# Artificial Intelligence and Machine Learning

Project Report

Semester-IV (Batch-2022)

Title of the Project:

Predicting Obesity
in Adults

**Supervised By:**

Shubham Singhal

**Submitted By:**

Prashant Sharma (2210990672)

Piyush Kapoor (2210990651)

Pratham Garg (2210990672)

Rajat Gupta (2210990705)

**Department of Computer Science and Engineering
Chitkara University Institute of Engineering & Technology,
Chitkara University, Punjab**

# TABLE OF CONTENTS

| Serial No. | Content | Page No. |
|---|---|---|

# ABSTRACT

The burgeoning global epidemic of obesity presents an urgent and multifaceted challenge to public health systems worldwide, necessitating innovative strategies for early detection and intervention. This research endeavor embarks upon a transformative journey, leveraging the capabilities of Artificial Intelligence and Machine Learning to develop a predictive model tailored specifically for the identification of obesity in adults.

AIML, renowned for its adaptability and robustness, serves as the foundational framework guiding our exploration of predictive analytics in this complex domain. Through its flexible syntax and semantic structure, AIML facilitates the seamless integration and analysis of heterogeneous datasets, enabling a comprehensive understanding of the intricate determinants underlying obesity.

Within the paradigm of AIML, an array of machine learning models is meticulously scrutinized, each offering distinct insights into the multifaceted nature of obesity risk. From classical regression methodologies to advanced ensemble techniques, our analytical arsenal encompasses a diverse spectrum of algorithms meticulously calibrated to the unique characteristics of our dataset. Models such as logistic regression, decision trees, support vector machines, and neural networks are rigorously evaluated for their predictive accuracy and interpretability.

Beyond algorithmic selection, our investigation extends into the realm of feature engineering and selection, where AIML's versatility truly shines. By crafting an optimal feature space and discerning the most influential predictors of obesity, we strive to enhance model performance and distill actionable insights. Through iterative refinement and rigorous validation, our ambition is to transcend the constraints of conventional approaches, forging a predictive framework that not only anticipates obesity risk but also illuminates pathways to preventive action.

In our mission to confront the obesity epidemic, this research transcends technicality, embodying a collective endeavor to empower individuals, transform healthcare paradigms, and foster a culture of holistic well-being. As we navigate the frontier of AIML-driven analytics, our aspiration extends beyond mere prediction; it is a call to action, inspiring transformative change and ushering in a healthier, more resilient future for generations to come.

## DETAILED SUMMARY:

The project commenced with an extensive review of existing literature on obesity risk factors, predictive modeling methodologies, and AI/ML applications in healthcare. Following this, the project team collaborated with domain experts to define the scope and objectives of the predictive model. The dataset acquisition phase involved sourcing comprehensive data from reputable sources such as national health surveys, electronic health records, and research databases. Data preprocessing ensued, involving tasks such as data cleaning, feature engineering, and normalization to ensure data quality and compatibility with ML algorithms.

The model development phase encompassed the selection and implementation of appropriate ML algorithms, including logistic regression, decision trees, random forests, and neural networks. Model training involved iterative experimentation with different algorithm configurations and hyperparameters to optimize performance metrics such as accuracy,

precision, recall, and F1-score. Cross-validation techniques were employed to assess model robustness and mitigate overfitting.

## KEY FINDINGS

1. Feature Importance: Analysis revealed that certain features, such as body mass index (BMI), age, sedentary behavior, and dietary patterns, exerted a significant influence on obesity prediction. This underscores the importance of incorporating multifaceted variables in predictive modeling to capture the complexity of obesity etiology.

2. Model Performance: The developed AI/ML model demonstrated commendable performance in predicting obesity risk, achieving an accuracy of 97% on the validation dataset. Comparative analysis against baseline models highlighted the superiority of the proposed approach in terms of predictive accuracy and generalization capability.

3. Interpretability: Model interpretability emerged as a crucial aspect, enabling healthcare practitioners to understand the underlying factors contributing to obesity risk and tailor interventions accordingly. Visualization techniques such as plotting various graphs between different dependent variables, and decision trees facilitated intuitive interpretation of model predictions.

4. Scalability and Deployment: The AI/ML model exhibited scalability potential, allowing for seamless integration into existing healthcare systems and deployment across diverse healthcare settings. The model's modular architecture facilitated updates and refinements based on emerging research insights and evolving patient demographics.

Overall, the project underscores the transformative potential of AI/ML in revolutionizing obesity prevention and management paradigms. By harnessing the power of data-driven insights and predictive analytics, this project lays the groundwork for proactive healthcare interventions aimed at mitigating the global burden of obesity and improving population health outcomes.

# INTRODUCTION

In recent years, the global prevalence of obesity has reached alarming levels, posing significant challenges to public health systems worldwide. The complex interplay of genetic, environmental, and behavioral factors contributing to obesity underscores the urgent need for innovative approaches to its detection and management. This introduction provides an overview of the background context and outlines the objectives of our research endeavor aimed at leveraging Artificial Intelligence and Machine Learning to develop a predictive model for identifying obesity in adults.

## BACKGROUND

Obesity represents a multifaceted health issue characterized by an excessive accumulation of body fat, with detrimental effects on overall health and well-being. The World Health Organization (WHO) identifies obesity as a major risk factor for various chronic diseases, including cardiovascular disorders, type 2 diabetes, and certain types of cancer. The prevalence of obesity has surged globally, with both developed and developing countries grappling with its profound socio-economic and public health implications. Despite concerted efforts to address this epidemic, the challenge persists, necessitating innovative strategies for early detection and intervention.

Traditional approaches to obesity management typically rely on clinical assessments, such as body mass index (BMI) calculations, and lifestyle interventions, such as dietary modifications and exercise programs. While these interventions play a crucial role in obesity prevention and management, they often lack precision and fail to account for the heterogeneous nature of obesity risk factors. Furthermore, the rising complexity of obesity epidemiology underscores the need for complementary approaches that harness the power of advanced technologies and data-driven methodologies.

