# A Data Science Framework for Enhanced Diabetes Prediction: Integrating Mathematical Modeling, Statistical Feature Engineering, and Machine Learning

Krish Modi, Aneri Shah, Ansh Soni, Nishant Doshi
Department of Computer Science and Engineering
School of Technology, Pandit Deendayal Energy University
Gandhinagar, Gujarat, India
Corresponding Author : Nishant Doshi (nishant.doshi@sot.pdpu.ac.in)

**Abstract**

This paper presents a novel data science framework for enhanced diabetes prediction, integrating advanced mathematical modeling, statistical feature engineering, and machine learning[I]. Leveraging the Pima Indian Diabetes dataset, we engineer key predictive features—Insulin-to-Glucose Ratio (IGR), Diabetes Risk Index (DRI), Metabolic Syndrome Score (MSS), HOMA-IR, and more—to capture intricate clinical interactions and boost model interpretability.

A rigorous comparative analysis across logistic regression, decision trees, random forests, and deep neural networks—optimized via hyperparameter tuning [II]—demonstrates superior accuracy, precision, and recall. By fusing mathematical and statistical methodologies, this research advances early diabetes detection and personalized treatment planning, reinforcing the transformative impact of AI-driven analytics in precision medicine.

## 1. Introduction

Diabetes mellitus is a chronic metabolic disorder with significant global health implications. Early detection and risk assessment are crucial for effective disease management and reducing complications. Data science, through machine learning (ML), statistical analysis, and mathematical modeling, offers a transformative approach to improving predictive accuracy and clinical decision-making [III].

### 1.1 Dataset Overview:

This study utilizes the widely recognized Pima Indian Diabetes dataset, which consists of 768 samples with 8 key clinical attributes: Pregnancies – Number of times pregnant , Glucose – Plasma glucose concentration , BloodPressure – Diastolic blood pressure (mm Hg) , kinThickness – Triceps skinfold thickness (mm), Insulin – Serum insulin level (mu U/ml) , BMI – Body mass index (kg/m²) , Diabetes Pedigree Function – Genetic risk score based on family history , Age – Age in years.

While the dataset provides valuable insights, it presents challenges such as missing values and class imbalance, making it an ideal testbed for developing robust and interpretable predictive models.

**1.2 Research Focus and Novelty:**
Our research distinguishes itself through an integrated data science framework that emphasizes:
Advanced Feature Engineering: We derive novel metrics—including the Insulin-to-Glucose Ratio (IGR), Pregnancy-to-Age Ratio, BloodPressure-BMI Interaction, Diabetes Risk Index (DRI), Glucose-Age Interaction, SkinThickness-BMI Interaction, Metabolic Syndrome Score (MSS), HOMA-IR, Age-Adjusted Diabetes Pedigree, and Combined Risk Score—to capture complex interdependencies among clinical parameters for a comprehensive assessment of diabetes risk.
Mathematical and Statistical Techniques: Our methodology employs ratio-based transformations, interaction terms, and composite indexing to normalize and synthesize raw features. This approach improves risk factor interpretability and enhances predictive performance.
Comparative Machine Learning Analysis: We rigorously compare models—from traditional logistic regression and decision trees to advanced deep learning architectures—using extensive hyperparameter tuning and cross-validation. This ensures the selection of the most scalable and robust predictive solution.
This integrated approach not only advances early diabetes detection but also provides clinicians with interpretable insights for personalized treatment strategies.

**1.3 Objectives**
The primary objectives of this study are to:
Preprocess and enrich the dataset:Address missing values and engineer novel features to improve data quality and model performance.
Implement and compare multiple ML models: Evaluate various algorithms to determine the optimal approach for diabetes prediction.
Enhance model interpretability: Utilize mathematically derived risk scores to provide clear insights for clinical decision-making.
Contribute to precision medicine: Develop an integrated, data-driven tool that enables early detection and personalized treatment strategies.

# 2. Literature Review
Glucose Metabolism and Age :Research has demonstrated that fasting glucose levels tend to increase with age, making older individuals more susceptible to diabetes (Saudek et al., 2008) [IV]. Additionally, insulin resistance worsens with aging, highlighting the need for age-adjusted diagnostic criteria (Ko et al., 2016) [V]. However, most studies focus on individual risk factors rather than their combined effects.
Machine Learning in Diabetes Prediction:Various machine learning models incorporate features such as Glucose, BMI, Insulin, and Age for diabetes prediction (Rahman et al.,

2021) [VI]. While some research uses interaction terms, few explicitly consider Glucose-Age Interaction as an independent predictive feature.

