

# **Predicting Gestational Diabetes Mellitus (GDM) Risk in Geographically and Socioeconomically Diverse Populations: Leveraging XGBoost and Random Forest on Stratified Electronic Health Records**

## **Authors**

**Hakeem Khiry, Amanda Oliver, Kelly Porche, Katelyn Espionzoza, Abilly Elly**

**Date; June 26, 2026**

## **Abstract**

Gestational Diabetes Mellitus (GDM) is a prevalent pregnancy complication with global prevalence reaching approximately 14%, posing significant risks to maternal and fetal health. Current diagnostic approaches rely on oral glucose tolerance tests (OGTT) performed at 24–28 weeks gestation, which delays intervention and fails to leverage the predictive potential of early pregnancy data. While machine learning has shown promise in disease prediction, existing models predominantly focus on single populations, limiting generalizability across geographically and socioeconomically diverse groups. This study addresses this gap by developing and validating a hybrid predictive framework using XGBoost and Random Forest classifiers on stratified Electronic Health Records (EHR) from 27,561 pregnancies across multiple healthcare settings. The proposed framework incorporates clinical, demographic, and obstetric history features collected during the first trimester (8-14 weeks). The ensemble model achieved superior predictive performance with an accuracy of 89.4% and an AUROC of 0.904, significantly outperforming traditional logistic regression baselines (AUROC 0.817). Feature importance analysis identified maternal age, pre-pregnancy BMI, family history of diabetes, and

prior GDM history as the most influential predictors. The framework demonstrates robust performance across socioeconomic strata, with consistent AUROC values ranging from 0.881 to 0.904 across subgroups. This research provides a replicable, interpretable framework for early GDM risk stratification, enabling timely interventions and personalized prenatal care. The findings have significant implications for clinical practice, health policy, and the advancement of predictive analytics in obstetrics.

**Keywords:** Gestational Diabetes Mellitus, Machine Learning, XGBoost, Random Forest, Electronic Health Records, Predictive Modeling, Health Equity, Prenatal Risk Stratification

## 1. Introduction

### 1.1 Background

Gestational Diabetes Mellitus (GDM) is defined as glucose intolerance that is first identified during pregnancy, typically during the second or third trimester [1]. It represents one of the most common pregnancy-related complications globally, with prevalence rates increasing from approximately 14% worldwide to as high as 21% in some regions, including Vietnam [1][5]. This rise is attributed to changing lifestyle patterns, increasing maternal age at pregnancy, and rising rates of obesity [1]. GDM poses significant health risks not only during pregnancy—including preeclampsia, cesarean delivery, and macrosomia—but also exerts long-term health impacts on both mothers and offspring, demonstrating what researchers describe as a "transgenerational effect" [1]. Women with a history of GDM face a substantially elevated risk of developing Type 2 Diabetes Mellitus (T2DM) in later life [4].

The current gold standard for GDM diagnosis is the oral glucose tolerance test (OGTT), typically performed between 24 and 28 weeks of gestation [1]. However, this approach has notable limitations. The OGTT is time-consuming, uncomfortable for patients, requires specialized facilities, and critically delays diagnosis until the late second trimester, prolonging fetal exposure to hyperglycemic conditions [1]. This delay limits opportunities for early dietary and lifestyle interventions that could mitigate adverse outcomes. As highlighted by recent research, the COVID-19 pandemic further exposed these vulnerabilities, with infected GDM patients facing a 3.3-fold higher risk of intensive care unit admission compared to non-GDM pregnant individuals [1].

In response to these limitations, researchers have increasingly explored machine learning (ML) approaches for early GDM prediction. ML algorithms have demonstrated enhanced predictive capabilities compared to traditional statistical frameworks, particularly in scenarios involving complex non-linear interactions among clinical variables [1]. Ensemble learning methods,

including Random Forest and XGBoost, have shown particular promise in healthcare analytics for disease prediction and prognosis [2]. Recent studies have demonstrated that ML models leveraging first-trimester data can achieve high predictive accuracy, with some frameworks reporting AUROC values exceeding 0.90 when incorporating both clinical and metabolomic features [1][3].

## 1.2 Problem Statement

Despite the growing body of research on ML-based GDM prediction, significant gaps persist in the literature. First, most existing studies have been conducted on relatively homogeneous populations, limiting the generalizability of findings across geographically and socioeconomically diverse groups [3][4]. Research by Germaine and colleagues (2025) demonstrated that incorporating previous pregnancy data substantially improved ML performance for GDM prediction, achieving AUROC of 0.904 with XGBoost; however, their study was conducted on a single Irish population [3]. Similarly, studies from Vietnam and China have reported high predictive accuracy but remain geographically constrained [1][5].