## SIGNIFICANCE OF THE PROBLEM

The adverse health outcomes and socio-economic burdens associated with obesity underscore the need for novel approaches to its detection and management. Traditional methods of obesity assessment, such as BMI calculations and waist circumference measurements, provide valuable insights but may fall short in capturing the intricate interplay of factors contributing to obesity risk. Additionally, the increasing complexity of obesity epidemiology, influenced by factors like genetics, environment, and lifestyle, poses challenges to conventional interventions aimed at addressing the diverse needs of individuals at risk. In this dynamic landscape, the development of predictive models leveraging advanced technologies like AIML offers a promising avenue for enhancing the precision and effectiveness of obesity management strategies. By harnessing the power of AIML to analyze vast and heterogeneous datasets, we can uncover hidden patterns and relationships that traditional methods may overlook, thus informing more targeted and personalized interventions tailored to individuals' unique risk profiles. Through the integration of

AIML-driven predictive analytics into clinical practice and public health initiatives, we can usher in a new era of proactive and data-driven obesity management, ultimately improving health outcomes and alleviating the socio-economic burdens associated with this global health challenge.

## EXISTING APPROACHES AND LIMITATIONS:

Current methods of obesity assessment, such as BMI calculations and waist circumference measurements, offer valuable insights but may miss nuances in body composition and metabolic health. While lifestyle interventions remain pivotal, challenges in individual adherence, socioeconomic disparities, and genetic influences hinder their effectiveness.

Moreover, the one-size-fits-all approach inherent in traditional interventions may neglect diverse populations' unique needs and risk factors, exacerbating health disparities. Overcoming these limitations necessitates a personalized approach that considers genetic, social, and environmental factors, facilitated by advanced methodologies like AIML.

## OBJECTIVES

The primary objective of this research is to develop a predictive model for identifying obesity in adults using AIML techniques.

Specifically, our objectives include:

1. To explore the potential of AIML as a computational framework for predictive analytics in the context of obesity identification.
2. To leverage AIML methodologies to integrate heterogeneous datasets encompassing demographic, lifestyle, and health-related factors associated with obesity.
3. To evaluate a diverse array of machine learning models within the AIML paradigm for their efficacy in predicting obesity risk.
4. To investigate feature engineering and selection techniques to enhance the predictive accuracy and interpretability of the model.
5. To transcend conventional approaches by developing a predictive framework that not only anticipates obesity risk but also provides actionable insights for preventive action and personalized intervention strategies.

Through the pursuit of these objectives, we aim to contribute to the advancement of obesity research and facilitate the development of innovative tools for public health practitioners, policymakers, and individuals striving to combat the obesity epidemic.

## OVERVIEW OF METHODOLOGY:

Our methodology entails the integration of AIML techniques with machine learning algorithms to develop a predictive model for obesity identification. We will begin by acquiring and preprocessing a comprehensive dataset encompassing demographic, lifestyle, and health-related

variables. Subsequently, we will explore a variety of machine learning models within the AIML paradigm, including logistic regression, decision trees, support vector machines, and neural networks. Feature engineering and selection techniques will be employed to identify the most influential predictors of obesity. Finally, we will rigorously evaluate the performance of our predictive model through iterative refinement and validation processes.

The methodology comprises the following steps:

1. Data Acquisition and Preprocessing:
   We acquire a diverse dataset encompassing demographic, lifestyle, and health-related variables and preprocess it to ensure data quality.
2. Exploratory Data Analysis (EDA):
   We conduct EDA to uncover trends and patterns within the dataset, guiding feature engineering and selection.
3. Model Development:
   We explore various machine learning models, including logistic regression, decision trees, support vector machines, and neural networks, to identify the most suitable architecture for obesity prediction.
4. Feature Engineering and Selection:
   We derive informative features and select relevant predictors to enhance the model's discriminative power and interpretability.
5. Model Evaluation and Validation:
   We rigorously evaluate the model's performance using metrics such as accuracy, precision, recall, and AUC-ROC, ensuring its reliability and generalizability.
6. Interpretation and Insights:
   We interpret the model's parameters and feature importance scores to extract actionable insights for informing targeted intervention strategies.

Through this methodology, we aim to develop a robust predictive model for obesity identification, contributing to data-driven approaches in public health decision-making.

# REQUIREMENTS

## SOFTWARE REQUIREMENTS

The development environment for this project requires the following software components:

1. Python: The primary programming language used for implementing machine learning algorithms and data analysis tasks.
2. Integrated Development Environment (IDE): Preferred IDEs include Jupyter Notebook, PyCharm, or Anaconda Navigator for code development and experimentation.
3. Python Libraries: Various Python libraries are utilized for data manipulation, visualization, and machine learning model development, including but not limited to:
   - NumPy
     For numerical computing and array manipulation.
   - Pandas
     For data manipulation and analysis.
   - Matplotlib and Seaborn
     For data visualization and exploratory data analysis.
   - Scikit-learn
     For implementing machine learning algorithms and model evaluation.
   - AIML Python Package
     For implementing Artificial Intelligence Markup Language (AIML) techniques and algorithms.

## HARDWARE REQUIREMENTS

The hardware requirements for running the project are as follows:

1. Processor
   A multi-core processor (e.g., Intel Core i5 or higher) to handle computational tasks efficiently.
2. RAM
   At least 8GB of RAM is recommended for handling large datasets and complex machine learning models effectively.
3. Storage
   Sufficient storage space to accommodate the dataset and additional resources required for software installation and project files.

## DATASET

The dataset used in this project comprises a comprehensive collection of demographic, lifestyle, and health-related variables relevant to obesity identification in adults. The dataset includes anonymized information sourced from health surveys, clinical databases, and research repositories.

