# BRAC UNIVERSITY

Inspiring Excellence

## CSE422
Artificial Intelligence Lab

## Project Title
Obesity Level Prediction Using Machine Learning Models

Section No. 28
Group No. 01
Semester: Summer_2025
Submitted Date: 04-09-2025
Submitted to - Mahjabin Chowdhury, Sadman Sakib Alif

## Group Members:

| | |
|---|---|
| Dihan Islam Dhrubo | 23301458 |
| Afnan Ul Islam Afif | 23301528 |

**Table Of Contents:**

# Introduction:

Our project focuses on predicting obesity levels based on demographic, dietary, and physical activity features. Obesity is a major global health issue, and early identification of risk factors can guide better health interventions. We use supervised machine learning models to classify individuals into different obesity categories, along with exploratory and unsupervised methods.
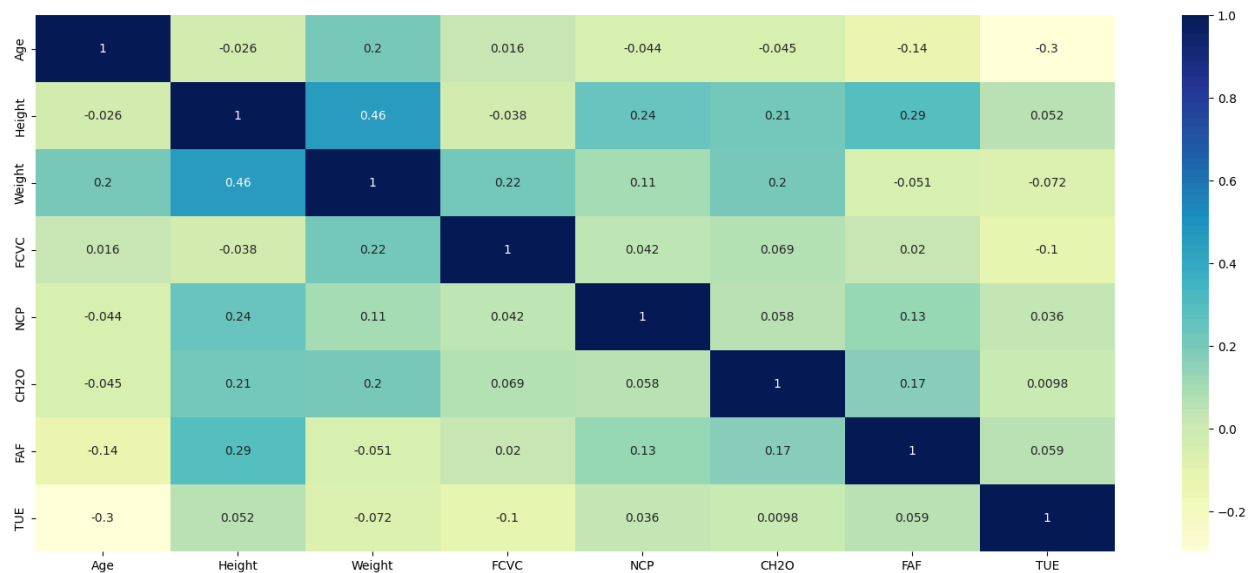
# Dataset-Description:

- **Total Features**: 16
- **Total Data Points:** 2111
- **Problem Type:**

  This is a **Classification problem** because the target variable, NObeyesdad, contains distinct categories. It classifies the result into categories rather than quantitative values.

- **Feature Type:**
  - **Quantitative Feature:**

    Age, Height, Weight, FCVC, NCP, CH2O, FAF, TUE.
  - **Categorical Features:**

    Gender, CALC, FAVC, SCC, SMOKE, family_history_with_overweight, CAEC, MTRANS.

**Dataset Heatmap before encoding:**



Only numeric features (Age, Height, Weight, FCVC, NCP, CH2O, FAF, TUE) are included. Categorical features (like Gender, CAEC, CALC, etc.) are still strings, so they are excluded from the correlation matrix because correlation can only be computed on numeric data. Interpretation is limited: you only see relationships among continuous variables.

- **Encoding Categorical Variables:**

  Yes, categorical values have been encoded here.
    - **Label Encoding** is used to encode Gender, CALC, FAVC, SCC, SMOKE, family_history_with_overweight, CAEC.
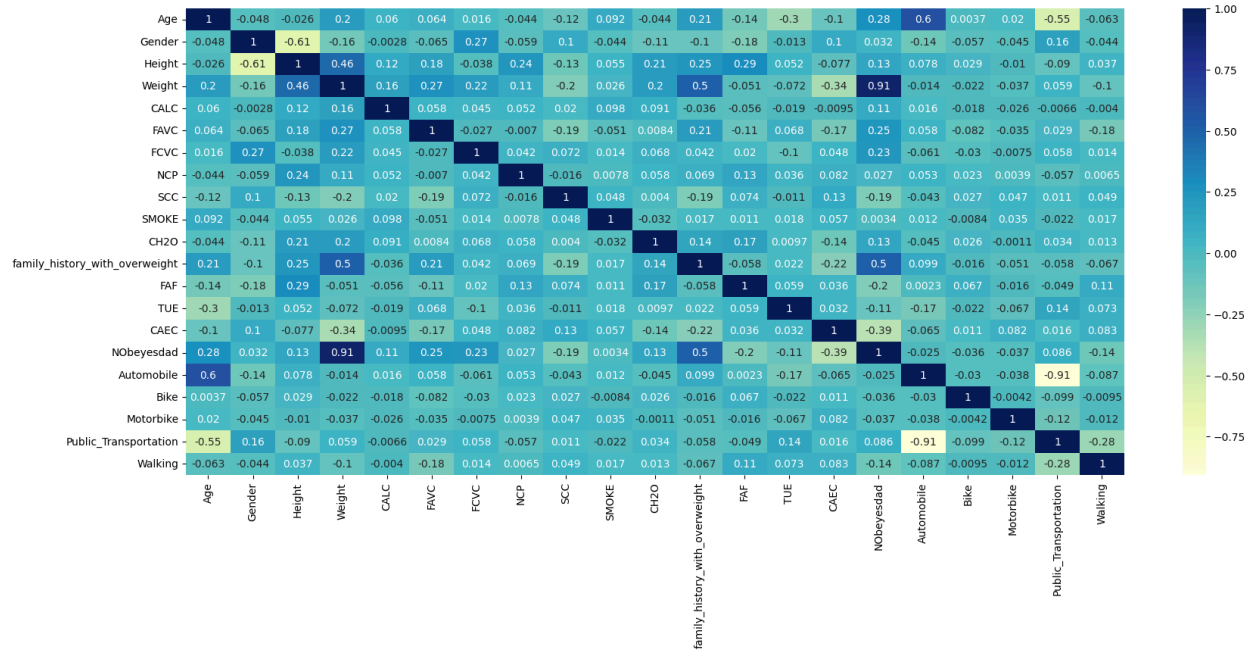    - **One-Hot Encoding** is used to encode MTRANS.

- **Correlation Insights:**
    - **Weight & Height (0.46 correlation)**
        - Taller individuals generally weigh more, as expected.
        - Both are strong predictors of obesity levels when combined into BMI.

- **Weight & Age (0.20 correlation)**
  - Older individuals show a slight trend of increased weight.
  - Indicates that obesity risk may rise with age.
- **Weight & Dietary/Activity features**
  - Weight has weak positive correlation with vegetable consumption(FCVC, 0.22) and water intake (CH2O, 0.20).
  - Suggests that even with better diet habits, higher weight persists in some individuals, highlighting complex interactions.
- **Height & Lifestyle Factors**
  - Height correlates moderately with NCP (number of main meals, 0.24), CH2O (water consumption, 0.21), and FAF (physical activity, 0.29).
  - Taller individuals tend to report slightly healthier lifestyle habits.
- **Negative Correlations**
  - TUE (time using technology) shows a negative correlation with Age (-0.30) and Weight (-0.07).
  - Indicates younger individuals spend more time on devices, potentially replacing active habits.
  - FAF (physical activity) correlates negatively with Age (-0.14), suggesting reduced activity in older participants.

