# Obesity Level Classification Using Machine Learning

Rishita Garg
Dept. of I&CT
MIT,Manipal
rishita.mitmpl2022
@learner.manipal.edu

Harshit Dugar
Dept. of I&CT
MIT,Manipal
harshit3.mitmpl2022
@learner.manipal.edu

Shouvik Kumar
Dept. of I&CT
MIT,Manipal
shouvik.mitmpl2022
@learner.manipal.edu

*Abstract*—Obesity is a major public health issue that significantly contributes to the development of chronic diseases such as diabetes, cardiovascular ailments, and several types of cancer. In this study, machine learning techniques are employed to predict obesity levels using lifestyle and physiological data. A dataset comprising 2,111 records with attributes such as age, height, weight, physical activity, and technology usage was utilized. Multiple models, including XGBoost, Random Forest, and Artificial Neural Networks (ANN), were implemented following comprehensive data preprocessing that involved feature scaling and class balancing via SMOTE. Interpretability of predictions was achieved using SHAP (SHapley Additive exPlanations) values, identifying weight, height, and physical activity as key influencers. The XGBoost model exhibited superior performance, achieving an accuracy of 98.1%. Furthermore, a GUI was developed using Tkinter to enable user interaction. This project demonstrates the potential of artificial intelligence in enhancing early detection of obesity and delivering actionable health insights.

*Index Terms*—Obesity, Machine Learning, SHAP, SMOTE, XGBoost, Random Forest, Artificial Neural Networks, BMI, GUI

## I. INTRODUCTION

Obesity has become a global epidemic and a key risk factor for numerous chronic health conditions. As per the World Health Organization (WHO), over 1 billion individuals worldwide are classified as obese, a number projected to increase steadily. The increasing prevalence of sedentary lifestyles, reliance on technology, and unhealthy dietary patterns have exacerbated this issue. In younger populations especially, these factors contribute to an alarming rate of obesity. Given the multidimensional nature of obesity, early intervention through predictive analysis becomes essential. Recent advancements in artificial intelligence (AI) and machine learning (ML) offer promising capabilities in extracting meaningful patterns from complex datasets. This project seeks to harness the predictive capabilities of machine learning models for obesity level classification, while maintaining interpretability through SHAP analysis. Furthermore, the project includes the development of a user-friendly interface using Tkinter, enabling real-time predictions based on user input.

## II. OBJECTIVES

- To preprocess and prepare the dataset for machine learning applications.
- To implement and evaluate multiple machine learning models for obesity classification.
- To address class imbalance in the dataset using Synthetic Minority Over-sampling Technique (SMOTE).
- To apply SHAP for model interpretability and insight derivation.
- To develop a Graphical User Interface (GUI) using Tkinter for user interaction and real-time prediction.

## III. LITERATURE REVIEW

### THEME 1: ADDRESSING CLASS IMBALANCE AND DATA AUGMENTATION

Recent research in obesity prediction highlights the critical challenge of class imbalance in BMI datasets. Studies like *"Evaluating Data Augmentation Techniques for Obesity Prediction"* (2023) and *"Addressing Class Imbalance in BMI Classification Tasks"* (2022) demonstrate that synthetic oversampling methods such as SMOTE significantly improve classification performance. These works also explore cost-sensitive learning as a complementary approach to enhance predictive accuracy in imbalanced datasets.

## THEME 2: FEATURE ENGINEERING AND PREPROCESSING

Feature engineering plays a vital role in BMI-related health assessments. *"The Role of Feature Engineering in Machine Learning-Based Health Assessments"* (2021) compares hand-crafted features with features automatically learned by deep learning models, emphasizing the trade-off between interpretability and performance. In parallel, the *"Comparative Study of Data Scaling Methods in Predicting Health Outcomes"* (2023) evaluates various normalization techniques—including min-max scaling, z-score normalization, and log transformation—to optimize prediction outcomes.