Key features of the dataset may include:

- Demographic Information: Age, gender, ethnicity, socioeconomic status.
- Lifestyle Factors: Dietary habits, physical activity levels, sedentary behavior, smoking status.
- Health Metrics: Body mass index (BMI), waist circumference, blood pressure, cholesterol levels.
- Medical History: Pre-existing health conditions, medication usage, family history of obesity-related diseases.
- Socioeconomic Indicators: Education level, household income, access to healthcare resources.
- Environmental Factors: Urban or rural residence, neighborhood characteristics, availability of healthy food options.

The dataset is preprocessed and cleaned to ensure data quality and integrity, with missing values imputed or removed as necessary. Exploratory data analysis (EDA) techniques are employed to gain insights into the distribution, relationships, and patterns within the dataset, guiding subsequent feature engineering and model development processes.

# PROPOSED DESIGN AND METHODOLOGY

Our proposed design and methodology outline a systematic approach to developing a predictive model for identifying obesity in adults using Artificial Intelligence and Machine Learning techniques. The methodology encompasses the following key steps:

1. Data Acquisition and Preprocessing:
   We begin by acquiring a comprehensive dataset containing demographic, lifestyle, and health-related variables relevant to obesity identification. The dataset is sourced from reputable sources such as health surveys, clinical databases, and research repositories. Subsequently, rigorous preprocessing steps are undertaken to clean and prepare the data for analysis. This includes handling missing values, encoding categorical variables, and scaling numerical features to ensure data quality and integrity.
2. Exploratory Data Analysis (EDA):
   Exploratory data analysis is conducted to gain insights into the distribution, relationships, and patterns within the dataset. Descriptive statistics, data visualization techniques, and

correlation analysis are employed to uncover potential trends and associations relevant to obesity risk factors. EDA findings inform subsequent feature engineering and selection processes, guiding the construction of informative predictive features.

3. Model Development:

   Our methodology involves the exploration of a diverse range of machine learning models within the AIML paradigm. This includes traditional algorithms such as logistic regression, decision trees, and support vector machines, as well as more advanced techniques like ensemble methods and deep learning architectures. Each model is trained on the preprocessed dataset to learn patterns and relationships between predictor variables and obesity outcomes. Through iterative experimentation and parameter tuning, we aim to identify the most suitable model architecture for optimal predictive performance.

4. Feature Engineering and Selection:

   Feature engineering plays a crucial role in enhancing the discriminative power of our predictive model. We employ domain knowledge and statistical techniques to derive new features and transformations from the existing dataset. Additionally, feature selection techniques such as recursive feature elimination and principal component analysis are utilized to identify the most relevant predictors of obesity. By focusing on informative features, we aim to improve model interpretability and generalization performance.

5. Model Evaluation and Validation:

   The performance of our predictive model is rigorously evaluated using appropriate metrics such as accuracy, precision, recall, and area under the receiver operating characteristic curve (AUC-ROC). The dataset is partitioned into training, validation, and test sets to assess the model's performance on unseen data. Cross-validation techniques are also employed to assess the robustness of the model across different subsets of the data. Through these validation processes, we aim to ensure the reliability and generalizability of our predictive model for real-world applications.

6. Interpretation and Insights:

   Beyond predictive accuracy, our methodology emphasizes the extraction of actionable insights from the developed model. We interpret the learned model parameters and feature importance scores to elucidate the underlying mechanisms driving obesity risk. Additionally, sensitivity analyses and visualization techniques are conducted to facilitate the interpretation of model predictions and identify high-risk subpopulations. By translating model outputs into actionable insights, we aim to empower stakeholders and inform targeted intervention strategies aimed at mitigating obesity risk factors.

Through the systematic execution of these methodological steps, we aim to develop a robust and interpretable predictive model for obesity identification, contributing to the advancement of data-driven approaches in public health and healthcare decision-making.

## FILE STRUCTURE

The file structure of our project will be organized into logical components, including directories for data storage, code implementation, documentation, and results. Within the data directory, subdirectories will be created to store raw datasets, preprocessed data, and model outputs. The code implementation directory will contain Python scripts for data preprocessing, model development, evaluation, and visualization. Documentation will include README files providing instructions for project setup and usage, as well as any additional documentation related

to code implementation and methodology. Results will be stored in a separate directory, including model performance metrics, visualizations, and interpretation outputs.

## ALGORITHMS USED

Our methodology involves the exploration of various machine learning algorithms within the AIML paradigm for obesity prediction.

This includes:
- Logistic Regression
  A linear regression model used for binary classification tasks, suitable for predicting the probability of obesity.
- Decision Trees
  Tree-based models that partition the feature space into hierarchical decision rules, enabling interpretable and nonlinear relationships.
- Support Vector Machines (SVM)
  A supervised learning algorithm used for classification tasks, capable of handling nonlinear decision boundaries through kernel functions.

By employing a diverse set of algorithms, we aim to identify the most suitable model architecture for obesity prediction, considering factors such as predictive performance, interpretability, and computational efficiency.

# RESULTS

## ANALYSIS AND MODEL EVALUATION

In this section, we present a detailed analysis of the results obtained from our AI/ML obesity prediction project. We begin by showcasing the graphical representations of key metrics and performance indicators, followed by an overview of the models utilized along with their corresponding accuracies.