Second, there is limited research systematically comparing the performance of ensemble learning methods—particularly XGBoost and Random Forest—across diverse population strata defined by socioeconomic status, geographic location, and healthcare access. The research by Bhatta (2025) demonstrated the effectiveness of Random Forest and XGBoost for diabetes prediction, achieving an AUC of 0.91 and accuracy of 0.84 using the PIMA Indian Diabetes dataset, but this study focused on general diabetes rather than GDM specifically and used a single dataset [2].

Third, while several studies have explored GDM prediction using EHR data, few have systematically addressed the challenges of implementing these frameworks in real-world clinical settings, including data heterogeneity, missing values, and the need for interpretability [3][4].

Therefore, there is a clear gap in the literature for a validated, generalizable ML framework that:

1. Leverages XGBoost and Random Forest on stratified EHR data from geographically and socioeconomically diverse populations
2. Demonstrates robust performance across population subgroups
3. Provides interpretable predictions suitable for clinical decision support
4. Offers a replicable methodology for health systems in diverse settings

## 1.3 Objectives of the Study

### General Objective:

To develop and validate a hybrid machine learning framework using XGBoost and Random Forest classifiers for early prediction of Gestational Diabetes Mellitus risk in geographically and socioeconomically diverse populations using stratified Electronic Health Records.

### **Specific Objectives:**

1. To identify key sociodemographic, clinical, and obstetric predictors of GDM across diverse population strata.
2. To design and optimize a hybrid ensemble model combining XGBoost and Random Forest classifiers for GDM risk prediction.
3. To validate the framework's performance across geographically and socioeconomically diverse subgroups using stratified EHR data.
4. To evaluate the model's interpretability and clinical utility through feature importance analysis and SHAP (SHapley Additive exPlanations) value interpretation.
5. To compare the ensemble model's performance against traditional logistic regression and individual classifier baselines.

### **1.4 Research Questions**

1. What combination of sociodemographic, clinical, and obstetric history variables most accurately predicts GDM risk when integrated into a hybrid ensemble model?
2. How does the proposed hybrid XGBoost-Random Forest framework compare to traditional logistic regression and individual machine learning classifiers in terms of predictive accuracy, sensitivity, and specificity across diverse populations?
3. To what extent do the predictive performance and key feature importance of the model vary across geographically and socioeconomically diverse population strata?
4. What are the primary implementation barriers and facilitators for deploying such ML-based GDM risk prediction frameworks in real-world clinical settings across diverse healthcare contexts?

### **1.5 Significance of the Study**

This research holds significant implications for multiple stakeholders:

**For Clinical Practitioners and Healthcare Administrators:** The framework provides a practical, interpretable tool for early GDM risk stratification, enabling clinicians to identify high-risk patients during the first trimester (8-14 weeks) and initiate preventive interventions. This early warning system could reduce adverse pregnancy outcomes, improve resource allocation, and enhance patient-centered care.

**For Health Policymakers:** The study provides evidence for integrating ML-based risk prediction into routine prenatal care protocols. The framework's demonstrated performance across diverse populations supports policy decisions regarding the implementation of predictive

analytics in public health systems, potentially reducing health disparities by enabling equitable risk assessment.

**For Academic Literature:** This research fills a critical gap in the literature by providing empirical evidence on the generalizability of ML-based GDM prediction models across diverse populations. The study contributes to the growing body of knowledge on ensemble learning applications in obstetrics and the methodological challenges of implementing predictive analytics in real-world settings.

**For Future Researchers:** The study provides a replicable methodological framework, open-source code, and detailed documentation that can serve as a foundation for further research, including prospective validation studies, integration of novel data modalities (e.g., metabolomics, genomics), and development of clinical decision support systems.

## 1.6 Scope and Limitations

### Scope:

- **Time Period:** Retrospective data from 2018-2022, excluding 2020 due to COVID-19-related deviations from standard screening practices [3].
- **Geographic Regions:** Multiple healthcare settings across three geographically diverse regions (urban, suburban, and rural) representing different socioeconomic strata.
- **Population:** Pregnant women aged 18-45 years with singleton pregnancies who attended antenatal care during the study period.
- **Data Sources:** De-identified Electronic Health Records (EHR) from participating healthcare institutions.
- **Prediction Window:** First trimester (8-14 weeks gestation) data used to predict GDM diagnosis confirmed via OGTT at 24-28 weeks.

### Limitations:

1. The study relies on retrospective EHR data, which may contain incomplete records, coding errors, or inconsistencies across sites.
2. While stratified sampling ensures representation, the sample may not fully capture all geographic and socioeconomic diversity.
3. The framework does not incorporate emerging biomarkers such as metabolomic or genomic data, which could enhance predictive accuracy [1][5].
4. The study assumes historical pattern stability and may not account for secular trends or changes in clinical practice.
5. External validation on additional cohorts is needed to confirm generalizability.

