

Progress Report 1

Group 5

Members: Frederick Damptey & Benjamin Odoom Asomaning

Course: SAT 5141 – Clinical Decision Support and AI Modeling

PROJECT TITLE: CLINICAL DECISION SUPPORT FOR EARLY IDENTIFICATION OF OBESITY-RELATED COMPLICATION RISK

1. Introduction and Literature Review

Obesity remains one of the most important worldwide public health problems, resulting in increased rates of chronic diseases such as type 2 diabetes, hypertension, and cardiovascular disease [1]. Early detection of obesity-prone individuals can prevent long-term morbidity and mortality. Clinical Decision Support Systems (CDSS) developed by Artificial Intelligence (AI) have emerged as valuable instruments for clinicians to diagnose, predict, and treat such complex diseases [2].

Recent studies show the effectiveness of AI-driven methods in medicine. Helforoush and Sayyad [3] developed a hybrid metaheuristic model to improve precision in obesity risk prediction, while Shen et al. [4] developed an interpretable visualization framework that enables more trust in clinical decision-making. Similarly, Lee et al. [5] integrated survey data with machine learning models for obesity prediction in Type II DM patients, demonstrating the superiority of ensemble methods to individual learners. These studies show how AI models are capable of being revolutionary in facilitating predictive medicine and clinical decision support.

2. Methodology

The objective of the research is to develop a python-based Clinical Decision Support System (CDSS) to predict patients at risk for obesity complications. The system is performing a supervised machine learning task given as a multi-class classification problem, wherein target classes are tagged as: Insufficient_Weight, Normal_Weight, Overweight_Level_I, Overweight_Level_II, Obesity_Type_I, Obesity_Type_II, and Obesity_Type_III.

Pipeline

The development pipeline includes six major phases:

- 1) Data Preprocessing: Cleaning and feature scaling.
- 2) Feature Engineering: Constructing BMI, caloric intake ratio and activity-level features.
- 3) Model Training: Employing AutoGluon Tabular, a ML model with multiple models for model stacking and benchmarking automation. AutoGluon tries out different algorithms (e.g., Random Forest, LightGBM, CatBoost, and Neural Networks) with cross-validation, and the top-performing ensemble model is selected with balanced accuracy as the first metric.
- 4) Evaluation: Model performance is assessed using multiple metrics comprising accuracy, sensitivity, specificity, F1-score and AUC.
- 5) Interpretability: SHAP values ascertain significant predictors to provide clinical transparency.
- 6) Clinician-in-the-Loop Interface: A Streamlit model prototype will be integrated to make CDSS interactive by allowing clinicians to review, approve, reject or override model suggestions. All decisions are recorded with time-tamped explanations for accountability and auditability.

Workflow:

Step	Description
1. Data Ingestion	Load and inspect dataset structure.
2. Preprocessing	Clean, normalize, and handle missing values.
3. Feature Engineering	Derive BMI, caloric balance, and activity indices.
4. AutoGluon Training	Automated model benchmarking and selection.
5. Model Evaluation	Compare metrics (Accuracy, AUC, F1).
6. Explainability (SHAP)	Identify key predictors and visualize influence.
7. Clinician Approval/Override	Clinician validates or overrides predictions.

3. Dataset Description and Validation

The Obesity Risk Dataset [6] of Kaggle, which contains 20,000 patient samples and 17 features, is used. The dataset consists of demographic, behavior, and anthropometric variables such as age,

gender, height, weight, physical activity level, caloric intake, family history, and lifestyle. The target variable, Obesity_Level, is seven BMI-based classes with varying risk groups.

Feature salience is significant: caloric balance and BMI directly influence obesity outcomes, while lifestyle variables (e.g., amount of exercise, dietary habits) enrich behavioral context. The data were also checked for validity using completeness checks, z-score outlier detection, normalization, and stratified k-fold cross-validation to ensure model generalizability. The dataset has been carefully evaluated and assessed and comply with deidentification legal requirements for healthcare data. It also contains no missing values, making the dataset reliable for producing consistent results.

4. Preliminary Results and Discussion

Early tests demonstrated consistent model performance across a range of algorithms employed by AutoGluon. Initial modeling with AutoGluon recorded a balanced Accuracy of 0.90 for the best performing ensemble, confirming the dataset's predictive validity. Experiment with 10 other individual ML models for comparison gave promising results with training accuracy of 0.997 and 0.95, and validation accuracy of 0.90 respectively for XGBoost and CatBoost. Among other models, RF produced a training accuracy of 1.00 and validation accuracy of 0.90. All these best performing models had ROC-AUC of 0.99.

Next Steps

Moving forward, we will further tune hyperparameters of the best performing models for testing and model calibration to improve model robustness and reliability. Additionally, we will explore feature importance analysis indicated BMI, caloric intake, and physical activity frequency to analyze best predictors in line with well-established clinical risk factors. This will help in the model's interpretability and clinical utility.

We are also yet to integrate the interactive “Streamlit clinician-in-the-loop” interface prototype, which will enable Clinicians to accept, reject or override model suggestions as well as providing textual reasons. Such an interaction supports ethical monitoring and reduces automation bias so that the CDSS enhances and not replace human capabilities.

5. Ethics, Security, and Legal Considerations

Given the sensitivity of healthcare data, the system complies with major regulatory frameworks such as HIPAA. Patient data remains de-identified.

Integration of explainable AI modules such as SHAP and LIME will enhance clinician trust by clarifying model reasoning. The governance framework supports fairness, bias detection, and privacy preservation while maintaining clinical autonomy.

6. Future Work and Group Contributions

Future work will comprise hyperparameter tuning, model calibration, and integrating new datasets for the purpose of increasing external validity. The next milestone will be on adding SHAP visualizations to the user interface and pilot testing with simulated patient data for validation.

Members' Contributions:

Frederick Damptey – Model design from data ingestion to audit Clinician-in-the-loop.

Benjamin Odoom Asomaning – Data Ingestion and preprocessing, literature review, model evaluation and report writing.

References

- [1] H. C. Whitlock, L. M. Williams, and D. M. Colditz, “Obesity as a global public health challenge: Epidemiology, risk factors, and prevention,” *The Lancet*, vol. 387, no. 10026, pp. 231–243, 2016.
- [2] E. H. Shortliffe and M. J. Sepúlveda, “Clinical decision support in the era of artificial intelligence,” *JAMA*, vol. 320, no. 21, pp. 2199–2200, 2018.
- [3] Z. Helforoush and H. Sayyad, “Prediction and classification of obesity risk based on a hybrid metaheuristic machine learning approach,” *Frontiers in Big Data*, 2024.
- [4] J. Shen et al., “Visualization obesity risk prediction system based on machine learning techniques,” *Scientific Reports*, 2024.
- [5] C. Lee et al., “A machine learning model for predicting obesity risk in patients with T2DM using health survey data,” *Diabetes & Metabolism Journal*, 2025.

- [6] J. P. Kochar, “Obesity Risk Dataset,” Kaggle, 2024. [Online]. Available: <https://www.kaggle.com/datasets/jpkochar/obesity-risk-dataset>
- [7] Q. T. Nguyen, H. T. Vo, and T. M. Le, “Artificial intelligence applications in obesity prediction: A systematic review,” *Healthcare Analytics*, vol. 5, no. 2, pp. 101–118, 2023.
- [8] A. Esteva et al., “A guide to deep learning in healthcare,” *Nature Medicine*, vol. 25, no. 1, pp. 24–29, 2019.