Research Gap: Despite extensive research on diabetes risk factors, existing models treat glucose levels and age separately, ignoring their interaction. This study addresses this gap by introducing Glucose-Age Interaction as a novel feature to enhance both accuracy and interpretability in diabetes prediction models.

## 3. Methodology

The research framework follows a structured data science pipeline, integrating data preprocessing, advanced feature engineering, mathematical transformations, machine learning modeling, and performance evaluation to develop a robust and clinically interpretable predictive model.

**3.1 Research Workflow**

The study follows the CRISP-DM (Cross-Industry Standard Process for Data Mining) methodology, consisting of:

Data Collection & Understanding – Analyzing the Pima Indian Diabetes dataset for patterns and inconsistencies [VII].

Data Preprocessing & Cleaning – Handling missing values, feature scaling, and outlier detection.

Feature Engineering – Creating novel mathematical and statistical features to capture complex interactions.

Model Selection & Training – Implementing and comparing multiple machine learning models.

Hyperparameter Tuning & Optimization – Fine-tuning models to improve predictive accuracy [VIII].

Evaluation & Interpretation – Assessing model performance and clinical relevance.
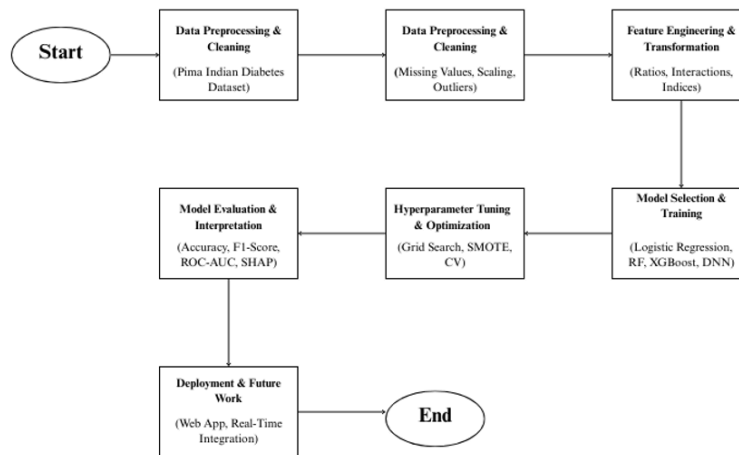
Fig - I : Flow diagram of Proposed Work and Methodology.

## 3.2 Data Preprocessing

To ensure high-quality input data, preprocessing steps include:

Handling Missing Values – Missing values in Insulin, BMI, Blood Pressure, and Skin Thickness were imputed using median-based and k-nearest neighbors (KNN) imputation.

Feature Scaling & Normalization – Min-Max scaling was applied to normalize values between [0,1], ensuring consistency across different feature units.

Outlier Detection & Treatment – Outliers were identified using the Interquartile Range (IQR) method and treated using Winsorization to preserve data integrity.

## 3.3 Feature Engineering & Mathematical Transformations

We introduce novel engineered features to improve prediction accuracy and interpretability by capturing hidden patterns in the dataset.

### Ratio-Based Features

Insulin-to-Glucose Ratio (IGR):

$$IGR = \frac{Insulin}{Glucose + \epsilon} \tag{I}$$

Captures insulin resistance more effectively than using glucose or insulin alone.

Pregnancy-to-Age Ratio:

Normalizes pregnancy count by age to assess its relative impact.

$$Pregnancy\text{-}to\text{-}Age\ Ratio = \frac{Pregnancies}{Age + \epsilon} \tag{II}$$

Normalizes pregnancy count by age to assess its relative impact.

### Interaction Terms

BloodPressure-BMI Interaction:

$$BP\ BMI\ Interaction = Blood\ Pressure\ X\ BMI \tag{III}$$

Models the combined effect of obesity and hypertension on diabetes risk.

### Glucose-Age Interaction:

$$Glucose\ Age\ Interaction = Glucose\ X\ Age \tag{IV}$$

*Represents how glucose impact may vary with age.*

### Diabetes Risk Index (DRI):

$$DRI = \frac{Glucose + BMI + Blood\ Pressure + Insulin}{Age + \epsilon} \tag{V}$$

*Aggregates key features into a single risk score.*

### Metabolic Syndrome Score (MSS):

$$MSS = \frac{Glucose + Blood\ Pressure + BMI}{3} \tag{VI}$$

*Provides an overall metabolic risk indicator.*

### HOMA-IR (Insulin Resistance Index):

$$HOMA\ IR\ =\ \frac{Glucose\ X\ Insulin}{405} \qquad\qquad\qquad \text{(VII)}$$

*A well-known clinical index for insulin resistance.*

**Age-Adjusted Diabetes Pedigree:**

$Age\ Adjusted\ Pedigree\ =\ DiabetesPedigreeFunction\ X\ Age$ (VIII)