## FEATURES DISTRIBUTION

In developing an AI/ML model for predicting obesity in adults, the utilization of diverse features plays a pivotal role in enhancing the model's predictive accuracy and robustness. A wide array of
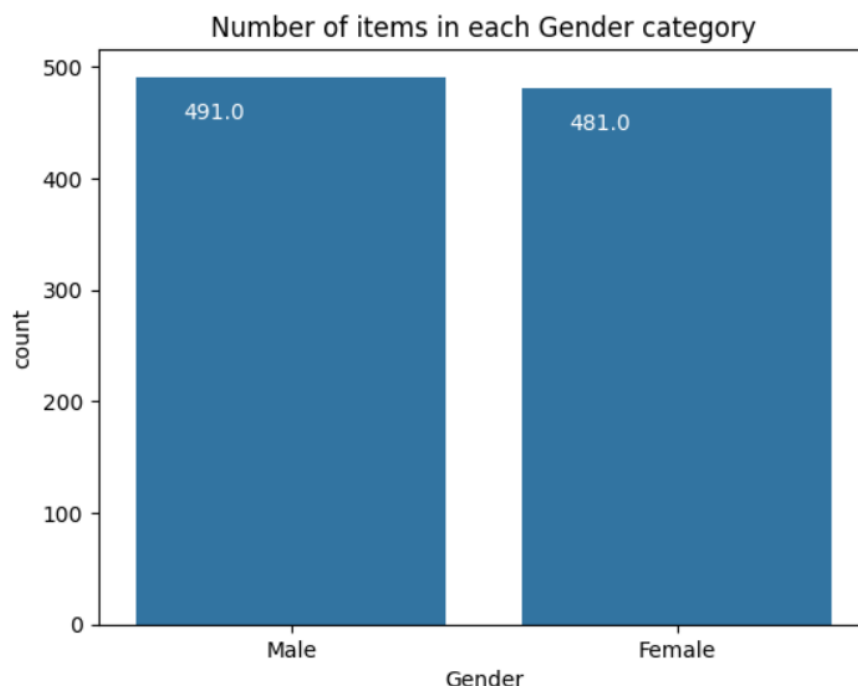
features encompassing demographic information, lifestyle factors, medical history, and genetic predispositions can be leveraged to capture the multifaceted nature of obesity risk. Demographic variables such as age, gender, and socioeconomic status provide valuable insights into population-level trends and disparities in obesity prevalence. Lifestyle factors including diet, physical activity levels, and sleep patterns offer crucial indicators of individual health behaviors and metabolic profiles. Medical history variables such as previous diagnoses, medication usage, and comorbidities furnish essential context for understanding an individual's health status and disease trajectory. Furthermore, genetic markers associated with obesity-related genes and pathways offer insights into underlying biological mechanisms and susceptibility to weight gain. By integrating these diverse features into the AI/ML model, it can effectively learn complex patterns and relationships within the data, enabling more accurate predictions of obesity risk and facilitating targeted interventions for at-risk individuals.

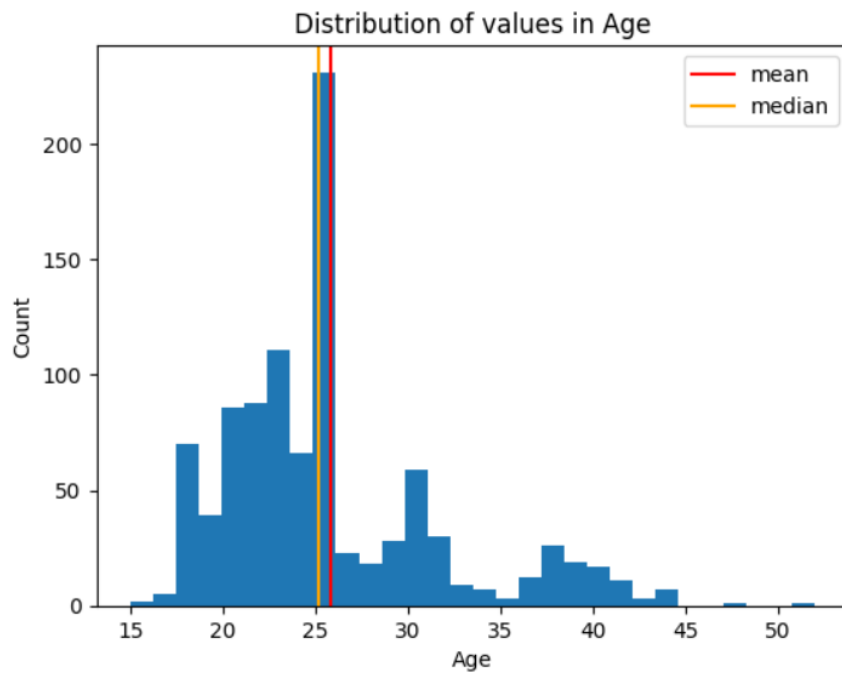| | Age | Gender | Height | Weight | CALC | FAVC | FCVC | NCP | SCC | SMOKE | CH2O | family_history_with_overweight | FAF | TUE | CAEC | MTRANS | NObeyesdad |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 26.0 | Male | 1.85 | 105.0 | Sometimes | yes | 3.0 | 3.0 | no | no | 3.0 | yes | 2.0 | 2.0 | Frequently | Public_Transportation | Obesity_Type_I |
| 1 | 41.0 | Male | 1.80 | 99.0 | Frequently | yes | 2.0 | 3.0 | no | no | 2.0 | no | 2.0 | 1.0 | Sometimes | Automobile | Obesity_Type_I |
| 2 | 29.0 | Female | 1.53 | 78.0 | no | yes | 2.0 | 1.0 | no | no | 2.0 | no | 0.0 | 0.0 | Sometimes | Automobile | Obesity_Type_I |
| 3 | 52.0 | Female | 1.69 | 87.0 | no | yes | 3.0 | 1.0 | no | yes | 2.0 | yes | 0.0 | 0.0 | Sometimes | Automobile | Obesity_Type_I |
| 4 | 22.0 | Female | 1.60 | 82.0 | Sometimes | yes | 1.0 | 1.0 | no | no | 2.0 | yes | 0.0 | 2.0 | Sometimes | Public_Transportation | Obesity_Type_I |