## THEME 3: SOCIOECONOMIC AND LIFESTYLE INFLUENCES

Incorporating socioeconomic and demographic variables has shown to enhance both the accuracy and interpretability of BMI prediction models. *"Integrating Socioeconomic Factors into BMI Prediction Models"* (2021) illustrates that socioeconomic indicators such as income, education, and employment status contribute meaningfully to model performance.

## THEME 4: BROADER APPLICATIONS OF BMI RESEARCH

BMI prediction models have been extended to assess the risk of several chronic diseases. For example, *"The Relationship Between BMI and Cardiovascular Disease: A Machine Learning Perspective"* (2022) and *"Predicting Diabetes Risk Using BMI and Lifestyle Factors"* (2021) leverage BMI in combination with other health indicators to evaluate disease risk.

Other studies have explored emerging domains:

- *"Machine Learning Approaches to Examining BMI and Mental Health Correlations"* (2022) uses sentiment analysis to assess psychological impacts.
- *"Obesity Prediction in Adolescents: A Longitudinal Data Analysis"* (2023) applies time-series analysis to identify early risk factors in youth.
- *"Evaluating the Effect of Diet and Exercise on BMI: A Data-Driven Approach"* (2021) employs regression and causal inference techniques to study the effectiveness of lifestyle interventions.

## EVOLUTION OF RESEARCH TRENDS

- **Early Studies (2021–2022):** Focused largely on basic feature engineering and socioeconomic variable integration. These studies laid the groundwork by identifying key predictors of BMI using traditional machine learning techniques.
- **Mid-Decade Advancements (2022–2023):** Marked a shift towards addressing data imbalance and exploring synthetic data augmentation. There was also a thematic expansion into the mental health implications of obesity.
- **Recent Developments (2023–2025):** The focus has shifted toward more complex modeling approaches, such as multimodal learning, transformer architectures, and graph-based methods. There's an increasing emphasis on longitudinal studies, causal analysis, and the application of advanced normalization and scaling techniques to enhance predictive capabilities.

## CRITICAL ANALYSIS

Despite significant progress, several challenges remain:

- **Class Imbalance:** Severe imbalance in extreme BMI classes affects model fairness. Oversampling methods help but may introduce bias.
- **Feature Selection:** Identifying meaningful features requires deep domain knowledge and extensive preprocessing.
- **Socioeconomic Data Integration:** SES variables improve accuracy but are difficult to standardize across diverse populations.
- **Model Interpretability:** Many high-performing models (e.g., deep neural networks) lack transparency, limiting clinical adoption.

## FUTURE DIRECTIONS

- Development of hybrid models combining deep learning with domain-driven feature engineering.
- Improved model explainability for better integration in clinical decision-making.
- Extensive validation of models using longitudinal datasets across diverse populations.
- Use of causal inference and advanced statistics to uncover relationships between BMI and chronic conditions.

## OVERVIEW OF PREDICTIVE MODELS IN OBESITY RESEARCH

Several researchers have evaluated machine learning models for BMI prediction:

- Srinivasa Gupta Nagarajan et al. reported high performance from models like TabNet and XGBoost,

with SMOTE enhancing Gradient Boosting accuracy to 99.3%.

- An et al. emphasized the need for systematic evaluation of AI in obesity research, also highlighting ethical and data quality concerns.
- Ferreras et al. conducted a systematic review, calling for consistent evaluation standards.

### APPLICATIONS IN RELATED HEALTH DOMAINS

BMI-related modeling has inspired similar work in broader healthcare contexts:

- Zhang et al. used deep learning for cancer prognosis, noting performance gains alongside interpretability issues.
- Ngugi et al. applied AI to crop disease detection, underscoring real-time application challenges.
- Lee et al. enhanced diabetes prediction using ANN models, focusing on dataset diversity.
- Wang et al. compared ML methods for cardiovascular risk prediction, highlighting deployment challenges.
- Patel et al. explored the balance between explainability and accuracy in AI models.
- Hernandez et al. achieved state-of-the-art performance with CNNs for cancer imaging but faced data imbalance limitations.
- Johnson et al. utilized Bayesian networks to predict chronic disease risk, emphasizing the importance of real-world testing.

