

# Predicting Diabetes Risk Using Machine Learning & Deep Learning Methods

---

**Course:** ALY6140 – Analytics System Technology

**Team: Immortals** (Sai Ashwin Anumula, Vamshi Krishna  
Korutla, Vyas Kadiyala, Xiaoxi Li)

**Date:** May 15, 2025

# Introduction

Diabetes remains one of the most serious global health concerns, impacting millions and placing substantial pressure on healthcare systems. The aim of this project is to identify and understand key demographic and lifestyle factors that contribute to diabetes and prediabetes risk.

Our key questions are:

- What lifestyle and demographic factors are most strongly associated with diabetes and prediabetes?
- How do income levels and healthcare access affect the likelihood of developing diabetes?
- Can we develop a reliable model that predicts diabetes status from health indicators?

To answer these questions, we applied advanced data analysis and predictive modeling techniques. Our methodology included data extraction, cleaning, exploratory data analysis (EDA), data visualization, and model building using logistic regression, random forest, XGBoost, SVM, and neural networks. These models were further optimized using SMOTE, undersampling, and class weighting.

## Data Extraction

We used the Diabetes Health Indicators Dataset, publicly available on Kaggle. This dataset includes 253,680 records with 22 features related to health status, lifestyle choices, demographics, and healthcare access. The primary target variable is Diabetes\_012, where:

- 0 = No diabetes
- 1 = Prediabetes
- 2 = Diabetes

The question we aimed to answer was: What variables most accurately predict the presence of diabetes, and can we build a predictive model to classify individuals accordingly?

## Data Cleanup

Initial data processing was done in Jupyter Notebook using Python libraries such as pandas, numpy, and sklearn. The cleanup included:

- Null value checks: The dataset had no missing values.
- Encoding check: All categorical data was already numerically encoded.
- Outlier detection and removal using the Interquartile Range (IQR) method.
- Balanced sampling techniques like SMOTE, undersampling, and class weighting to address class imbalance.

## Exploratory Data Analysis (EDA)

We conducted exploratory data analysis (EDA) to understand feature distributions and correlations.

Key insights:

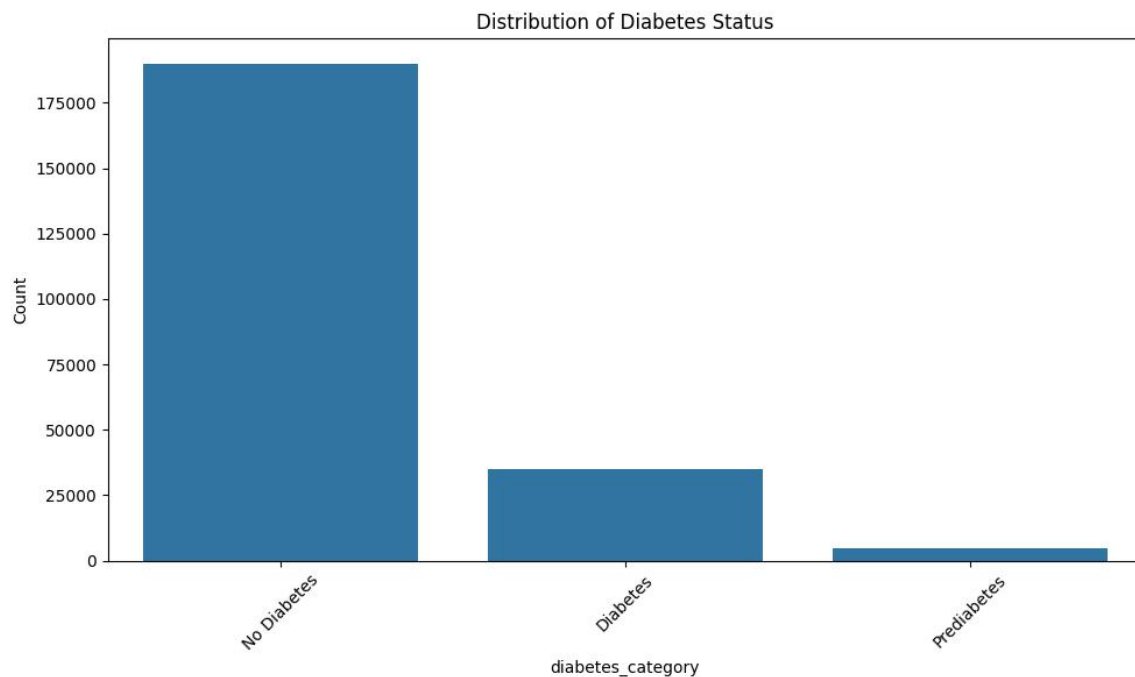
- Most participants had no diabetes (~175,000), followed by ~35,000 with diabetes, and ~5,000 prediabetics.
- Obesity Class III had the highest percentage of diabetes cases.
- Individuals with 3–5 risk factors were significantly more likely to have diabetes.
- People with healthier lifestyle scores (3–5) had lower diabetes prevalence.
- Poor general health was strongly correlated with diabetes.

We also used correlation heatmaps to visualize feature relationships with diabetes, identifying high BP and BMI as notable predictors.

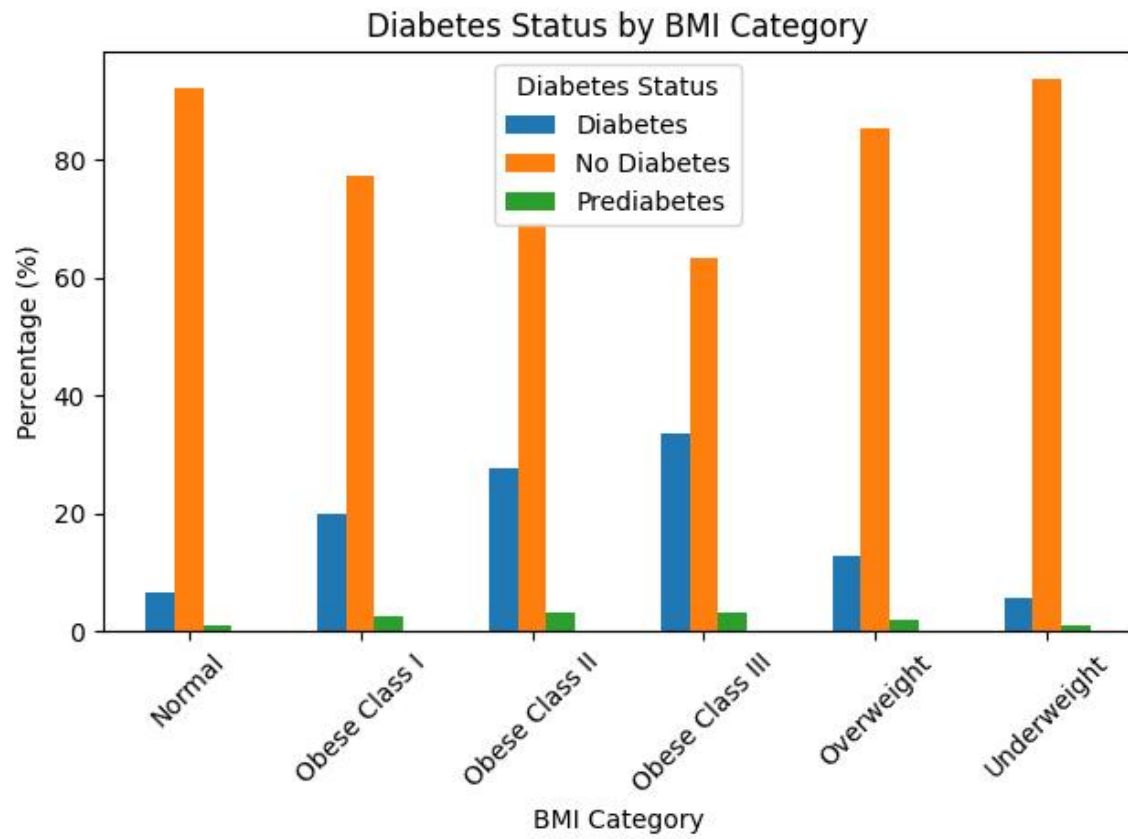
## Data Visualization

We used matplotlib and seaborn to explore:

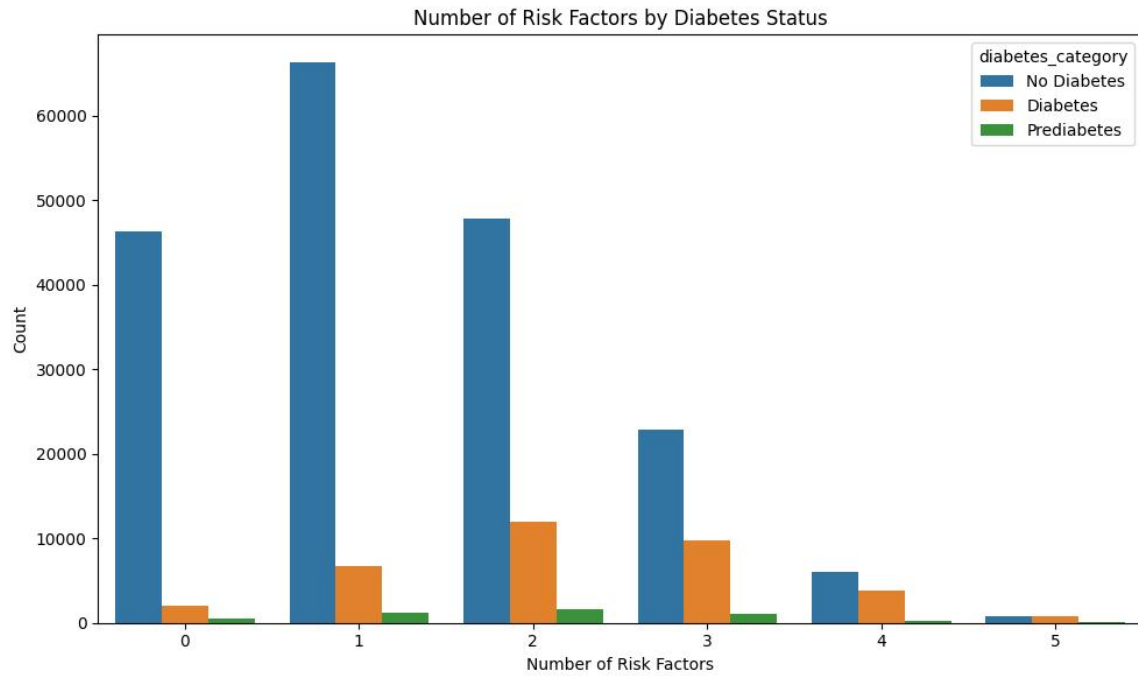
- Diabetes status distribution



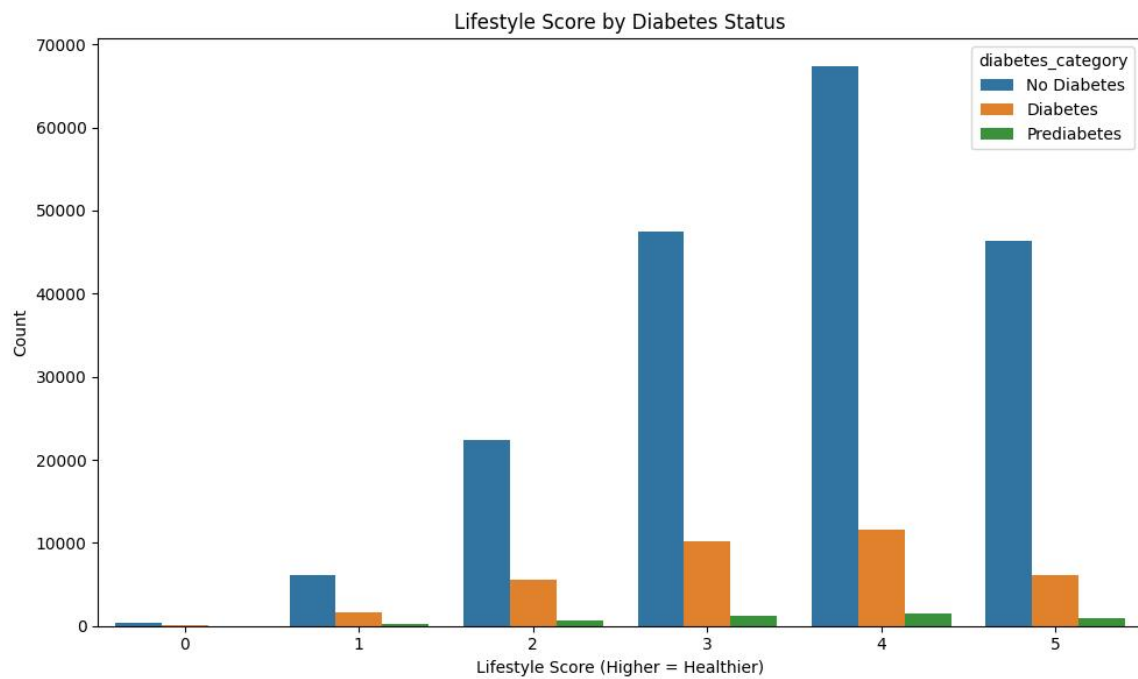
- BMI vs. Diabetes



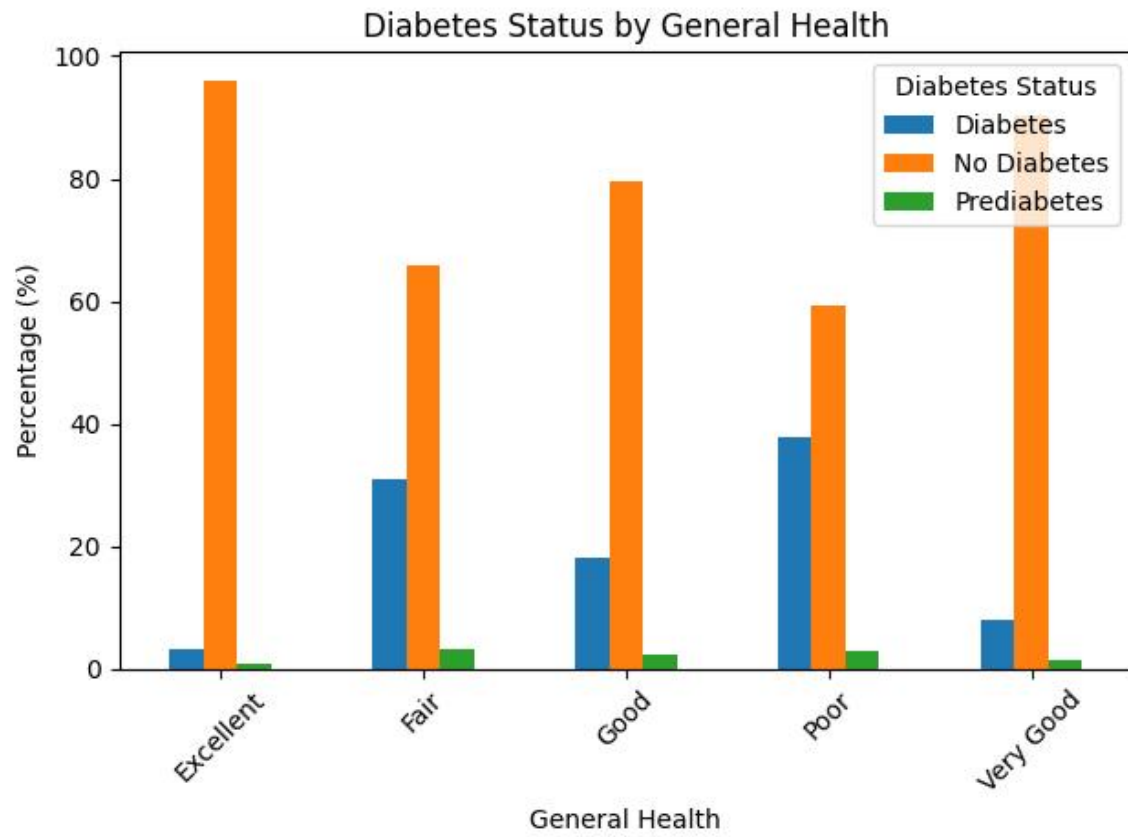
- Risk factors