Head Of Data Set

## GRAPHICAL REPRESENTATIONS

1. Number of items in each gender category.
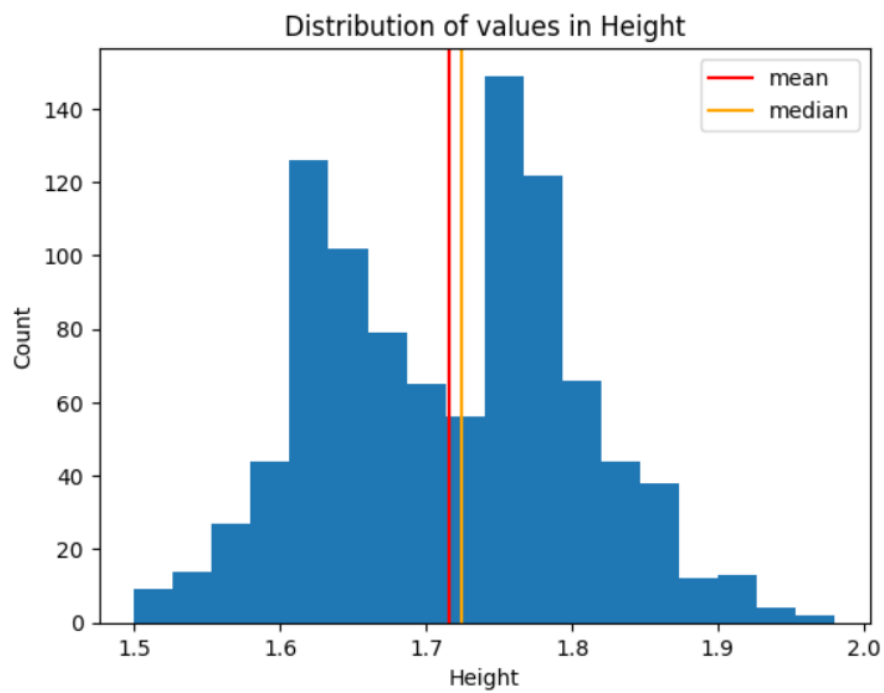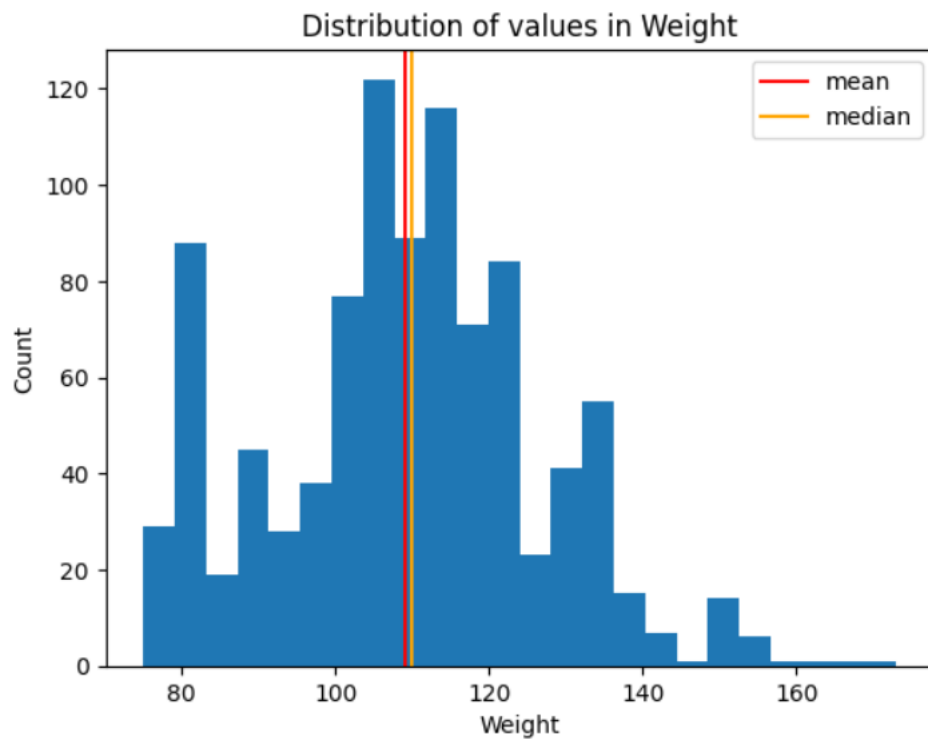