# Correlations heatmap of All Features:

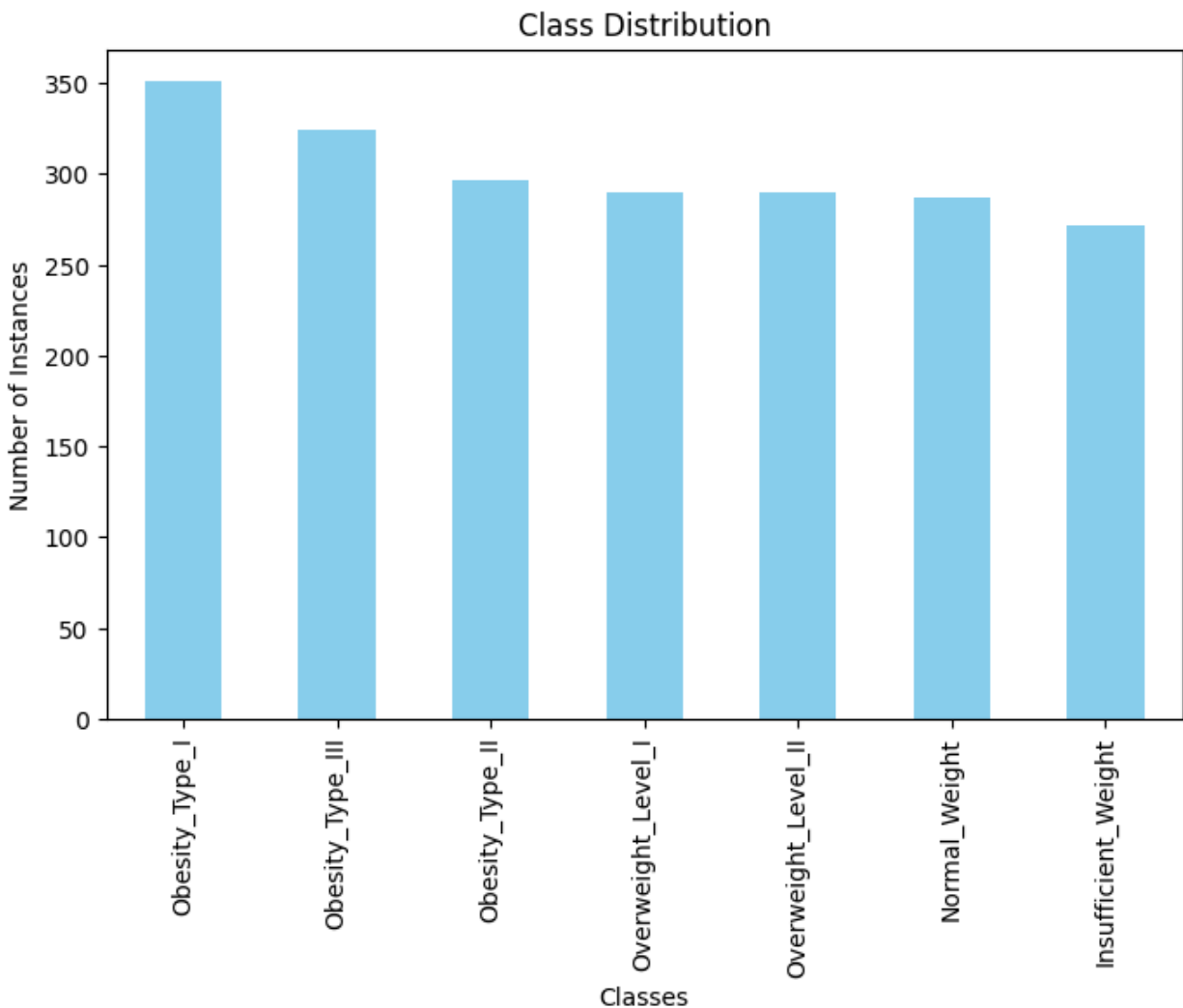| | Age | Gender | Height | Weight | CALC | FAVC | FCVC | NCP | SCC | SMOKE | CH2O | family_history_with_overweight | FAF | TUE | CAEC | NObeyesdad | Automobile | Bike | Motorbike | Public_Transportation | Walking |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Age | 1 | -0.048 | -0.026 | 0.2 | 0.06 | 0.064 | 0.016 | -0.044 | -0.12 | 0.092 | -0.044 | 0.21 | -0.14 | -0.3 | -0.1 | 0.28 | 0.6 | 0.0037 | 0.02 | -0.55 | -0.063 |
| Gender | -0.048 | 1 | -0.61 | -0.16 | -0.0028 | -0.065 | 0.27 | -0.059 | 0.1 | -0.044 | -0.11 | -0.1 | -0.18 | -0.013 | 0.1 | 0.032 | -0.14 | -0.057 | -0.045 | 0.16 | -0.044 |
| Height | -0.026 | -0.61 | 1 | 0.46 | 0.12 | 0.18 | -0.038 | 0.24 | -0.13 | 0.055 | 0.21 | 0.25 | 0.29 | 0.052 | -0.077 | 0.13 | 0.078 | 0.029 | -0.01 | -0.09 | 0.037 |
| Weight | 0.2 | -0.16 | 0.46 | 1 | 0.16 | 0.27 | 0.22 | 0.11 | -0.2 | 0.026 | 0.2 | 0.5 | -0.051 | -0.072 | -0.34 | 0.91 | -0.014 | -0.022 | -0.037 | 0.059 | -0.1 |
| CALC | 0.06 | -0.0028 | 0.12 | 0.16 | 1 | 0.058 | 0.045 | 0.052 | 0.02 | 0.098 | 0.091 | -0.036 | -0.056 | -0.019 | -0.0095 | 0.11 | 0.016 | -0.018 | -0.026 | -0.0066 | -0.004 |
| FAVC | 0.064 | -0.065 | 0.18 | 0.27 | 0.058 | 1 | -0.027 | -0.007 | -0.19 | -0.051 | 0.0084 | 0.21 | -0.11 | 0.068 | -0.17 | 0.25 | 0.058 | -0.082 | -0.035 | 0.029 | -0.18 |
| FCVC | 0.016 | 0.27 | -0.038 | 0.22 | 0.045 | -0.027 | 1 | 0.042 | 0.072 | 0.014 | 0.068 | 0.042 | 0.02 | -0.1 | 0.048 | 0.23 | -0.061 | -0.03 | -0.0075 | 0.058 | 0.014 |
| NCP | -0.044 | -0.059 | 0.24 | 0.11 | 0.052 | -0.007 | 0.042 | 1 | -0.016 | 0.0078 | 0.058 | 0.069 | 0.13 | 0.036 | 0.082 | 0.027 | 0.053 | 0.023 | 0.0039 | -0.057 | 0.0065 |
| SCC | -0.12 | 0.1 | -0.13 | -0.2 | 0.02 | -0.19 | 0.072 | -0.016 | 1 | 0.048 | 0.004 | -0.19 | 0.074 | -0.011 | 0.13 | -0.19 | -0.043 | 0.027 | 0.047 | 0.011 | 0.049 |
| SMOKE | 0.092 | -0.044 | 0.055 | 0.026 | 0.098 | -0.051 | 0.014 | 0.0078 | 0.048 | 1 | -0.032 | 0.017 | 0.011 | 0.018 | 0.057 | 0.0034 | 0.012 | -0.0084 | 0.035 | -0.022 | 0.017 |
| CH2O | -0.044 | -0.11 | 0.21 | 0.2 | 0.091 | 0.0084 | 0.068 | 0.058 | 0.004 | -0.032 | 1 | 0.14 | 0.17 | 0.0097 | -0.14 | 0.13 | -0.045 | 0.026 | -0.0011 | 0.034 | 0.013 |
| family_history_with_overweight | 0.21 | -0.1 | 0.25 | 0.5 | -0.036 | 0.21 | 0.042 | 0.069 | -0.19 | 0.017 | 0.14 | 1 | -0.058 | 0.022 | -0.22 | 0.5 | 0.099 | -0.016 | -0.051 | -0.058 | -0.067 |
| FAF | -0.14 | -0.18 | 0.29 | -0.051 | -0.056 | -0.11 | 0.02 | 0.13 | 0.074 | 0.011 | 0.17 | -0.058 | 1 | 0.059 | 0.036 | -0.2 | 0.0023 | 0.067 | -0.016 | -0.049 | 0.11 |
| TUE | -0.3 | -0.013 | 0.052 | -0.072 | -0.019 | 0.068 | -0.1 | 0.036 | -0.011 | 0.018 | 0.0097 | 0.022 | 0.059 | 1 | 0.032 | -0.11 | -0.17 | -0.022 | -0.067 | 0.14 | 0.073 |
| CAEC | -0.1 | 0.1 | -0.077 | -0.34 | -0.0095 | -0.17 | 0.048 | 0.082 | 0.13 | 0.057 | -0.14 | -0.22 | 0.036 | 0.032 | 1 | -0.39 | -0.065 | 0.011 | 0.082 | 0.016 | 0.083 |
| NObeyesdad | 0.28 | 0.032 | 0.13 | 0.91 | 0.11 | 0.25 | 0.23 | 0.027 | -0.19 | 0.0034 | 0.13 | 0.5 | -0.2 | -0.11 | -0.39 | 1 | -0.025 | -0.036 | -0.037 | 0.086 | -0.14 |
| Automobile | 0.6 | -0.14 | 0.078 | -0.014 | 0.016 | 0.058 | -0.061 | 0.053 | -0.043 | 0.012 | -0.045 | 0.099 | 0.0023 | -0.17 | -0.065 | -0.025 | 1 | -0.03 | -0.038 | -0.91 | -0.087 |
| Bike | 0.0037 | -0.057 | 0.029 | -0.022 | -0.018 | -0.082 | -0.03 | 0.023 | 0.027 | -0.0084 | 0.026 | -0.016 | 0.067 | -0.022 | 0.011 | -0.036 | -0.03 | 1 | -0.0042 | -0.099 | -0.0095 |
| Motorbike | 0.02 | -0.045 | -0.01 | -0.037 | -0.026 | -0.035 | -0.0075 | 0.0039 | 0.047 | 0.035 | -0.0011 | -0.051 | -0.016 | -0.067 | 0.082 | -0.037 | -0.038 | -0.0042 | 1 | -0.12 | -0.012 |
| Public_Transportation | -0.55 | 0.16 | -0.09 | 0.059 | -0.0066 | 0.029 | 0.058 | -0.057 | 0.011 | -0.022 | 0.034 | -0.058 | -0.049 | 0.14 | 0.016 | 0.086 | -0.91 | -0.099 | -0.12 | 1 | -0.28 |
| Walking | -0.063 | -0.044 | 0.037 | -0.1 | -0.004 | -0.18 | 0.014 | 0.0065 | 0.049 | 0.017 | 0.013 | -0.067 | 0.11 | 0.073 | 0.083 | -0.14 | -0.087 | -0.0095 | -0.012 | -0.28 | 1 |

