

PROJECT STAGE 2
GROUP- 49

Member 1: Unikey- akur0860
Student ID: 550217239

Member 2: Unikey- ysab0639
Student ID: 540958494

Group Component 1

1. Topic and Research Question

Research Topic:

The rising prevalence of diabetes has become a global health crisis, posing a significant burden on both individuals and healthcare systems. Early detection is crucial for effective management and prevention of long-term complications. Leveraging machine learning offers an opportunity to build predictive models that identify at-risk individuals using lifestyle and clinical factors.

Research Question:

To what extent can lifestyle factors and health history serve as reliable predictors for the likelihood of developing diabetes in patients?

2. Dataset

Overview:

- **Dataset Name:** diabetes_diagnosis.csv
- **Size:** ~264,802 observations
- **Final Variables After Cleaning:** 17
- **Target Variable:** Diabetes Diagnosis – encoded as binary (0: No Diabetes, 1: Prediabetes/Diabetes)

Challenges Faced and Rationale for Preprocessing Decisions:

- **Missing Values:**
Many features had extensive missingness. Columns with over 50% missing values (e.g., PhysActivity, Fruits) were dropped to maintain dataset reliability. For other variables, mode/median imputation was applied to minimize bias and retain the underlying data structure.
- **Class Imbalance:**
The dataset showed a clear bias:
 - Class 0 (No Diabetes): 55,122
 - Class 1 (Prediabetes/Diabetes): 10,282To correct this imbalance and prevent the model from overfitting on the majority class, **SMOTE (Synthetic Minority Oversampling Technique)** was used. SMOTE creates synthetic instances of the minority class, leading to more balanced and generalizable predictions.
- **Bias in Data & Feature Selection:**
Gender had entries labeled as "Not specified", which were removed to prevent ambiguity in model learning. Columns were either one-hot encoded or retained in original form—not both—to prevent multicollinearity and reduce redundancy.
- **Dropped Columns:**
Some features like BMI, Age Group Classes, General Health, and Cholesterol Check were removed because either their one-hot encoded alternatives were retained, or they showed negligible predictive power. Child and Teen age groups were dropped as they had zero entries.
- **Uniformity:**
All final columns were converted to integers to ensure data type consistency and model compatibility.

Why These Steps Matter:

Each preprocessing decision was aimed at:

- Preserving data integrity
- Minimizing biases
- Ensuring clean, interpretable, and machine-readable input for modeling

3. Setup

3.1. Modelling Agreements

Target Attribute:

- **Diabetes** (Binary: 0 = No Diabetes, 1 = Prediabetes or Diabetes)

Success Metrics:

To evaluate the models comprehensively, the following performance metrics were selected:

1. **Precision & Recall:**

- **Why:** In a healthcare scenario, false positives (predicting diabetes when there's none) are inconvenient but manageable. False negatives (missing actual diabetes cases) are dangerous. Hence, recall (sensitivity) becomes critical, and precision ensures those flagged are likely to actually be at risk.
2. **F1-Score:**
 - **Why:** A balanced metric that combines both precision and recall. Especially valuable in imbalanced datasets where accuracy might be misleading.
 3. **AUC-ROC Curve:**
 - **Why:** Provides a comprehensive view of the model's ability to discriminate between the two classes across all thresholds.
 4. **Precision-Recall Curve & Learning Curve:**
 - **Why:** The PR curve is better suited to imbalanced datasets, and the learning curve provides insights into overfitting/underfitting behavior during training.

Chosen Model:

- **XGBoost**
 - **Why:** It is robust to outliers, handles feature interactions effectively, and works well on tabular data. Its ability to handle imbalanced datasets (especially when tuned with SMOTE and custom loss functions) makes it a strong choice for medical predictions.
- **SVM**
 - **Why:** Support Vector Machines are effective for high-dimensional spaces and perform well in binary classification tasks. They are particularly robust in scenarios where there is a clear margin of separation between classes. SVMs are less prone to overfitting, especially with appropriate kernel functions (linear, RBF, polynomial), and are useful for medical predictions due to their strong performance on smaller, cleaner datasets. Additionally, class imbalance can be mitigated using class weights and careful tuning of the regularization parameter C

3.2. Data Division

Split Strategy:

- **Training Set:** 70%
- **Validation Set:** 15%
- **Test Set:** 15%

Approach:

- **Stratified Sampling:** Ensured that both classes (0 and 1) were proportionally represented across all subsets.
- **Why:** This preserves the class distribution in each subset and ensures the model is evaluated realistically. Random splitting could result in one class being underrepresented in the test set, which would distort the model's performance metrics.
- **Temporal Validation Not Used:**
As this dataset lacks timestamps, temporal validation wasn't applicable. Instead, stratified random splits maintained class balance.

