


Diabetes Risk: The Prediction of Diabetes Risk Through Unsupervised Learning

By: Katherine Nguyen



Agenda

- Outline/Agenda
- Introduction
- Approach
 - Data
 - Rundown
 - Data Cleaning
 - Models
 - Model Analysis
- Results
- Conclusion



Introduction

- **Problem**

- Diabetes is prevalent and being able to predict whether one has diabetes is useful
- How can we predict diabetes based on potential risk factors of diabetes?

- **Purpose**

- To understand the diabetes risk factors that can help predict whether an individual will acquire diabetes based on those risk factors in the future

- **Why Is it Important?**

- To raise awareness and potentially address prevention strategies for earlier stages before diabetes



Approach: Data

- **Dataset:** "Diabetes Risk Prediction"
 - Diabetes Database: <https://www.kaggle.com/datasets/rcratos/diabetes-risk-prediction>
 - Categorical Data; Binary

Approach: Data

[illegible]



Approach: Rundown

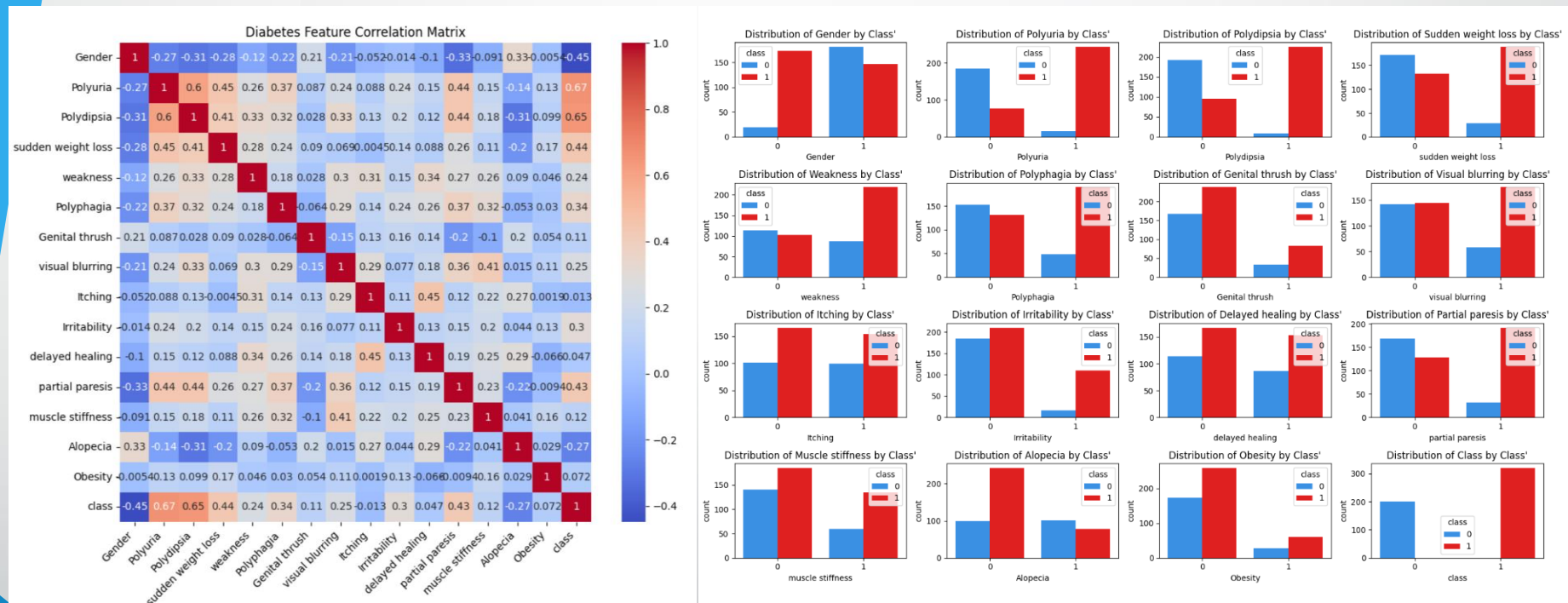
- **Data Cleaning :**

- Simplified the Data
 - Convert labels into binary values (e.g. Yes = 1, No = 0)
 - Removed unnecessary features for prediction

