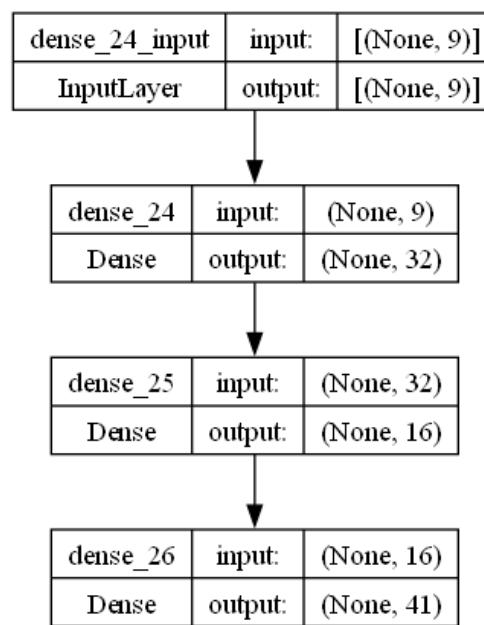
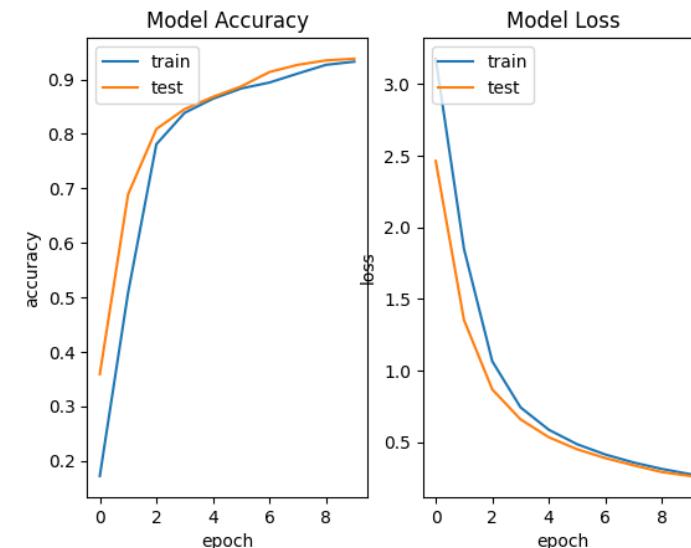
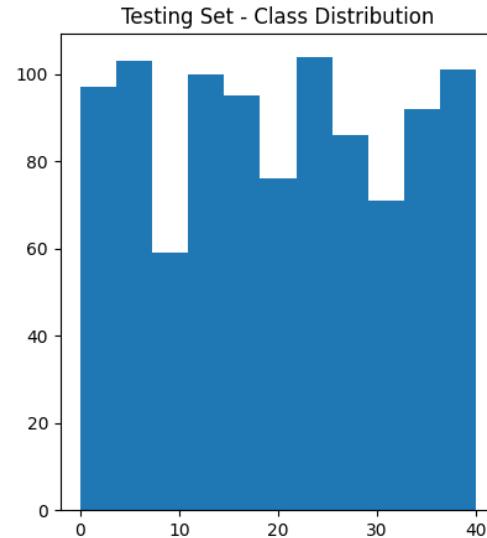
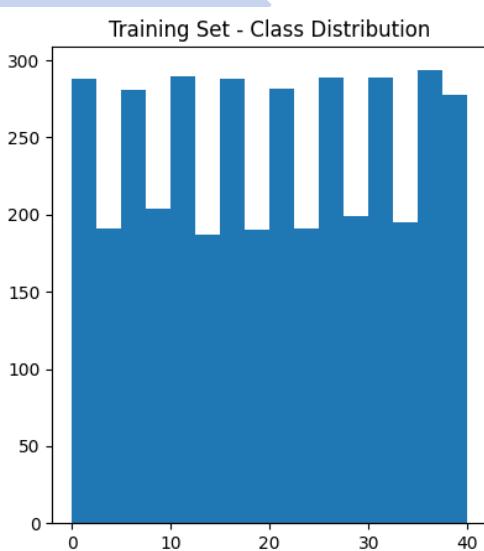
Random forest





ENSEMBLE

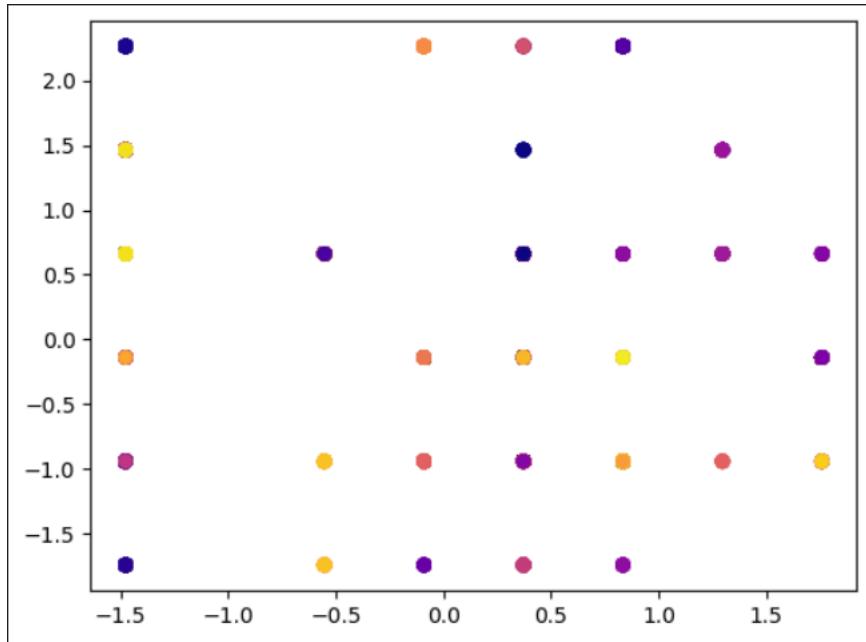
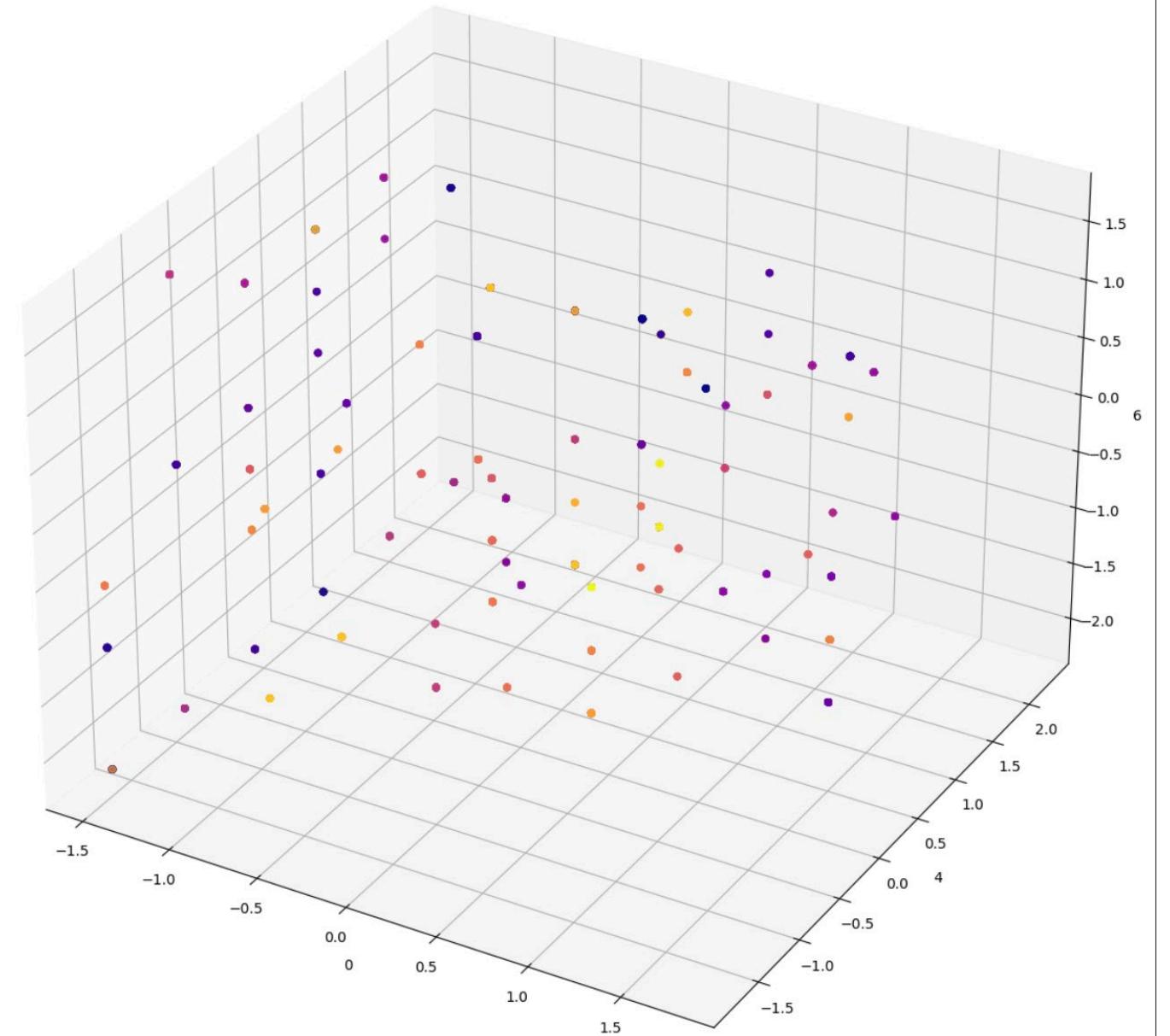


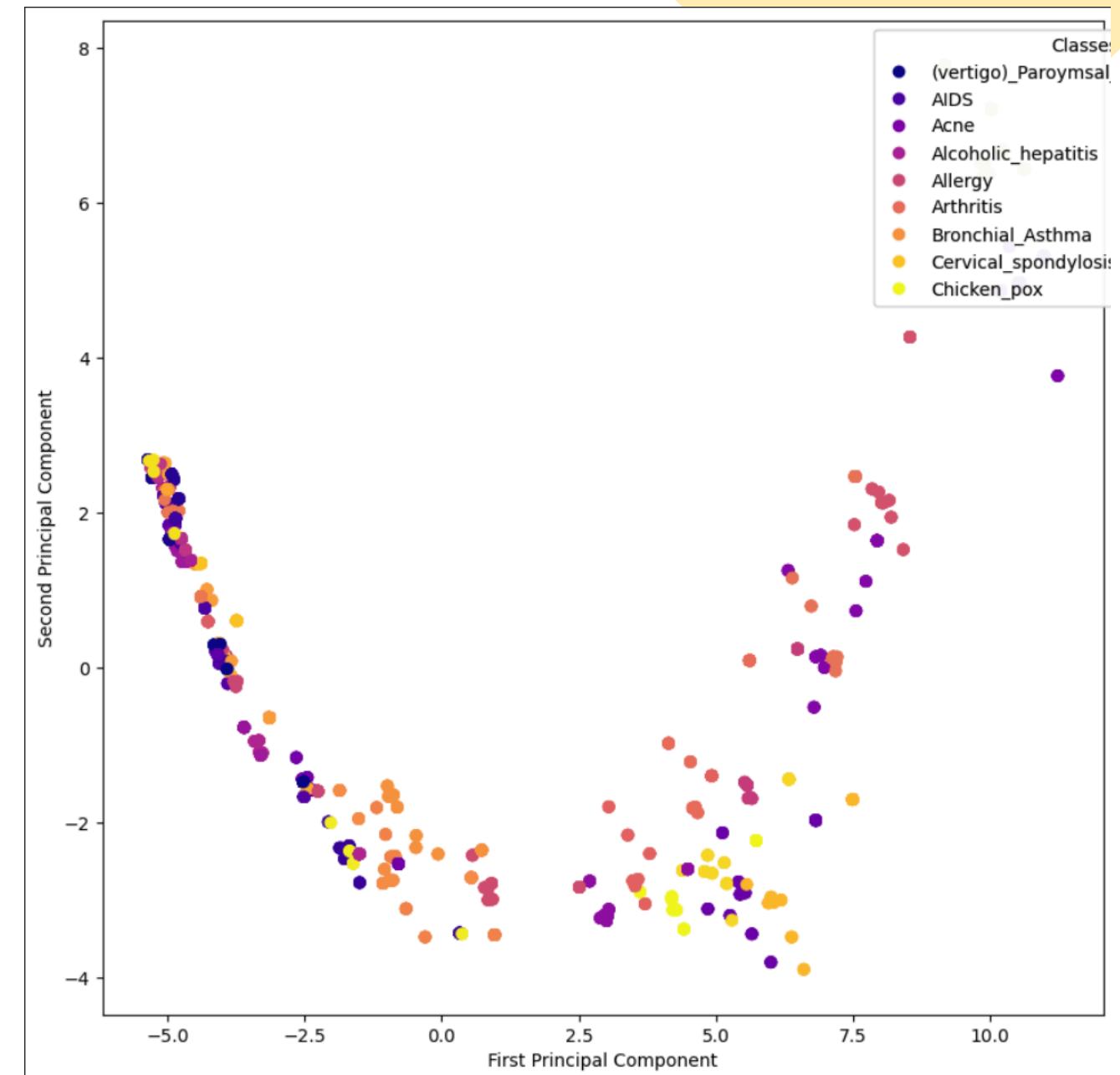
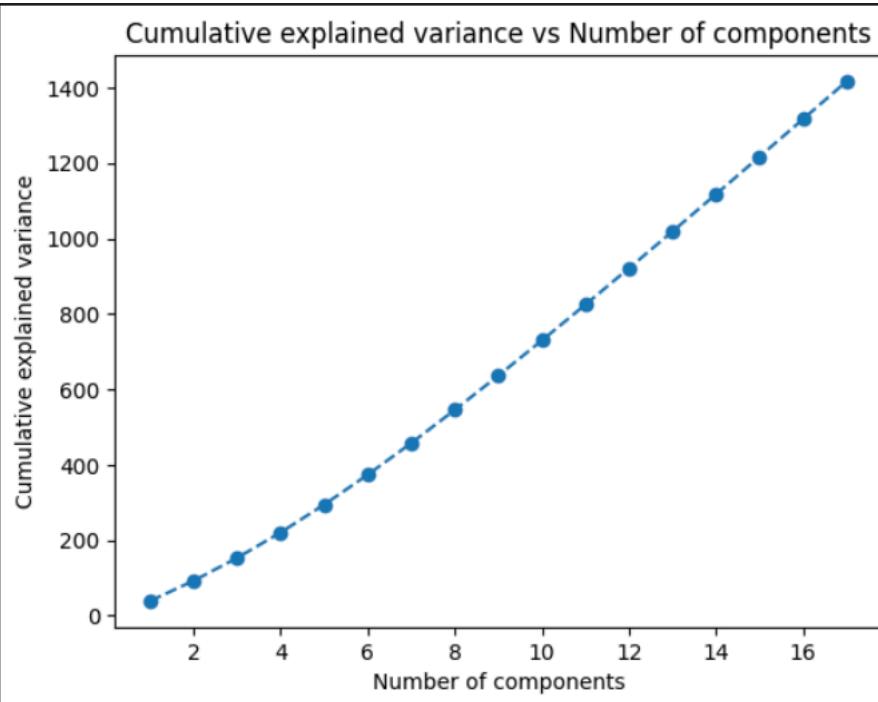


A 3D visualization of a network graph. The nodes are represented by small 3D cubes of various sizes, some with a glowing interior. They are interconnected by a complex web of thin red lines forming a mesh. The background is dark, making the red lines and the glowing nodes stand out.