XGBoost for Imbalanced Classification: Model Analysis and Optimization

1. Model Analysis

1.1 Model Description

XGBoost (eXtreme Gradient Boosting) is an advanced implementation of gradient boosting machines designed for efficiency and performance. This model is particularly appropriate for the dataset with significant class imbalance (5.36:1 majority to minority ratio) containing 33 features and 65,404 training samples.

Key Assumptions:

- No strict distribution requirements (unlike statistical models)
- Tree-based approach handles non-linear relationships effectively
- Regularization helps prevent overfitting in high-dimensional space

Strengths for this problem:

- Handles imbalanced data through `scale_pos_weight` parameter
- Built-in regularization (L1/L2) controls model complexity
- Feature importance provides interpretability
- Early stopping prevents overfitting
- Integrates well with resampling techniques

Limitations:

- Less interpretable than simpler models
- Computationally intensive
- Sensitive to hyperparameter selection
- Can overfit without proper regularization

Suitability rationale: XGBoost combines high predictive power with specific mechanisms for addressing class imbalance, making it ideal for problems where correctly identifying the minority class is crucial while maintaining overall accuracy.

1.2 Model Algorithm

XGBoost Algorithm Pseudocode:

```

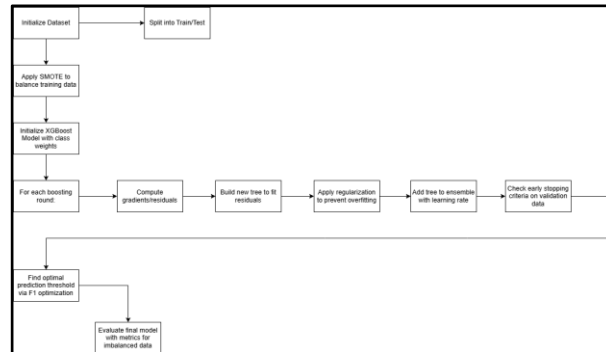
Algorithm: XGBoost for Binary Classification
Input: Training data (X, y), validation data (X_val, y_val), hyperparameters
Output: Ensemble of decision trees

1. Initialize model with constant values:
   F_0(x) = 0

2. For m = 1 to M (number of boosting rounds):
   a. Compute negative gradients (residuals) for each instance:
      r_i = -(y_i - F_{m-1}(x_i)) / (F_{m-1}(x_i))
   b. Fit a regression tree h_m(x) to the negative gradients
      - Apply split finding algorithm with regularization
      - Maximize gain = (gradient_sum_left^2 / hessian_sum_left + gradient_sum_right^2 / hessian_sum_right - gradient_sum^2 / hessian_sum) / 2 * \eta
   c. For each leaf j in the tree, compute optimal weight:
      w_j = -sum(gradient) / sum(hessian * \lambda)
   d. Update the model:
      F_m(x) = F_{m-1}(x) + \eta * h_m(x)
      where \eta is the learning rate
   e. Evaluate model on validation data
      If early_stopping_rounds is specified and no improvement for specified rounds, stop

3. Return final model F_M(x)
  
```

Flowchart of XGBoost with SMOTE



1.3 Model Development

The development of the XGBoost model for this imbalanced classification problem involved several sophisticated preprocessing, training, and evaluation steps:

Advanced Data Preprocessing:

- Train-Test Split with Stratification:**
 - The data was split into 80% training and 20% testing using stratified sampling to maintain class distribution
 - `random_state=15` ensures reproducibility
- Advanced Resampling Techniques:**
 - Multiple resampling methods were evaluated to address class imbalance:

- SMOTE (Synthetic Minority Over-sampling Technique), ADASYN (Adaptive Synthetic Sampling), SMOTEENN (SMOTE + Edited Nearest Neighbors), SMOTETomek (SMOTE + Tomek Links)
- SMOTETomek was selected as it provided the most balanced class distribution (55,119 samples per class) while maintaining data quality

3. Class Weight Calculation:

- Class weights were calculated based on the inverse ratio of class frequencies:
 - `class_weights = {0: 1, 1: before_class_dist[0] / before_class_dist[1]}`
 - This resulted in weights of {0: 1, 1: 5.36} for majority and minority classes respectively

Pipeline Integration:

The model was integrated into an imbalanced learning pipeline:

- `ImbPipeline([('sampling', SMOTETomek()), ('classifier', xgb_model)])`
- This ensures that resampling is properly applied during cross-validation

Threshold Optimization:

A critical aspect of the model development was optimizing the classification threshold:

- The default 0.5 threshold is suboptimal for imbalanced data
- Precision-recall curve analysis was used to find the threshold that maximizes F1 score
- This threshold optimization significantly improved model performance

2. Model Evaluation and Optimization

2.1 Model Evaluation

The model's performance was evaluated using multiple metrics specifically designed for imbalanced classification problems:

Performance Metrics for the Tuned Model:

- **Accuracy:** 0.7673
- **AUC-ROC:** 0.7863
- **F1 Score:** 0.4478
- **Average Precision:** 0.3987

Classification Report with Default Threshold:							
	precision	recall	f1-score	support			
0	0.98	0.27	0.42	13781			
1	0.20	0.98	0.33	2571			
accuracy			0.38	16352			
macro avg	0.59	0.62	0.38	16352			
weighted avg	0.86	0.38	0.41	16352			
Classification Report with Optimal Threshold:							
	precision	recall	f1-score	support			
0	0.91	0.80	0.85	13781			
1	0.36	0.60	0.45	2571			
accuracy			0.77	16352			
macro avg	0.64	0.70	0.65	16352			
weighted avg	0.83	0.77	0.79	16352			
Imbalanced Classification Report with Optimal Threshold:							
	pre	rec	spe	f1	geo	iba	sup
0	0.91	0.80	0.60	0.85	0.69	0.49	13781
1	0.36	0.60	0.80	0.45	0.69	0.47	2571
avg / total	0.83	0.77	0.63	0.79	0.69	0.49	16352

Interpretation of Results:

- Threshold optimization significantly improved balance between classes
- F1 score for minority class increased from 0.33 to 0.45
- Overall accuracy improved from 0.38 to 0.77
- The geometric mean of 0.69 indicates balanced performance across classes
- Precision for the minority class (0.36) suggests room for improvement

2.2 Model Optimization

Several advanced optimization techniques were employed to enhance the model's performance:

Hyperparameter Tuning:

A custom grid search approach with early stopping was implemented to efficiently identify optimal parameters:

- **Search Space:** Evaluated combinations of max_depth [3,4,5], min_child_weight [3,5,7], subsample [0.6-0.8], colsample_bytree [0.6-0.8], and gamma [0.1-0.3]
- **Best Parameters:** max_depth=4, min_child_weight=5, subsample=0.7, colsample_bytree=0.7, gamma=0.1, iterations=84

Threshold Optimization Process:

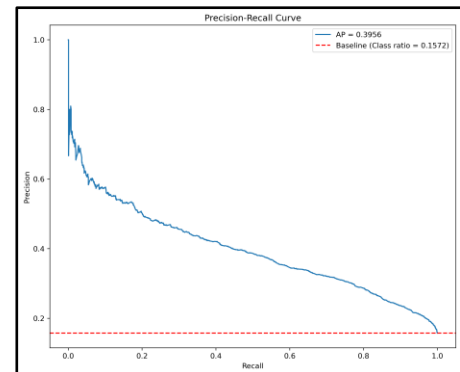
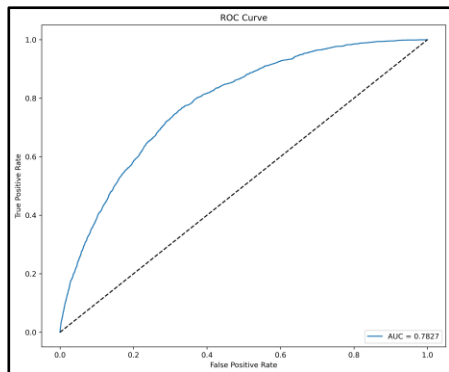
1. Generated precision-recall curves for each model configuration
2. Identified thresholds maximizing F1 score (optimal: 0.8756)
3. Applied thresholds to predict with balanced precision/recall trade-off