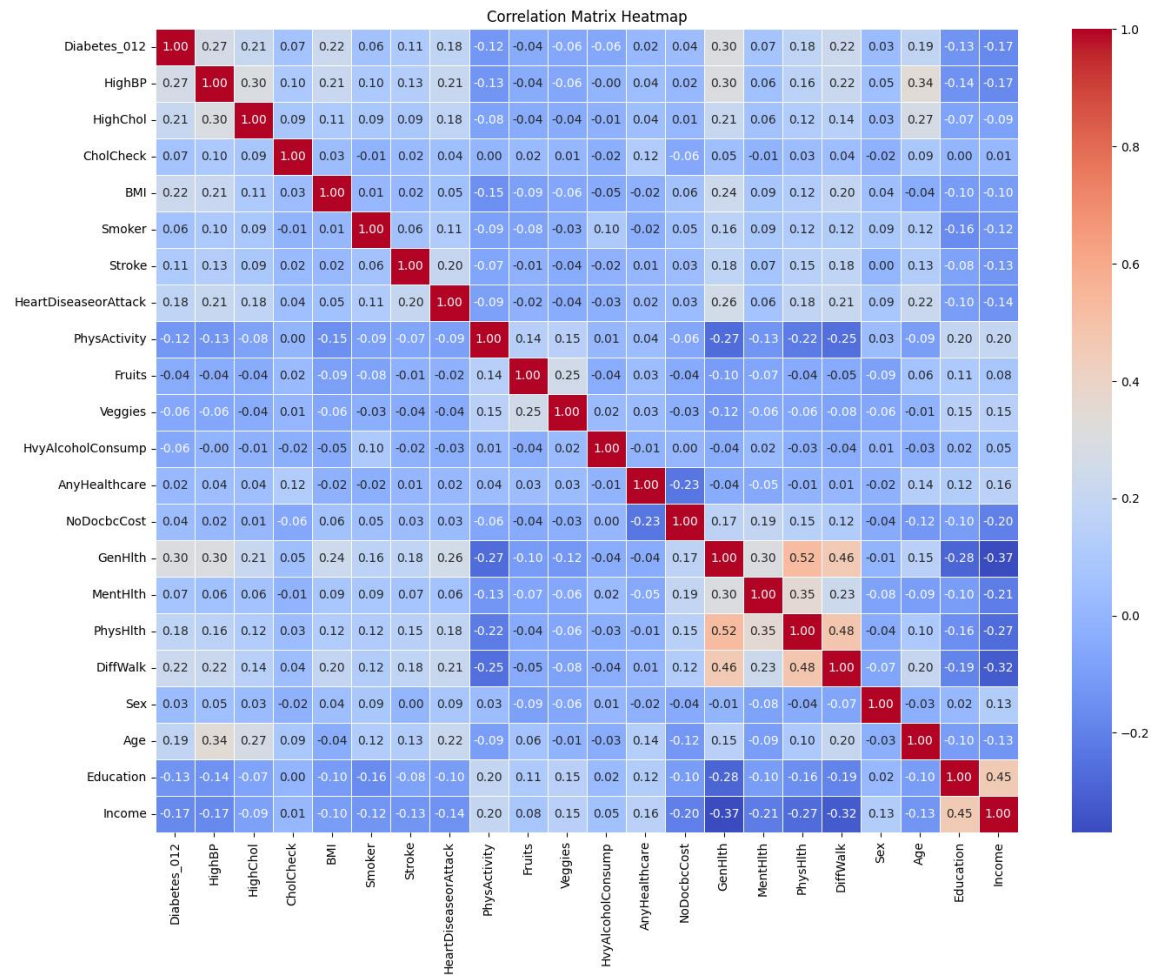
- Lifestyle scores



- General health



- Correlation heatmaps



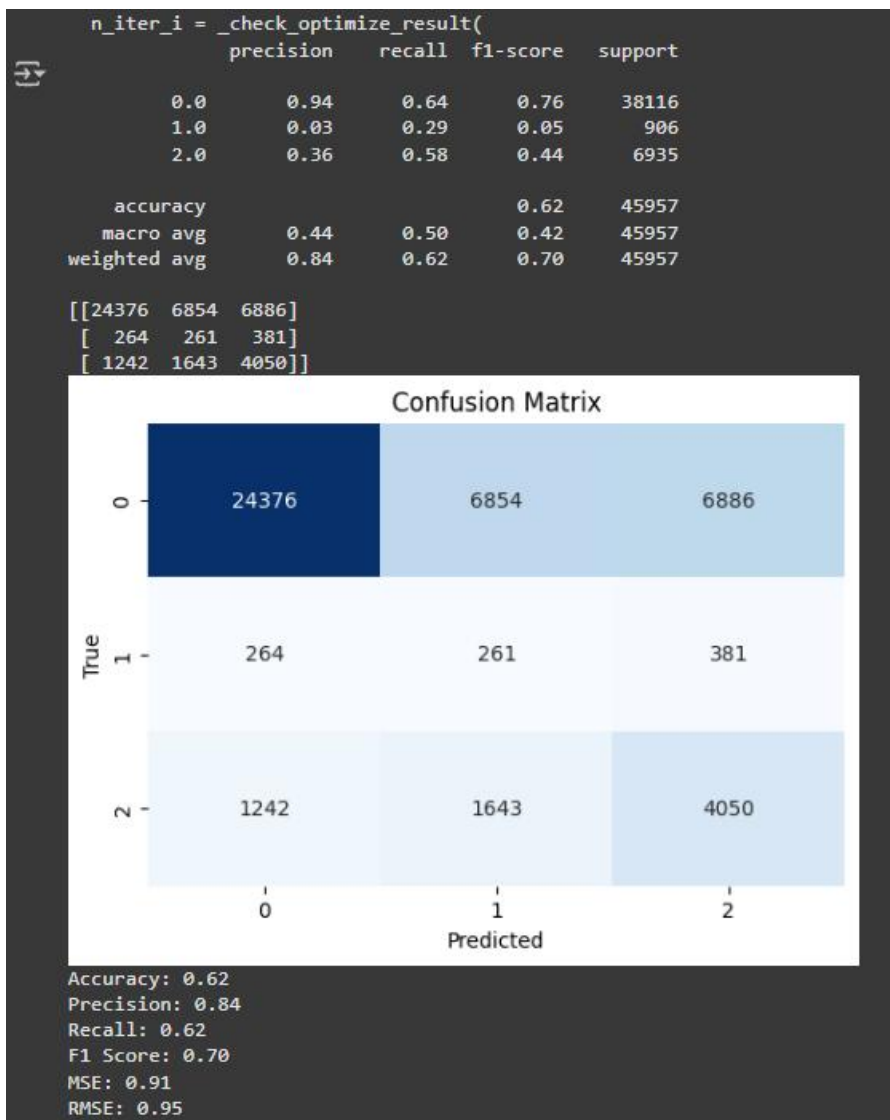
## Predictive Models (Statistical/Predictive Analysis)