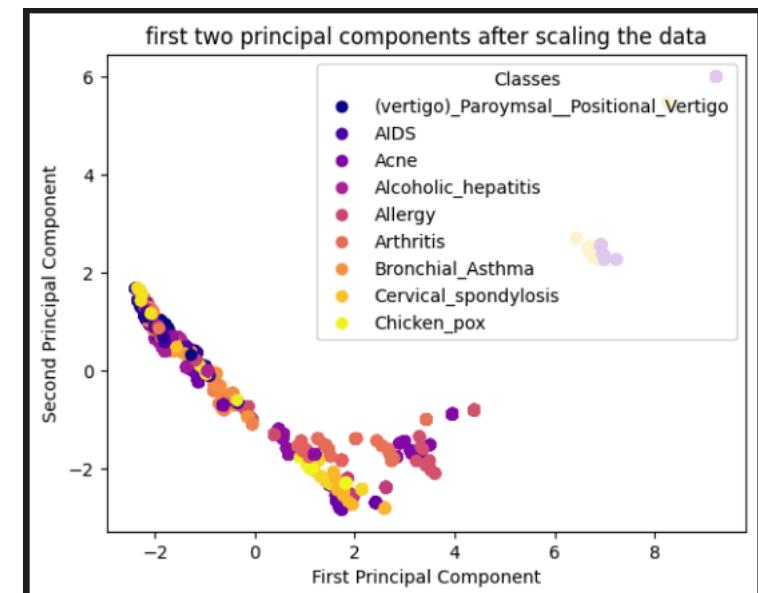
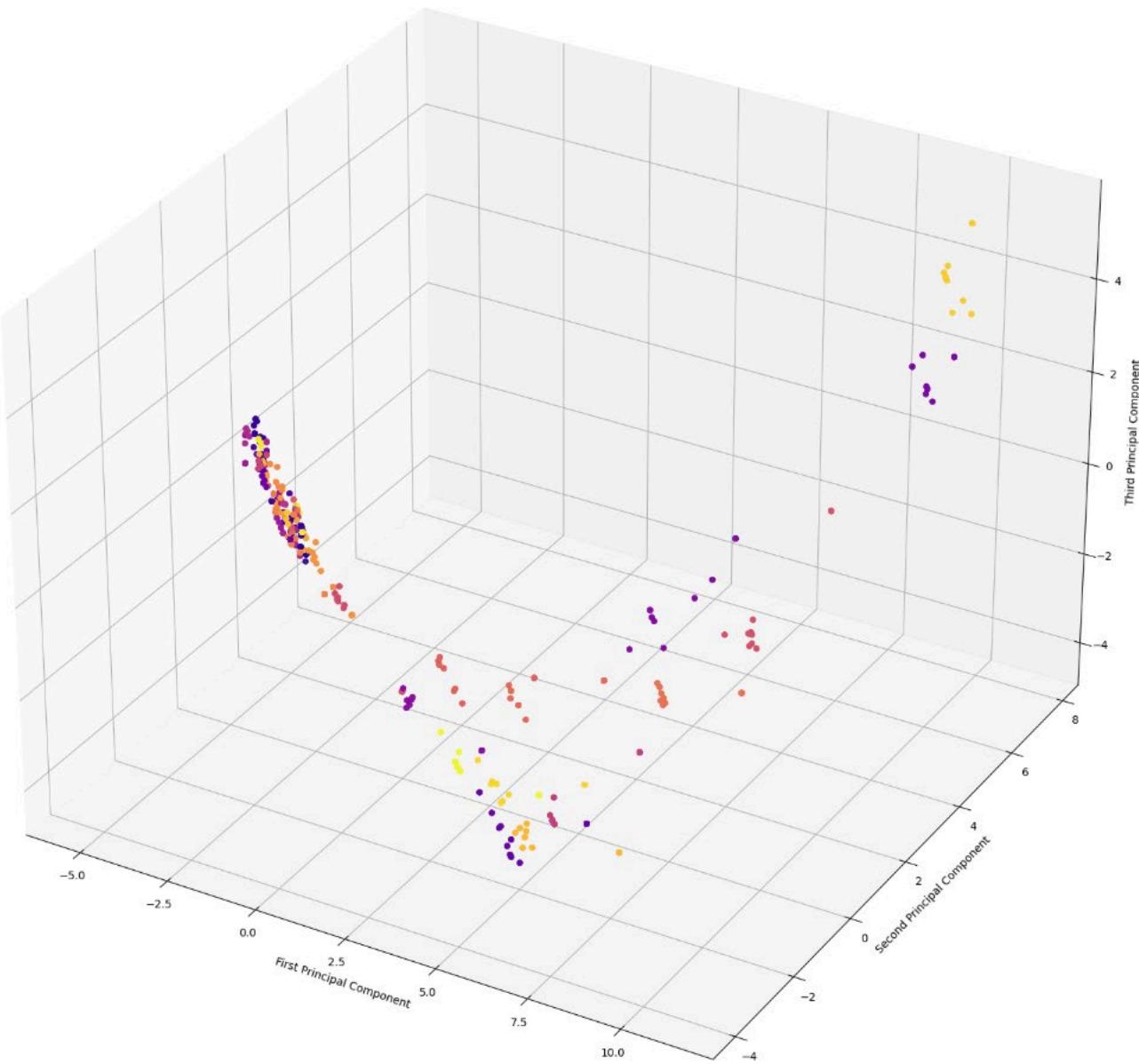
PCA

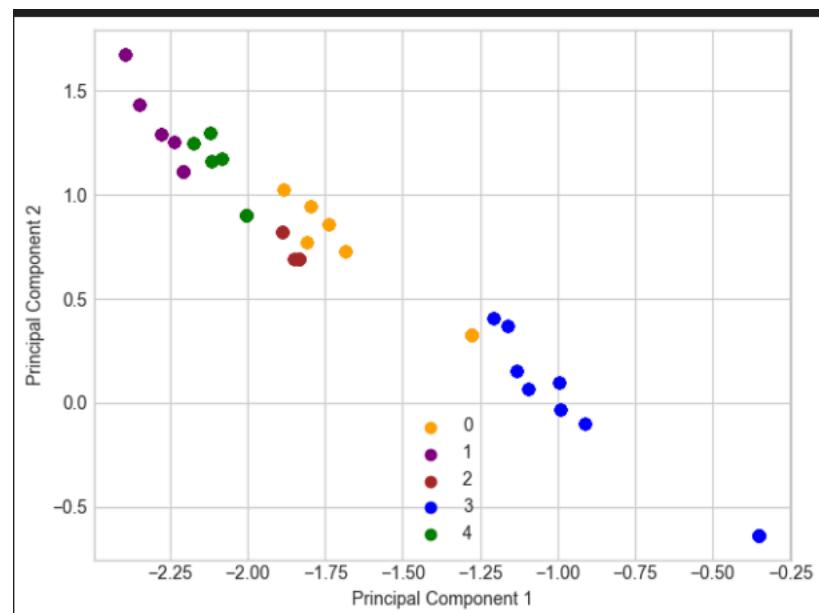
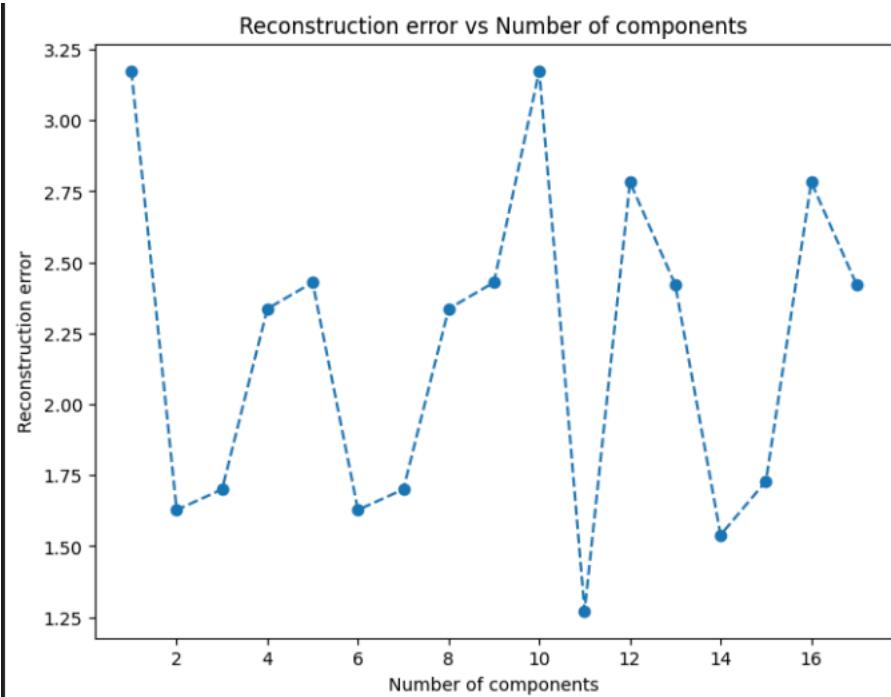
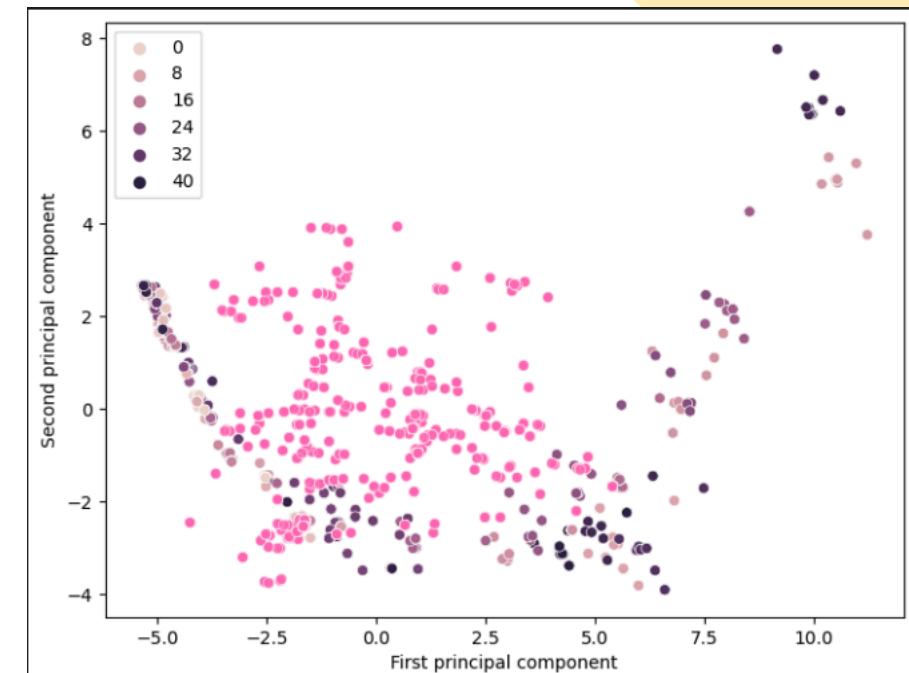
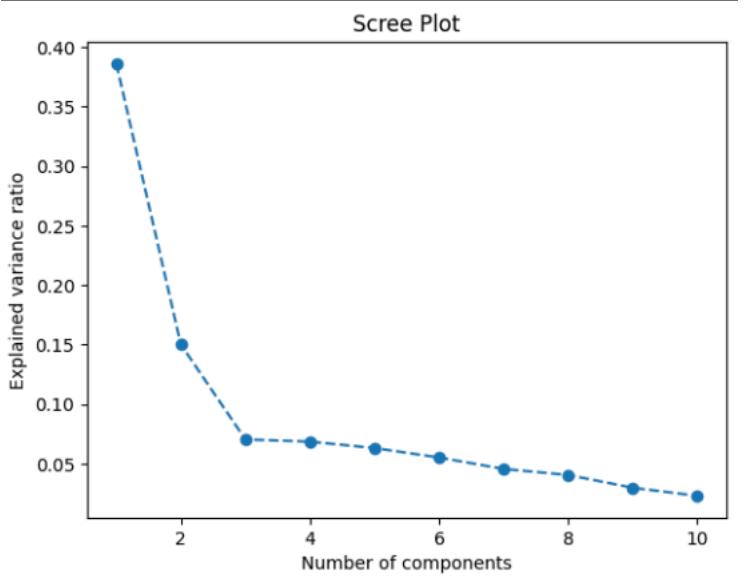
3D Scatter Plot of 3 particular symptoms