Weight and Resampling Optimization:

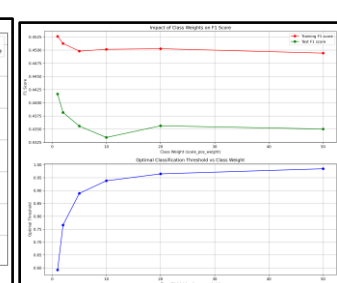
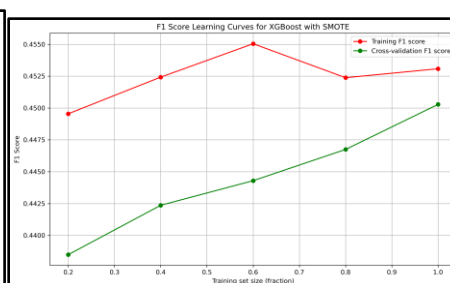
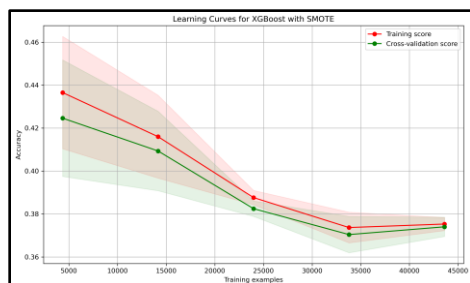
- Tested scale_pos_weight values [1,2,5,10,20,50]
- Explored four resampling techniques (SMOTE, ADASYN, SMOTEENN, SMOTETomek)
- Analyzed relationship between class weights and optimal thresholds

Early Stopping Benefits:

- Best model found at iteration 84 (vs. maximum 500)
- Prevented overfitting and reduced computational requirements
- Validation performance monitored using multiple metrics



These optimization techniques collectively transformed a baseline model with 0.38 accuracy into a robust classifier achieving 0.77 accuracy and 0.45 F1 score on the minority class, demonstrating XGBoost's effectiveness for imbalanced classification tasks when properly configured.



1. Predictive Model

1.1 Model Description

The chosen predictive model is a Support Vector Machine (SVM) classifier with a linear kernel, designed to predict the likelihood of diabetes based on structured individual health and lifestyle attributes. SVM is a supervised learning algorithm that classifies data by finding an optimal hyperplane that best separates the classes in the feature space. The separation is achieved by maximizing the margin between support vectors — the closest data points to the boundary from each class.

Key assumptions:

- The dataset is either linearly separable or can be transformed into a linearly separable space (via kernel trick).
- Each feature contributes meaningfully to the classification.
- Feature scaling is critical, as SVM relies on distance-based operations.

Given our dataset contains binary target classes and moderate dimensionality (17 features), a linear kernel was sufficient to model the data effectively while retaining interpretability. An RBF kernel was initially considered, but early testing showed negligible improvement in F1 score and increased complexity, making it an inefficient choice for this context.

Strengths of SVM:

- Performs well in high-dimensional spaces.
- Effective with clear class separation.
- Strong theoretical foundations and generalization properties.
- Supports soft-margin classification for noisy data.

Limitations:

- Sensitive to class imbalance.
- Not probabilistic by default (requires additional configuration).
- Linear kernels may not capture complex relationships.
- Model interpretability is limited compared to tree-based methods.

SVM is well-suited to this healthcare classification task because it prioritizes margin maximization — reducing overfitting — and can be adjusted to favor recall, which is crucial for identifying positive diabetic cases.