We built and compared multiple classification models:

### 1. Logistic Regression:

- Interpretable model, decent performance using SMOTE.

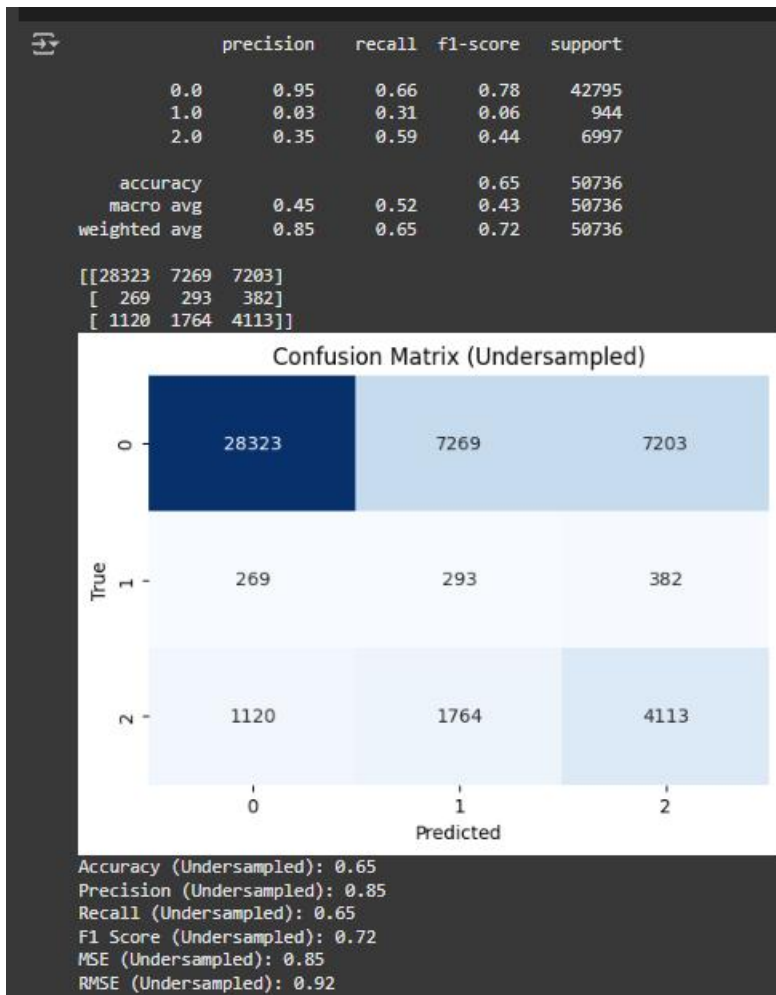
Using Smote Sampling for Logistic Regression:



The confusion matrix shows the model identifies the "No Diabetes" group well, with high precision and decent recall. However, it performs poorly on the "Diabetes" group and only moderately on "Prediabetes." With 62% accuracy and a 0.70 F1 score, the model favors the majority class, likely due to class imbalance.

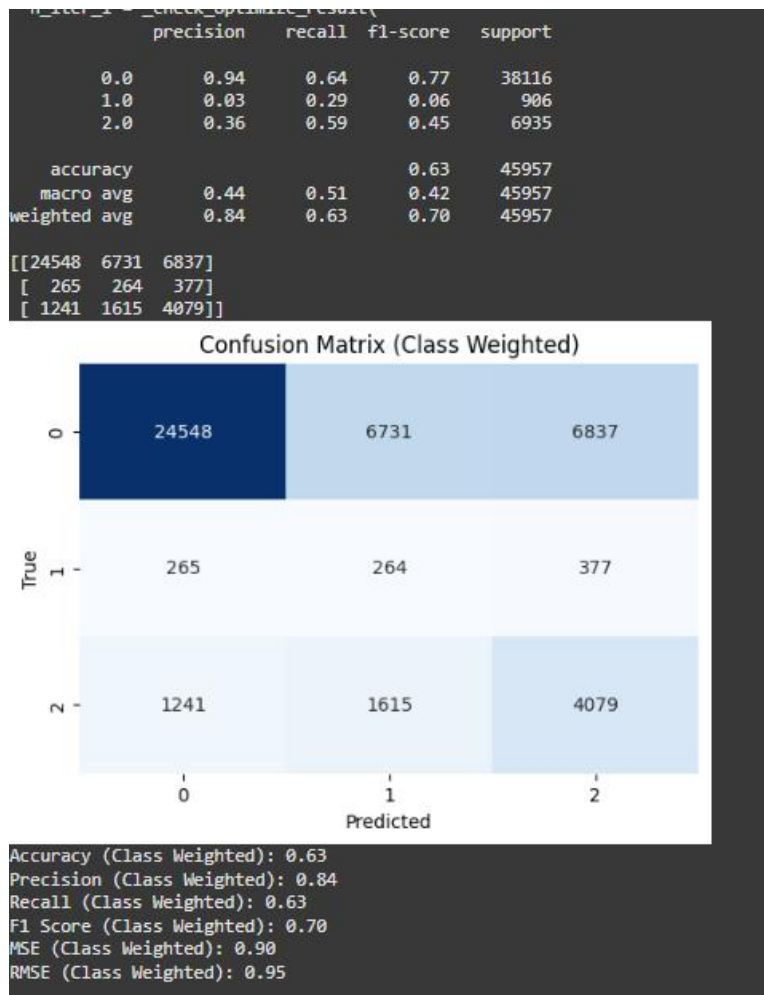
Using Undersampling for Logistic regression:





After undersampling, the model shows slightly better balance, with accuracy rising to 65% and an F1 score of 0.72. It still predicts "No Diabetes" well, while recall for "Diabetes" and "Prediabetes" improves slightly. However, detecting minority classes remains difficult.