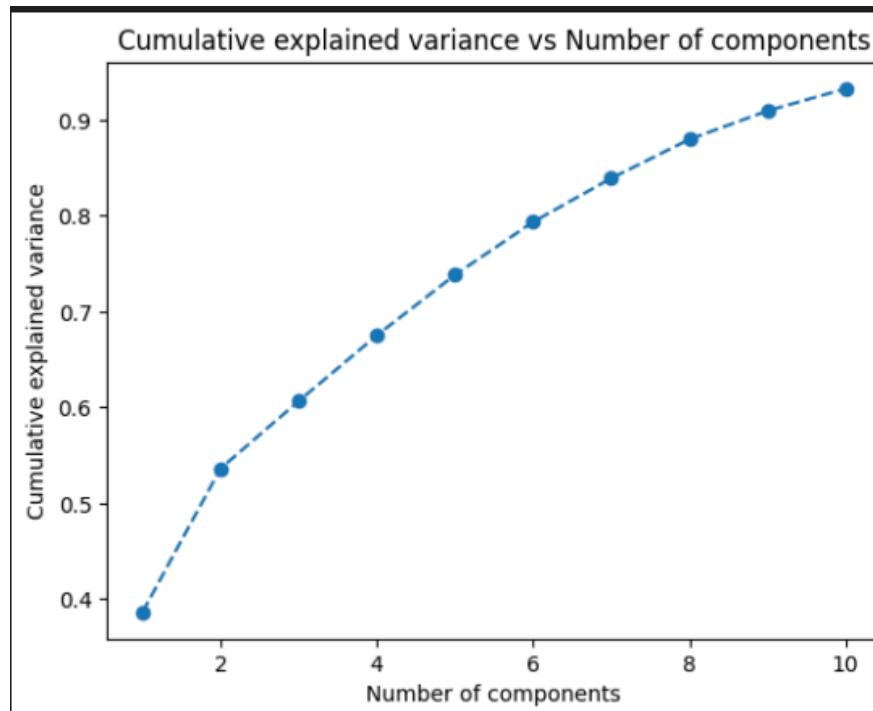
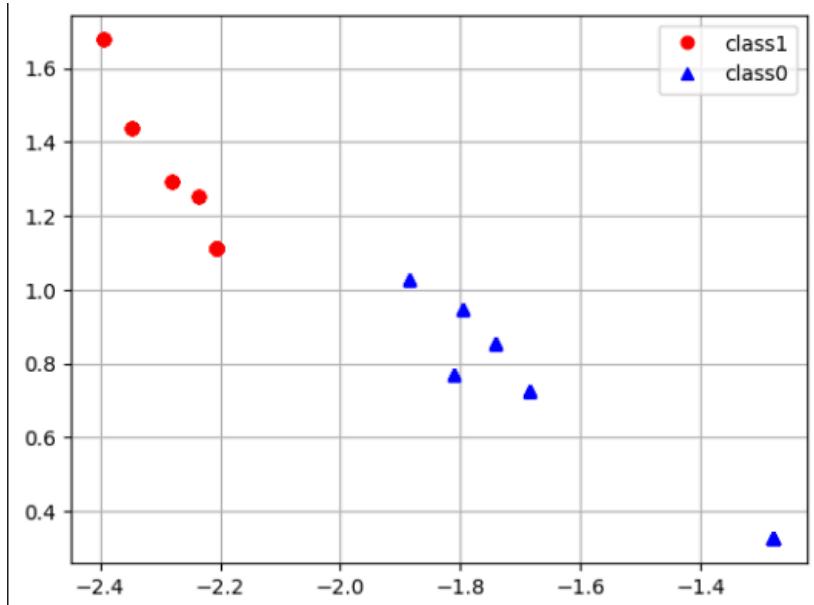


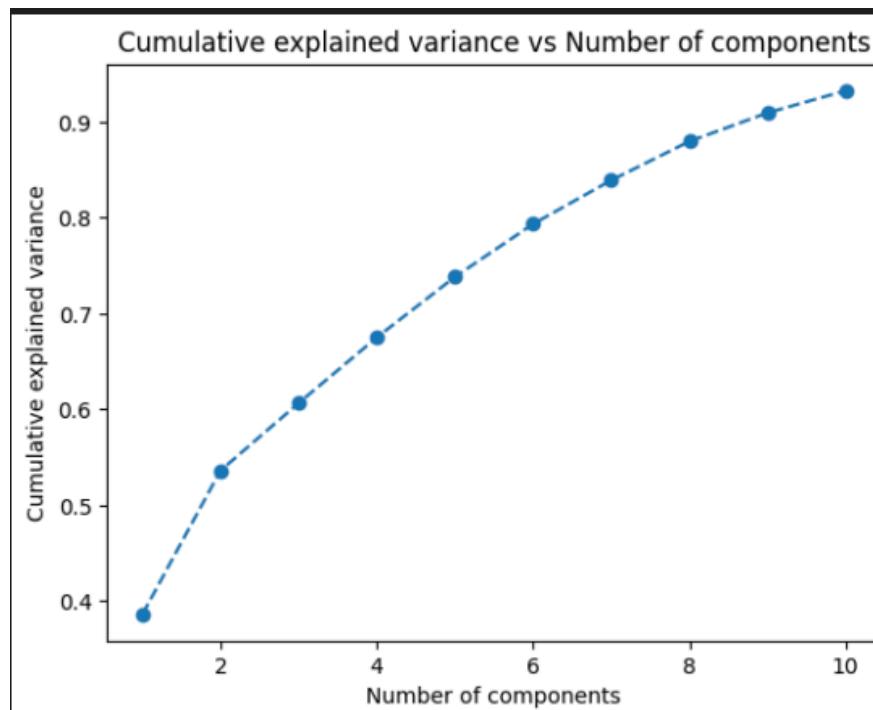
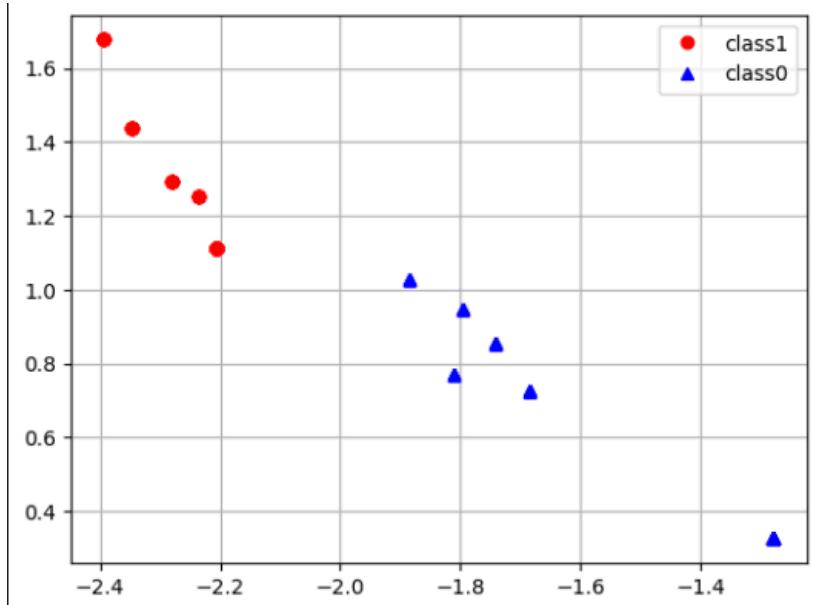


3D Scatter Plot of 3 Principal Components



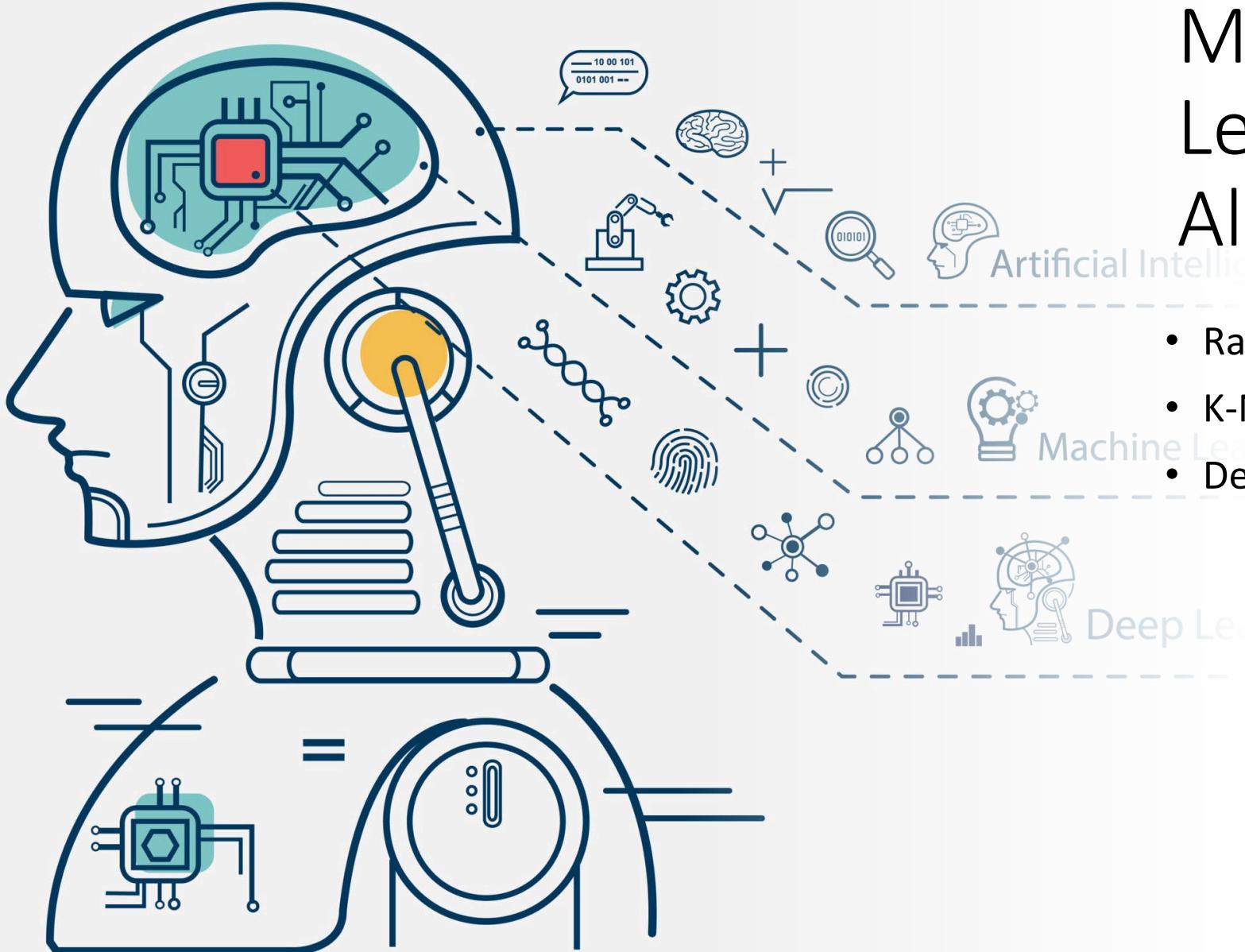


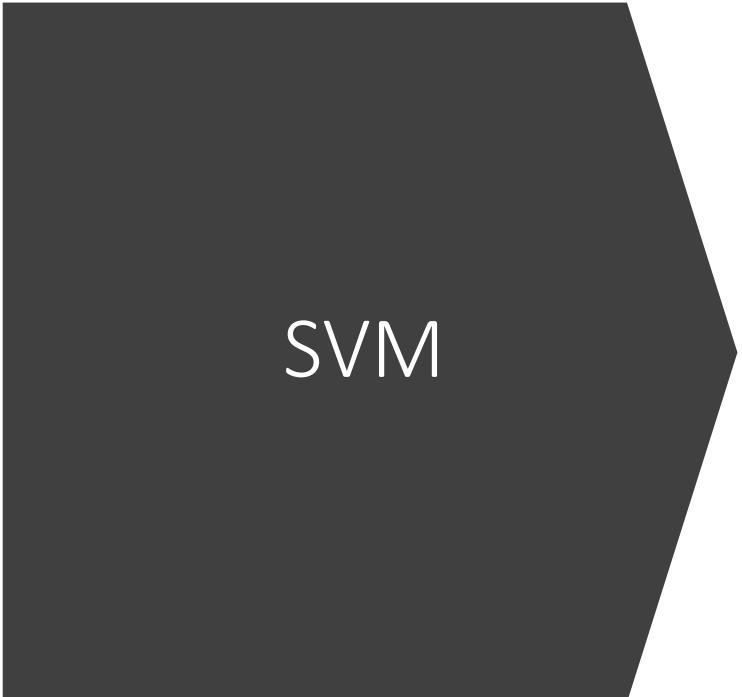




Machine Learning Algorithms Used

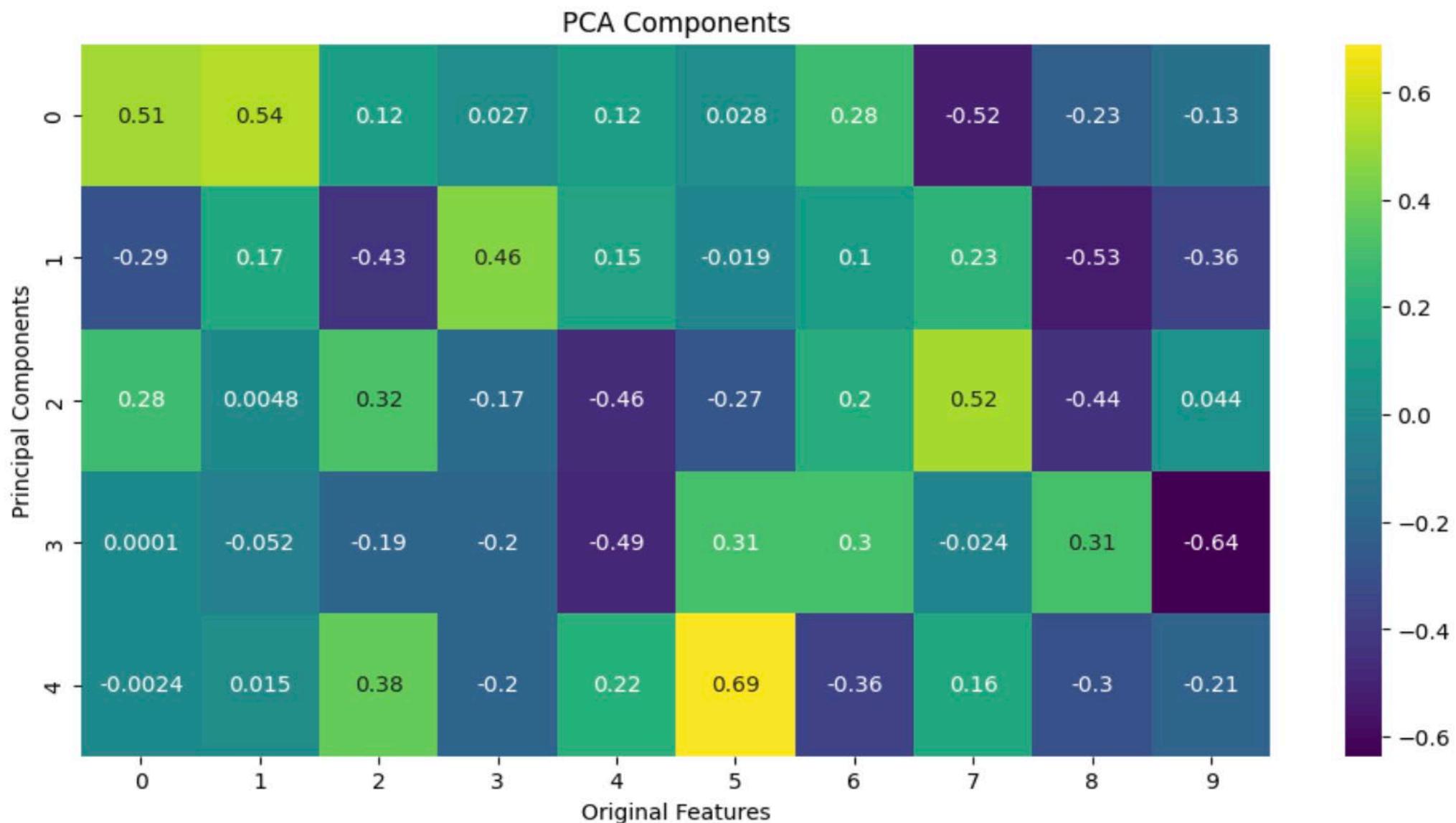
- Random Forest Classifier
- K-Nearest Neighbour Classifier
- Decision Tree Classifier