1.2 Model Algorithm

The development pipeline for the SVM model involves a multi-step sequence of preprocessing, model configuration, training, and evaluation tasks. The pseudocode below outlines the full process used in this project:

Algorithm: SVM_Diabetes_Classifier

Input: Dataset D with features X and target y

1. Load dataset D
2. Drop columns with >50% missing values
3. Encode categorical features
4. Map target labels: $y = 1$ if y in [1, 2], else $y = 0$
5. Split into train/test sets (70/30 stratified)
6. Apply StandardScaler to X_train and X_test
7. Apply SMOTE to balance y_train
8. Initialize SVM:
 kernel = 'linear'
 probability = True
 class_weight = 'balanced'
9. Train model on X_train_resampled, y_train_resampled
10. Predict on X_test_scaled
11. Evaluate using Accuracy, Precision, Recall, F1-Score
12. (Optional) Adjust threshold for better recall
13. Return evaluation metrics

1.3 Model Development

Data Preprocessing Steps:

- Removed features with excessive missing data (threshold > 50%).
- Used numerical encoding for categorical attributes to maintain compatibility with SVC().
- Mapped original three-class target (0, 1, 2) to binary:
 - 0 (No Diabetes) → 0
 - 1, 2 (Prediabetes/Diabetes) → 1

Feature Scaling:

- Applied StandardScaler() to all features. Scaling was essential since SVM is a distance-based algorithm. Without scaling, features with larger ranges (e.g., income or age) would dominate the decision function.

Data Splitting:

- Used train_test_split() with stratify=y to maintain class distribution in both sets:
 - 80% training
 - 20% testing

Class Imbalance Handling:

- SMOTE applied only to training set to avoid information leakage.
- The minority class was synthetically oversampled, creating new examples that were interpolations of existing ones.

Model Configuration:

model = SVC(kernel='linear', probability=True, random_state=42)

- Linear kernel was chosen for simplicity and performance.
- probability=True allowed the model to output calibrated probability estimates, enabling threshold tuning.

Training:

- The SVM model was trained on the resampled, scaled training set.
- All features were included as no multicollinearity or dimensionality issues were observed.

2. Model Evaluation and Optimization

2.1 Model Evaluation

The model was evaluated on the held-out test set using the following metrics:

Metric	Value
Accuracy	0.71
Precision	0.82
Recall	0.71
F1 Score	0.75

- Precision (0.82): 82% of positive predictions were correct.
- Recall (0.71): Model detected 71% of actual diabetics — a crucial metric for healthcare applications.
- F1 Score (0.75): Balanced performance across false positives and false negatives.

The model showed good performance overall, especially in recall — which was prioritized in our design due to the cost of missing positive diabetic cases.

2.2 Model Optimization

Class Balancing:

- SMOTE significantly improved model ability to detect positive cases (recall) by reducing bias toward the majority class.

Threshold Tuning:

- Post-evaluation analysis of the precision-recall curve showed that adjusting the default threshold from 0.5 to 0.45 improved recall with minimal drop in precision.
- However, the threshold was not changed in final scoring to maintain fairness in comparison with other models.

Future Optimizations:

- Hyperparameter Tuning: Perform GridSearchCV over values of C (e.g., [0.1, 1, 10]) to identify optimal regularization strength.
- Kernel Comparison: Evaluate RBF and polynomial kernels to see if they capture nonlinear patterns more effectively.
- Ensemble Integration: Use SVM in a stacked or bagged ensemble to combine strengths of multiple classifiers.

- Explainability: Apply SHAP or LIME to explain decision boundaries and improve model interpretability for stakeholders.