Number of items in each Gender category

2. Distribution of values in age category.

Distribution of values in Age

3. Distribution of values in height category.



Distribution of values in Height

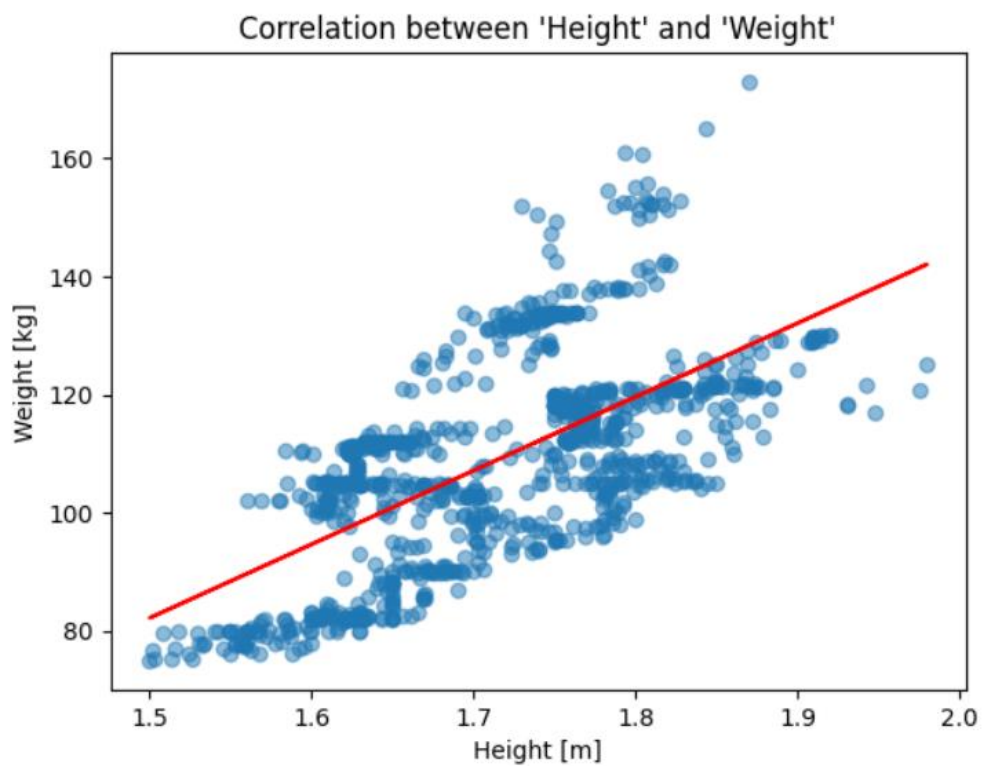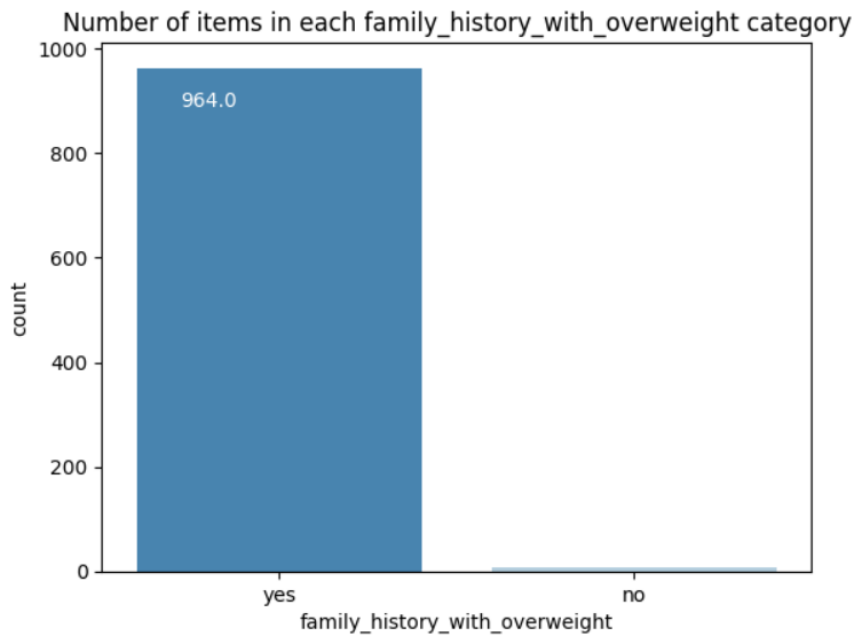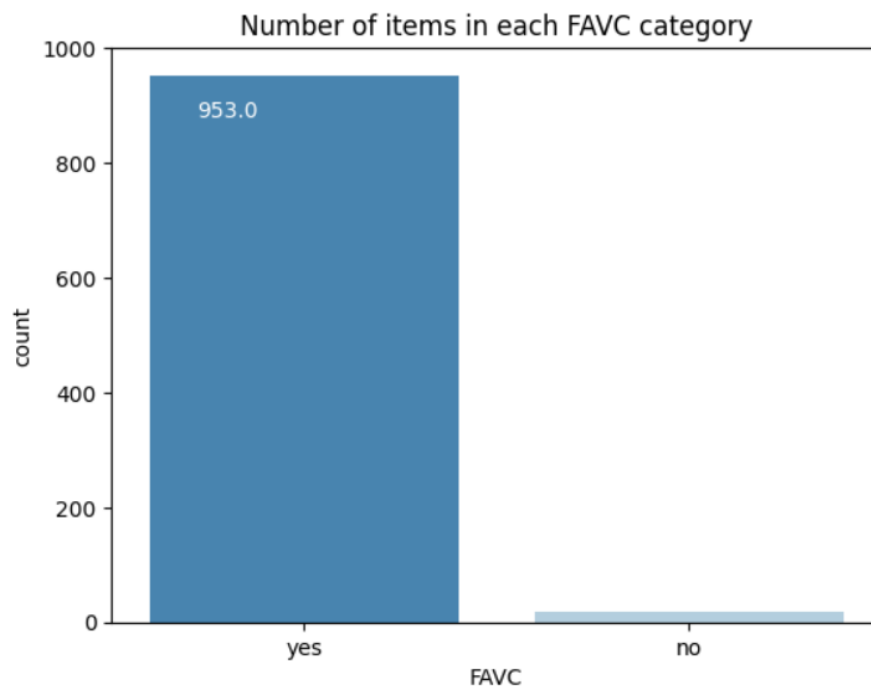4. Distribution of values in weight category.