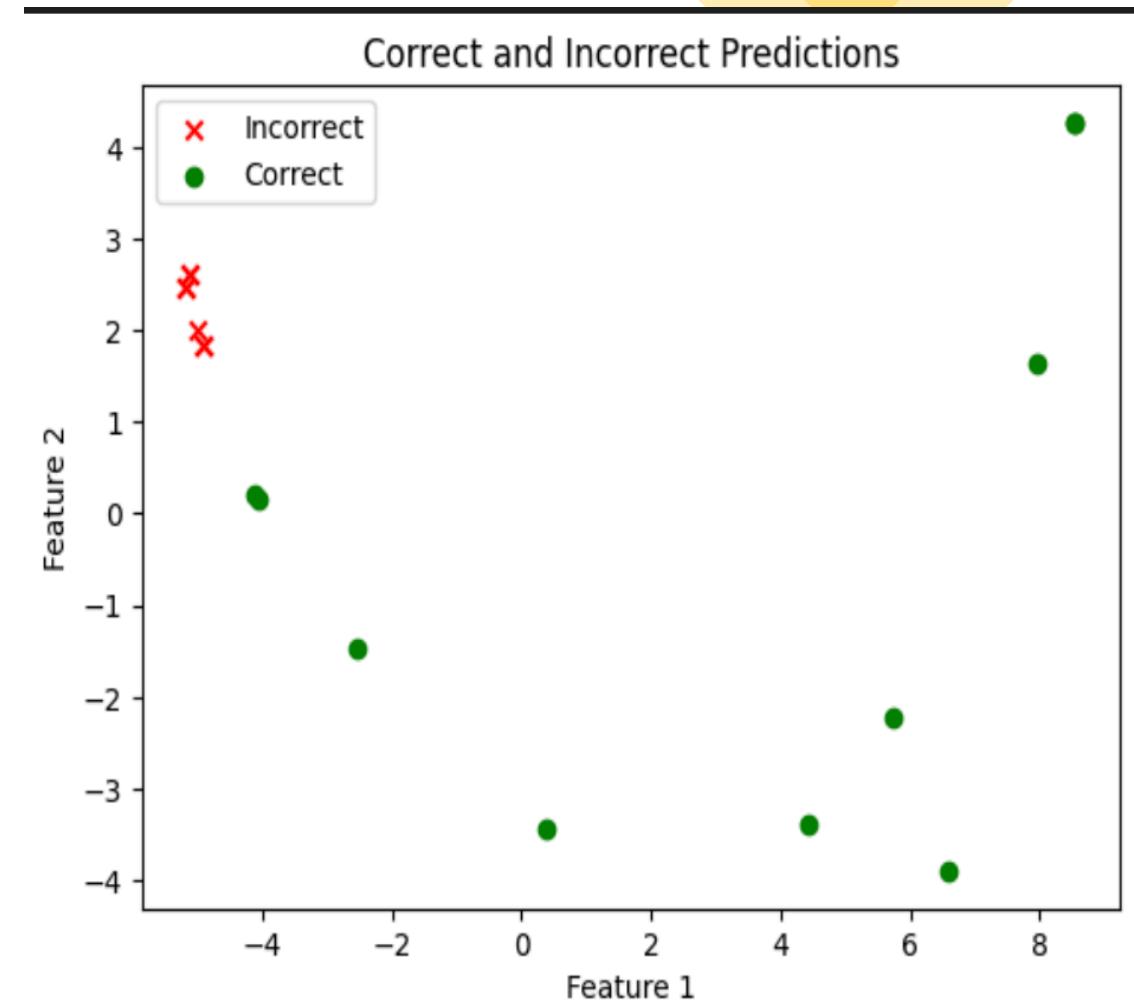
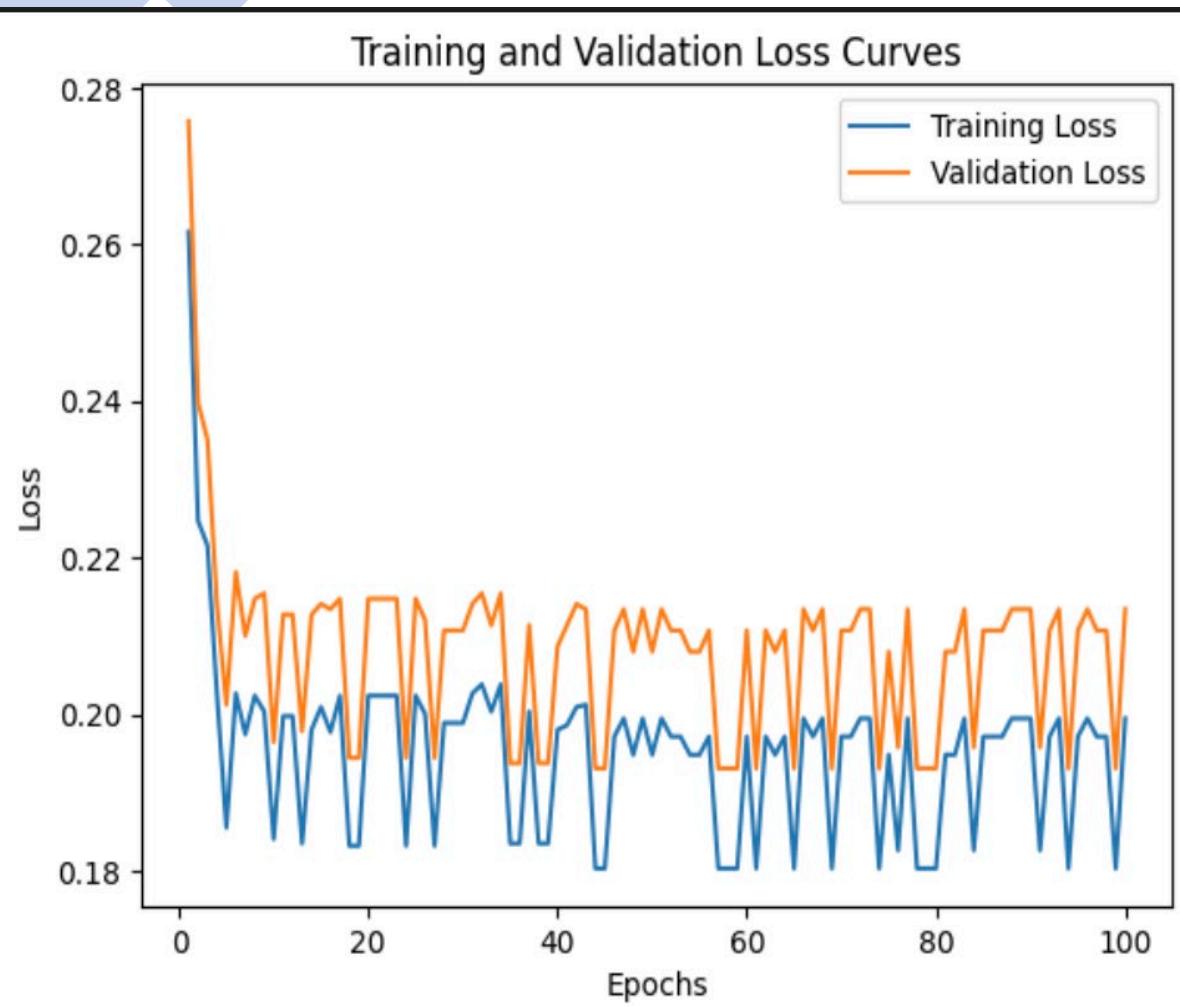




SVM

Metrics	
Accuracy	0.9932249322493225
Precision	0.9937446854296159
Recall	0.9932249322493225





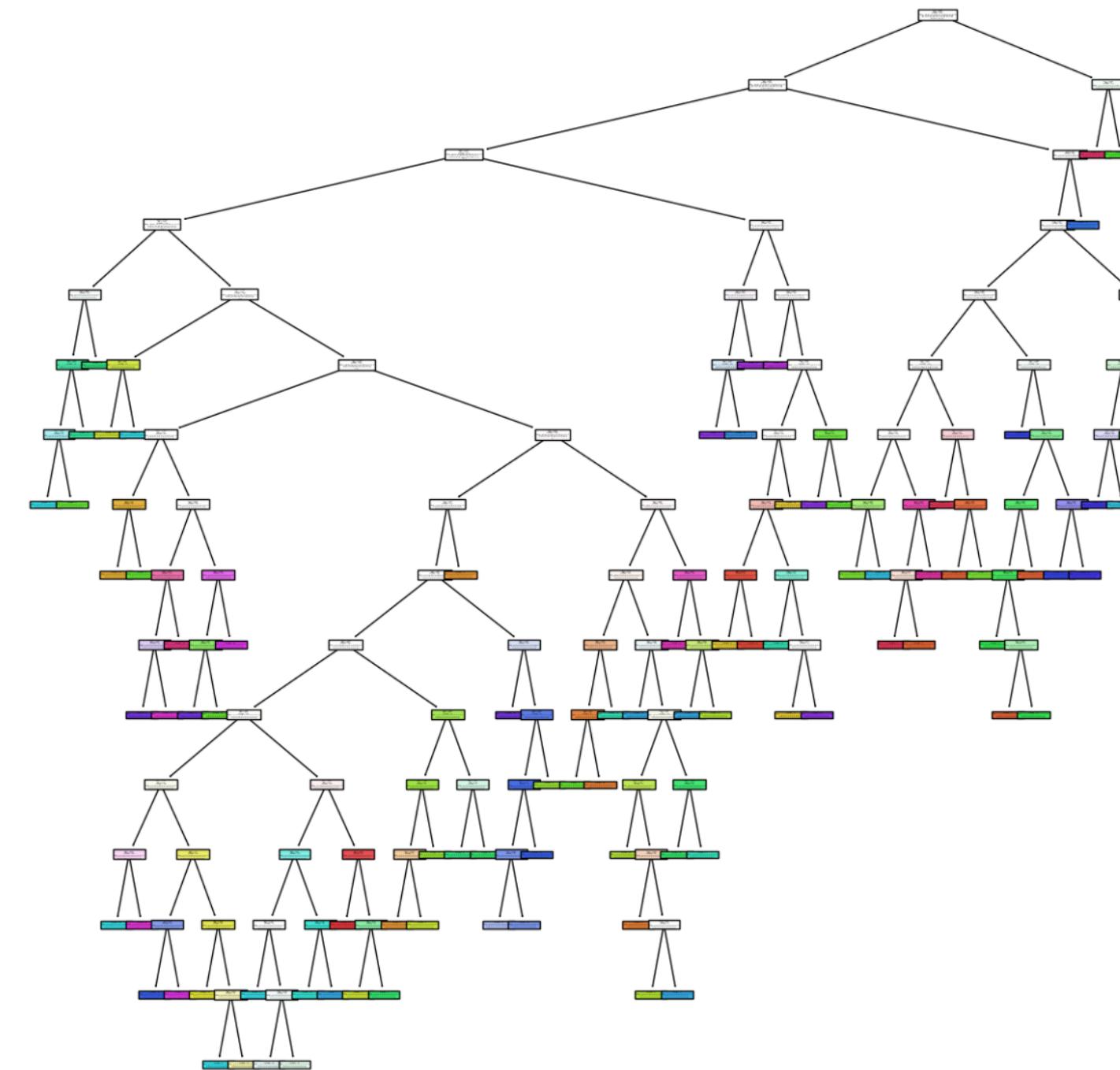
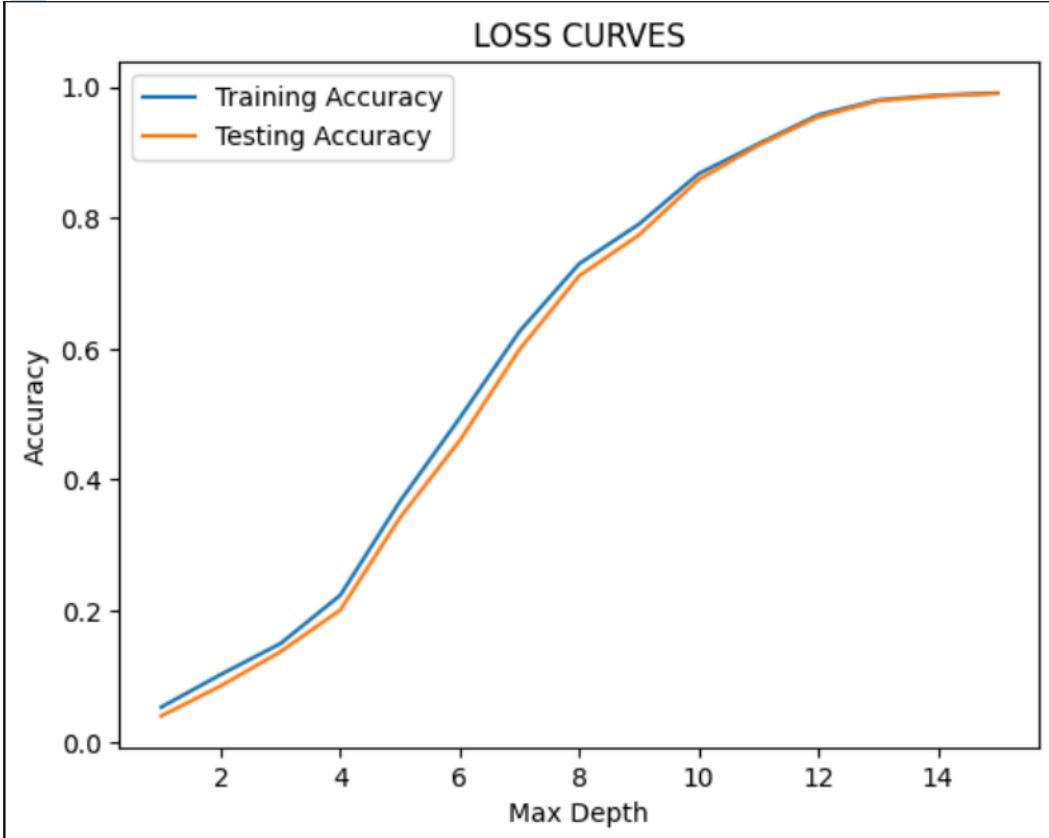
DECISION TREE