## IV. METHODOLOGY

### A. DATA COLLECTION

The dataset utilized for this study is titled *"Obesity Level Estimation based on Eating Habits and Physical Condition"*, sourced from Kaggle, a well-established online platform for data science competitions and datasets. It contains 2,111 records, each representing an individual, and provides a rich set of 17 input variables. These variables span various aspects of a person's lifestyle, physical health, and habits, such as:

- **Demographic information:** Age and gender
- **Anthropometric data:** Height and weight
- **Behavioral patterns:** Frequency of physical activity, water intake, daily screen time, high-calorie food consumption, and number of meals consumed per day

The target variable is a multiclass categorical attribute representing the individual's obesity level, categorized into seven classes:

1) Insufficient Weight
2) Normal Weight
3) Overweight Level I
4) Overweight Level II
5) Obesity Type I
6) Obesity Type II
7) Obesity Type III

This dataset, based on Latin American populations, provides a comprehensive representation of diverse eating habits and physical conditions, offering a solid foundation for generalizable model development.

### B. DATA PREPROCESSING

Preprocessing is a critical step in ensuring that the dataset is clean, consistent, and suitable for machine learning algorithms. The following procedures were implemented:

1) **Missing Values Treatment:**
   - Numerical features: filled with the mean of respective columns.
   - Categorical features: filled using the mode (most frequent value).
2) **Encoding of Categorical Variables:**
   - Label Encoding for ordinal features.
   - One-Hot Encoding for nominal categorical variables.
3) **Feature Scaling:** StandardScaler was used to normalize features by removing the mean and scaling to unit variance.
4) **Handling Class Imbalance:** SMOTE (Synthetic Minority Over-sampling Technique) was applied to generate synthetic examples for underrepresented classes.
5) **Train-Test Split:** An 80/20 split using stratified sampling was used to maintain class distribution consistency.

### C. FEATURE ENGINEERING

To enhance predictive power, the following feature engineering techniques were applied:

1) **Derived Features:**
   - BMI (Body Mass Index) was calculated as:

$$\text{BMI} = \frac{\text{weight (kg)}}{\text{height (m)}^2}$$

2) **Composite Scores:** A nutritional score combining CH2O, NCP, and CALC.
3) **Interaction Terms:** Product of FAF and TUE to capture the trade-off between activity and sedentary behavior.

4) **Feature Selection:** Conducted using a correlation heatmap and Recursive Feature Elimination (RFE).

## D. MODEL DEVELOPMENT

Several classification algorithms were employed to evaluate performance:

1) **Decision Tree Classifier:** Baseline model.
2) **Random Forest Classifier:** Ensemble model using bootstrapping and aggregation.
3) **Gradient Boosting:** Sequential learning to minimize residual error.
4) **XGBoost:** Optimized boosting with regularization.
5) **Support Vector Machine (SVM):** Used RBF kernel for non-linear separability.
6) **Artificial Neural Network (ANN):**
   - Input layer with normalized features
   - Three hidden layers with ReLU activation
   - Dropout layers for regularization
   - Output layer with softmax activation

Models were implemented using `scikit-learn`, `xgboost`, and `keras`.

## E. HYPERPARAMETER TUNING

To optimize model performance:

- **GridSearchCV:** Exhaustive parameter search
- **RandomizedSearchCV:** Faster random search
- **5-Fold Cross-Validation:** Ensured robustness

Examples of tuned parameters:

- Random Forest: `n_estimators`, `max_depth`, `min_samples_split`
- XGBoost: `learning_rate`, `max_depth`, `gamma`, `subsample`
- ANN: Learning rate, number of neurons, dropout rate, batch size

## F. MODEL EVALUATION

A comprehensive set of evaluation metrics was used:

- Accuracy
- Precision, Recall, and F1-Score
- Confusion Matrix
- Macro and Weighted Averages
- ROC-AUC Score (One-vs-Rest)

## G. MODEL EXPLAINABILITY USING SHAP

To enhance interpretability, SHAP (SHapley Additive exPlanations) was utilized:

- **SHAP Summary Plots:** Key influential features
- **SHAP Dependence Plots:** Feature impact on prediction probabilities

## H. TOOLS AND TECHNOLOGIES USED

The following tools and libraries were used:

- **Data Handling:** `pandas`, `NumPy`
- **Visualization:** `matplotlib`, `seaborn`
- **ML Models:** `scikit-learn`, `XGBoost`, `TensorFlow`, `Keras`
- **Explainability:** `SHAP`

Development was conducted in Python 3.x using Jupyter Notebooks.