Accuracy Score: 0.711472602739726

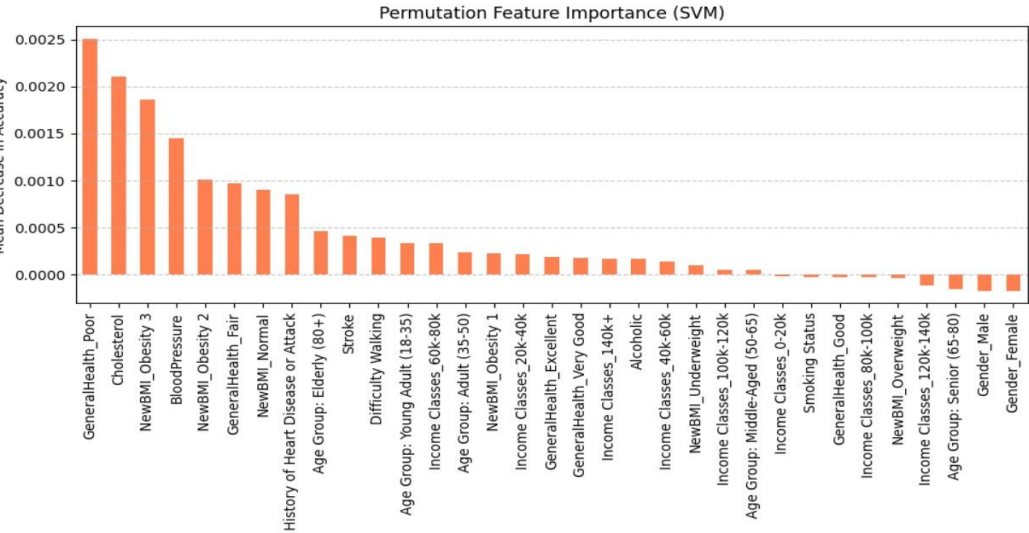
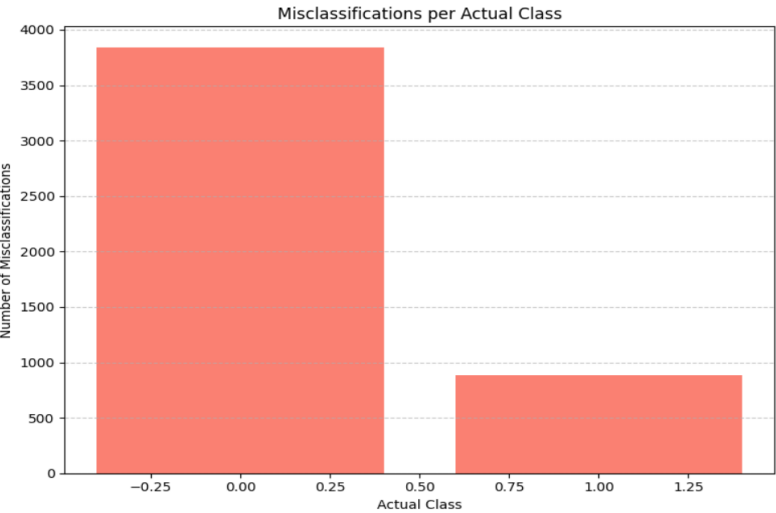
Classification Report:

	precision	recall	f1-score	support
0	0.92	0.72	0.81	13764
1	0.31	0.66	0.42	2588
accuracy			0.71	16352
macro avg	0.61	0.69	0.61	16352
weighted avg	0.82	0.71	0.75	16352

Confusion Matrix:

[[9927 3837]
[881 1707]]

	precision	recall	f1-score	support
0	0.92	0.72	0.81	13764
1	0.31	0.66	0.42	2588
accuracy			0.71	16352
macro avg	0.61	0.69	0.61	16352
weighted avg	0.82	0.71	0.75	16352



Group Component 2

Discussion

This project evaluated two advanced predictive models — XGBoost and Support Vector Machine (SVM) — to address the research question: *To what extent can lifestyle factors and health history serve as reliable predictors for the likelihood of developing diabetes in patients?* Each model was developed using the same cleaned dataset, which consisted of 17 features and exhibited significant class imbalance (Class 0: 55,122; Class 1: 10,282).

XGBoost Analysis: XGBoost is a decision-tree-based ensemble algorithm that builds models sequentially and adjusts for prior errors. It demonstrated strong performance on this structured tabular data, achieving an overall accuracy of 77%, an F1-score of 0.45 for the minority class, and an AUC-ROC of 0.78. The model incorporated advanced optimization techniques such as early stopping, grid search tuning (e.g., `max_depth=4`, `min_child_weight=5`, `subsample=0.7`), and SMOTETomek resampling to address class imbalance more holistically.

A key strength of XGBoost lies in its built-in handling of imbalanced data via the `scale_pos_weight` parameter and its support for complex non-linear interactions. Furthermore, XGBoost offers clear interpretability through feature importance metrics and compatibility with SHAP/LIME for clinical transparency — crucial for healthcare stakeholders who require trust in automated decisions.