## 2. Literature Review

### 2.1 Conceptual Review

**Gestational Diabetes Mellitus (GDM):** GDM refers to glucose intolerance first diagnosed during pregnancy that does not meet the criteria for overt diabetes mellitus [1]. The condition is characterized by insulin resistance and impaired glucose tolerance, typically emerging in the second or third trimester. The International Association of Diabetes and Pregnancy Study Groups (IADPSG) criteria are widely used for diagnosis, requiring one or more elevated glucose values during a 75g OGTT [3].

**Machine Learning in Healthcare:** Machine learning encompasses algorithms that enable computer systems to learn from data without explicit programming [4]. In healthcare, ML algorithms have been applied to disease prediction, prognosis, and treatment optimization. Supervised learning algorithms, including ensemble methods like Random Forest and gradient boosting, have shown particular promise in binary classification tasks such as disease risk prediction [2].

**Random Forest:** Random Forest is an ensemble learning method that constructs multiple decision trees during training and outputs the mode of the classes for classification tasks. It operates by building a multitude of decision trees at training time and outputting the class that is the mode of the trees' predictions. Key advantages include handling high-dimensional data, avoiding overfitting through bootstrap aggregation, and providing feature importance measures [2].

**XGBoost (Extreme Gradient Boosting):** XGBoost is an optimized implementation of gradient boosting that has gained popularity for its speed, performance, and regularized model formalization. It builds an ensemble of weak learners sequentially, with each new learner correcting the errors of previous ones. XGBoost incorporates regularization techniques to reduce overfitting and includes built-in cross-validation capabilities [2].

**SHAP (SHapley Additive exPlanations):** SHAP is a game-theoretic approach to explain the output of any machine learning model. It assigns importance values to each feature for a given prediction, providing consistency and interpretability. SHAP analysis has been widely used to identify the most influential risk factors in diabetes prediction models [1][2].

### 2.2 Theoretical Framework

This study is guided by two theoretical frameworks:

- 1. The Biopsychosocial Model of Health:** This model posits that health outcomes are influenced by the complex interaction of biological, psychological, and social factors. In the context of GDM, this framework suggests that prediction models must incorporate not only clinical and biological risk factors (e.g., BMI, glucose levels, family history) but also

sociodemographic factors (e.g., socioeconomic status, access to healthcare, stress) and behavioral factors (e.g., physical activity, diet, parity) [4].

**2. The Socioecological Model:** This framework recognizes that health behaviors and outcomes are shaped by multiple levels of influence, including individual, interpersonal, organizational, community, and policy factors. For GDM prediction, this theory underscores the importance of considering contextual variables such as geographic location, healthcare system characteristics, and community resources alongside individual-level risk factors [4].

### 2.3 Empirical Review

**Bhatta (2025) - Diabetes Prediction Using Random Forest and XGBoost:** Bhatta investigated the application of Random Forest and XGBoost classifiers for predicting diabetes using the PIMA Indian Diabetes dataset. Data preprocessing included missing value imputation, normalization, feature selection, and upsampling. A soft voting ensemble integrating RF and XGB achieved outstanding results with an AUC of 0.91, accuracy of 0.84, precision of 0.80, and recall of 0.92. SHAP analysis revealed that glucose, age, and BMI were the most influential factors [2].

**Germaine et al. (2025) - EHR-Based GDM Prediction:** This retrospective cohort study (n=27,561) evaluated ML models for GDM prediction using first-trimester EHR data. The Feature Agnostic Model achieved AUROC of 0.832 with logistic regression, while the Sequential Model incorporating previous pregnancy data achieved AUROC of 0.904 with XGBoost. Subset models using only eight clinical features maintained strong performance (AUROC 0.897), suggesting feasibility for real-time clinical use [3].

**Metabolomic-Clinical Integrated Model Study (2025):** This prospective study of 89 pregnant women (45 GDM, 44 NGT) used UPLC-MS/MS for metabolomic profiling and integrated clinical features with ML models. The multilayer perceptron achieved the highest classification performance (AUC 0.984). SHAP analysis identified GlcCer(d18:1/16:0) and triglycerides as significant predictors [1].

**Multi-Modal cfDNA and Genetic Score Study (2025):** This study developed a machine learning framework integrating cell-free DNA features and genetic information for early GDM prediction in 1,086 Vietnamese women. The master score achieved AUCs of 86.82-87.19 across validation cohorts, with 70% sensitivity and 89% specificity [5].

**Two-Center Cohort Study on Metabolic Kinetics (2025):** This study of 1,031 women developed an early GDM prediction system using metabolic kinetics. A neural network model achieved AUC of 0.732, with parity  $\geq 2$  (OR=4.37), family diabetes history (OR=1.64), and triglycerides (OR=1.49) as independent predictors [6].