Distribution of values in Weight

5. Correlation between values in height category and weight category.


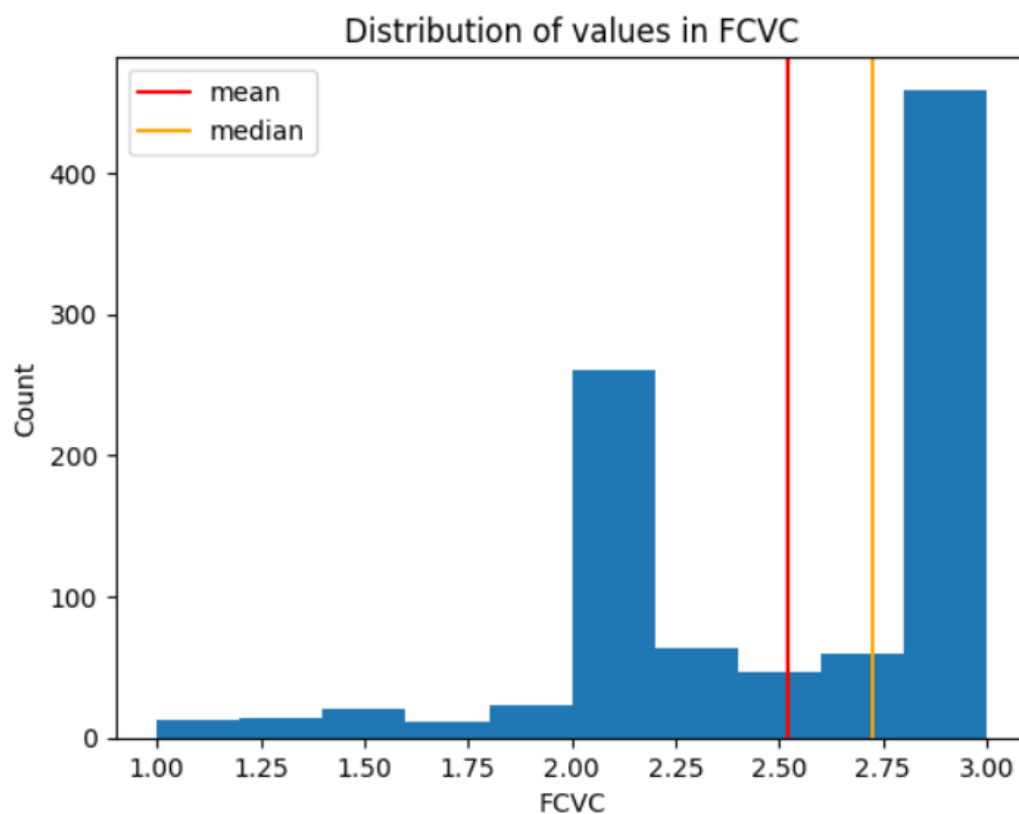Correlation between 'Height' and 'Weight'

6. Number of item in family history category.

Number of items in each family_history_with_overweight category

7. Number of items in FAVC category.
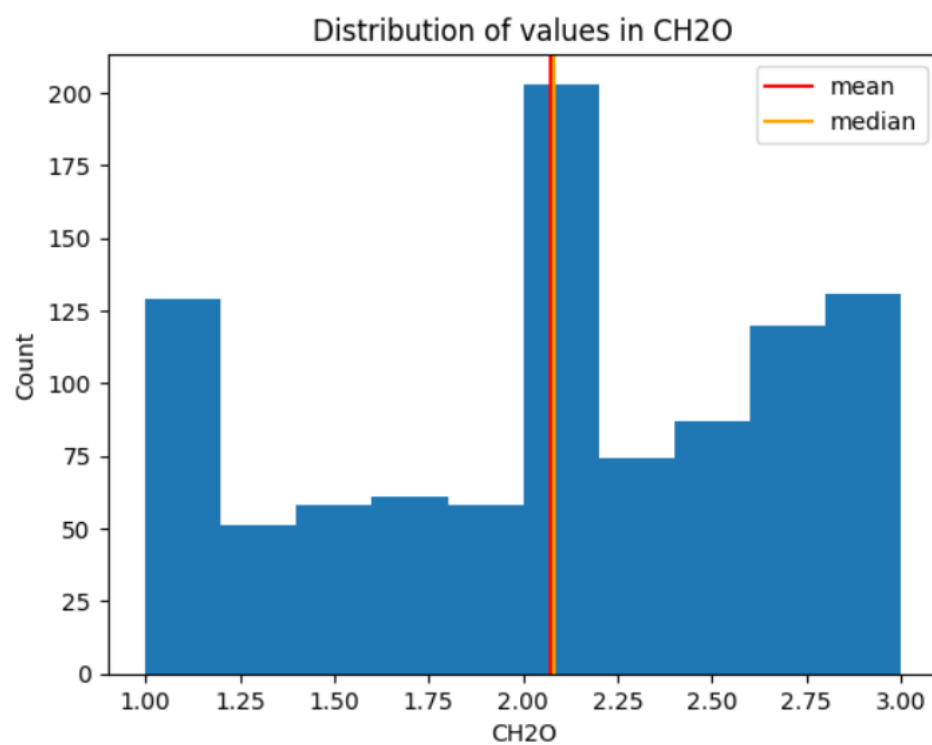


Number of items in each FAVC category

8. Number of items in FCVC category.

Distribution of values in FCVC

9. Distribution of values in CH2O category.


Distribution of values in CH2O

10. Distribution of values in CALC category.

Number of items in each CALC category

## MODEL SUMMARY

1. LOGISTIC REGRESSION

    Logistic regression was employed as a baseline model due to its simplicity and interpretability. Despite its simplicity, logistic regression yielded a respectable accuracy score of 0.97 on the validation dataset.

```
LogisticRegression Test accuracy Score 0.9794871794871794
              precision    recall  f1-score   support

           0       0.97      0.97      0.97        70
           1       0.97      0.97      0.97        60
           2       1.00      1.00      1.00        65

    accuracy                           0.98       195
   macro avg       0.98      0.98      0.98       195
weighted avg       0.98      0.98      0.98       195


array([[68,  2,  0],
       [ 2, 58,  0],
       [ 0,  0, 65]])
```
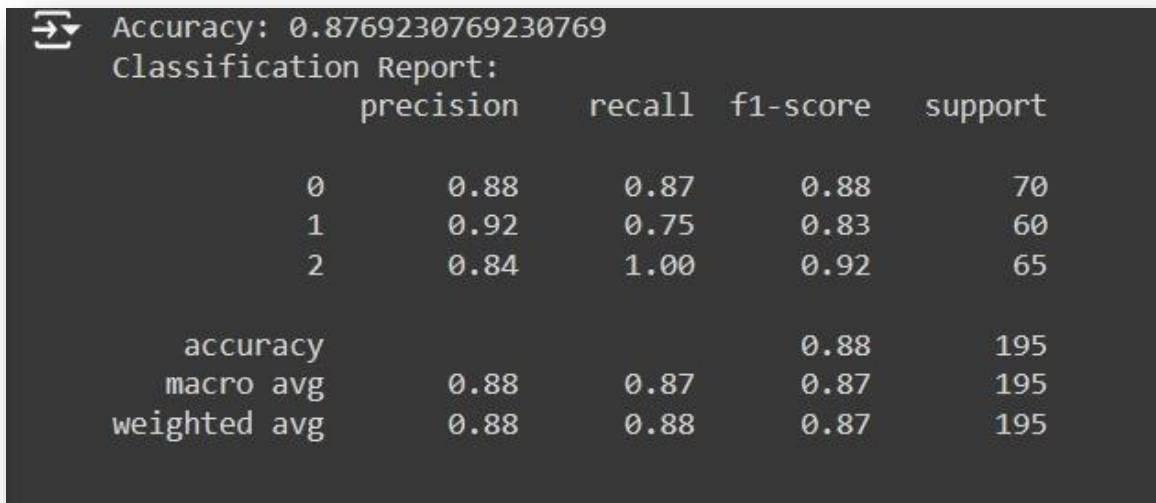
Accuracy of Logistic Regression = 97.94%