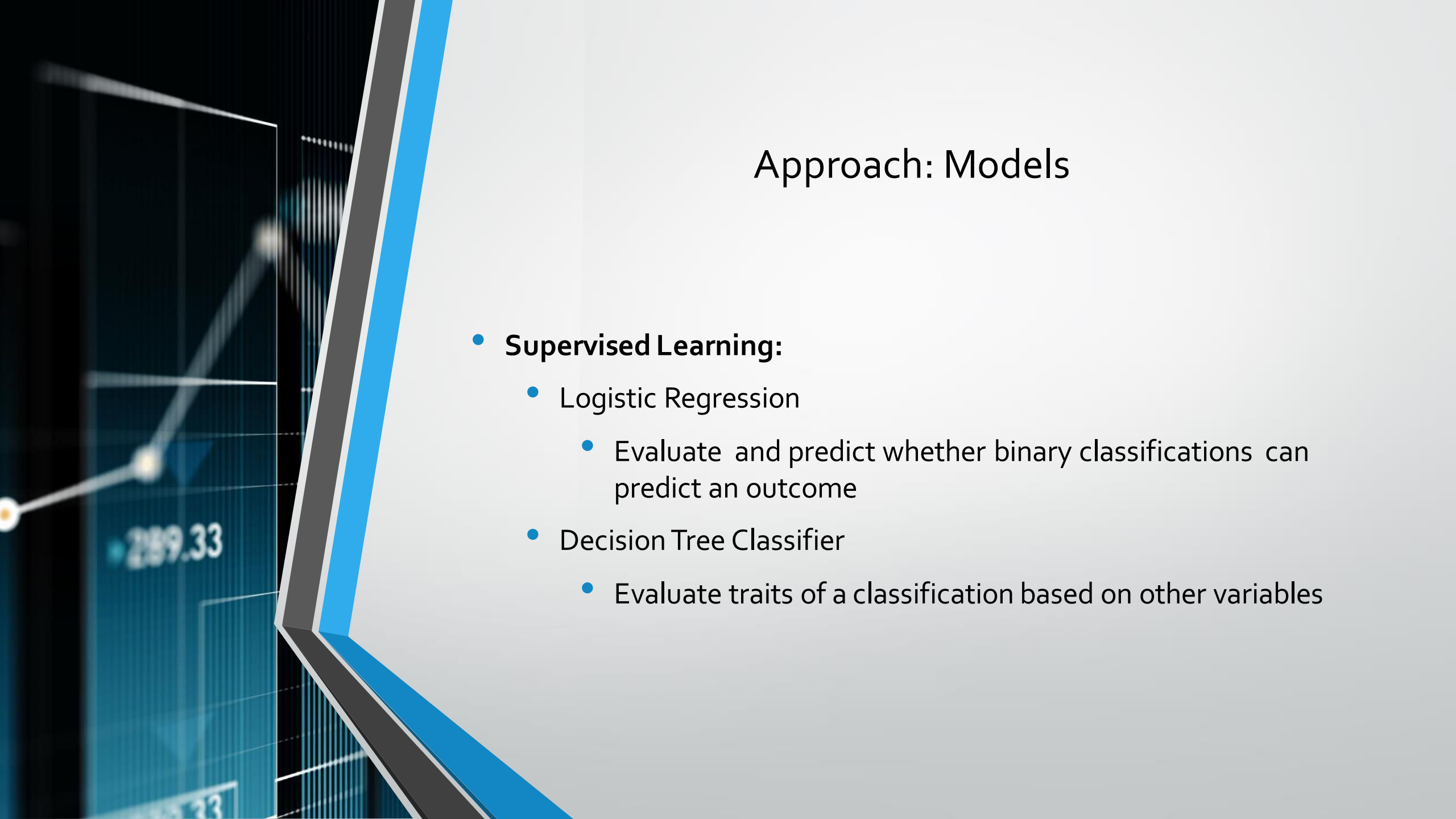
- **Exploratory Data Analysis :**

- Used Bar Plots to show classifications of diabetes in respect to other features

Approach: Data Cleaning



[illegible]