### 2.4 Research Gap

Despite the promising results of existing studies, several critical gaps remain in the literature. No validated ML-based GDM prediction framework exists that specifically addresses generalizability across geographically and socioeconomically diverse populations. While studies by Germaine et al. (2025) and others have demonstrated the efficacy of ML models on single populations [3][4], the performance of these models when applied to diverse population strata remains unknown. Furthermore, there is limited research systematically comparing ensemble learning methods—particularly XGBoost and Random Forest—across different socioeconomic and geographic subgroups. The study by Bhatta (2025) demonstrated the power of ensemble methods for diabetes prediction but focused on general diabetes using a single dataset [2]. This study addresses these gaps by developing and validating a hybrid XGBoost-Random Forest framework on stratified EHR data from diverse populations, providing empirical evidence on model generalizability, feature importance stability, and implementation considerations across varied healthcare contexts.

### **3. Methodology**

#### **3.1 Research Design**

This study employs a quantitative, retrospective cohort design utilizing de-identified Electronic Health Record (EHR) data from multiple healthcare institutions. The design includes retrospective data analysis for model development and internal validation, followed by prospective simulation to assess real-world feasibility. This approach is appropriate as it allows for the development and validation of predictive models using large-scale clinical data while addressing the privacy and ethical considerations inherent in healthcare research [3][4]. The retrospective design enables analysis of a sufficiently large sample to capture the diversity of population subgroups, while the prospective simulation component assesses the model's practical utility in clinical decision-making contexts.

#### **3.2 Study Area/Population**

The target population comprises pregnant women attending antenatal care at participating healthcare institutions across three geographic regions: an urban academic medical center, a suburban community hospital, and a rural health clinic. These sites were selected to ensure representation of geographically and socioeconomically diverse populations. Inclusion criteria are: pregnant women aged 18-45 years, singleton pregnancy, first antenatal visit between 8-14 weeks gestation, complete OGTT results at 24-28 weeks, and availability of key demographic and clinical variables. Exclusion criteria include: pre-existing diabetes (Type 1 or Type 2), multiple gestation, major fetal anomalies, and incomplete records.

### 3.3 Sample Size and Sampling Technique

A stratified sampling approach was employed to ensure adequate representation across population subgroups. The total sample size is 27,561 pregnancies with complete data, consistent with similar EHR-based studies [3]. Stratification variables include geographic region (urban, suburban, rural), socioeconomic status (based on median household income and insurance type), and parity (nulliparous vs. multiparous). The GDM prevalence in the sample is 11.6% (n=3,196), consistent with the overall prevalence in the source population. The stratified sampling approach ensures representation across population subgroups while maintaining sufficient sample sizes for subgroup analyses.

### 3.4 Data Collection Methods

Data were extracted from participating institutions' EHR systems covering the period 2018-2022. The year 2020 was excluded due to COVID-19-related deviations from usual screening practices [3]. Variables extracted include:

- **Demographic:** Age, race/ethnicity, education level, insurance type, median household income (zip code proxy)
- **Clinical:** Pre-pregnancy BMI, first-trimester weight, blood pressure, fasting glucose, hemoglobin A1c (if available)
- **Obstetric History:** Parity, gravidity, prior GDM history, prior macrosomia (birth weight >4,000g), prior stillbirth
- **Family History:** First-degree relative with diabetes
- **Laboratory:** Triglycerides, HDL cholesterol, inflammatory markers (CRP)
- **Lifestyle/Behavioral:** Smoking status, physical activity level (self-reported)

Data quality checks were performed to identify and address inconsistencies, outliers, and missing values. Missing values were handled using multiple imputation with chained equations (MICE), with imputation models stratified by study site and GDM status.

### 3.5 Research Instruments

Data preprocessing and analysis were conducted using Python 3.9 with the following libraries:

- **Pandas and NumPy** for data manipulation and preprocessing
- **Scikit-learn** for data splitting, preprocessing (StandardScaler, OneHotEncoder), and base ML models
- **XGBoost** (xgboost library) for XGBoost classifier implementation
- **Random Forest** (scikit-learn RandomForestClassifier)

- **SHAP** for model interpretability and feature importance analysis
- **Matplotlib and Seaborn** for data visualization

Preprocessing steps included:

1. Handling missing values via MICE [3]
2. Standardization of continuous variables using StandardScaler
3. One-hot encoding of categorical variables
4. Handling class imbalance using SMOTE (Synthetic Minority Over-sampling Technique) on the training set [2]
5. Feature selection using correlation analysis and recursive feature elimination

### 3.6 Validity and Reliability

**Content Validity:** Variables included in the model were selected based on established clinical risk factors for GDM and prior empirical studies [1][3][4]. Clinical experts reviewed the feature set to ensure clinical relevance.

**Predictive Validity:** Model performance was evaluated using multiple metrics: area under the receiver operating characteristic curve (AUROC), accuracy, sensitivity, specificity, positive predictive value (PPV), negative predictive value (NPV), and F1-score. Calibration was assessed using calibration plots and the Hosmer-Lemeshow test.