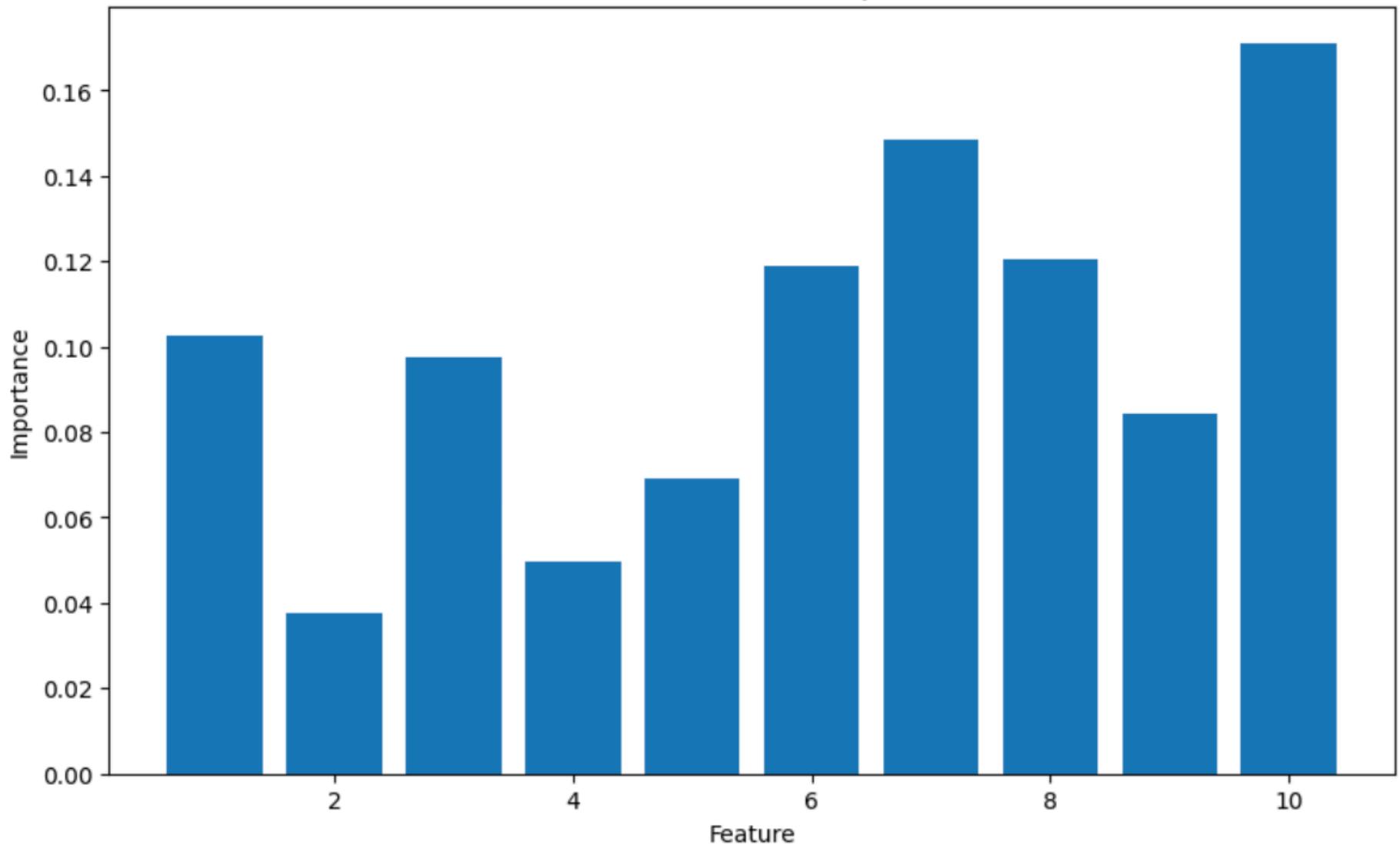


Decision Tree

Metrics	Accuracy
Accuracy	0.9932249322493225
Precision	0.9937446854296159
Recall	0.9932249322493225



Decision Tree Feature Importance

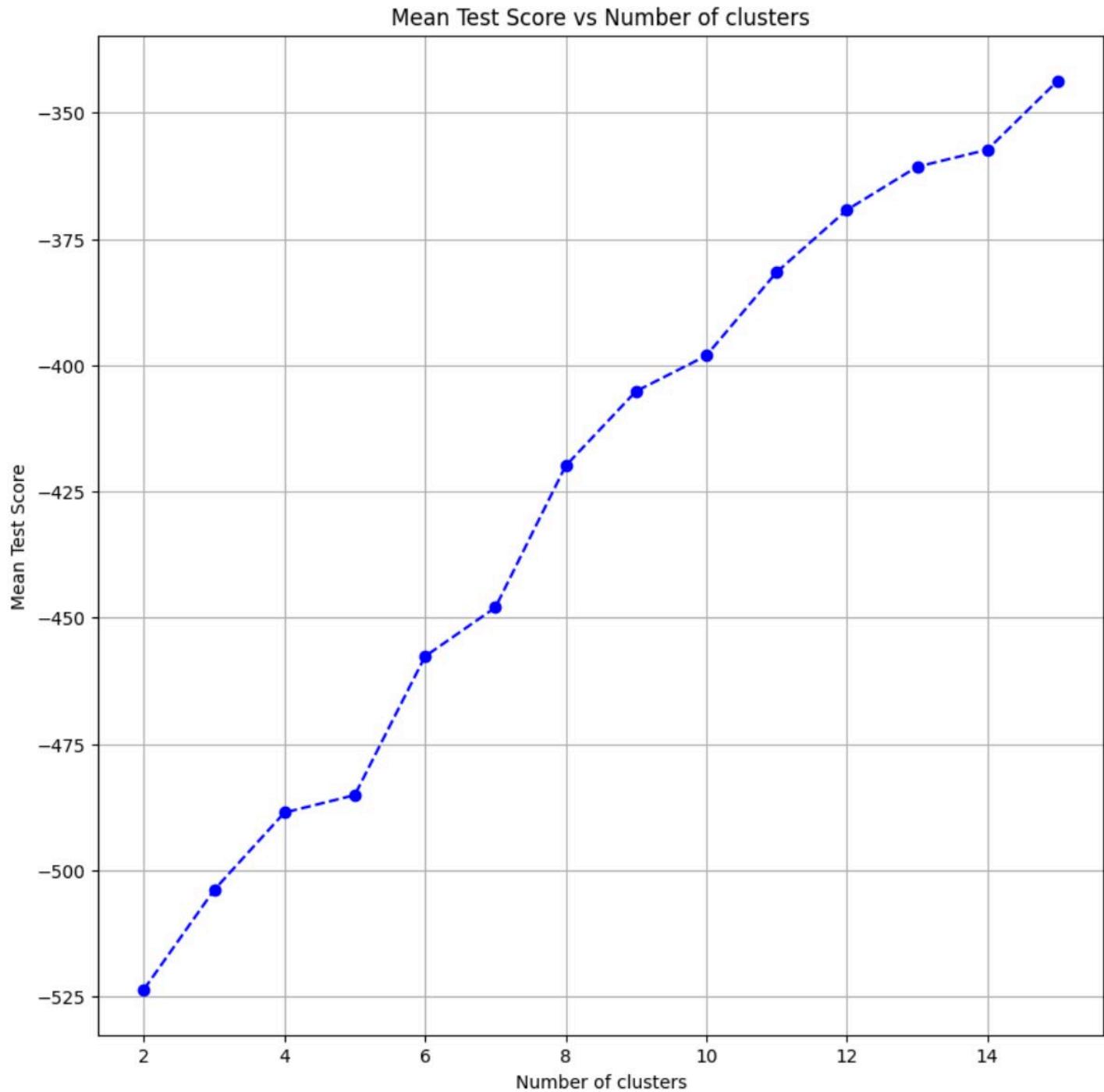


CLUSTERING

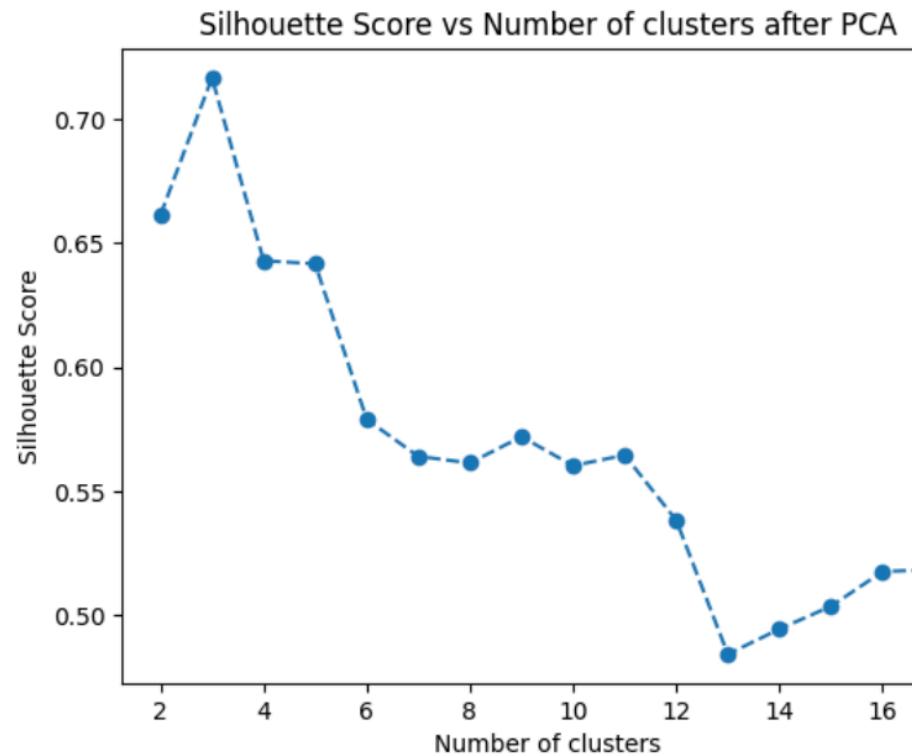
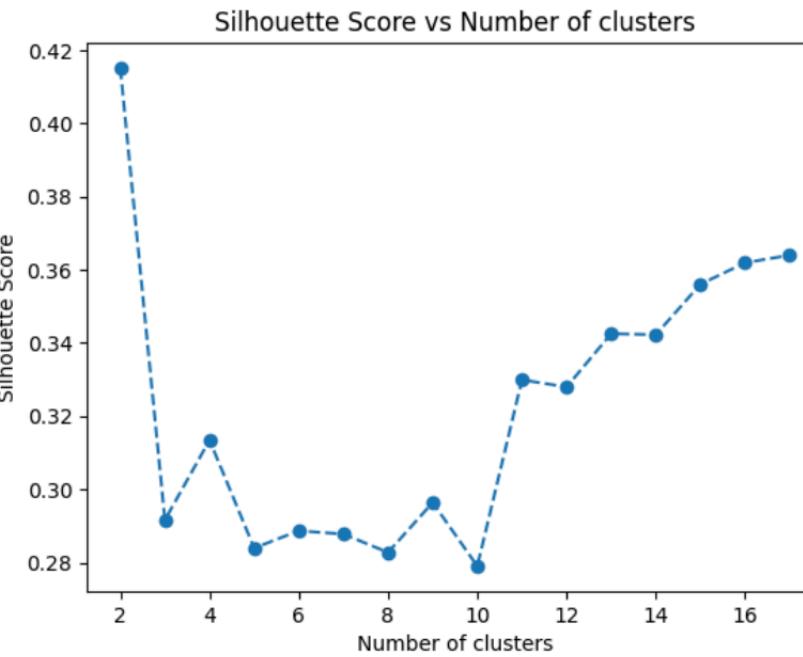
- K-Means Clustering
- Agglomerative Clustering



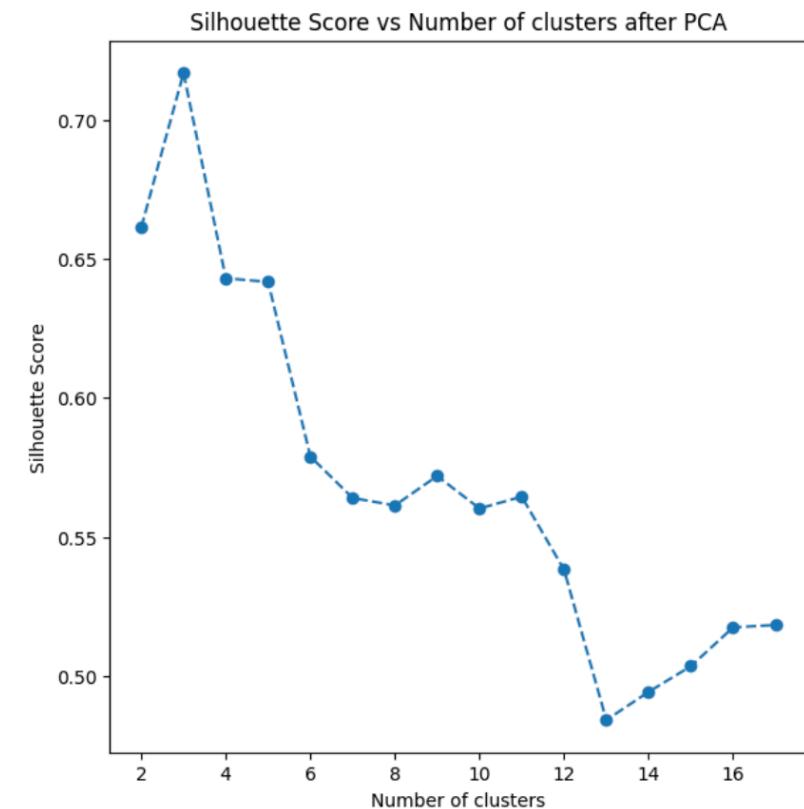
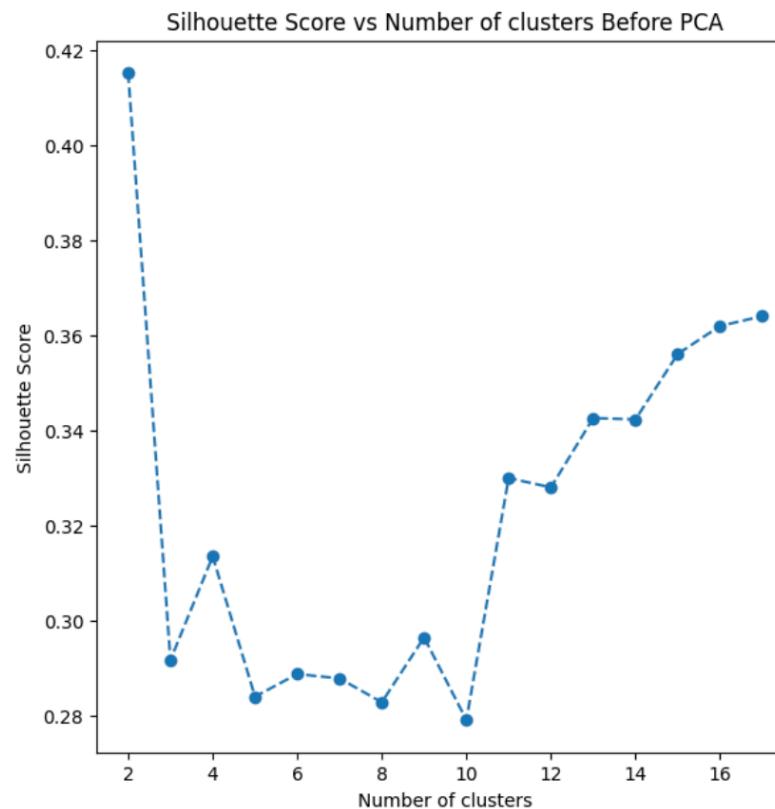
Metrics for K-Means