## ● Interpretation of Correlation Test:

Weight and Height are the most influential features in relation to obesity classification. Lifestyle factors (diet, exercise, screen time) show weaker but meaningful correlations, reinforcing that obesity is multi-factorial. The heatmap suggests that while physical measurements dominate prediction, behavioral factors refine class separation.
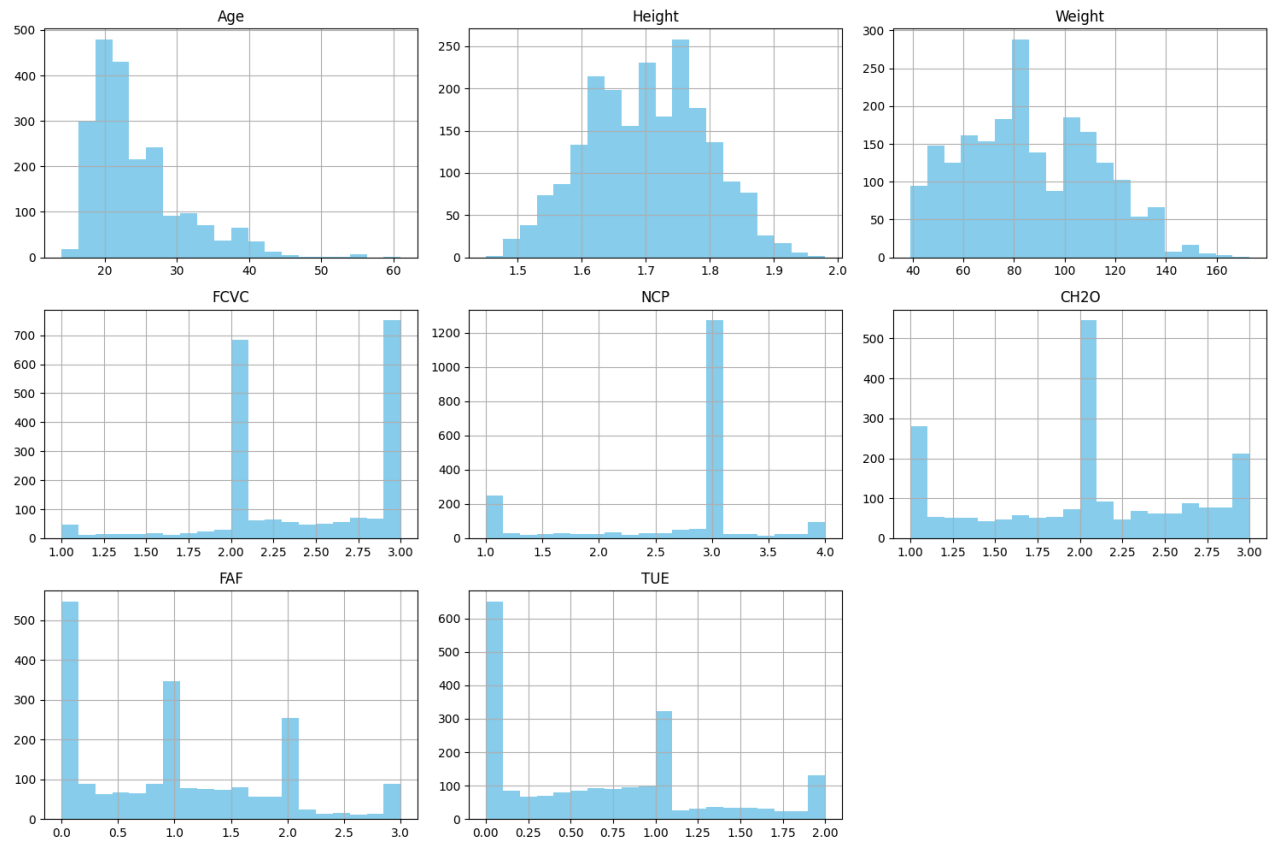
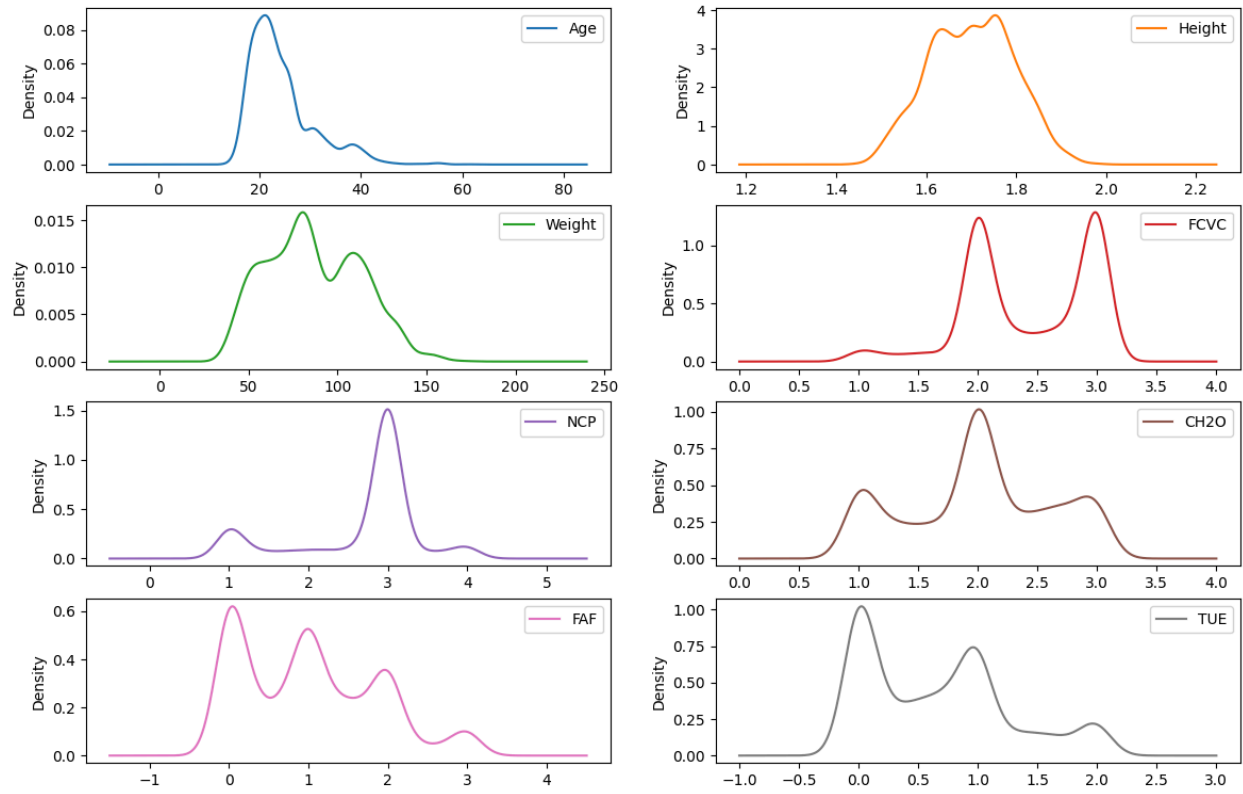# Imbalanced Dataset