Using Class Weighing Sampling for Logistic regression:

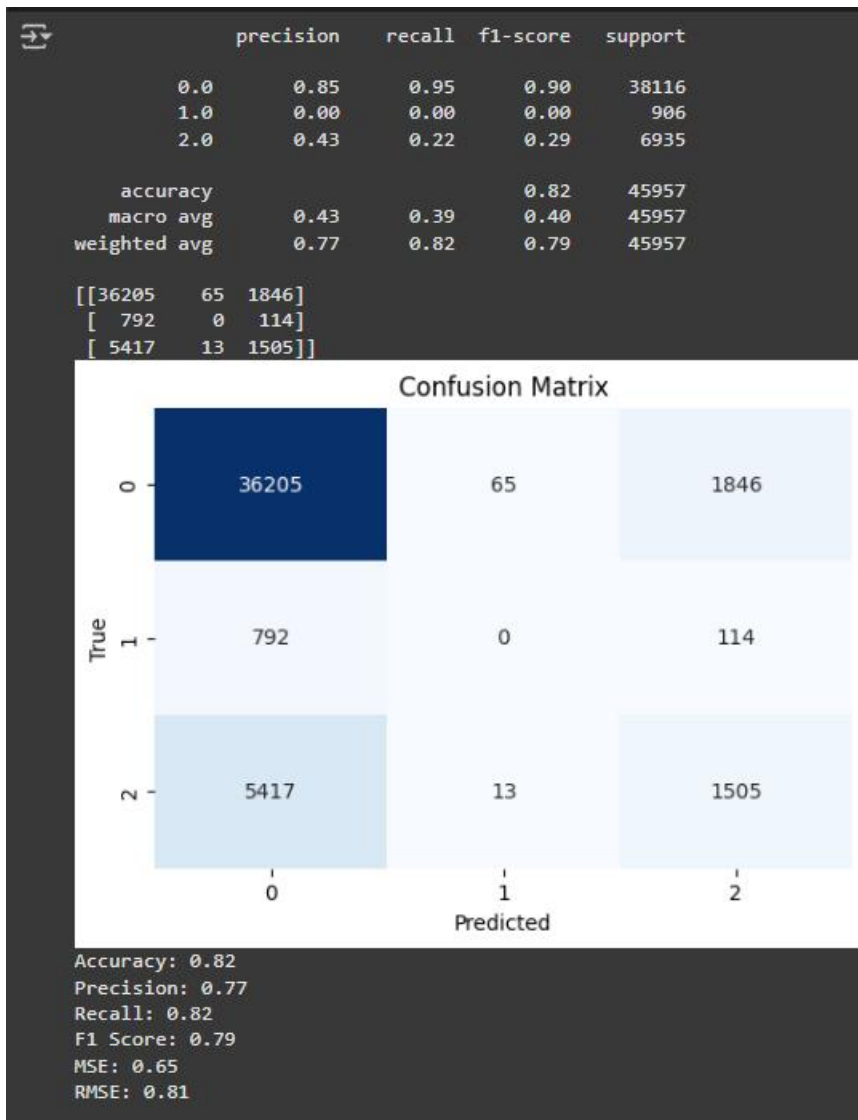


The class-weighted model reaches 63% accuracy and a 0.70 F1 score, similar to earlier models. It predicts "No Diabetes" well but continues to struggle with "Diabetes" due to low precision and recall. "Prediabetes" shows moderate improvement. Class weighting helps overall balance but doesn't fully fix the issue of poor minority class detection.

## 2. Random Forest:

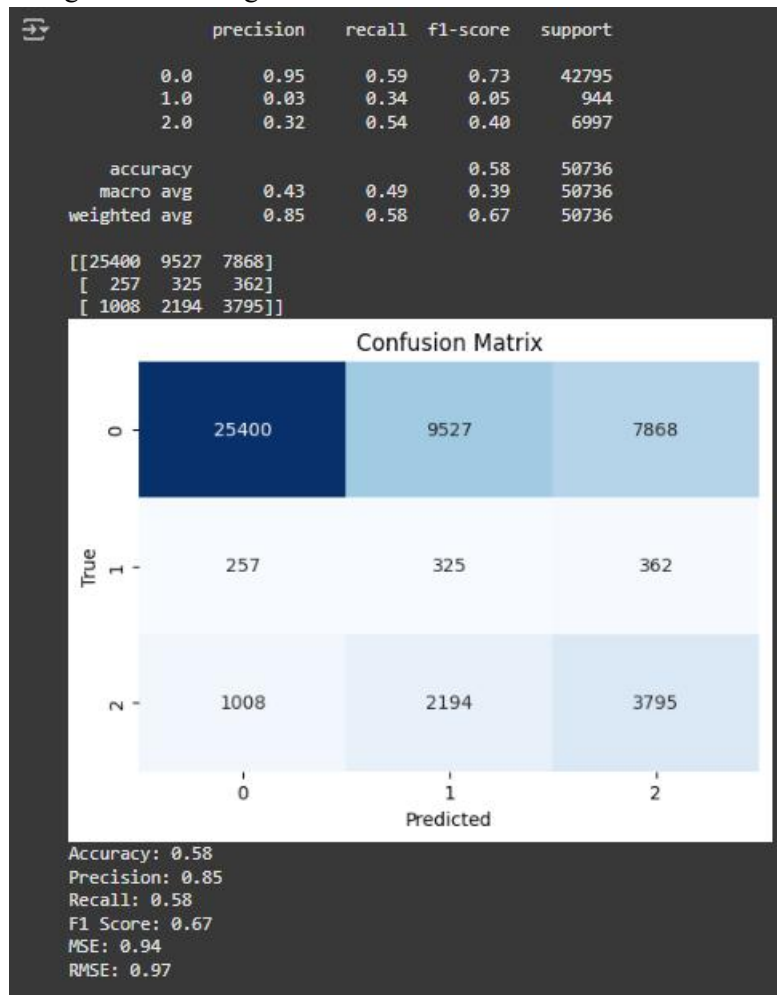
- Outperformed other models, especially with class weighting. Highlighted key features like HighBP, BMI, and Age.

Using Smote sampling for RandomForest:



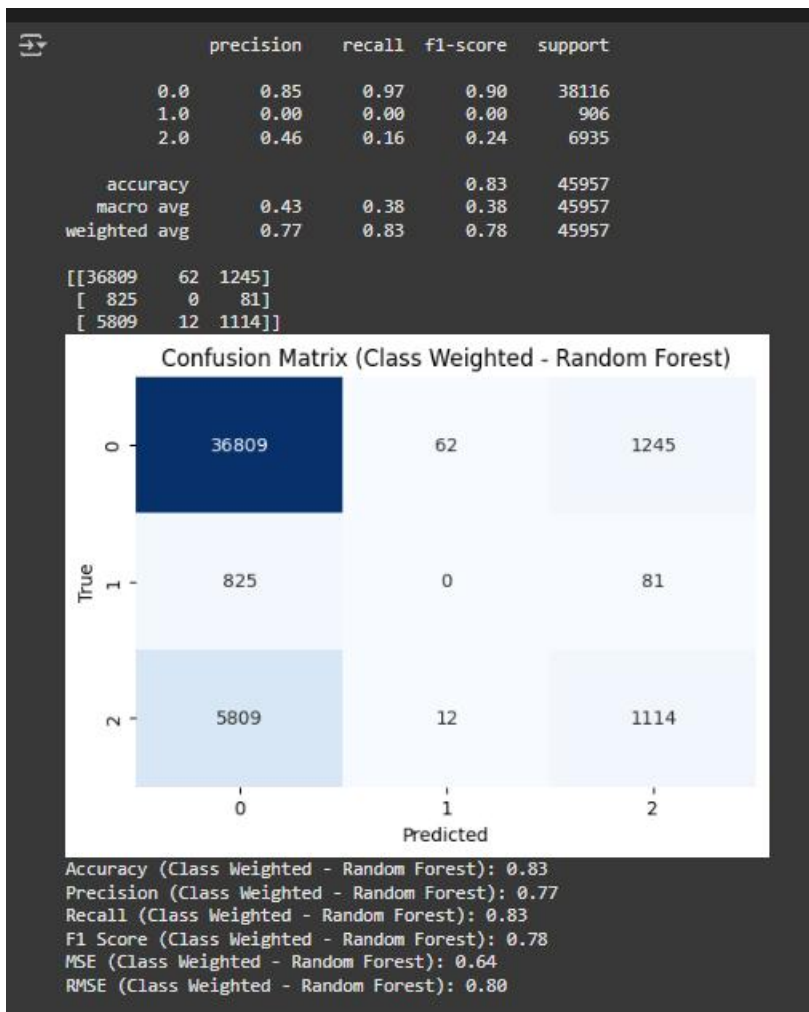
This model achieves 82% accuracy and a 0.79 F1 score, mainly due to strong performance on the "No Diabetes" group with high recall. However, it completely fails to detect "Diabetes" cases and performs poorly on "Prediabetes" as well. Despite strong overall metrics, the model heavily favors the majority class, limiting its usefulness in real-world healthcare.