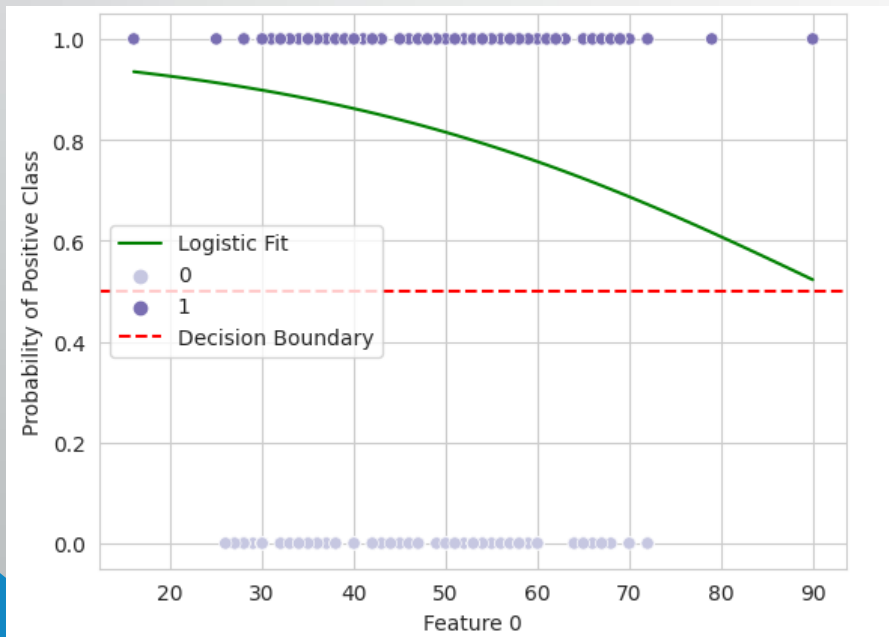


Approach: Models

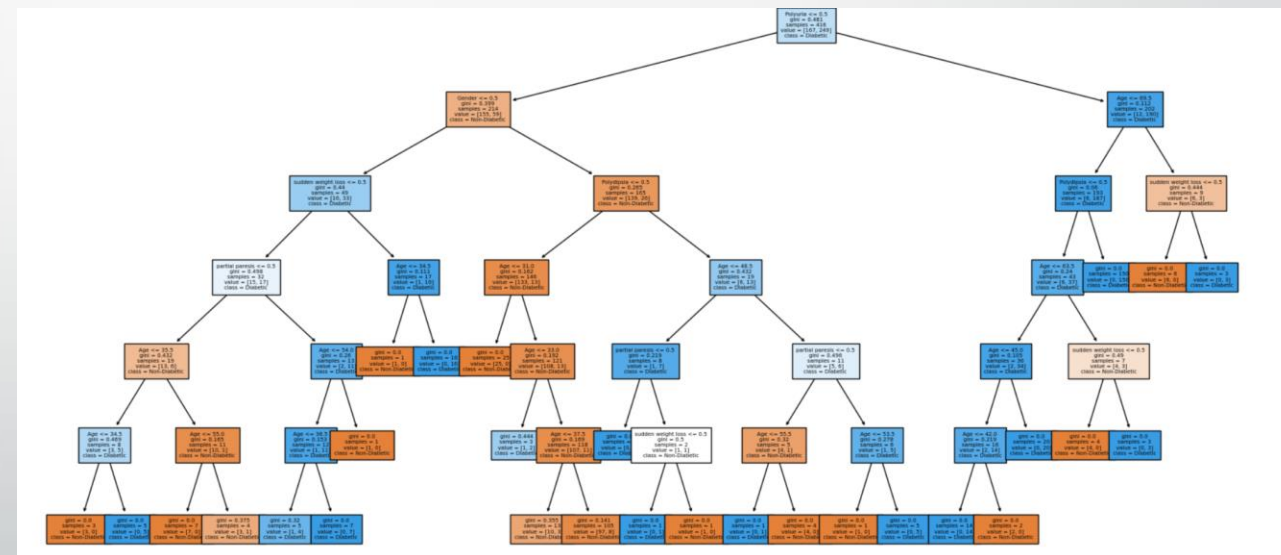
- **Unsupervised Learning:**
 - K-Means Clustering
 - Predicts where assumed clusters will be in data
 - Hierarchical Clustering
 - Predicts clusters based on similarities and branches to next decision