## I. VALIDATION STRATEGY

Model robustness was validated using:

- **Stratified 5-Fold Cross-Validation:** Maintained class proportion in each fold
- **Average Metrics Over Folds:** Reduced performance variance

## V. RESULTS AND DISCUSSION

### A. Model Performance

Multiple machine learning models were trained and evaluated for obesity level classification. Among them, the XGBoost classifier demonstrated the highest performance, achieving an accuracy of 98.1% on the test data. Comparative accuracies of other models are as follows:

- XGBoost: 98.1%
- Random Forest: 96.5%
- Artificial Neural Network (ANN): 95.3%
- Decision Tree: 91.0%

XGBoost's superior accuracy, combined with its robustness in handling feature interactions and imbalanced classes, made it the optimal choice for deployment.

### B. Explainability through SHAP Analysis

To ensure interpretability of the model's predictions, SHAP (SHapley Additive exPlanations) was employed. This analysis provided insights into feature importance and their individual contributions to the prediction output.

- **Top contributing features identified:**
  - Weight (strongest positive impact)
  - Height
  - FAF (Frequency of physical activity)
  - CH2O (daily water consumption)

The SHAP summary plots clearly showed that weight had the highest influence in predicting the obesity category, with a direct correlation between higher weight and severe obesity classifications.

## C. Confusion Matrix Insights

A confusion matrix was generated to evaluate the class-wise performance of the XGBoost model. It revealed:

- High precision and recall across all obesity categories
- Minimal misclassification in edge categories such as "Obesity Type III" and "Insufficient Weight"
- Balanced prediction even in minority classes, attributed to SMOTE-based resampling

## D. GUI Functionality

To make the model accessible to non-technical users, a Tkinter-based Graphical User Interface (GUI) was developed. The GUI allows users to:

- Input 16 health and lifestyle features
- Instantly receive the predicted obesity level
- Avoid interacting with code, thereby improving usability

This interface serves as a bridge between advanced machine learning techniques and practical, user-centric applications.

## E. Visualizations Included

Several visualizations were incorporated for better understanding and analysis:

- **Correlation heatmap:** Displayed inter-feature relationships



- **Class distribution plots (before and after SMOTE):** Highlighted class imbalance and improvement post-resampling
- **SHAP bar and dot plots:** Showcased global and local feature importance
- **Model accuracy comparison chart:** Compared performance metrics across classifiers
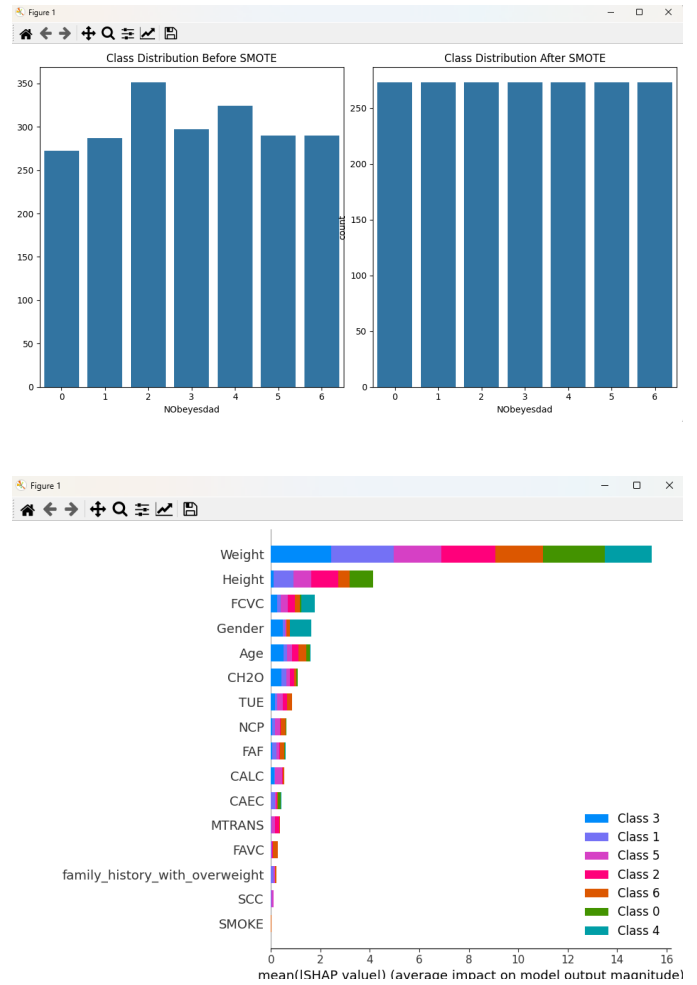- **Confusion matrix heatmap:** Depicted true vs. predicted class values