**The dataset is imbalanced:**

not all unique classes of the target feature "NObeyesdad" have an equal number of instances. Some obesity categories (e.g., Obesity_Type_I with 351 samples) are more frequent, while others (e.g., Insufficient_Weight with 272 samples) are less represented.
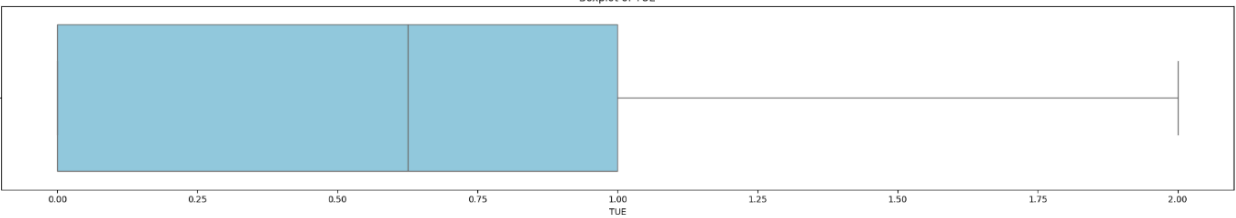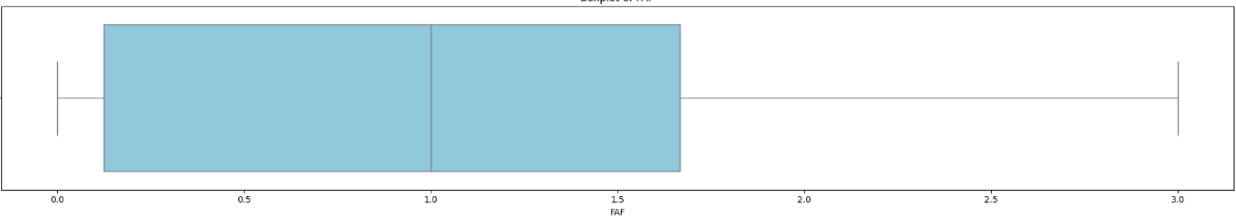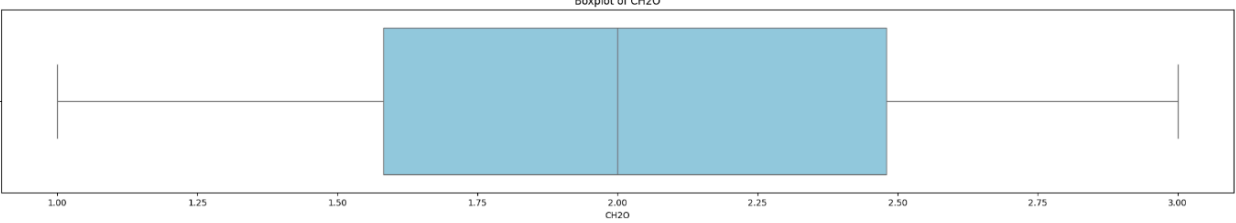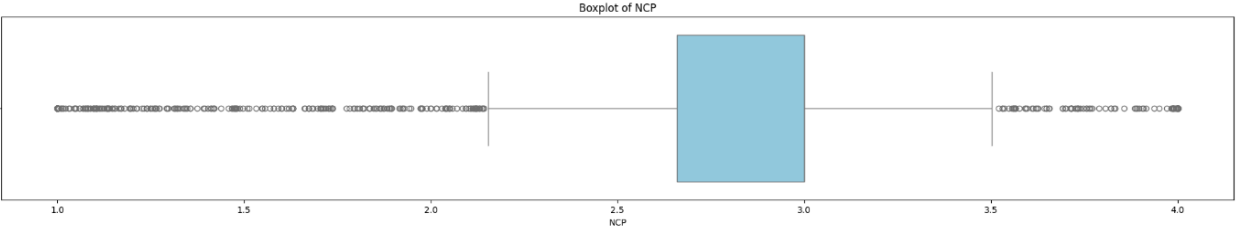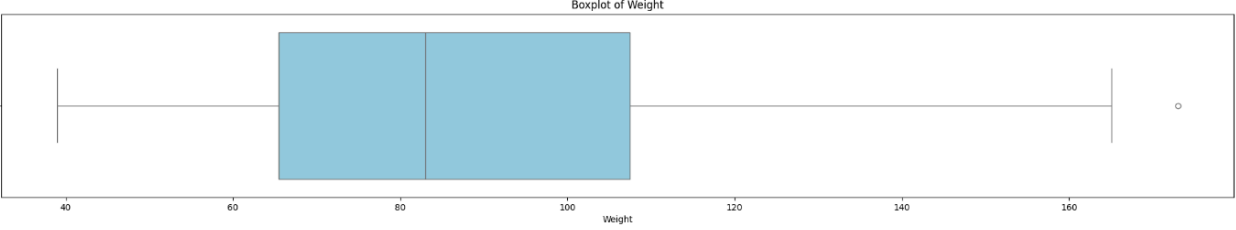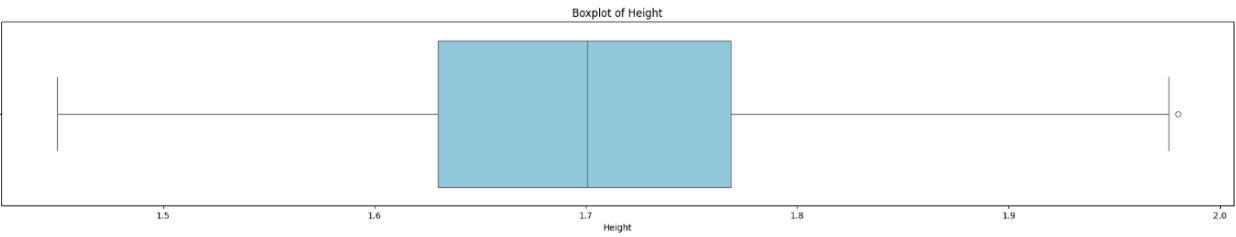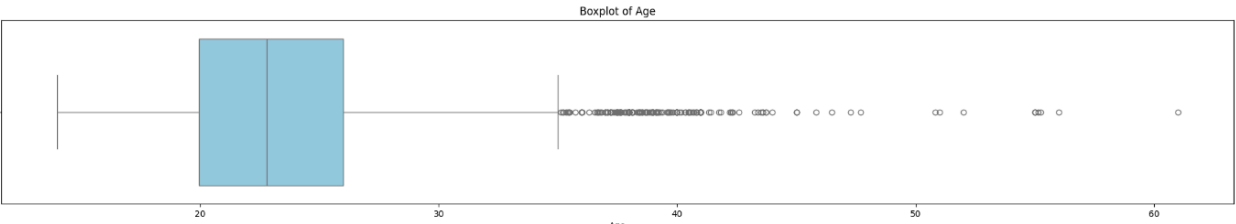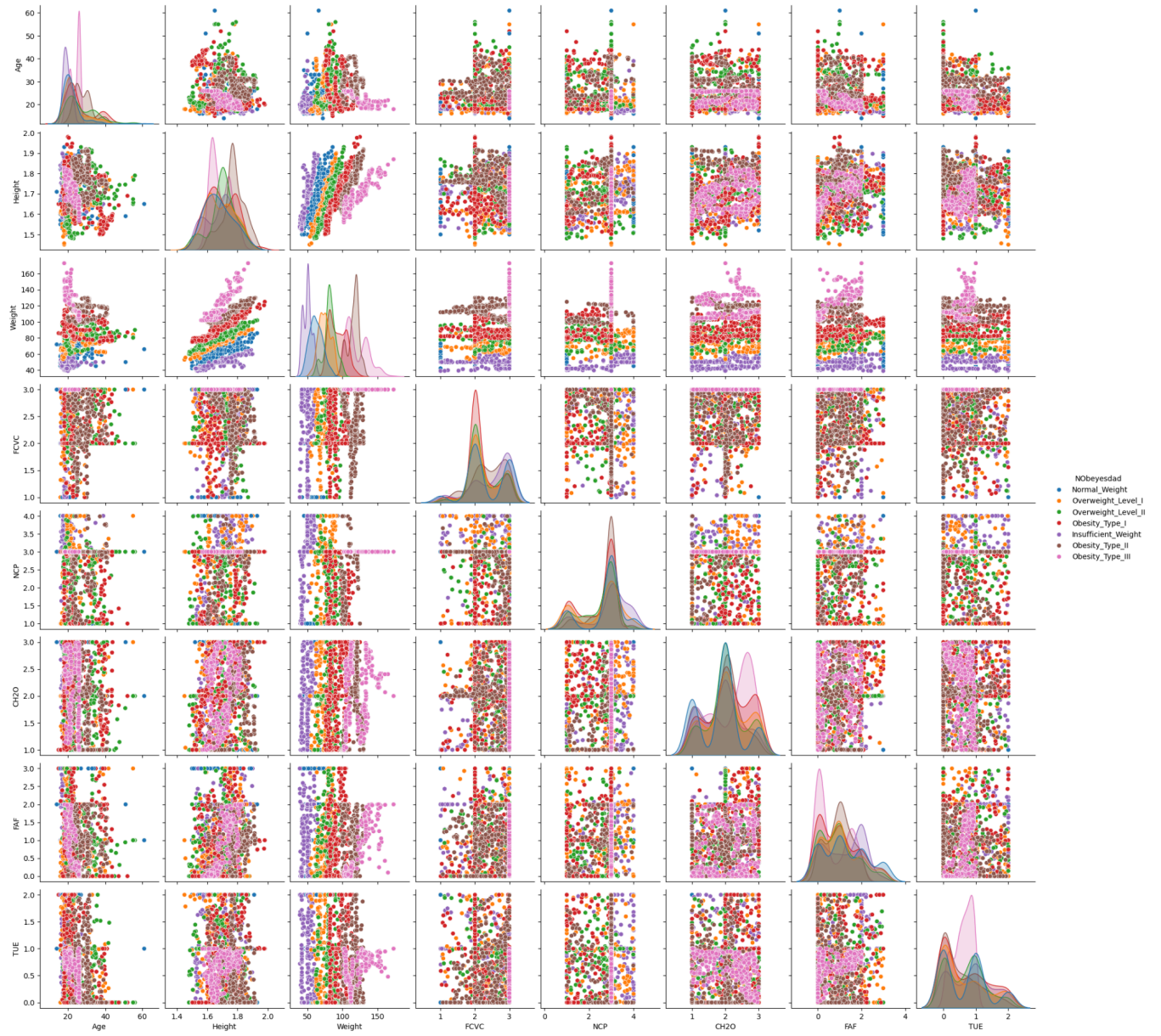


Class Distribution

# Exploratory Data Analysis (EDA)

Density plot of Numerical features

Boxplot of Age


Boxplot of Height


Boxplot of Weight


Boxplot of FCVC


Boxplot of NCP


Boxplot of CH2O


Boxplot of FAF


Boxplot of TUE

Distribution of Gender

Distribution of CALC

Distribution of FAVC

Distribution of SCC

Distribution of SMOKE

Distribution of family_history_with_overweight

Distribution of CAEC



Distribution of MTRANS



Distribution of NObeyesdad

# Dataset Pre-Processing

- **Faults:**

### Null Values:

```
Missing Values:
 Age                              0
Gender                          11
Height                           0
Weight                           0
CALC                             0
FAVC                             0
FCVC                             0
NCP                              0
SCC                              0
SMOKE                            0
CH2O                            21
family_history_with_overweight  21
FAF                              0
TUE                              0
CAEC                             0
MTRANS                           0
NObeyesdad                       0
```

### Categorical Values:

```
Feature Types:
 Age                             float64
Gender                           object
Height                          float64
Weight                          float64
CALC                             object
FAVC                             object
FCVC                            float64
NCP                             float64
SCC                              object
SMOKE                            object
CH2O                            float64
family_history_with_overweight   object
FAF                             float64
TUE                             float64
CAEC                             object
MTRANS                           object
NObeyesdad                       object
dtype: object
```

# <u>Feature Scaling</u>

1. **Robust Scaler:** Values with too many outliers.
2. **Standard Scaler:** Values with lower number of outliers.

## ● Solution:

- ○ **Imputing Values:**
  - ■ **Mode Imputation:** The null values in the features with less outliers are suitable for using mode imputation. Because these values don't have a numeric average. So the most frequent values are the ones chosen as a substitute for the null values. E.g.: Gender, CALC , CAEC , FAVC , SMOKE, SCC , MTRANS, family_history_with_overweight.
  - ■ **Median Imputation:** Median is robust to outliers, so it was applied to skewed numerical features. E.g.: Age, Weight, Height, FCVC, CH2O, FAF, TUE.

- ○ **Encoding:**
  - ■ **One-Hot Encoding:** The MTRANS feature (mode of transportation) was handled using One-Hot Encoding, since it is a categorical feature with no intrinsic ordering.
  - ■ **Label-Encoding:** Gender, CALC, Sometimes, FAVC, SCC, SMOKE, family_history_with_overweight, CAEC, these features are with no intrinsic values and mostly True False. That is why label encoding is used.
- ○ **Scaling:**

  Feature scaling is used different scalers depending on the distribution of the data
  - ■ **StandardScaler:** For normally distributed features such as Height, Weight, FCVC, CH2O, FAF, TUE.
  - ■ **RobustScaler:** For skewed or outlier-prone features such asAge, NCP, CALC, CAEC.

# Dataset Splitting

To evaluate model performance fairly, the dataset is divided into training and testing subsets.

- **Stratified Split**:

  Since the target variable (NObeyesdad) has 7 classes with imbalanced distribution, we used a stratified split. This ensured that each subset maintained the same class proportion as the original dataset, avoiding bias toward majority classes.

- **Train/Test Ratio**:

  The dataset was split into:
  - 70% Training set – used to train the models.
  - 30% Testing set – used to evaluate the trained models on unseen data.

- **Validation**:

  For the **Neural Network**, we also enabled **early stopping**, which internally uses a portion of the training set as a validation set to prevent overfitting.

This splitting strategy ensured that the models learned from sufficient data while being tested on a representative and unbiased subset.