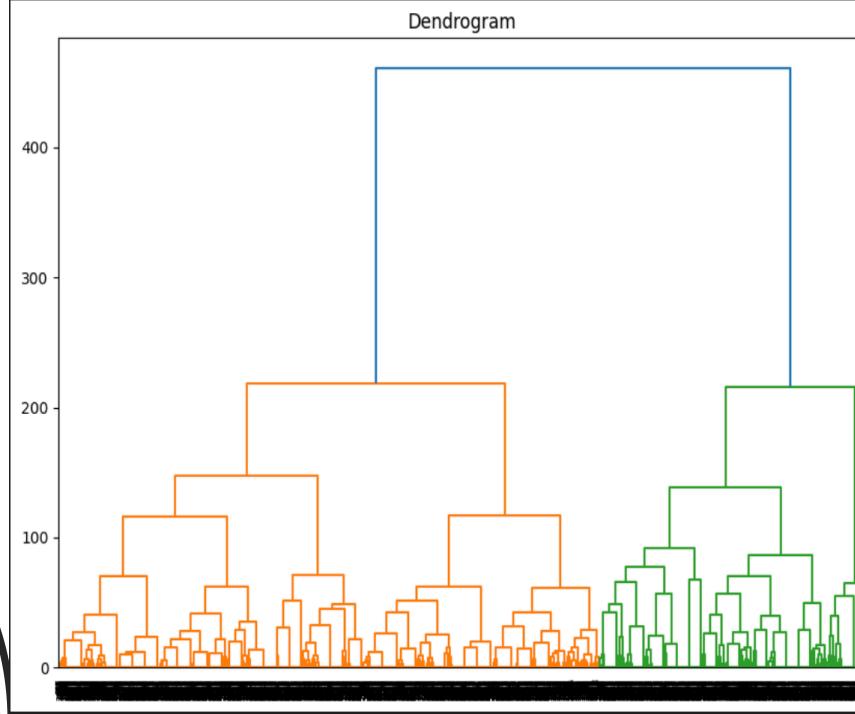
Silhouette Score



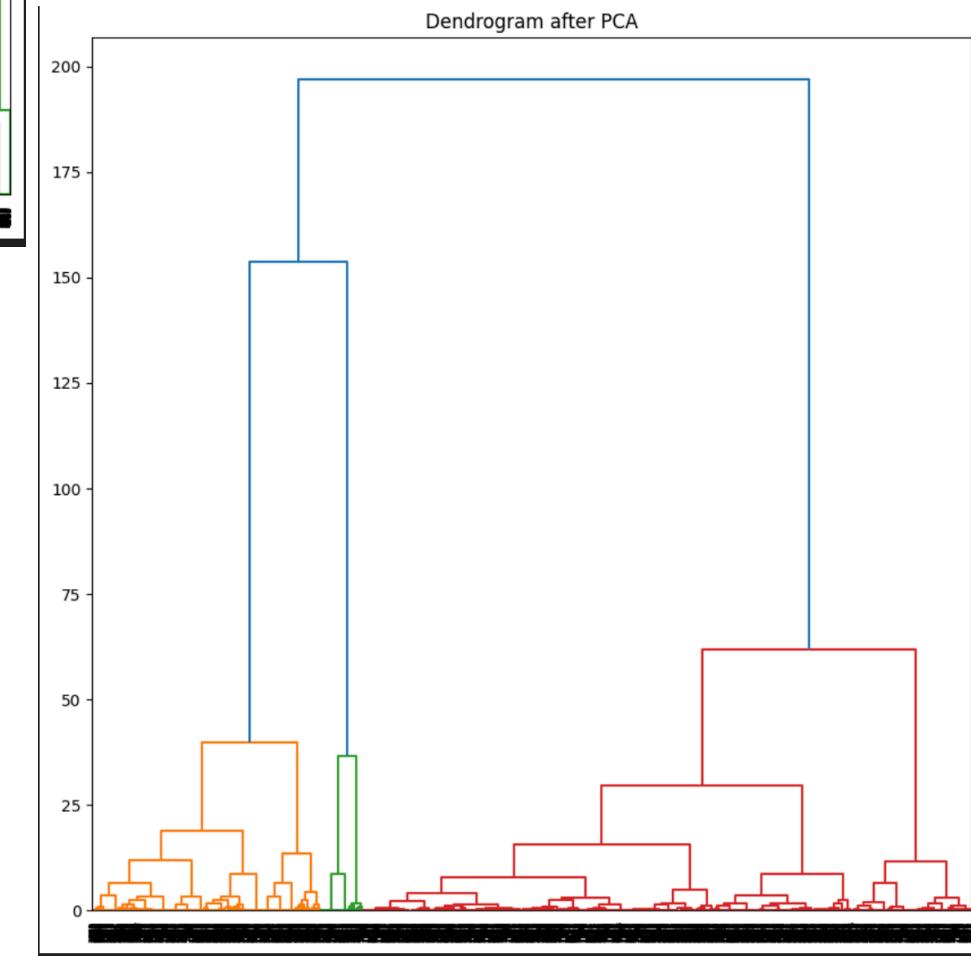
Silhouette Score Before Vs After PCA



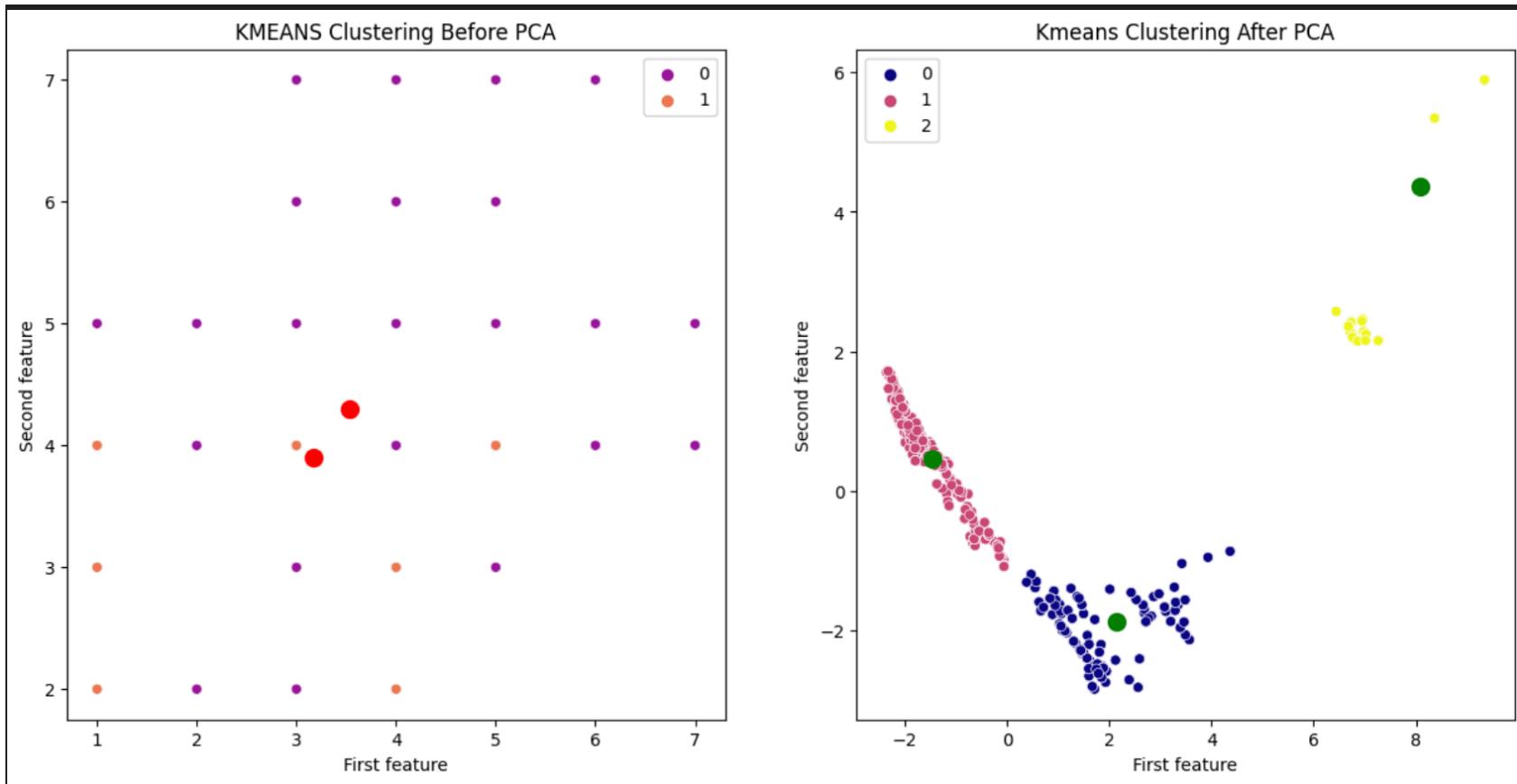
Dendrogram



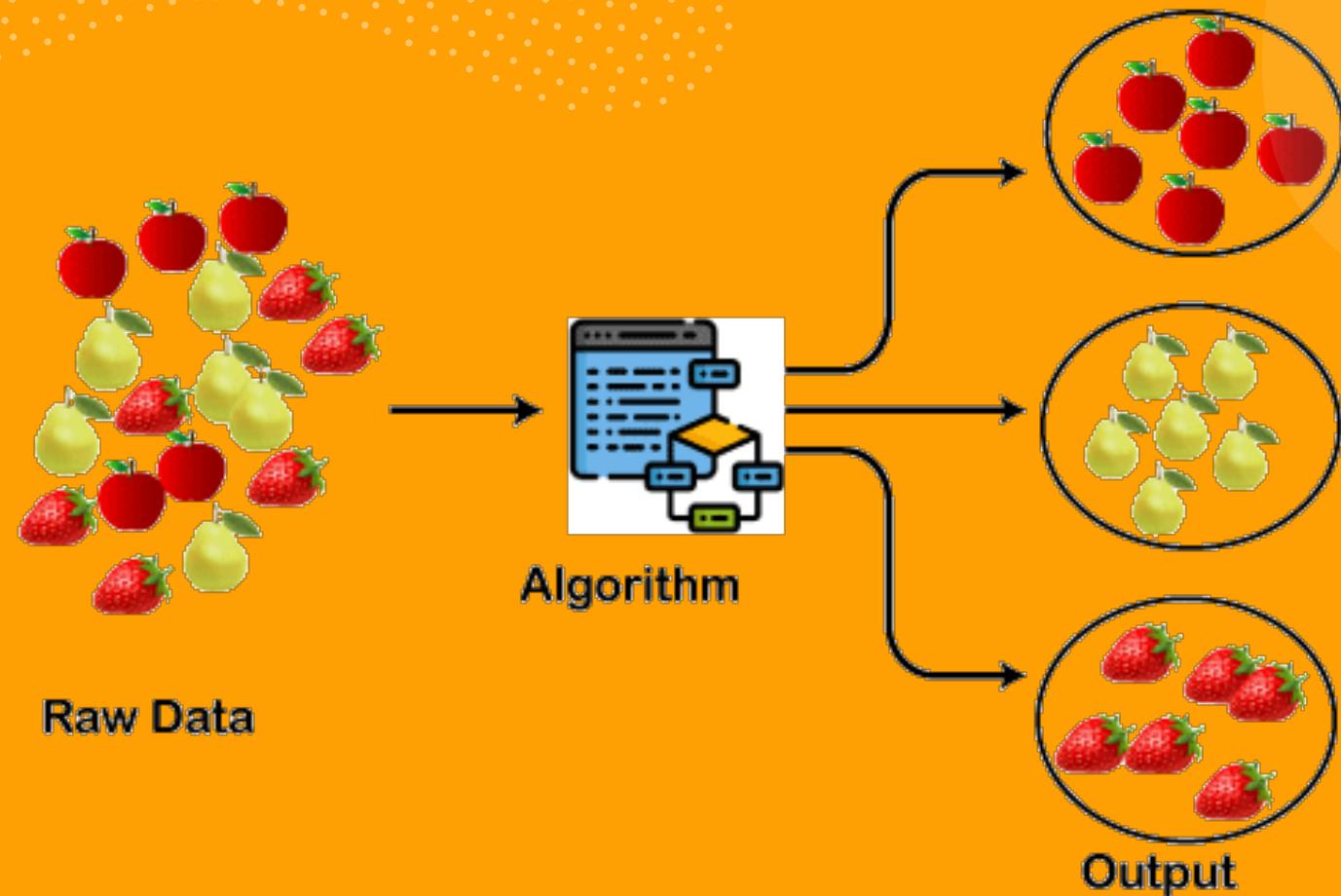
Dendrogram after PCA



K-MEANS Before Vs After PCA