## Using Understanding for Random Forest



This model has a lower accuracy of 58% but still predicts the "No Diabetes" class well, with high precision but moderate recall. It struggles with the "Diabetes" class, showing very low precision and recall, while "Prediabetes" has slightly better recall but poor precision. Overall, the model favors the majority class and remains weak at detecting minority groups, limiting its real-world reliability.

Class Weighing for Randomforest:

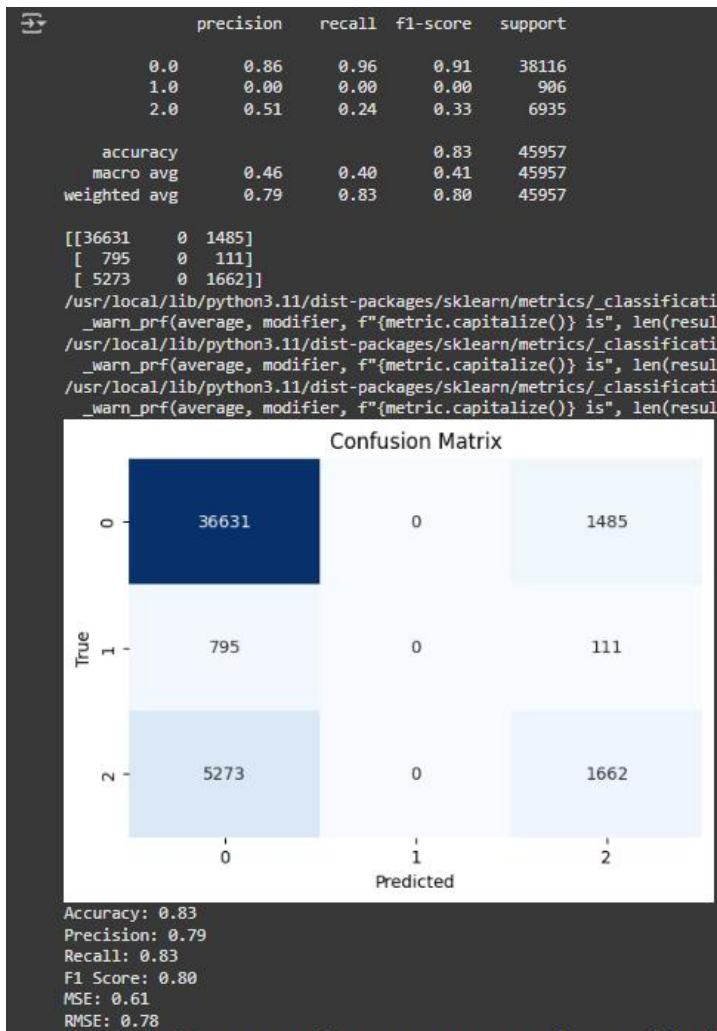


This class-weighted Random Forest model reaches 83% accuracy and a 0.78 F1 score, mainly due to strong recall for "No Diabetes." However, it fails to detect any "Diabetes" cases and performs poorly on "Prediabetes" as well. Despite good overall metrics, the model is heavily biased toward the majority class, making it less effective for balanced healthcare predictions.

### 3. XGBoost:

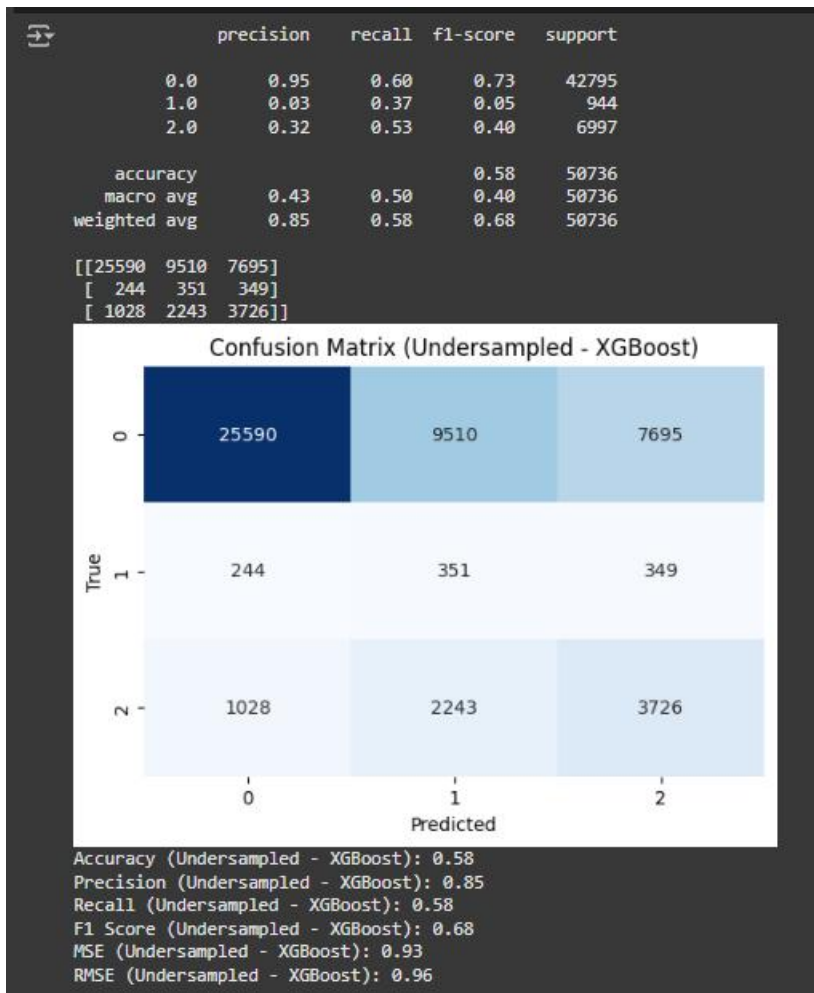
- High-performing model with good precision and recall using class weighting.