# Model Training and Testing

We trained and evaluated several supervised learning models, along with an unsupervised approach, on the preprocessed obesity dataset. Each model was implemented using scikit-learn.

**1. K-Nearest Neighbors (KNN)**

- **Implementation:** Tested for values of k from 1 to 30.
- **Best k:** Selected based on maximum accuracy on the test set.
- **Observation:** Performance depended strongly on the choice of k.
- **Result:** Best value of K is 1. Accuracy of the model for the value of K mentioned above: 86.12%

## 2. Logistic Regression

- **Implementation:** Solver = "lbfgs", max_iter = 1000.
- **Strengths:** Stable linear baseline.
- **Limitation:** Limited when classes are not linearly separable.
- **Result:** Logistic Regression prediction accuracy: 88.01%

## 3. Naive Bayes (GaussianNB)

- **Implementation:** Trained on scaled features.
- **Strengths:** Simple and fast.
- **Limitation:** Assumes feature independence, which is not fully valid in this dataset, leading to lower accuracy.
- **Result:** Naive Bayes Prediction accuracy 54.42%

## 4. Decision Tree

- **Implementation:** Criterion = "entropy", random_state = 42.
- **Strengths:** Interpretable, captures non-linear relationships.
- **Limitation:** Without pruning, risk of overfitting.
- **Result:** Decision Tree prediction accuracy: 94.01%

## 5. Neural Network (MLPClassifier)

- **Architecture:** Hidden layers (128, 64), activation = ReLU.
- **Training:** Solver = Adam, learning rate = adaptive, early stopping enabled.
- **Strengths:** Achieved the **highest accuracy**, captured complex feature interactions.
- **Limitation:** Requires more computation and careful tuning.
- **Result:** Test accuracy is 95.27%

## 6. KMeans Clustering (Unsupervised)

- **Implementation:** Applied KMeans to uncover natural groupings in the dataset.
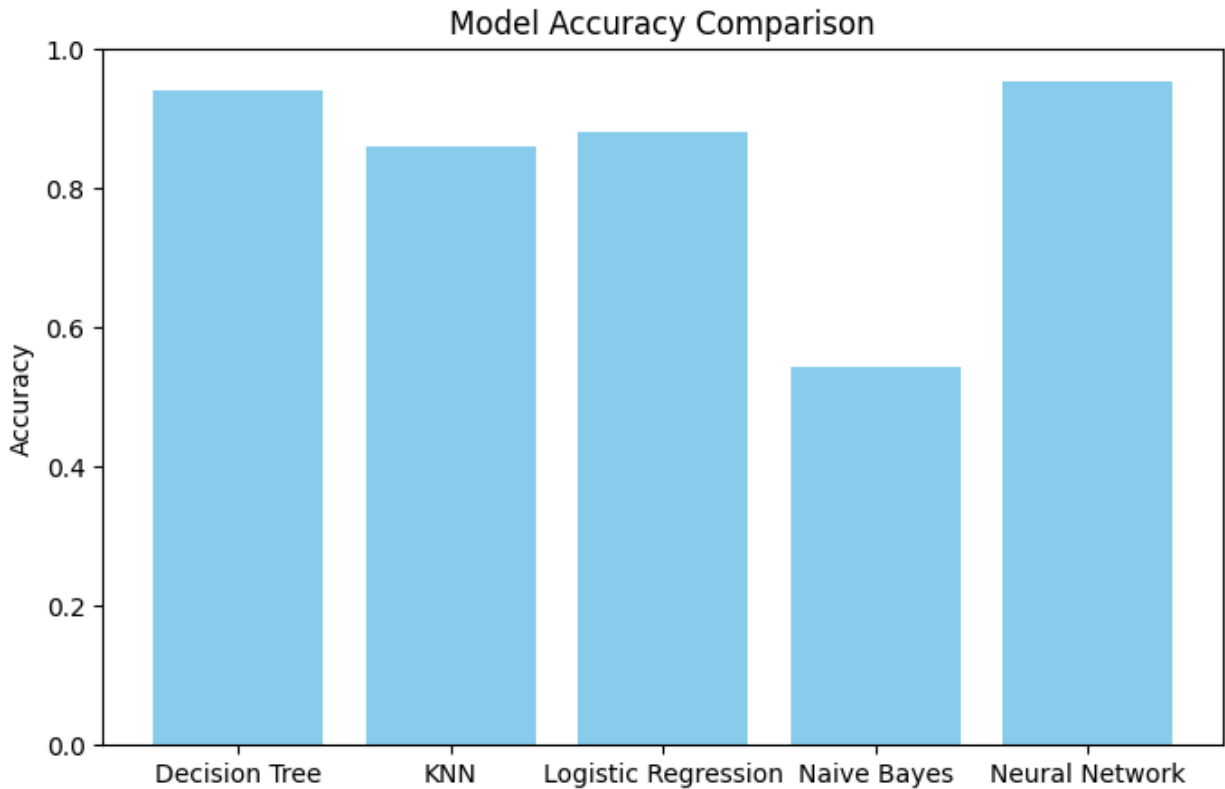- **Evaluation:** Measured with silhouette score.

- **Observation:** Some clusters aligned with obesity levels, though unsupervised performance was weaker than supervised methods.

# **Model Selection / Comparison Analysis**

To evaluate the performance of all models, we used accuracy, precision, recall, confusion matrix, and ROC-AUC curves. Since this is a classification problem, regression metrics ($R^2$, loss) were not applicable.
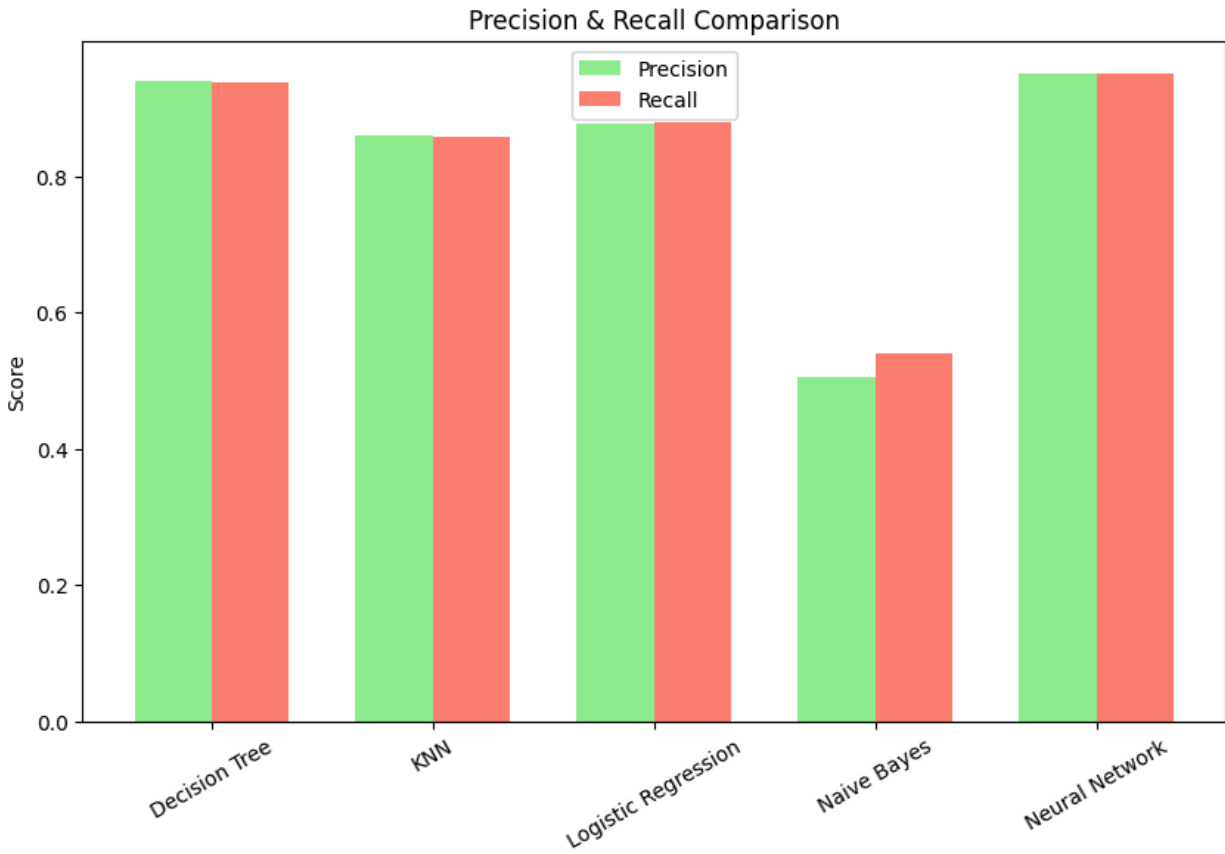
**Accuracy Comparison:**
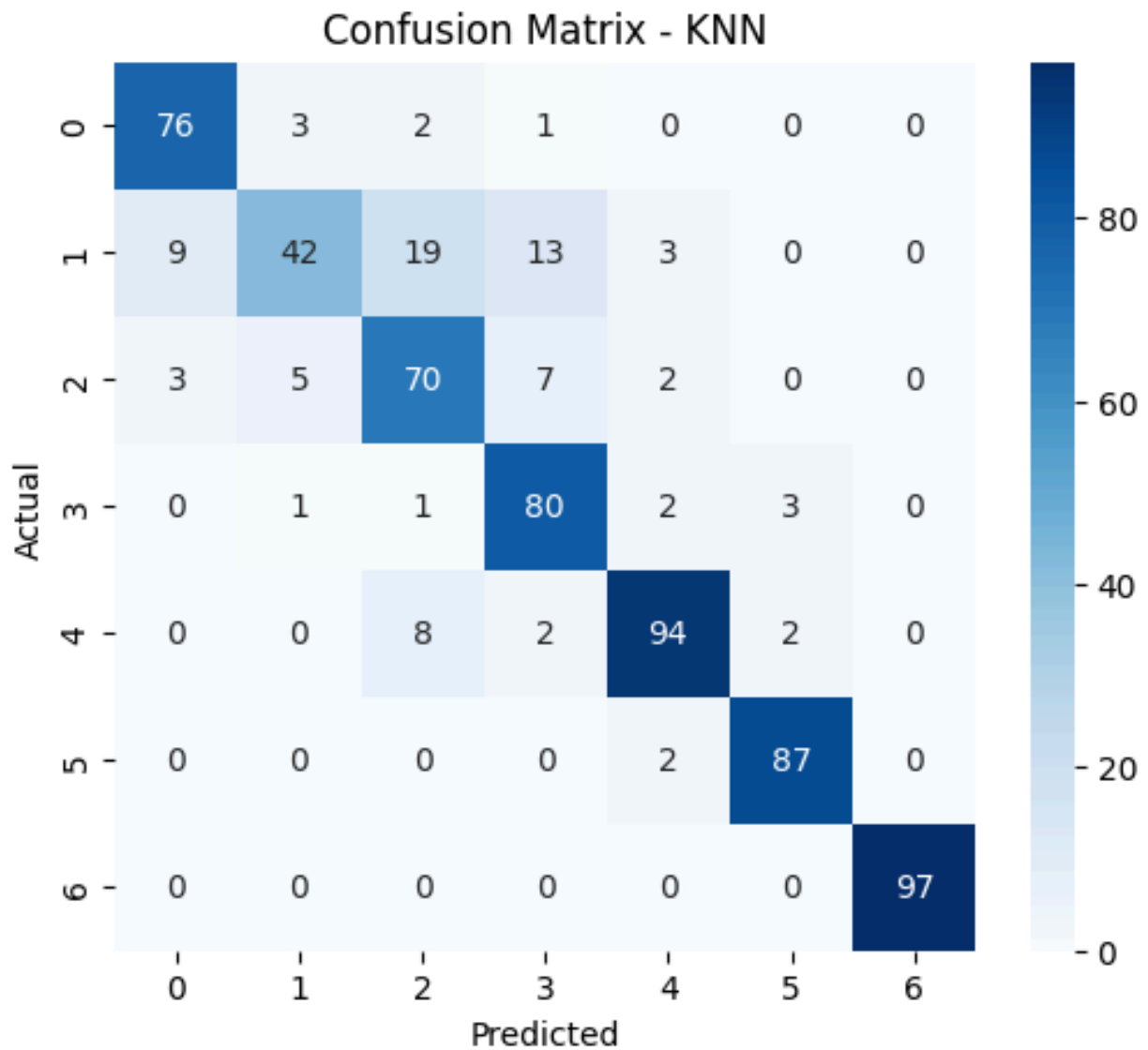
A bar chart was plotted comparing the accuracy of each supervised model.