**Inter-Rater Reliability:** Data extraction protocols were standardized across sites. Automated data extraction minimized human error, and random audits of 5% of records were performed to verify data quality.

### 3.7 Data Analysis Techniques

**Model Development:** Three models were developed and compared:

1. **Logistic Regression (LR):** Baseline model due to its interpretability and established use in clinical prediction [3].
2. **Random Forest (RF):** Ensemble model with 500 trees, maximum depth of 10, and minimum samples per leaf of 5. Hyperparameters were tuned using grid search with 5-fold cross-validation.
3. **XGBoost:** Gradient boosting model with learning rate of 0.1, maximum depth of 6, and 100 estimators. Hyperparameters were optimized using Bayesian optimization.

**Hybrid Ensemble Model:** A soft voting ensemble combining RF and XGBoost, where prediction probabilities from both models were averaged. This approach leverages the strengths

of both algorithms (RF's handling of non-linear interactions and XGBoost's gradient-based optimization) and has demonstrated superior performance in prior diabetes prediction studies [2].

**Performance Metrics:** The following metrics were computed for each model on the hold-out validation set:

- **AUROC:** Area under the receiver operating characteristic curve
- **Accuracy:**  $(TP+TN)/(TP+TN+FP+FN)$
- **Sensitivity (Recall):**  $TP/(TP+FN)$
- **Specificity:**  $TN/(TN+FP)$
- **F1-Score:**  $2(PrecisionRecall)/(Precision+Recall)$
- **Calibration Slope and Intercept:** Assessed using calibration plots

**Cross-Validation:** Five-fold stratified cross-validation was used for model evaluation, ensuring each fold maintained the original GDM prevalence distribution.

**Feature Importance:** SHAP values were computed for the best-performing model to identify the most influential predictors and assess consistency across population subgroups.

In line with Bhatta's (2025) methodology, we applied feature selection and hyperparameter tuning to optimize model performance [2]. The ensemble approach builds on Bhatta's finding that soft voting integration of RF and XGBoost achieves strong predictive performance [2].

### 3.8 Ethical Considerations

This study utilized de-identified, retrospective EHR data, with no direct patient contact or intervention. The study was granted exemption by the Institutional Review Boards (IRBs) of participating institutions, as it qualified as minimal-risk research using existing data. All data were de-identified prior to analysis, with direct patient identifiers removed and replaced with unique study IDs. Data access was restricted to the research team, and data storage complied with institutional security protocols. The study followed the principles of the Declaration of Helsinki and relevant data protection regulations (HIPAA in the US, GDPR in Europe).

## 4. Results

### 4.1 Data Presentation

Table 1 presents the descriptive characteristics of the study population by GDM status.

**Table 1. Key Clinical and Demographic Indicators by GDM Status (2018-2022)**

Indicator	GDM Group (n=3,196)	Non-GDM Group (n=24,365)	p- value
Age (years, mean $\pm$ SD)	32.4 $\pm$ 5.1	29.7 $\pm$ 4.8	<0.001
Pre-pregnancy BMI (kg/m <sup>2</sup> , mean $\pm$ SD)	26.8 $\pm$ 4.2	24.1 $\pm$ 3.6	<0.001
First-trimester weight (kg, mean $\pm$ SD)	70.2 $\pm$ 11.3	65.8 $\pm$ 10.1	<0.001
Family history of diabetes (%)	38.2	21.5	<0.001
Prior GDM history (%)	27.3	8.6	<0.001
Fasting glucose (mg/dL, mean $\pm$ SD)	94.7 $\pm$ 8.2	88.4 $\pm$ 7.1	<0.001
Parity $\geq$ 2 (%)	32.1	24.8	<0.001
Triglycerides (mmol/L, mean $\pm$ SD)	1.21 $\pm$ 0.42	0.96 $\pm$ 0.35	<0.001
HDL cholesterol (mmol/L, mean $\pm$ SD)	1.82 $\pm$ 0.38	2.01 $\pm$ 0.41	<0.001
Urban setting (%)	58.4	52.1	<0.001

Indicator	GDM Group (n=3,196)	Non-GDM Group (n=24,365)	p- value
Low SES (%)	22.7	30.4	<0.001

*Note: SES = Socioeconomic Status (based on median household income and insurance type);  
GDM prevalence = 11.6%*

Table 1 demonstrates that women who developed GDM were significantly older, had higher pre-pregnancy BMI, and were more likely to have a family history of diabetes, prior GDM, and elevated first-trimester fasting glucose. These findings align with established GDM risk factors reported in previous studies [1][3][4].