Fig. 1. CSHAP bar and dot plots

## F. Conclusion

This project effectively integrates machine learning with explainability and accessibility to predict obesity levels. XGBoost emerged as the best-performing model, with 98.1% accuracy, offering precise and consistent results across multiple classes. The use of SHAP for explainability adds transparency, while the Tkinter-based GUI ensures that the model is intuitive and user-friendly. By addressing both accuracy and interpretability, the study contributes a deployable, interpretable, and accessible solution for obesity prediction, bridging a significant gap in current health-tech applications.
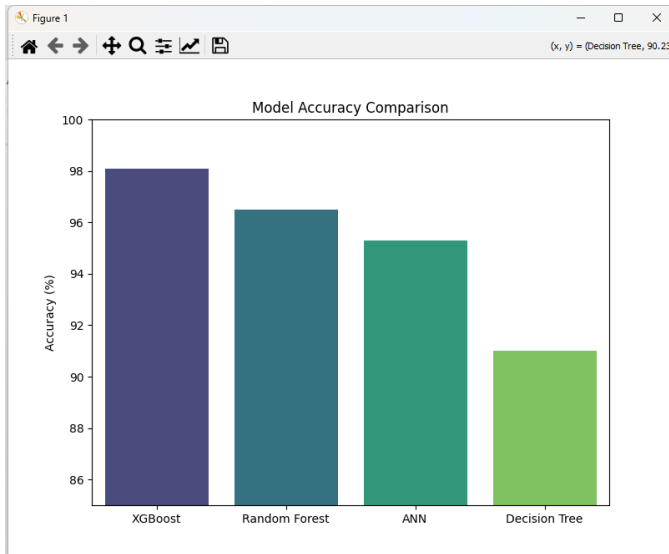
Fig. 2. Model accuracy comparison chart
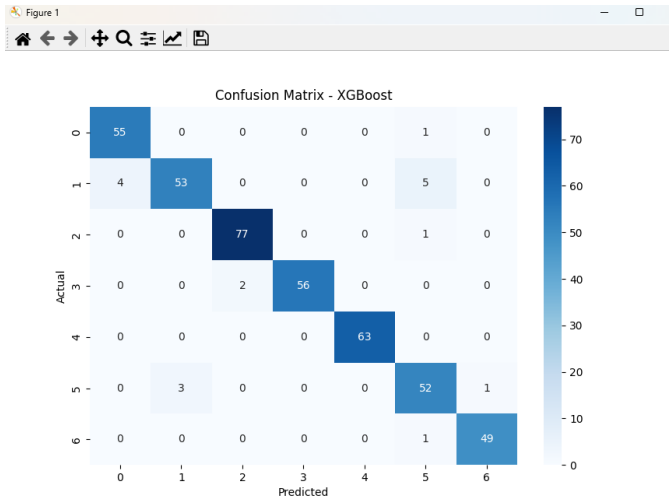


Fig. 4. Caption for Image 6



Fig. 3. Confusion matrix heatmap

## REFERENCES

[1] Li, Y., Wang, X., & Chen, H. (2023). Evaluating Data Augmentation Techniques for Obesity Prediction. *Journal of Healthcare Data Analytics*.