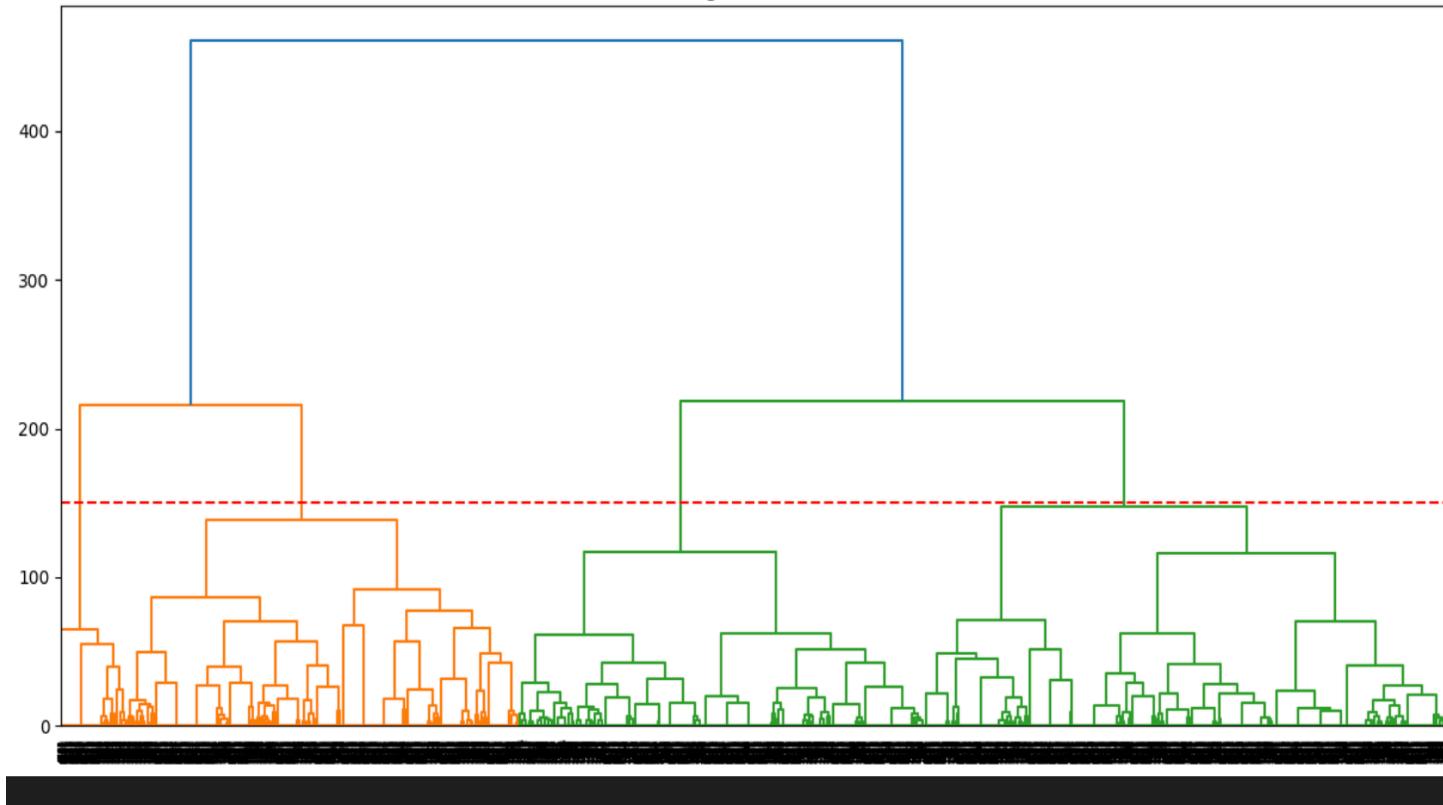


Hierarchical clustering

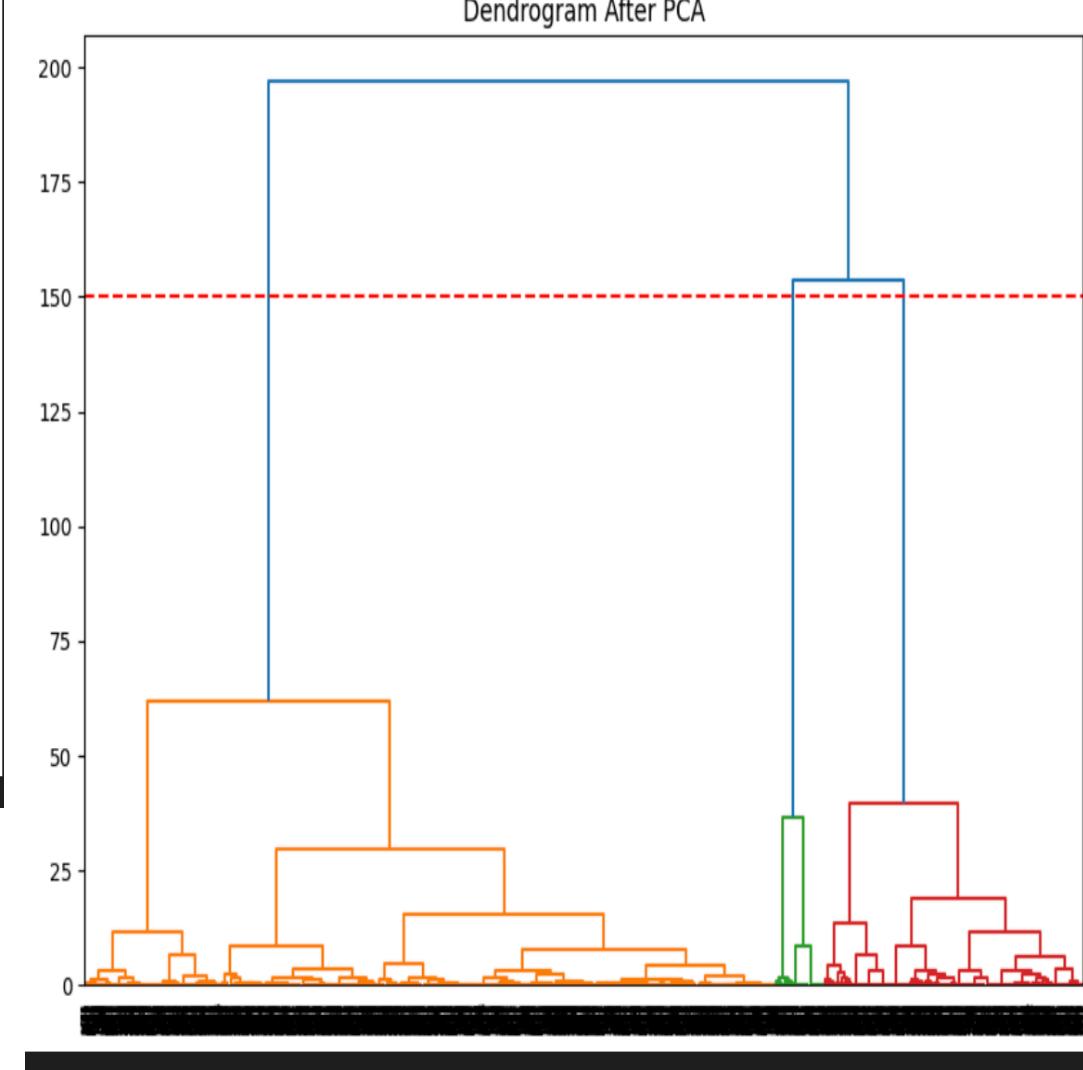


Hierarchial Clustering

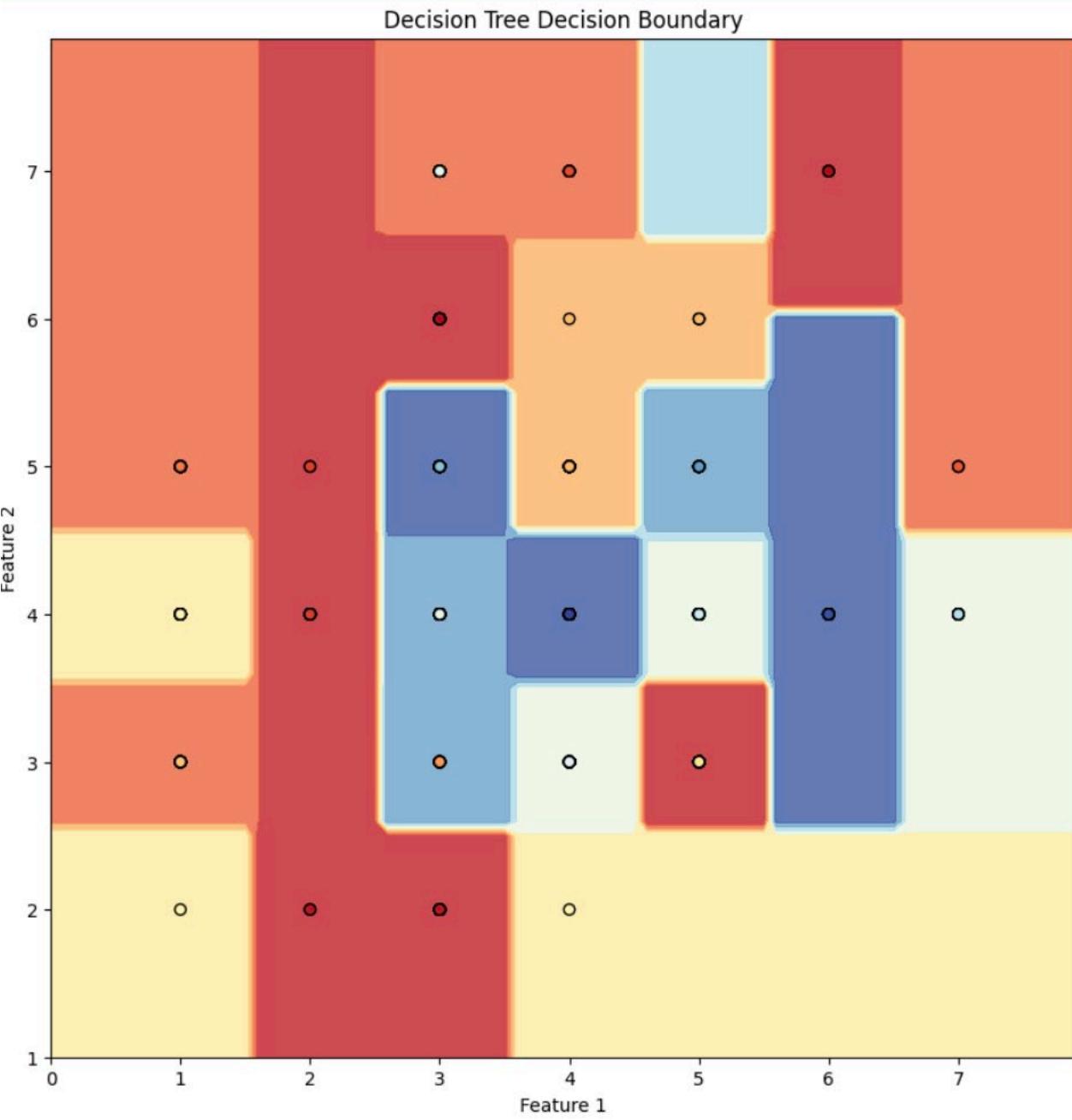
Dendrogram Before PCA



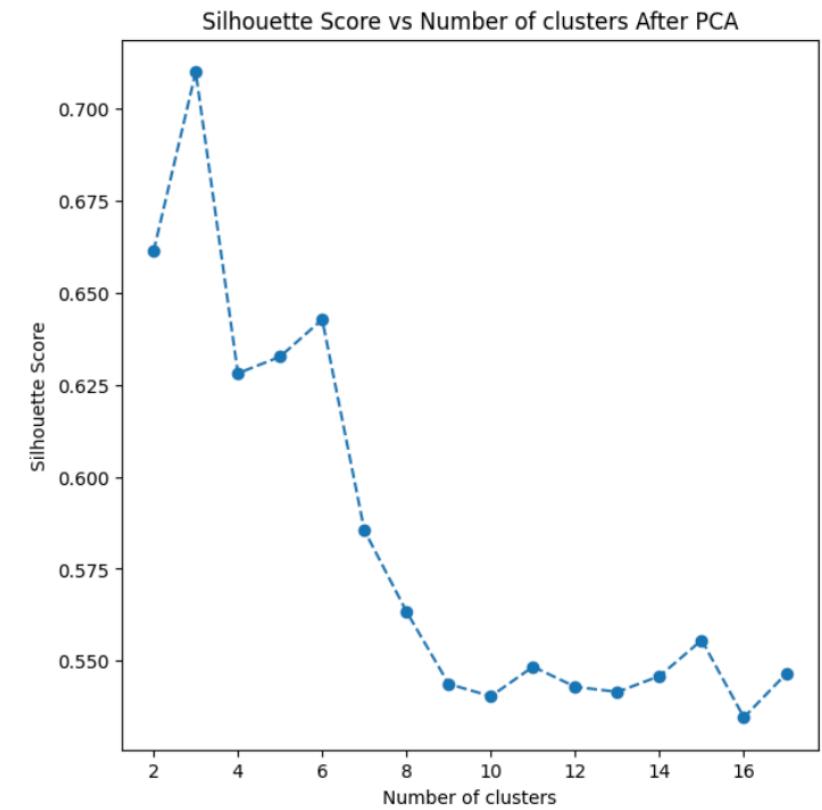
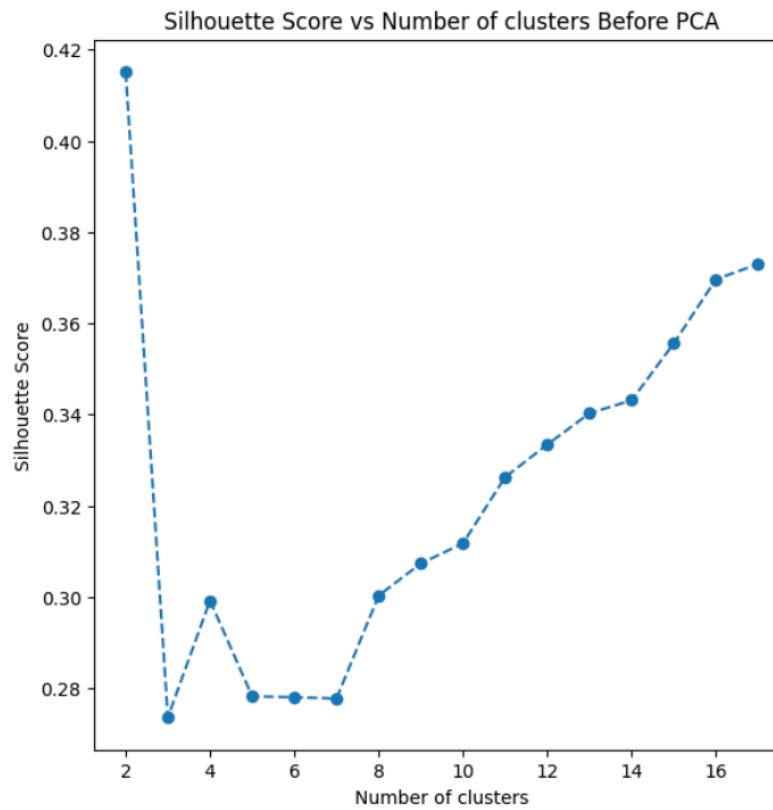
Dendrogram After PCA