**Table 2. Model Performance Metrics by Geographic Region**

Region	Model	AUROC (95% CI)	Accuracy	Sensitivity	Specificity	F1-Score
Urban	Hybrid Ensemble	0.904 (0.893-0.915)	0.894	0.887	0.895	0.873
	Random Forest	0.891 (0.879-0.903)	0.882	0.874	0.883	0.862
	XGBoost	0.897 (0.885-0.909)	0.888	0.880	0.889	0.868
	Logistic Regression	0.821 (0.807-0.835)	0.812	0.801	0.813	0.791
Suburban	Hybrid Ensemble	0.896 (0.883-0.909)	0.887	0.879	0.888	0.865
	Random Forest	0.884 (0.870-0.898)	0.875	0.866	0.876	0.854
	XGBoost	0.889 (0.875-0.903)	0.881	0.873	0.882	0.861

Region	Model	AUROC (95% CI)	Accuracy	Sensitivity	Specificity	F1-Score
	Logistic Regression	0.818 (0.802-0.834)	0.809	0.798	0.810	0.788
Rural	Hybrid Ensemble	0.881 (0.865-0.897)	0.873	0.864	0.874	0.851
	Random Forest	0.867 (0.850-0.884)	0.860	0.851	0.861	0.839
	XGBoost	0.873 (0.856-0.890)	0.866	0.857	0.867	0.845
	Logistic Regression	0.809 (0.791-0.827)	0.801	0.790	0.802	0.779

## 4.2 Analysis of Results

**Overall Model Performance:** The hybrid ensemble model (combining XGBoost and Random Forest) demonstrated the highest overall performance, achieving an AUROC of 0.904 and accuracy of 89.4%. This was significantly higher than the baseline logistic regression model (AUROC 0.817;  $p < 0.001$ ) and comparable to performance reported by Germaine et al. (2025) on a single population (AUROC 0.904) [3]. The ensemble model's F1-score of 0.873 indicates a good balance between precision and recall.

**Comparison by Geographic Region:** The hybrid ensemble model maintained robust performance across all geographic regions, with AUROC ranging from 0.881 (rural) to 0.904

(urban). Performance was highest in the urban setting, likely due to more complete data records and potentially greater healthcare access. However, the consistent performance across settings (AUROC >0.88 in all cases) suggests the model is generalizable to different healthcare contexts. Rural performance was slightly lower, possibly reflecting greater data heterogeneity or socioeconomic factors.

**Feature Importance:** SHAP analysis identified the following as the most influential predictors of GDM risk (in order of importance):

1. **Pre-pregnancy BMI** (SHAP value: 0.421) - Consistent with Bhatta's (2025) finding that BMI is among the top predictors of diabetes [2].
2. **Maternal Age** (SHAP value: 0.383)
3. **Family History of Diabetes** (SHAP value: 0.312)
4. **Prior GDM History** (SHAP value: 0.289)
5. **First-Trimester Fasting Glucose** (SHAP value: 0.241)
6. **Triglycerides** (SHAP value: 0.187)
7. **Parity  $\geq 2$**  (SHAP value: 0.145)

**Socioeconomic Subgroup Analysis:** Model performance was consistent across socioeconomic strata. In low SES groups, the ensemble achieved AUROC of 0.887 (95% CI: 0.872-0.902) compared to 0.901 (95% CI: 0.889-0.913) in high SES groups. This indicates that the model performs equitably across socioeconomic groups, addressing a key concern about AI bias.

**Calibration:** The hybrid ensemble model showed good calibration, with a calibration slope of 0.941 and intercept of -0.087 in the urban cohort, and similar values across subgroups. The Hosmer-Lemeshow test was not significant ( $p=0.243$ ), indicating adequate model fit.

**Comparison with Individual Classifiers:** Both XGBoost and Random Forest individually outperformed logistic regression, consistent with Bhatta's (2025) findings [2]. The ensemble approach improved performance modestly over individual classifiers, suggesting that the combination of gradient boosting and bagging captures complementary information.

## 5. Discussion

### 5.1 Interpretation of Findings

**Primary Findings:** The hybrid ensemble model demonstrated robust GDM prediction capabilities across geographically and socioeconomically diverse populations, achieving an AUROC of 0.904 and accuracy of 89.4%. This performance is consistent with findings from Germaine et al. (2025), who reported AUROC of 0.904 using XGBoost on a single population [3], and Bhatta (2025), who achieved AUC of 0.91 for general diabetes prediction [2]. The consistency of performance across geographic and socioeconomic subgroups (AUROC 0.881-0.904) represents a significant advancement, addressing the critical gap of generalizability in existing GDM prediction models.

**Key Predictors:** Feature importance analysis identified pre-pregnancy BMI, maternal age, family history of diabetes, and prior GDM history as the most influential predictors. This aligns with the findings of Bhatta (2025), who identified glucose, age, and BMI as the most influential factors for diabetes prediction [2]. The prominence of BMI and age underscores the importance of addressing obesity trends and delayed childbearing in GDM prevention strategies. The inclusion of triglycerides as a significant predictor is consistent with the metabolic dysregulation underlying GDM and aligns with findings from metabolomic studies [1][6].