Using Smote Sampling for XG boost:



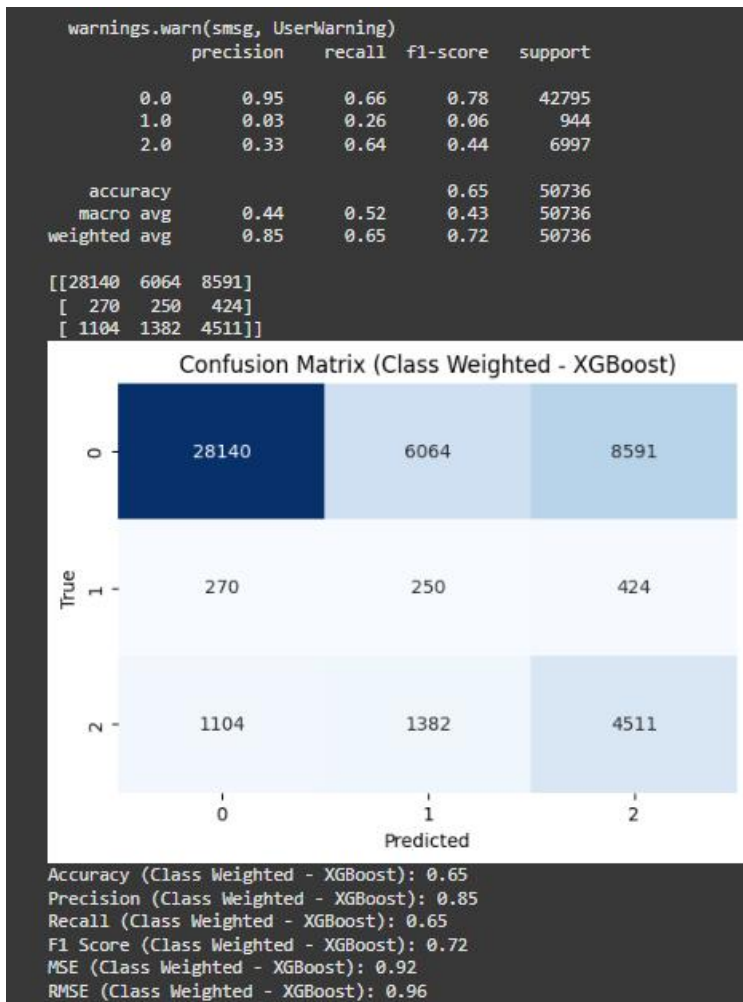
This model shows 83% accuracy and a strong F1 score of 0.80, driven by high recall for the "No Diabetes" class. However, it fails to detect any "Diabetes" cases and performs poorly on "Prediabetes," with low recall. The confusion matrix highlights a strong bias toward the majority class, making the model unreliable for real-world healthcare use where detecting all conditions accurately is essential.

Undersampling for XG boosting:



The undersampled XGBoost model reaches 58% accuracy and an F1 score of 0.68. It predicts the "No Diabetes" group well but only slightly improves recall for "Diabetes" and "Prediabetes." Precision for these classes remains low. Undersampling improves balance somewhat, but the model still favors the majority class and needs further refinement for reliable detection across all groups.

Class Weighing for XG boosting:



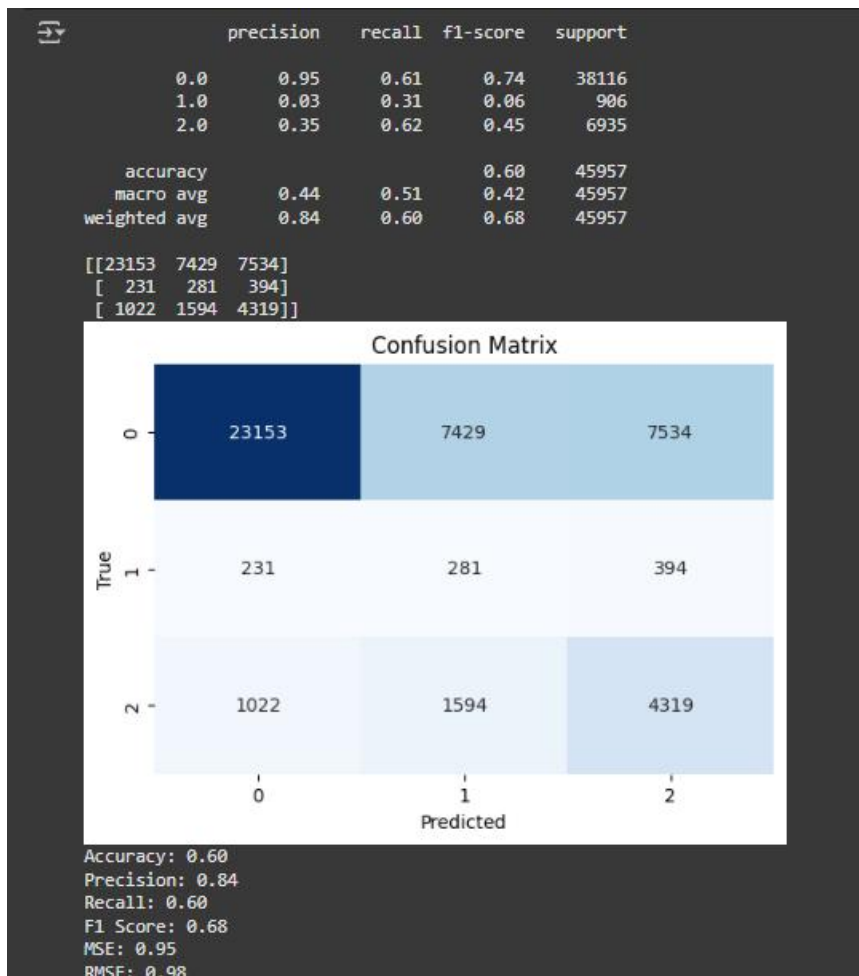
The class-weighted XGBoost model shows 65% accuracy and an F1 score of 0.72, with more balanced results across classes. It maintains strong precision for "No Diabetes" and improves recall for "Prediabetes" to 0.64. However, performance on "Diabetes" remains weak, with low recall and precision. While class weighting helps, the model still favors the majority class, and further adjustments are needed for better overall balance.

## 4. SVM:

- Effective on balanced data, but resource-intensive.

Using smote sampling for SVM:



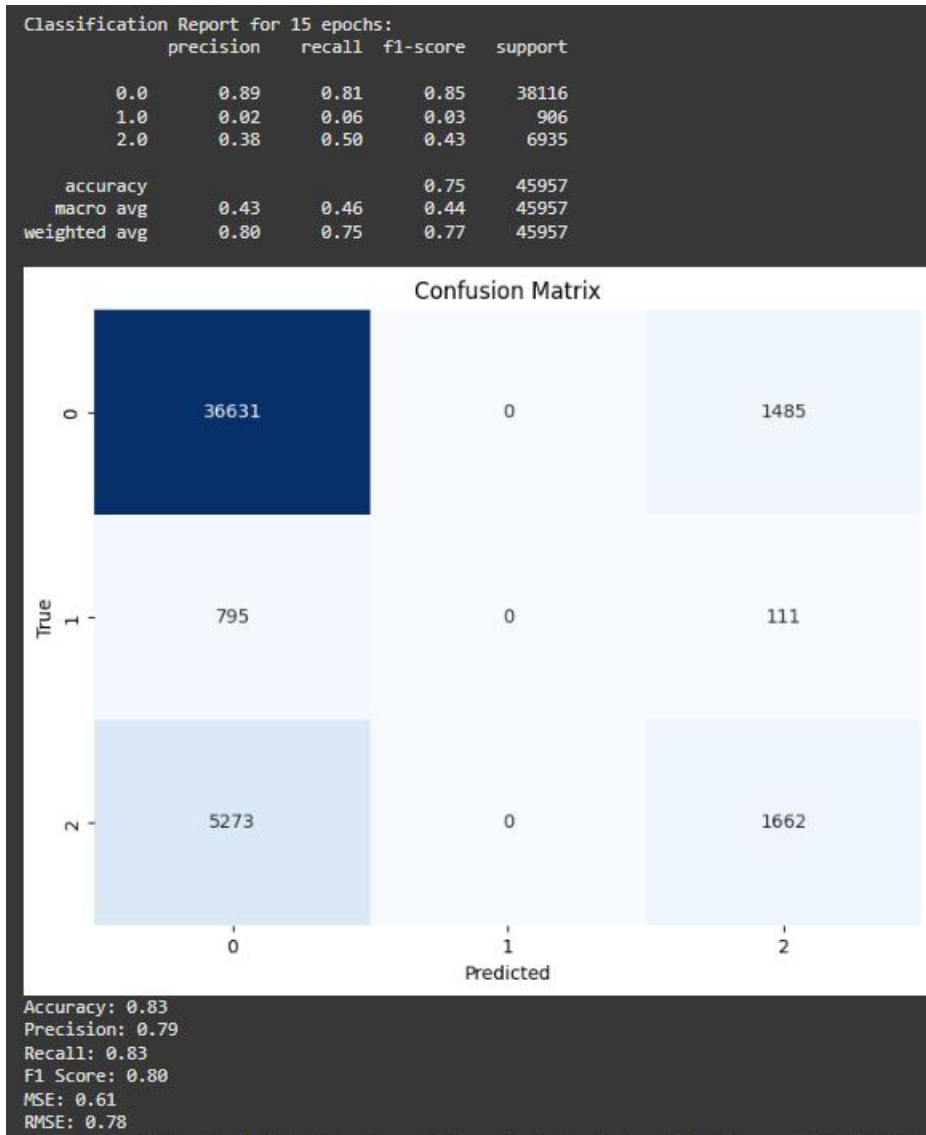


This model reaches 60% accuracy and an F1 score of 0.68. It performs well on "No Diabetes" with high precision but moderate recall. "Prediabetes" shows fair recall but low precision, while "Diabetes" remains poorly detected with very low scores. Overall, the model handles the majority class well but continues to underperform on minority class detection.

## 5. Neural Networks:

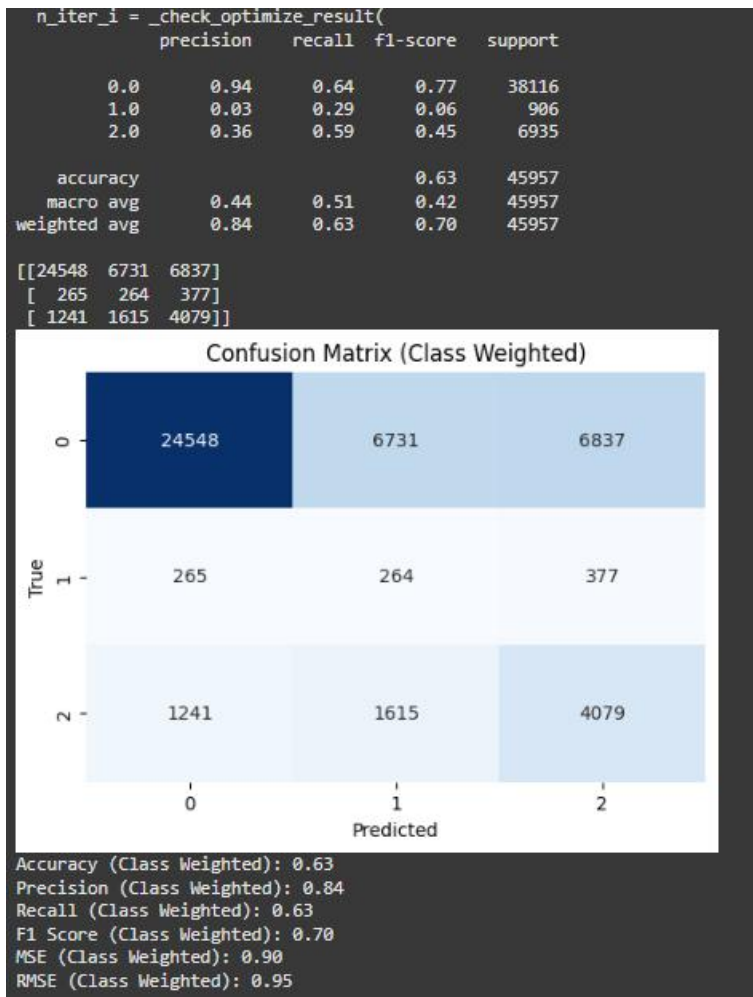
- Used Keras, performed reasonably well but required more resources and was less interpretable.

Smote Sampling for neural network:



This model delivers 83% accuracy and a strong F1 score of 0.80, largely due to its high performance on the "No Diabetes" group. However, it fails to detect "Diabetes" cases, with nearly zero precision and recall. The "Prediabetes" class performs slightly better but remains weak. Despite strong overall metrics, the model is still biased toward the majority class and needs improvement in identifying minority groups.

Class Weighing for Neural Network:



This class-weighted model reaches 63% accuracy and a 0.70 F1 score. It performs well on the "No Diabetes" group and shows moderate improvement for "Prediabetes" with better recall. However, it continues to struggle with detecting "Diabetes" due to low precision and recall. While class weighting helps a bit, the model still leans toward the majority class and underperforms on minority groups.

Evaluation metrics included accuracy, precision, recall, F1 score, and ROC-AUC where applicable. The best model was Random Forest with class weighting. It achieved strong overall accuracy and high performance in predicting the majority class. However, like many others, it still struggled to correctly identify minority classes, particularly individuals with diabetes. While class weighting helped improve balance slightly, most models, including XGBoost and neural networks, showed a clear bias toward the "No Diabetes" group. This highlights the need for further tuning or advanced techniques to boost minority class detection without sacrificing overall performance.

## Interpretation and Conclusions

Our analysis confirmed that demographic, lifestyle, and healthcare access variables significantly influence diabetes risk. High BMI, poor general health, and high blood pressure are the strongest indicators, while higher income and education offer protection.

Random Forest with class weighting proved to be the best-performing model. This can be used to support early intervention programs.

Based on the findings, several recommendations can help improve outcomes. First, launching preventive programs that promote healthy lifestyles may reduce health risks over time. Second, making healthcare more accessible can ensure earlier detection and better management of conditions. Lastly, including more detailed information, like genetic data and long-term health history, could improve the accuracy of future predictions and analyses.

## References

1. Teboul, A. (2022). Diabetes Health Indicators Dataset. Kaggle.  
<https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset>
2. American Diabetes Association. (2023). Standards of Medical Care in Diabetes—2023. Diabetes Care, 46(Supplement 1), S1–S272.
3. World Health Organization. (2021). Global report on diabetes. WHO.