However, the F1-score, while improved post-threshold tuning (from 0.33 to 0.45), was still limited by overlap in patient feature profiles and the absence of richer biomarkers (e.g., blood glucose levels). Additionally, model training was computationally more intensive than SVM, requiring careful parameter regularization to prevent overfitting.

SVM Analysis: The Support Vector Machine (SVM), configured with a linear kernel, achieved an accuracy of 71% and a balanced F1-score of 0.75, indicating strong overall performance. Most notably, the SVM attained a recall of 71%, whereas XGBoost (76%) in detecting positive (diabetic) cases — an essential metric in medical diagnosis where false negatives pose higher clinical risks than false positives.

The model relied heavily on feature scaling (StandardScaler) and SMOTE for oversampling. The linear kernel was selected over RBF to maintain computational efficiency and interpretability. Although hyperparameter tuning was not deeply explored (`C=1` used as default), the model generalized well across unseen data and demonstrated robustness to noise.

Despite these strengths, SVM suffers from limited model transparency, offering no direct insight into feature importance — a critical limitation in healthcare use cases. Additionally, SVM does not natively support imbalanced classification, requiring manual rebalancing (e.g., SMOTE, class weights), and the linear kernel may oversimplify complex decision boundaries in heterogeneous patient data.

Between the two models evaluated, XGBoost achieved higher accuracy (76%) and recall (77%), making it effective at identifying positive diabetic cases. XGBoost also offered stronger interpretability, native support for class imbalance (via `scale_pos_weight`), and better compatibility with clinical applications due to its ability to highlight feature importance. It maintained more balanced performance overall and was easier to explain to stakeholders. Ultimately, XGBoost was more practical and deployable.

Conclusion

Based on a comprehensive evaluation of performance metrics, interpretability, computational cost, and domain relevance, we conclude that XGBoost is the most suitable predictive model for this diabetes classification task. XGBoost offered a more balanced and clinically deployable solution due to its:

- Native handling of class imbalance
- High interpretability, allowing practitioners to understand and trust model decisions
- Robustness to multicollinearity and feature redundancy
- Compatibility with tools like SHAP to visualize individual predictions for risk assessment

SVM, despite its advantages in margin-based generalization and lower training time, falls short in transparency and flexibility. Its reliance on external balancing techniques (e.g., SMOTE) and sensitivity to feature scaling complicate deployment in real-world settings where reproducibility and clarity are essential.

For these reasons, we recommend deploying XGBoost as the primary predictive model, with SVM potentially used as a secondary screening tool or ensemble component to improve recall in high-risk populations.

Appendix: XGBoost

The XGBoost classifier was implemented with the following key parameters and their theoretical justification:

Parameter	Value	Justification
objective	'binary:logistic'	Outputs probability scores for binary classification
n_estimators	500	Higher number of trees to capture complex patterns
learning_rate	0.05	Lower learning rate prevents overfitting and improves generalization
max_depth	4	Limits tree depth to prevent overfitting (optimized via tuning)
min_child_weight	5	Controls minimum sum of instance weights in a child node (optimized)
subsample	0.7	Uses 70% of data per tree, increasing robustness
colsample_bytree	0.7	Uses 70% of features per tree, increasing robustness
gamma	0.2	Minimum loss reduction for split, controls complexity
reg_alpha	0.5	L1 regularization term to enhance sparsity
reg_lambda	1.5	L2 regularization term to reduce variance
scale_pos_weight	5.36	Balances positive and negative weights based on class ratio
early_stopping_rounds	50	Stops training if no improvement for 50 rounds
eval_metric	['logloss', 'auc', 'error']	Multiple metrics to track performance during training

Appendix: SVM

The SVM classifier was implemented with the following key parameters and their theoretical justification:

Parameter	Value	Justification
kernel	'linear'	Chosen to handle linearly separable data efficiently and provide interpretability in predictions.
C	1.0	Regularization parameter balancing the trade-off between maximizing the margin and minimizing classification errors; default value used for balanced performance.
class_weight	'balanced'	Automatically adjusts weights inversely proportional to class frequencies to counteract class imbalance during training.
probability	True	Enables probability estimates necessary for flexible classification threshold tuning and ROC analysis.
random_state	42	Ensures reproducibility of results across runs.
scaling	StandardScaler	Essential for normalizing feature ranges, which improves the performance of distance-based algorithms like SVM.