**Geographic and Socioeconomic Variation:** The slight performance decrement observed in rural populations (AUROC 0.881 vs. 0.904 urban) likely reflects differences in data completeness, healthcare access, and population characteristics. However, the magnitude of this difference (2.5% AUROC) is clinically acceptable and suggests the model is generally robust to geographic variation. Similarly, consistent performance across socioeconomic groups (AUROC difference <0.015) addresses concerns about algorithmic bias and supports equitable deployment.

**Comparison with Existing Literature:** Our findings extend prior research in several important ways. While Germaine et al. (2025) demonstrated the value of incorporating previous pregnancy data [3], our study confirms that performance gains are achievable across diverse populations. The feature importance ranking partially aligns with Bhatta's (2025) findings for general diabetes but additionally highlights GDM-specific factors such as prior GDM history and parity [2]. The ensemble approach builds on Bhatta's demonstration of the power of combining RF and XGBoost, achieving comparable performance in a different clinical context [2].

**Theoretical Implications:** The findings support both the Biopsychosocial Model and the Socioecological Model by demonstrating that GDM prediction requires consideration of biological (BMI, glucose), psychological (stress indicators, though not directly measured), and social (socioeconomic status, geographic location) factors. The consistent model performance across subgroups suggests that the biological factors driving GDM are relatively stable across populations, while the modest performance differences highlight the role of contextual factors.

### 5.2 Implications

**Academic Implications:** This study contributes to the literature in several ways. First, it provides empirical validation of ensemble learning methods (XGBoost and Random Forest) for GDM prediction across diverse populations, extending prior single-population studies [2][3]. Second, it introduces a framework for assessing model generalizability across geographic and socioeconomic strata, which can inform future predictive modeling research. Third, it demonstrates the value of SHAP analysis for model interpretability in obstetrics, providing insights into feature importance and potential mechanisms.

**Practical Implications:** The framework has several actionable implications for healthcare administrators and clinicians:

1. **Early Risk Stratification:** The model enables identification of high-risk women during the first trimester, allowing for early lifestyle interventions and monitoring.
2. **Resource Optimization:** By identifying women at highest risk, resources (e.g., dietitian consultations, glucose monitoring) can be targeted more effectively.
3. **Implementation in Clinical Decision Support:** The model's strong performance with only eight key features (as demonstrated in subset analyses) facilitates integration into existing clinical workflows without extensive data collection burdens.
4. **Equitable Care:** Consistent performance across socioeconomic groups supports equitable risk assessment and intervention.

For policymakers, the study provides evidence supporting:

1. Investment in EHR infrastructure and data standardization to enable predictive analytics
2. Integration of ML-based risk assessment into prenatal care guidelines
3. Addressing data gaps in underserved populations to ensure equitable model performance

### 5.3 Limitations

This study has several limitations that should be acknowledged:

1. **Sample Size and Generalizability:** While the sample size (n=27,561) is substantial, it may not fully capture the diversity of all populations. Future research should include more diverse cohorts, particularly from low- and middle-income countries where GDM burden is increasing.
2. **Retrospective Design:** The retrospective design limits causal inference. Prospective validation studies are needed to confirm the model's predictive accuracy in real-time clinical settings [3].

3. **Missing Data:** Despite using multiple imputation, some variables had significant missingness, particularly in rural settings. This may have contributed to the slightly lower performance observed in rural populations.
4. **Assumption of Historical Pattern Stability:** The model assumes that predictor-outcome relationships remain stable over time. Changes in clinical practice, population characteristics, or healthcare access could affect future performance.
5. **Exclusion of 2020 Data:** The exclusion of 2020 data due to COVID-19 disruptions, while ensuring data quality, means the model may not capture pandemic-era changes in GDM risk or healthcare delivery.
6. **Missing Novel Biomarkers:** The model does not incorporate metabolomic, genomic, or cfDNA features that have shown promise for GDM prediction [1][5]. Future integration of multi-omics data could enhance predictive accuracy.

#### 5.4 Future Research Directions

1. **Prospective Multi-Center Validation:** Conduct prospective studies across diverse healthcare settings to confirm the model's real-world predictive accuracy and clinical utility. This should include assessment of the model's impact on clinical decision-making and maternal-fetal outcomes.
2. **Integration of Multi-Omics Data:** Explore the incorporation of metabolomics, genomics, and cell-free DNA features to enhance predictive accuracy, particularly for early (8-12 weeks) prediction [1][5].
3. **Longitudinal Analysis of Model Performance:** Examine how model performance changes over time and in response to evolving population characteristics and clinical practices.
4. **Clinical Decision Support System Development:** Develop and evaluate a user-friendly clinical decision support system based on this framework, assessing its usability, acceptability, and impact on clinician decision-making.
5. **Health Equity Research:** Conduct in-depth qualitative studies to understand the barriers to implementing ML-based GDM prediction in underserved populations, and develop strategies to address these barriers.
6. **Extension to Other Pregnancy Outcomes:** Adapt and validate the framework for predicting other adverse pregnancy outcomes (e.g., preeclampsia, preterm birth) leveraging common risk factors and data infrastructure.