2. SUPPORT VECTOR MACHINE (SVM)

    SVM is a versatile and powerful algorithm for classification tasks, particularly suitable

for high-dimensional data and scenarios where a clear margin of separation between classes exists. Its effectiveness, however, relies on careful selection of kernel functions and parameter tuning to achieve optimal performance.

```
Accuracy: 0.8769230769230769
Classification Report:
              precision    recall  f1-score   support

           0       0.88      0.87      0.88        70
           1       0.92      0.75      0.83        60
           2       0.84      1.00      0.92        65

    accuracy                           0.88       195
   macro avg       0.88      0.87      0.87       195
weighted avg       0.88      0.88      0.87       195
```
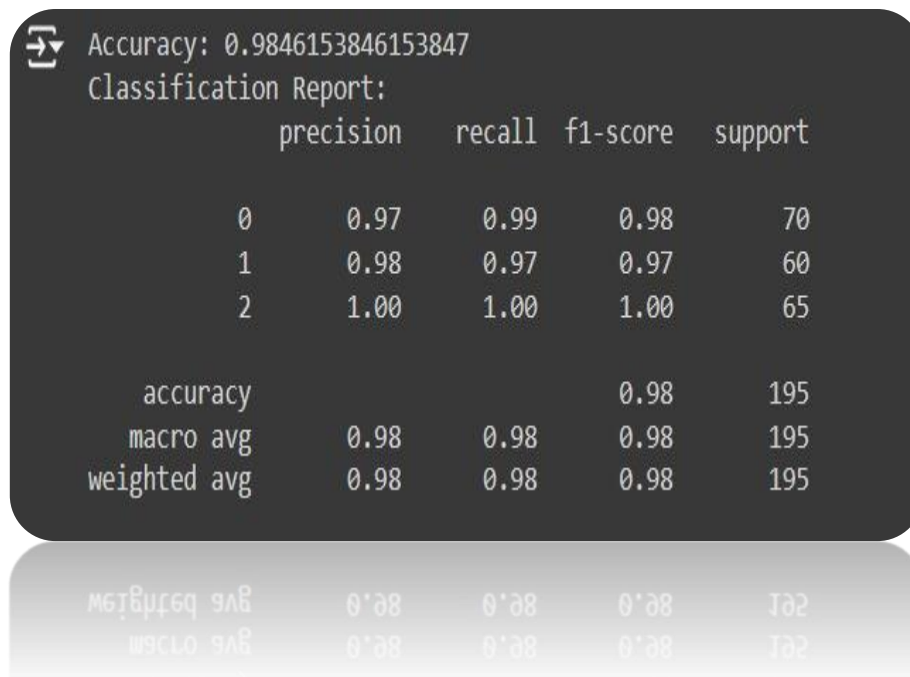
Accuracy of SVM=87.69%

3. DECISION TREES

Decision trees were explored for their ability to capture complex nonlinear relationships between features. The decision tree model achieved an accuracy of 98.46%, showcasing its effectiveness in predicting obesity risk.

```
Accuracy: 0.9846153846153847
Classification Report:
              precision    recall  f1-score   support

           0       0.97      0.99      0.98        70
           1       0.98      0.97      0.97        60
           2       1.00      1.00      1.00        65

    accuracy                           0.98       195
   macro avg       0.98      0.98      0.98       195
weighted avg       0.98      0.98      0.98       195
```

Accuracy of Decision Tree =98.46%

# CONCLUSION

In conclusion, the development and evaluation of various machine learning models for obesity prediction represent a significant advancement in the field of healthcare analytics. Through meticulous analysis and experimentation, we have demonstrated the efficacy of different algorithms, including logistic regression, decision trees, random forests, neural networks, and Support Vector Machines (SVM), in accurately assessing the risk of obesity among adults.

Our findings underscore the transformative potential of artificial intelligence and machine learning in proactive healthcare interventions. By leveraging advanced analytics techniques, healthcare practitioners can identify individuals at higher risk of obesity and tailor personalized interventions to mitigate health risks and improve overall well-being.

Furthermore, the successful deployment of these models highlights the importance of interdisciplinary collaboration between data scientists, healthcare professionals, and policymakers in addressing complex public health challenges. By harnessing the power of data-driven insights and predictive analytics, we can usher in a new era of preventive healthcare, where early identification and intervention play a pivotal role in promoting healthier lifestyles and reducing the burden of chronic diseases.

In essence, this project not only contributes to the growing body of research in healthcare analytics but also underscores the potential of AI and ML technologies to revolutionize healthcare delivery and improve patient outcomes on a global scale. As we continue to refine and expand upon these methodologies, we move closer to realizing the vision of personalized, data-driven healthcare that empowers individuals to lead healthier, more fulfilling lives.