Approach: Models

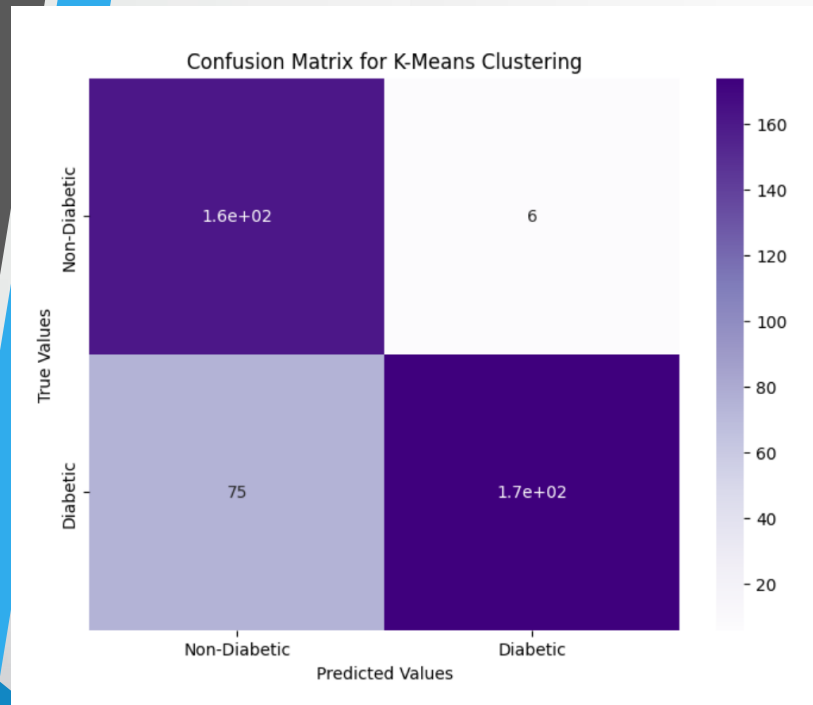
Logistic Regression



Decision Tree Classifier

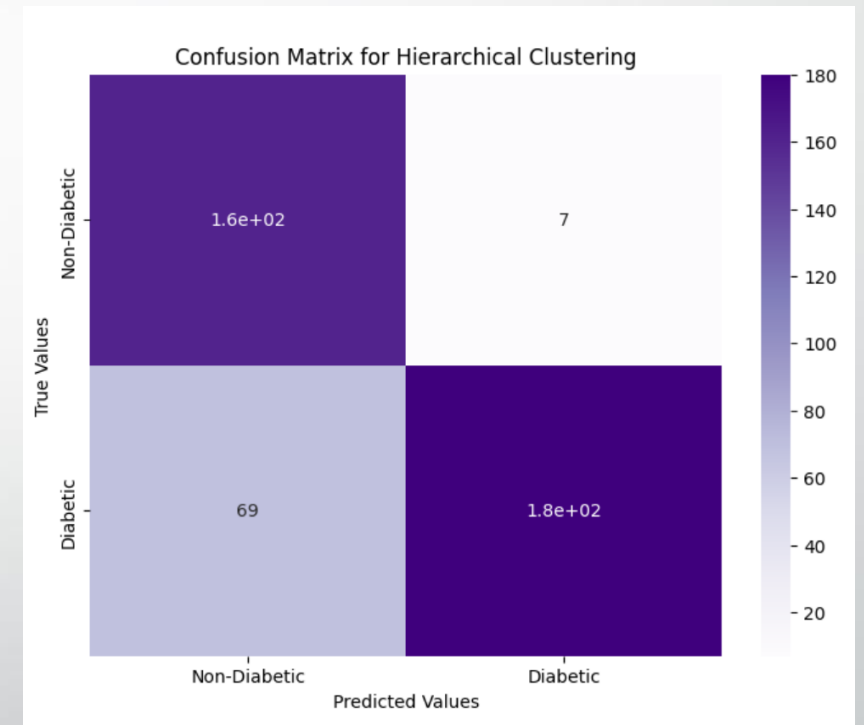


Results: Model Analysis



K-Means Clustering Label Ordering: (0, 1)
K-Means Clustering Accuracy: 0.8052884615384616
K-Means Clustering Precision: 0.9666666666666667

Silhouette Score: 0.44072574501920675



Hierarchical Clustering Label Ordering: (0, 1)
Hierarchical Clustering Accuracy: 0.8173076923076923
Hierarchical Clustering Precision: 0.9625668449197861

Silhouette Score: 0.37396613190254996



Results

- **Evaluation:**

- K-Means Cluster
 - Somewhat accurate and precise; Decent predictor model
 - Not ideal based on silhouette score; lacks similarity in clusters
 - Has potential with different parameters and less features
- Decision Tree classifier
 - Somewhat accurate and precise; Decent predictor model
 - Not ideal based on silhouette score; lacks similarity in clusters
 - Worse than K-Means Cluster
 - Has potential with different parameters and less features



Conclusion

- **K-Means Clustering** is the better model and has potential, but it fails to cluster by similarity, probably due to too many features
- **Hierarchical Clustering** has potential, but it fails to cluster by similarity, probably due to too many features and hyperparameters