Decision Tree Decision Boundary



Agglomerative Clustering

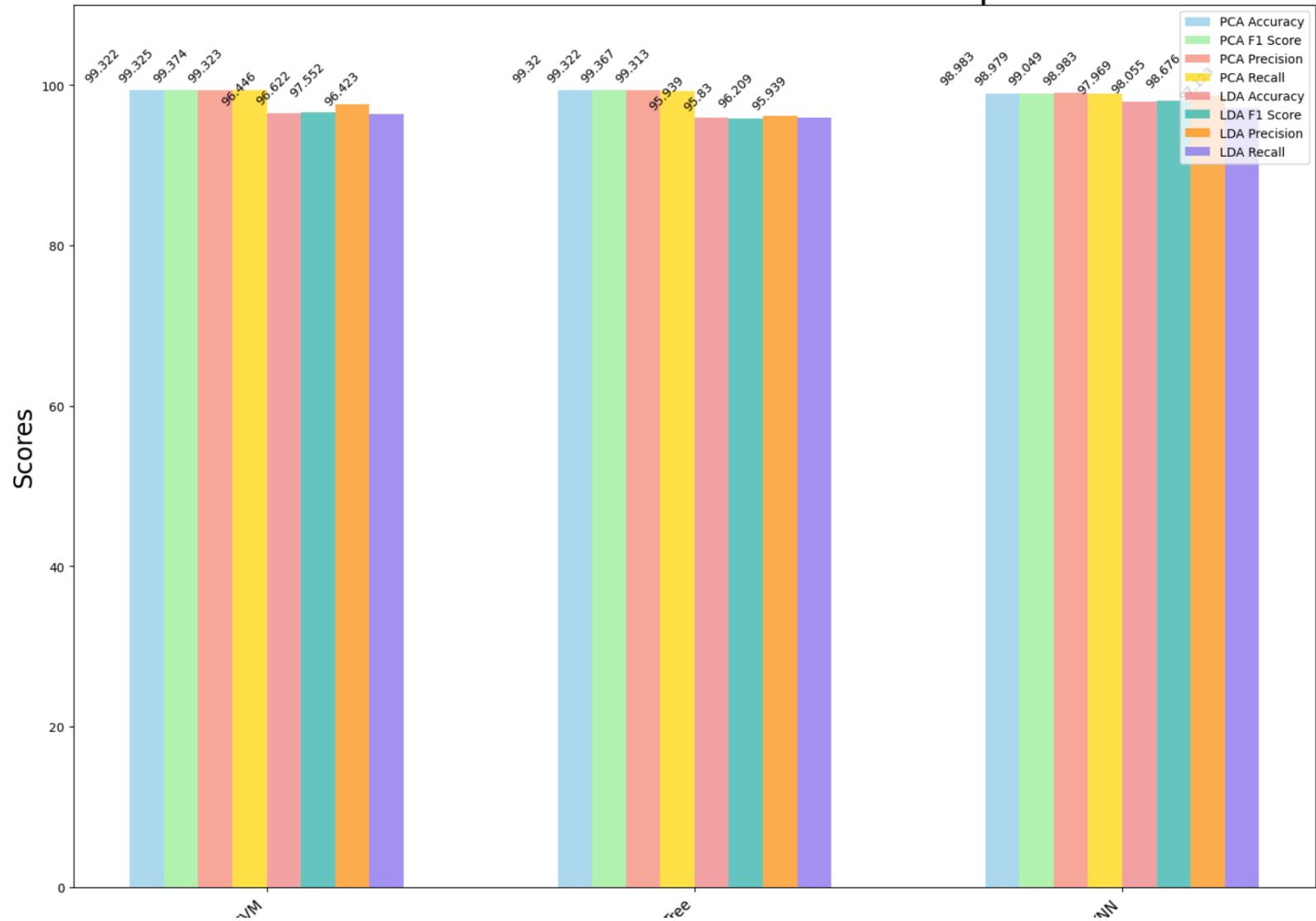


S.NO	MODEL	ACCURACY	F1 SCORE	PRECISION	RECALL
1	KNN-PCA	99%	99.2	99.3	99.18
2	DECISION TREE-PCA	99.18%	99.09	99.37	99.17
3	RANDOM FOREST	99.58%	99.14	99.28	99.09
4	SVM-Kernel-PCA				
5	DECISION-TREE-clustering	100%	93.91	94.816	93.8
6	SVM-RBF KERNEL clustering		97.5	97.523	97.26
7	KNN-CLUSTERING	99.28%	99.18	99.25	99.18

PCA and LDA Model Performance Comparison

Legend:

- PCA Accuracy
- PCA F1 Score
- PCA Precision
- PCA Recall
- LDA Accuracy
- LDA F1 Score
- LDA Precision
- LDA Recall

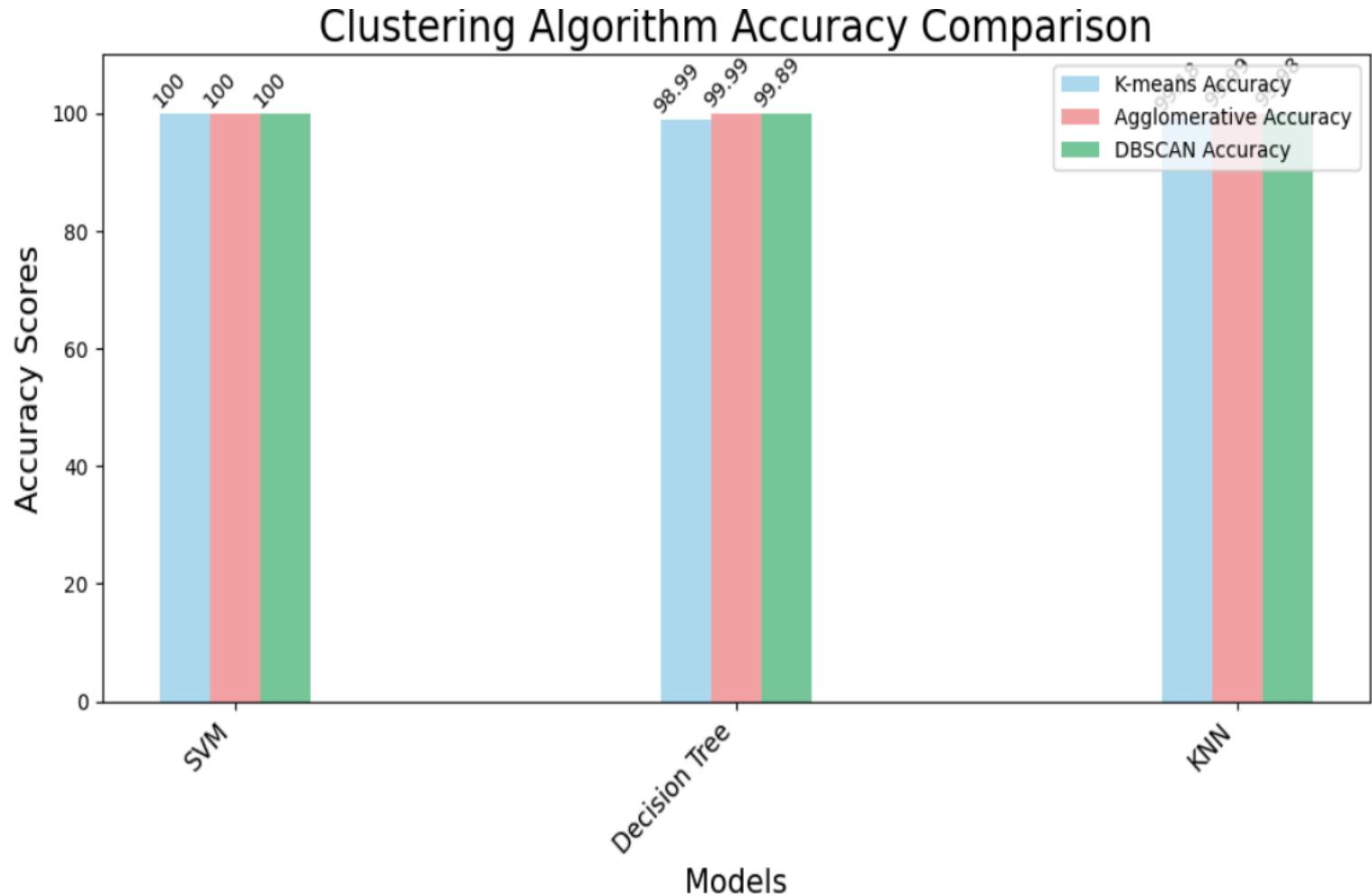


Performance Analysis
for PCA AND LDA

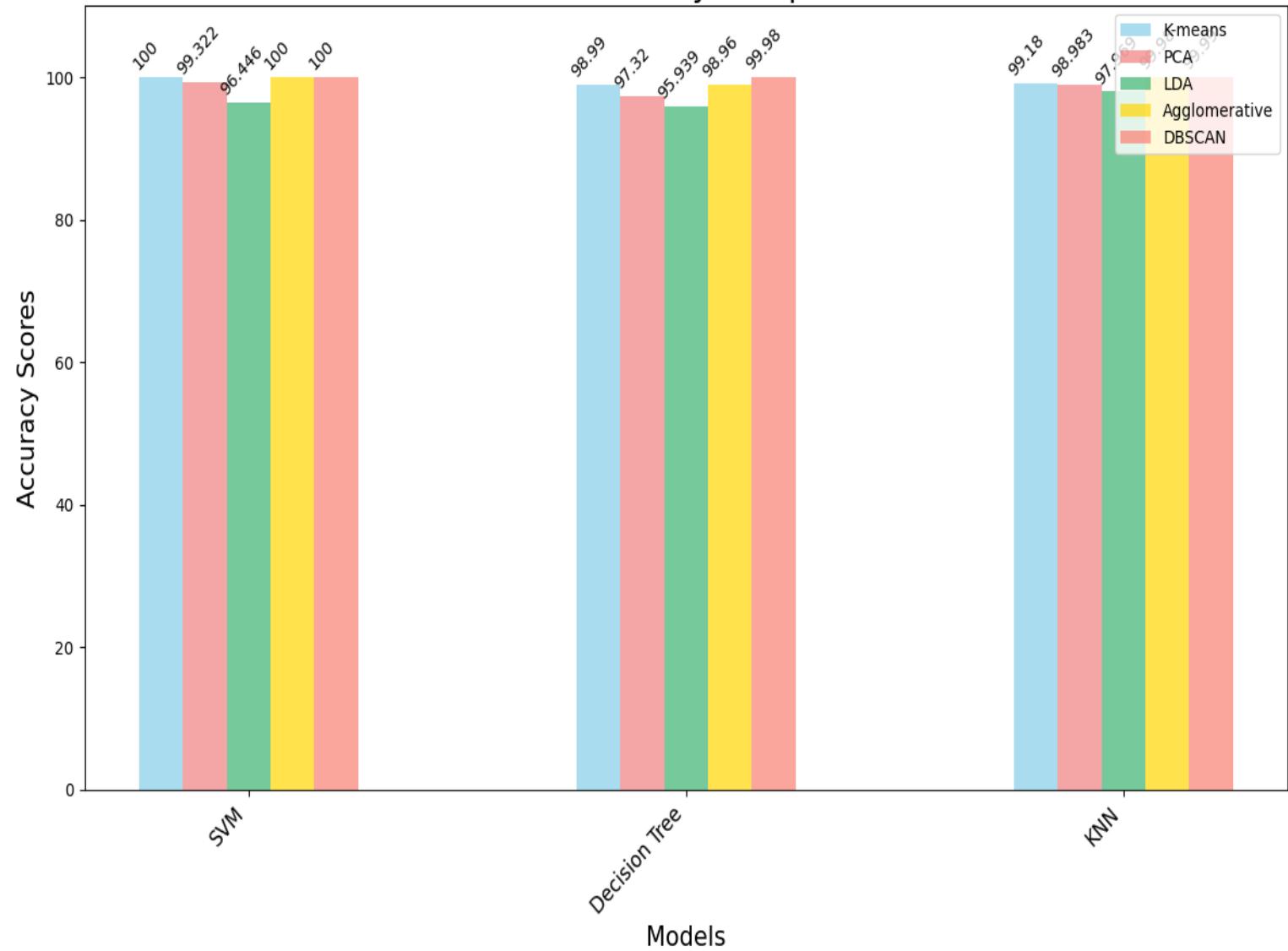
PERFORMANCE ANALYSIS FOR K MEANS



ACCURACY ANALYSIS FOR CLUSTERING ALGORITHMS



Model Accuracy Comparison



INFERENCE

In conclusion, the Support Vector Machine (SVM) model consistently demonstrated superior performance across clustering and Principal Component Analysis (PCA) tasks. Its accuracy, precision, recall, and F1 score outperformed alternative models, including Decision Trees and K-Nearest Neighbors, indicating its robustness in capturing complex patterns within the data. The SVM's exceptional performance suggests its suitability for both clustering and dimensionality reduction tasks, highlighting its versatility as a reliable and effective machine learning algorithm for the given dataset. These findings are integral to the model selection process, providing a strong basis for recommending SVM as the preferred choice in diverse analytical scenarios.

GUI

Top 3 Predicted Diseases

Selected Symptoms: phlegm, redness_of_eyes, sinus_pressure, runny_nose, puffy_face_and_eyes

Common Cold - Probability: 0.44

Drug Reaction - Probability: 0.2449967164151191

Pneumonia - Probability: 0.08

TEAM MEMBERS

KRISHNA MURTHY - CB.EN.U4CSE21016
S MEENAKSHI - CB.EN.U4CSE21035
LIKITHA PICHERI - CB.EN.U4CSE21044
RAGALA TEJDEEP - CB.EN.U4CSE21046

19CSE305-Machine Learning

[NOTE: This app is meant for demo purposes only. Please consult a Doctor if you have any symptoms.]

Select a symptom from the list :

- back_pain constipation abdominal_pain diarrhoea mild_fever
- yellow_urine yellowing_of_eyes acute_liver_failure fluid_overload swelling_of_stomach
- swelled_lymph_nodes malaise blurred_and_distorted_vision phlegm throat_irritation
- redness_of_eyes sinus_pressure runny_nose congestion chest_pain
- weakness_in_limbs fast_heart_rate pain_during_bowel_movements pain_in_anal_region bloody_stool
- irritation_in_anus neck_pain dizziness cramps bruising
- obesity swollen_legs swollen_blood_vessels puffy_face_and_eyes enlarged_thyroid
- brittle_nails swollen_extremities excessive_hunger extra_marital_contacts drying_and_tingling_lips
- slurred_speech knee_pain hip_joint_pain muscle_weakness stiff_neck
- swelling_joints movement_stiffness spinning_movements loss_of_balance unsteadiness
- weakness_of_one_body_side loss_of_smell bladder_discomfort foul_smell_of_urine continuous_feel_of_urine
- passage_of_gases internal_itching toxic_look_typhos depression irritability
- muscle_pain altered_sensorium red_spots_over_body belly_pain abnormal_menstruation
- dischromic_patches watering_from_eyes increased_appetite polyuria family_history
- mucoid_sputum rusty_sputum lack_of_concentration visual_disturbances receiving_blood_transfusion
- receiving_unsterile_injections coma stomach_bleeding distention_of_abdomen history_of_alcohol_consumption
- fluid_overload blood_in_sputum prominent_veins_on_calf palpitations painful_walking
- pus_filled_pimples blackheads scurving skin_peeling silver_like_dusting
- small_dents_in_nails inflammatory_nails blister red_sore_around_nose yellow_crust_ooze

Choose symptoms and click the button to predict possible diseases.

Submit

Clear