Model Accuracy Comparison

- **Neural Network (MLP)** achieved the highest accuracy, due to its ability to learn complex non-linear patterns.
- **KNN** and Logistic Regression performed moderately well, serving as reliable baselines.
- **Decision Tree** captured non-linear relationships but tended to overfit, resulting in unstable performance.
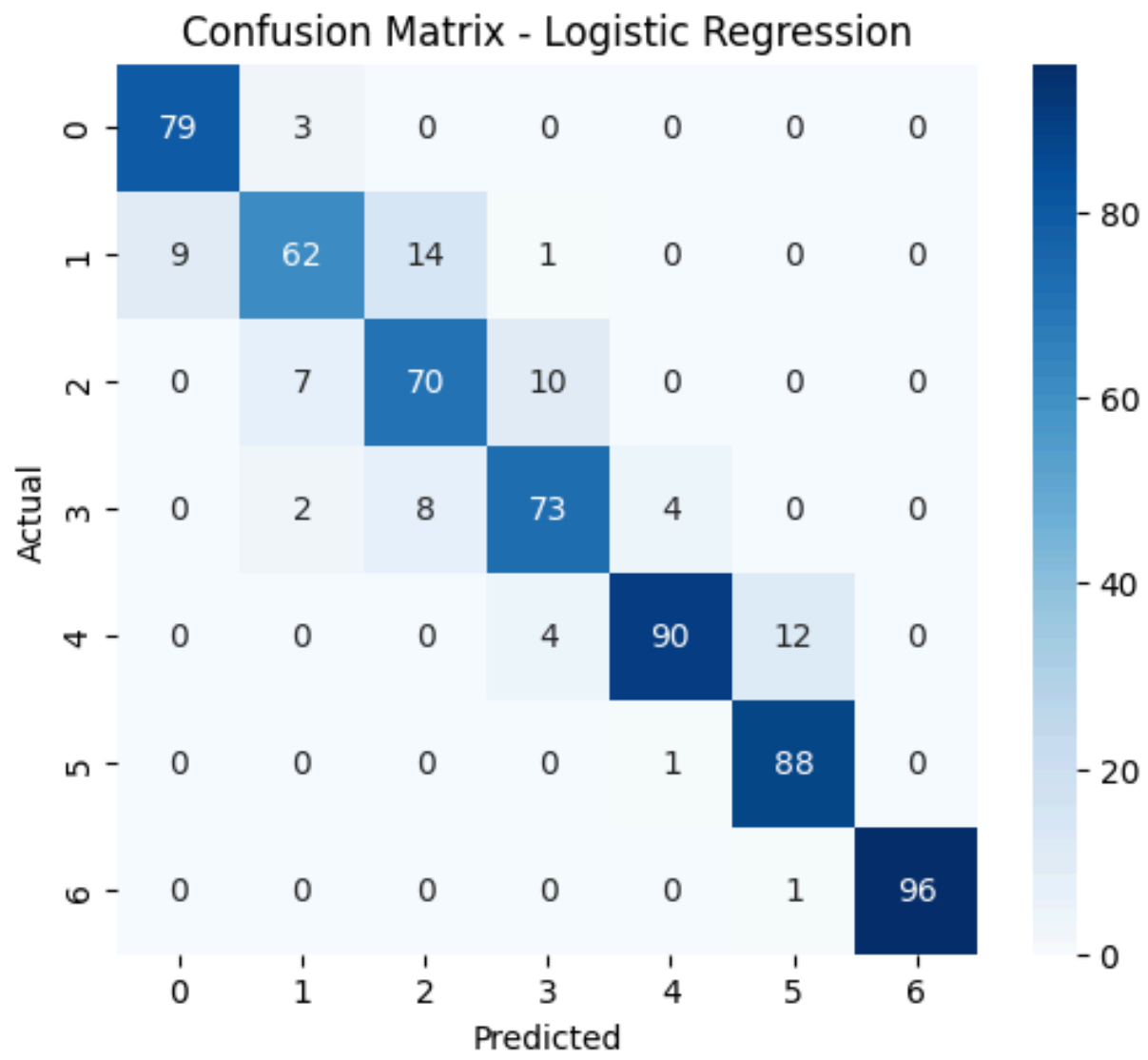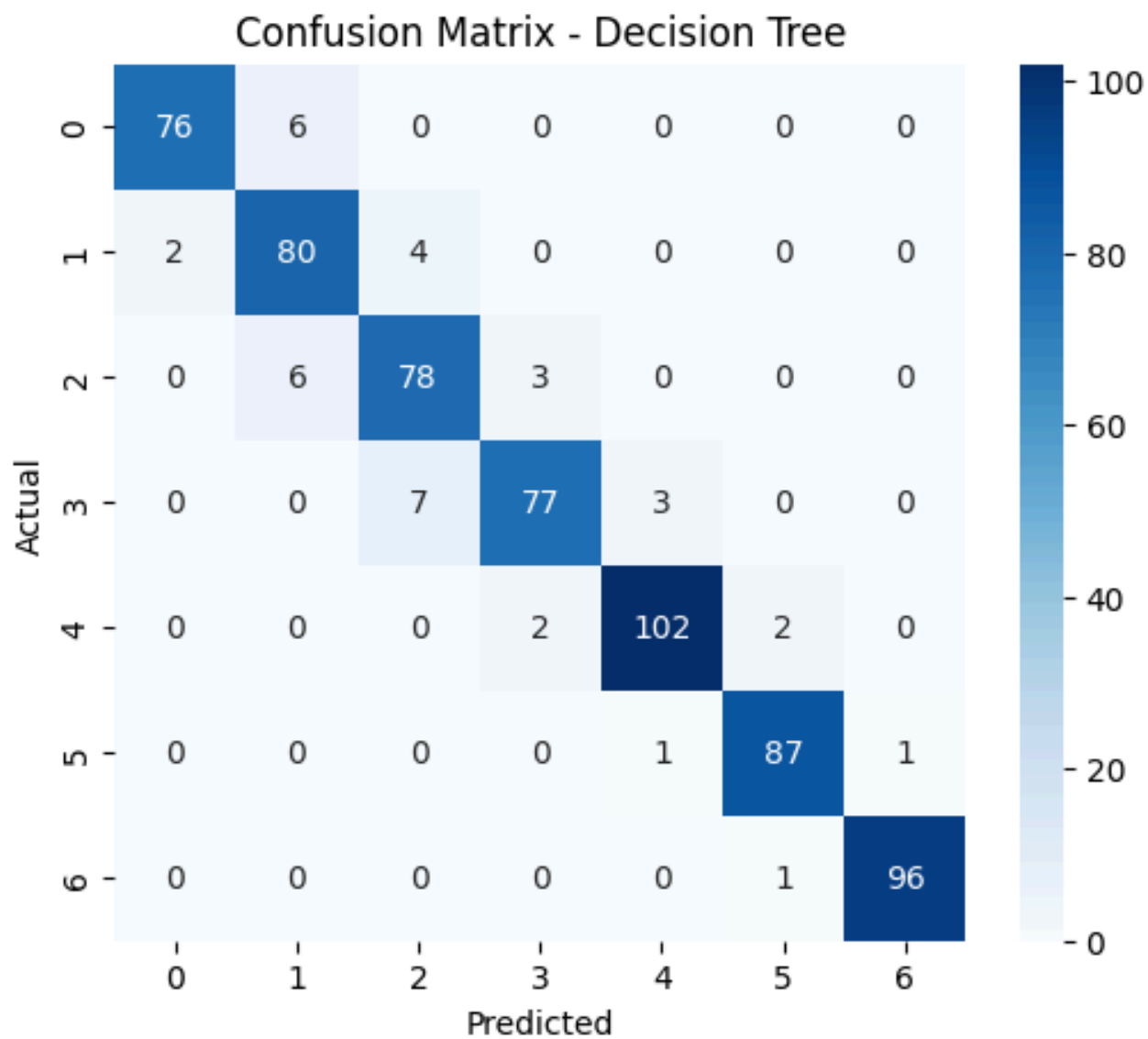- **Naive Bayes** was the weakest performer, as it assumes independence among features, which is unrealistic for this dataset.
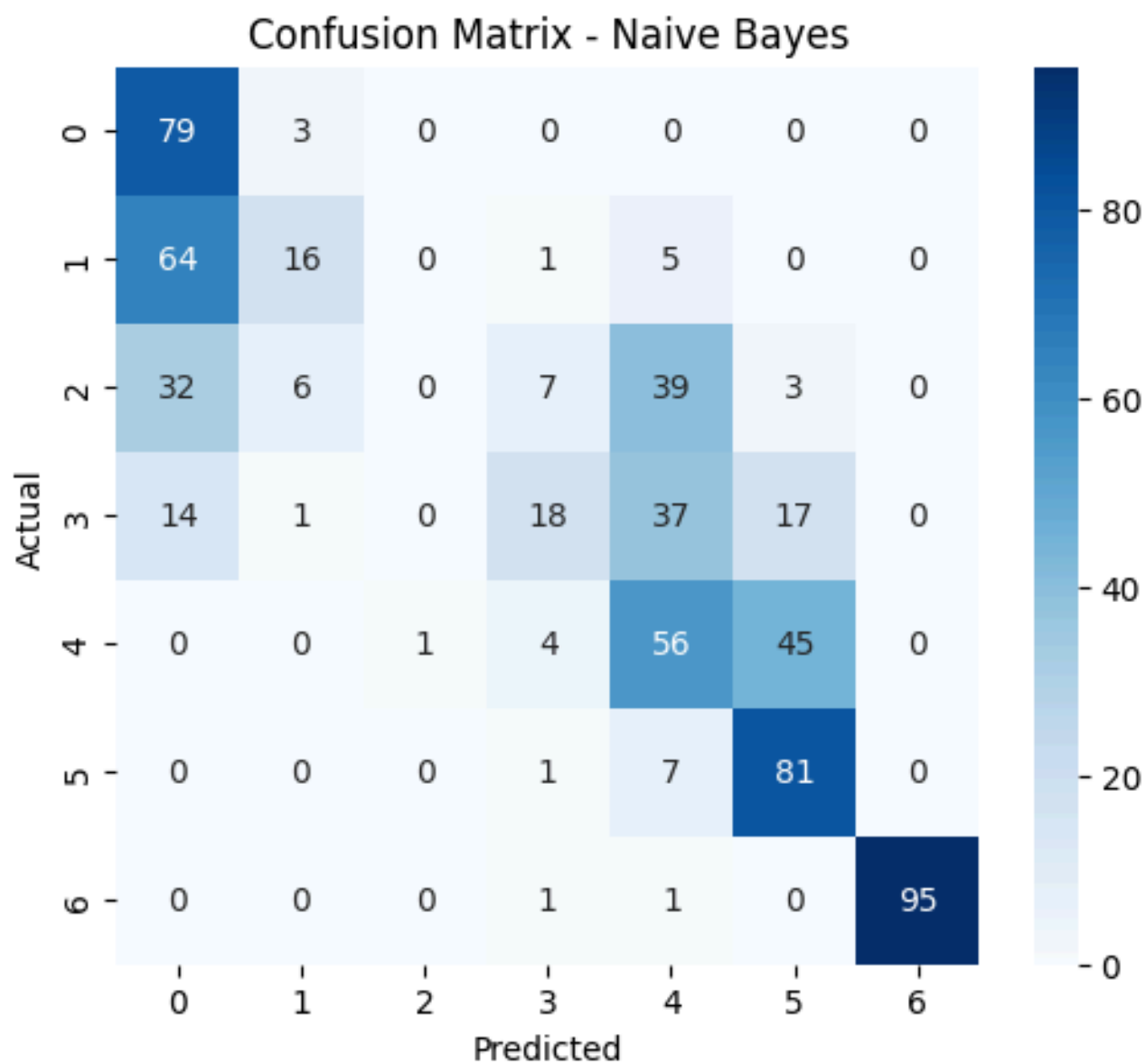
**Precision & Recall:**