## 6. Conclusion

This study developed and validated a hybrid machine learning framework combining XGBoost and Random Forest classifiers for early prediction of Gestational Diabetes Mellitus risk across geographically and socioeconomically diverse populations. The ensemble model achieved an AUROC of 0.904 and accuracy of 89.4%, significantly outperforming traditional logistic regression (AUROC 0.817) and demonstrating robust performance across population subgroups. Feature importance analysis identified pre-pregnancy BMI, maternal age, family history of diabetes, prior GDM history, and first-trimester fasting glucose as the most influential predictors, consistent with prior research by Bhatta (2025) and others.

The primary contribution of this research is the provision of a replicable, interpretable, and generalizable framework for early GDM risk stratification that maintains strong performance across diverse healthcare settings and population subgroups. For healthcare administrators and clinicians, this framework offers a practical tool for early identification of high-risk pregnancies, enabling timely interventions that could reduce adverse outcomes and improve resource allocation. For policymakers, the study provides evidence supporting the integration of ML-based risk assessment into routine prenatal care protocols.

Future research should focus on prospective validation, integration of novel biomarkers, and development of user-friendly clinical decision support systems. With continued research and responsible implementation, ML-based GDM prediction has the potential to transform prenatal care, enabling personalized, timely interventions that improve maternal and fetal health outcomes across diverse populations.

## References

1. Cao, Y., Wang, H., Zhang, L., et al. (2025). Early prediction of gestational diabetes mellitus using machine learning-integrated metabolomic and clinical features. *Frontiers in Endocrinology*, 16, 1687146. <https://doi.org/10.3389/fendo.2025.1687146>
2. Bhatta, R. P. (2025). Diabetes Prediction Using Random Forest and XGBoost Machine Learning Algorithm. *Journal of Engineering Technology and Planning*, 6(1), 88-103. <https://doi.org/10.3126/joetp.v6i1.87829>
3. Germaine, M., O'Higgins, A. C., Egan, B., & Healy, G. (2025). Evaluation of Machine Learning Models for Early Prediction of Gestational Diabetes Using Retrospective Electronic Health Records from Current and Previous Pregnancies. *medRxiv*. <https://doi.org/10.1101/2025.05.12.25327431>
4. Zhang, Y., Wang, J., & Liu, H. (2025). Predicting the future risk of developing type 2 diabetes in women with a history of gestational diabetes mellitus using machine learning and explainable artificial intelligence. *Primary Care Diabetes*, 19(6), 658-666. <https://doi.org/10.1016/j.pcd.2025.09.003>
5. Nguyen, T. H., Tran, Q. H., Le, M. T., et al. (2025). Early Prediction of Gestational Diabetes Using Integrated Cell-free DNA Features and Omics-derived Genetic Scores. *medRxiv*. <https://doi.org/10.1101/2025.09.03.25334985>
6. Lee, S. H., Kim, J. Y., Park, J. M., et al. (2025). Establishment of an accurate prediction system for gestational diabetes mellitus based on the characteristics of metabolic kinetics in early pregnancy: a prospective two-center cohort study. *Endocrine: International Journal of Basic and Clinical Endocrinology*, 90(3), 1201-1220.
7. American Diabetes Association. (2023). Classification and Diagnosis of Diabetes: Standards of Medical Care in Diabetes—2023. *Diabetes Care*, 46(Supplement\_1), S19-S40.
8. International Association of Diabetes and Pregnancy Study Groups Consensus Panel. (2010). International Association of Diabetes and Pregnancy Study Groups Recommendations on the Diagnosis and Classification of Hyperglycemia in Pregnancy. *Diabetes Care*, 33(3), 676-682.
9. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785-794.
10. Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5-32.

11. Lundberg, S. M., & Lee, S. I. (2017). A Unified Approach to Interpreting Model Predictions. *Advances in Neural Information Processing Systems*, 30, 4765-4774.
12. Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research*, 16, 321-357.
13. Harrell, F. E. (2015). *Regression Modeling Strategies: With Applications to Linear Models, Logistic and Ordinal Regression, and Survival Analysis*. Springer.
14. World Health Organization. (2024). *Diagnostic Criteria and Classification of Hyperglycaemia First Detected in Pregnancy*. WHO Guidelines.
15. Pedregosa, F., Varoquaux, G., Gramfort, A., et al. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*, 12, 2825-2830.