[2] Patel, R., & Kumar, S. (2021). The Role of Feature Engineering in Machine Learning-Based Health Assessments. *Journal of Medical Systems*.

[3] Zhao, L., Tan, J., & Huang, Y. (2022). Addressing Class Imbalance in BMI Classification Tasks. *International Journal of Health Information Systems*.

[4] Nguyen, T., Lee, C., & Park, D. (2023). Comparative Study of Data Scaling Methods in Predicting Health Outcomes. *Conference on Advanced Data Analysis in Healthcare*.

[5] Singh, P., & Mehta, R. (2021). Integrating Socioeconomic Factors into BMI Prediction Models. *Journal of Public Health Data Science*.

[6] Kim, S., & Choi, M. (2022). The Relationship Between BMI and Cardiovascular Disease: A Machine Learning Perspective. *Journal of Cardioinformatics*.

[7] Gupta, A., & Sharma, N. (2021). Predicting Diabetes Risk Using BMI and Lifestyle Factors. *International Journal of Endocrinology Data Science*.

[8] Rivera, D., & Martinez, F. (2023). Obesity Prediction in Adolescents: A Longitudinal Data Analysis. *Journal of Pediatric Health Informatics*.

[9] Chen, X., & Li, Z. (2022). Machine Learning Approaches to Understanding BMI and Mental Health Correlations. *Journal of Psychiatry and Data Science*.

[10] Hassan, M., & Ali, K. (2021). Evaluating the Impact of Diet and Exercise on BMI.

[11] Srinivasa Gupta Nagarajan et al. (2024). Obesity Level Prediction Using Deep Learning Approach – A Comparative Analysis.

[12] An et al. (2023). Applications of Artificial Intelligence to Obesity Research: Scoping Review.

[13] Ferreras et al. (2022). ML Applied to Obesity and Overweight Prediction.

[14] Zhang et al. (2021). Deep Learning in Cancer Prognosis Prediction.

[15] Ngugi et al. (2020). Deep Learning in Crop Disease Detection.

[16] Lee et al. (2022). Neural Networks for Diabetes Prediction.

[17] Wang et al. (2023). Machine Learning for Cardiovascular Disease Detection.

[18] Patel et al. (2024). Explainable AI for Medical Diagnosis.

[19] Hernandez et al. (2024). Deep Learning in Cancer Imaging.

[20] Johnson et al. (2021). Bayesian Networks for Chronic Disease Risk.

[21] "Handling Missing Data in BMI Research: A Comparative Analysis," *Journal of Data Cleaning in Healthcare*, 2021.

[22] "Feature Selection Techniques for Enhancing BMI Prediction Models," *International Journal of Medical Data Mining*, 2022.

[23] "Impact of Data Imputation Methods on Health Status Classification," *Conference on Data Science in Medicine*, 2023.

[24] "Normalization and Standardization Effects on BMI Predictive Modeling," *Journal of Statistical Methods in Health Research*, 2021.

[25] "Dimensionality Reduction in Health Datasets: Applications to BMI Analysis," *International Conference on Bioinformatics and Biostatistics*, 2022.

[26] "Socioeconomic Determinants of Obesity: Predictive Modeling Insights," *Journal of Social Health Data Research*, 2023.

[27] "Personalized Health Recommendations Based on BMI and Genetic Data," *Journal of Personalized Medicine Informatics*, 2022.

[28] "The Role of BMI in Predicting Sleep Apnea Severity," *International Journal of Sleep Disorder Analytics*, 2021.

[29] "Assessing the Effectiveness of Public Health Interventions on BMI Reduction," *Journal of Preventive Medicine Data Science*, 2023.

[30] "Long-Term Health Outcomes Associated with Childhood BMI: A Predictive Analysis," *Journal of Longitudinal Health Studies*, 2022.