Precision & Recall Comparison

- **Neural Network** achieved strong precision and recall across most obesity classes.

- **KNN** and Logistic Regression showed moderate balance between precision and recall.

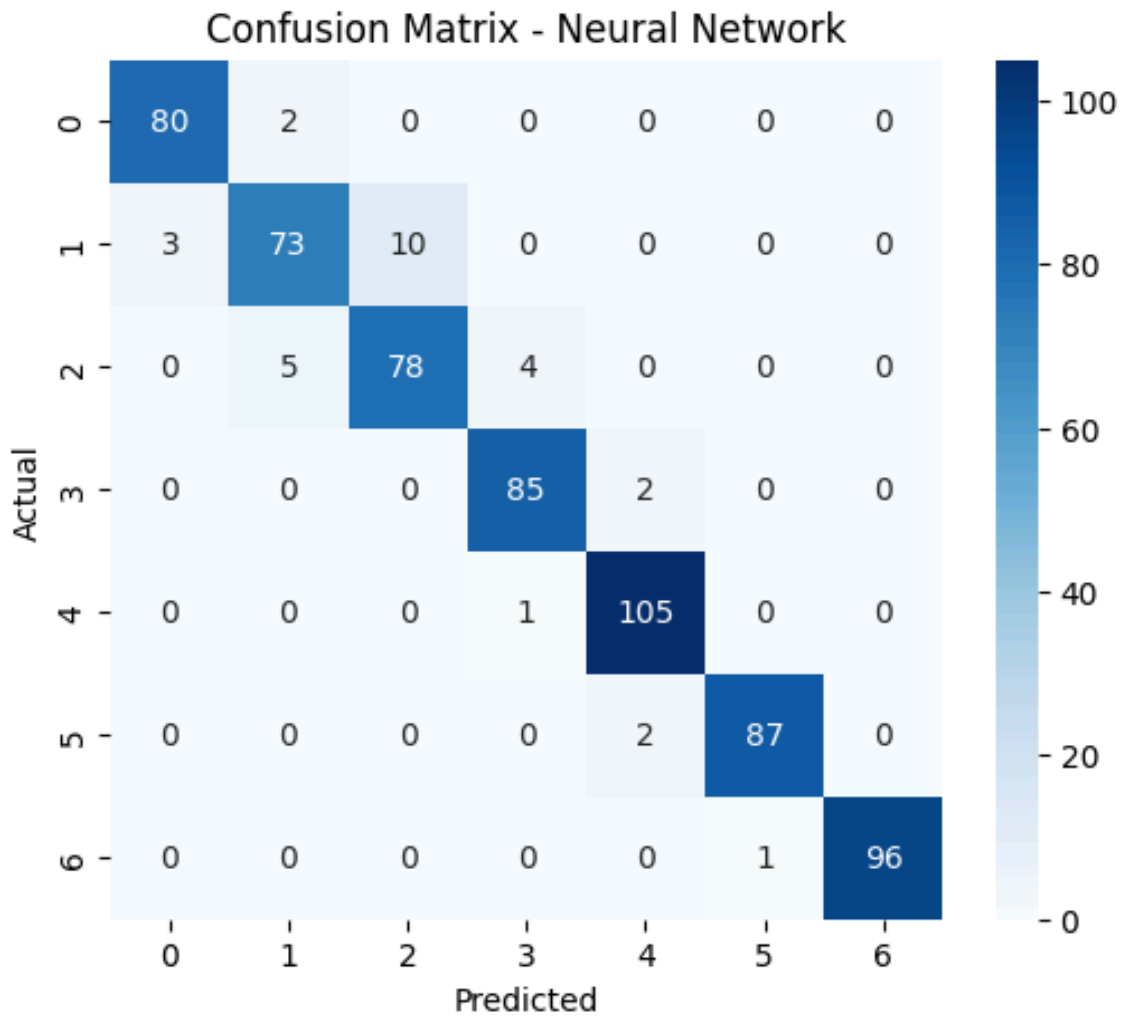- **Naive Bayes** suffered from low recall, especially for minority classes, misclassifying many instances.
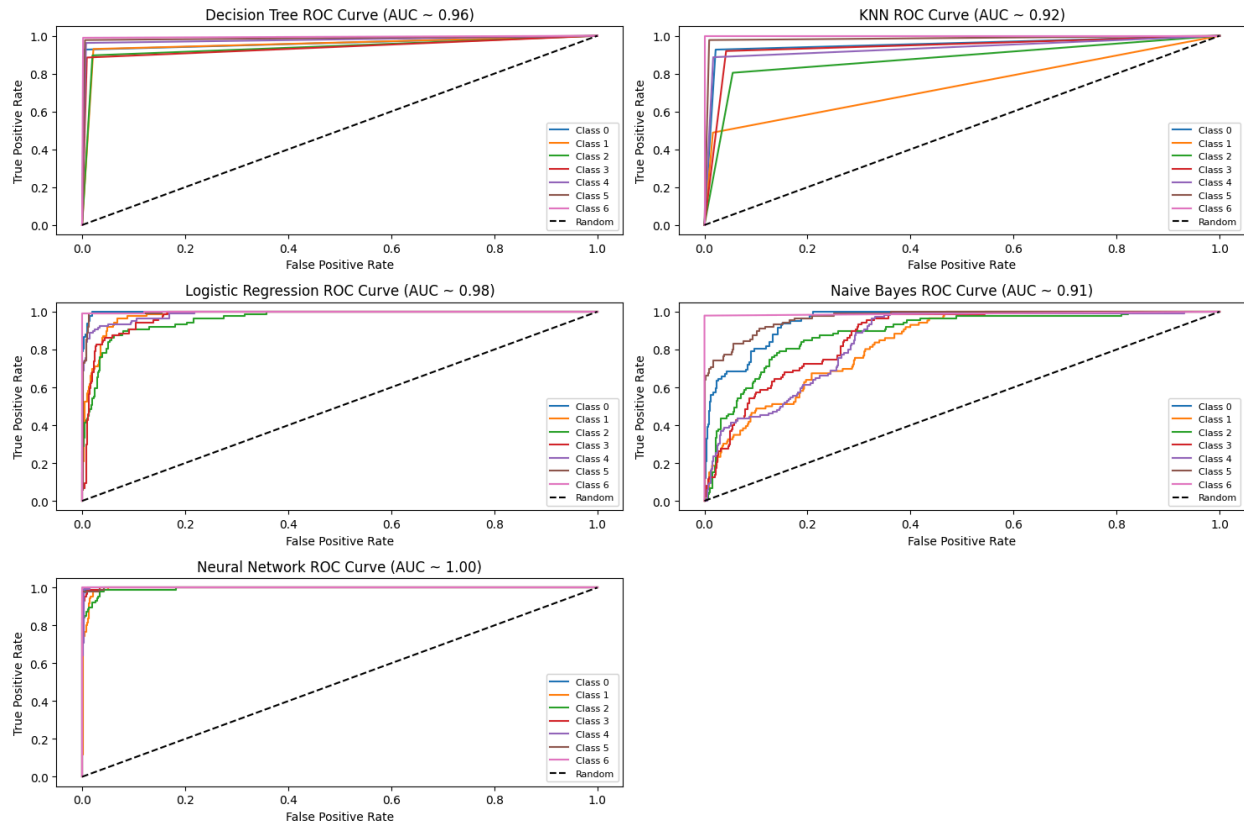
# Confusion Matrix:



Confusion Matrix - KNN

Confusion Matrix - Logistic Regression

Confusion Matrix - Decision Tree

Confusion Matrix - Naive Bayes

## Confusion Matrix - Neural Network

|  | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| **0** | 80 | 2 | 0 | 0 | 0 | 0 | 0 |
| **1** | 3 | 73 | 10 | 0 | 0 | 0 | 0 |
| **2** | 0 | 5 | 78 | 4 | 0 | 0 | 0 |
| **3** | 0 | 0 | 0 | 85 | 2 | 0 | 0 |
| **4** | 0 | 0 | 0 | 1 | 105 | 0 | 0 |
| **5** | 0 | 0 | 0 | 0 | 2 | 87 | 0 |
| **6** | 0 | 0 | 0 | 0 | 0 | 1 | 96 |

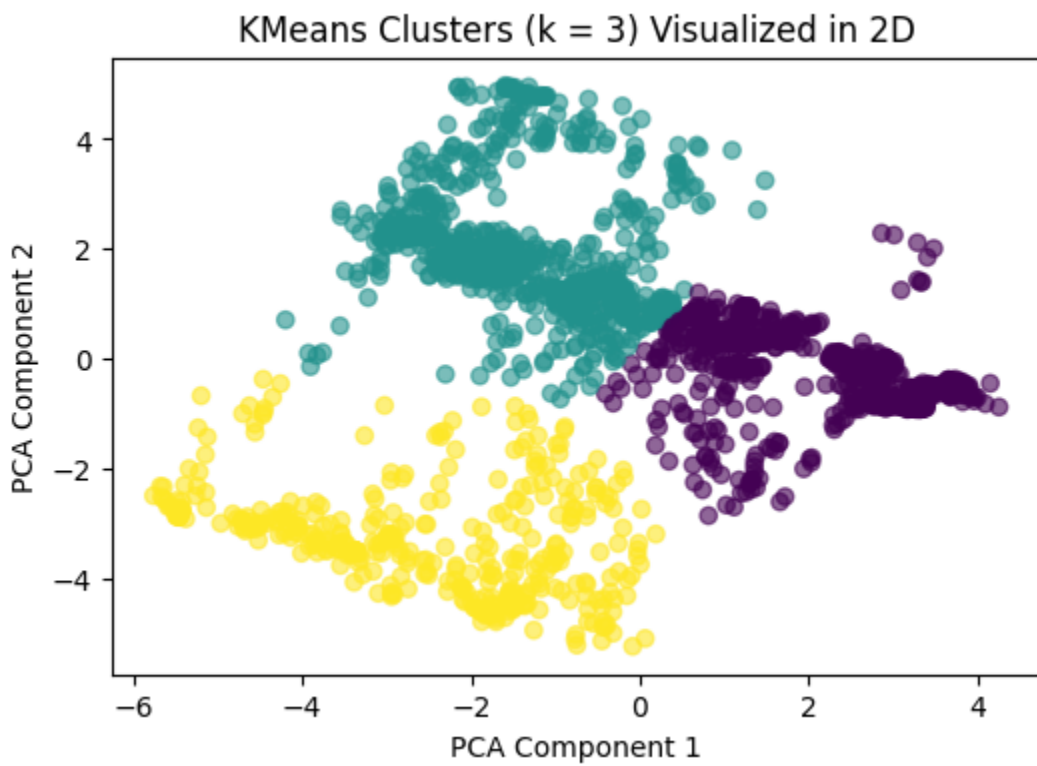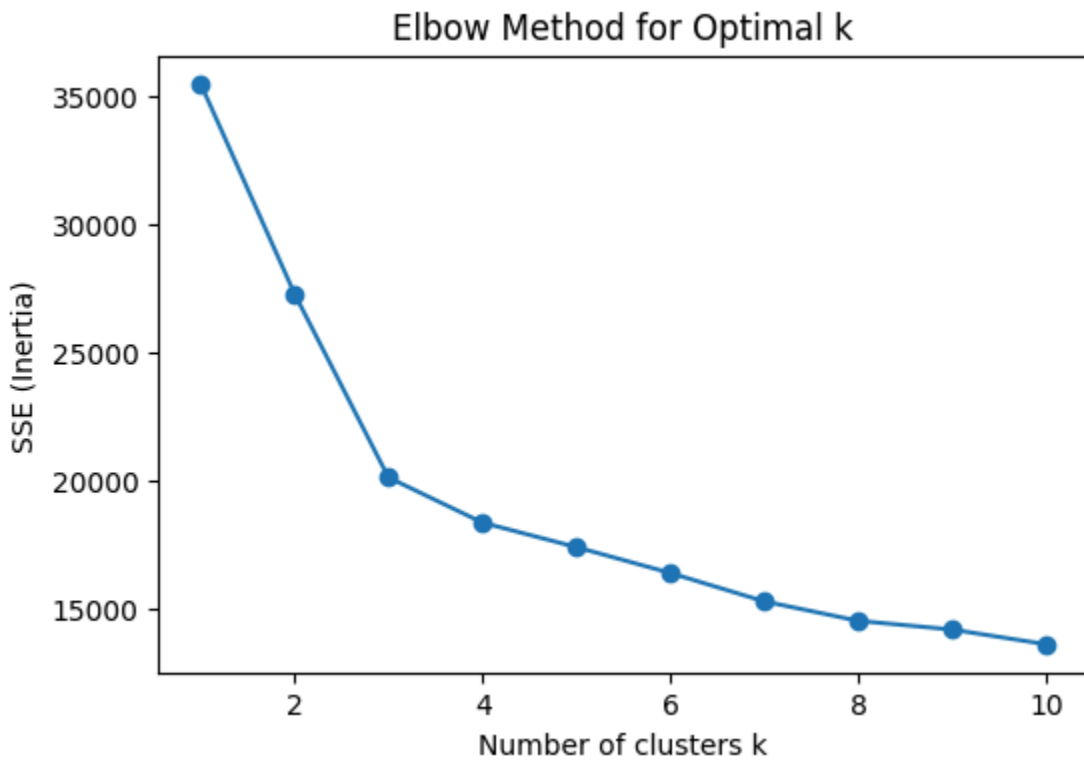Actual (rows) / Predicted (columns)

- Neural Network confusion matrix showed the least misclassification, with better class balance.
- KNN misclassified some minority classes into nearby categories (e.g., "Obesity I" confused with "Overweight").
- Decision Tree overfitted some classes, leading to many correct predictions but also severe misclassifications in other classes.

**ROC & AUC Score:**



- Neural Network had the highest ROC-AUC, showing it is the most discriminative model.
- Logistic Regression also had stable ROC curves.
- Naive Bayes produced weaker ROC curves, reflecting lower predictive power.

**KMeans (Unsupervised):**



Elbow Method for Optimal k



KMeans Clusters (k = 3) Visualized in 2D

We also applied KMeans clustering to treat the problem in an unsupervised setting.

- **Evaluation:** Silhouette score used to measure cluster quality.
- **Observation:** Some clusters partially aligned with obesity categories, but separation was not as strong as in supervised methods.

# **Conclusion**

From the experiments, the following insights were obtained:

- **Best Model:** The Neural Network (MLPClassifier) consistently outperformed other models, showing robustness across all evaluation metrics.
- **Reliable Alternatives:** KNN and Logistic Regression also gave good performance and can serve as interpretable baseline models.
- **Overfitting Risks:** Decision Tree performed well on training data but overfitted without pruning.
- **Weak Performer:** Naive Bayes was not well-suited due to correlated features and class imbalance.
- **Unsupervised Learning:** KMeans clustering confirmed some grouping patterns but was less effective compared to supervised methods.

**Why these results?**

- Neural Networks capture complex, non-linear relationships and adapt well with proper scaling and early stopping.
- KNN benefited from feature scaling but struggled with imbalanced classes.
- Logistic Regression is limited by linearity assumptions.
- Naive Bayes underperformed because independence assumptions don't hold in real-world obesity features.
- Decision Tree overfitted due to unlimited growth without maximum depth or pruning, capturing noise in the dataset.

**Challenges faced:**

- Handling imbalanced class distribution across 7 obesity categories.
- Choosing correct scalers for different feature distributions.
- Avoiding overfitting in models like Decision Trees and Neural Networks.
- Computational cost of Neural Network training.