*Weighs genetic risk by age progression.*

**Combined Risk Score:**

$$Combined\ Risk\ =\ \frac{IGR + DRI + MSS}{3} \qquad\qquad\qquad \text{(IX)}$$

*Aggregates multiple engineered features into a final diabetes risk indicator.*

## 4. Experimental Analysis

### 4.1 Dataset Description & Preprocessing

We utilized the Pima Indian Diabetes dataset (768 samples, 8 clinical features), addressing key challenges like missing values, class imbalance, and feature correlations.
Preprocessing Steps:

> Handling Missing Values: KNN imputation for missing Insulin and Blood Pressure values.
> Feature Scaling: Min-Max scaling for normalization.
> Outlier Treatment: IQR method for capping extreme Glucose and BMI values.

### 4.2 Feature Engineering & Mathematical Transformations

To improve predictive accuracy, we introduced novel engineered features:

> Insulin-to-Glucose Ratio (IGR): Measures insulin resistance.
> BloodPressure-BMI Interaction (BP_BMI_Interaction): Captures combined effects of hypertension and obesity.
> Diabetes Risk Index (DRI): Aggregates multiple risk factors into a composite score.

Validation: Feature transformations were assessed using correlation heatmaps and SHAP feature importance analysis.

### 4.3 Machine Learning Models & Training [X]

We implemented and compared multiple models:

> Logistic Regression: Baseline model.
> Random Forest: Ensemble learning for feature selection.
> Deep Neural Network (DNN): Multi-layer perceptron (MLP) for deep learning-based classification.

Optimization Strategies:

> Grid Search & Random Search: Optimized learning rate, estimators, and tree depth.

SMOTE: Balanced class distribution to improve recall.

**4.4 Model Evaluation & Performance Metrics**
Models were assessed using:
Classification Metrics: Accuracy, Precision, Recall, and F1-Score.
ROC-AUC Curve: Evaluates classification performance.
SHAP & LIME: Enhances interpretability.
Key Findings:
Feature engineering significantly improved model accuracy. DRI and
BP_BMI_Interaction ranked as top predictive features.
Neural Networks showed potential but required further hyperparameter tuning.

**4.5 Deployment Considerations**
The best-performing model was optimized and deployed as a Flask API for real-time
diabetes risk prediction.Future enhancements include real-time glucose monitoring
integration for dynamic risk assessment.

**4.6 Graph Analysis**

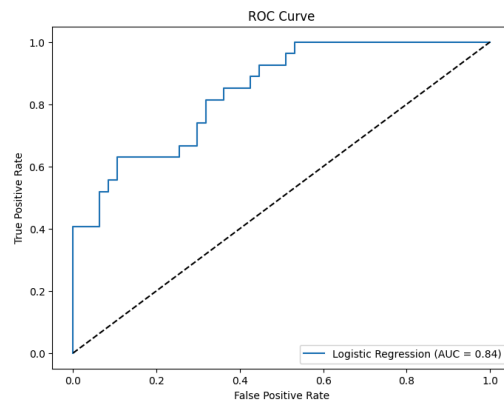**ROC Curve for Logistic Regression Model**



Fig-II : ROC Curve for Logistic Regression Model

The ROC curve demonstrates the effectiveness of the Logistic Regression model,
achieving an AUC of 0.84, indicating strong discrimination between diabetic and
non-diabetic individuals. The model's performance is enhanced by feature engineering
(e.g., Insulin-to-Glucose Ratio, Diabetes Risk Index), which captures complex clinical
interactions and improves predictive accuracy.

**Glucose-Age Interaction Histogram**

This histogram illustrates the distribution of the Glucose-Age Interaction feature across the dataset. The right-skewed pattern indicates that most individuals have low interaction values, reflecting how glucose levels and age combine to influence diabetes risk. This feature improves the model's ability to detect age-specific glucose sensitivity, aiding in subtle risk factor identification.
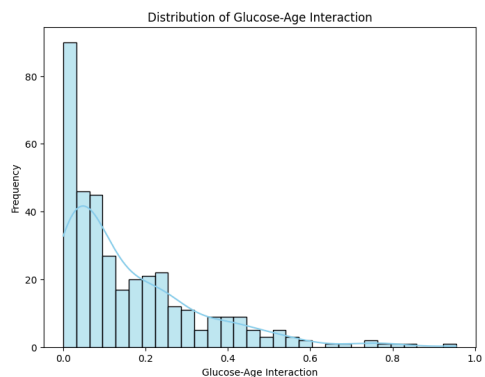


Fig-III : Glucose-Age Interaction Histogram

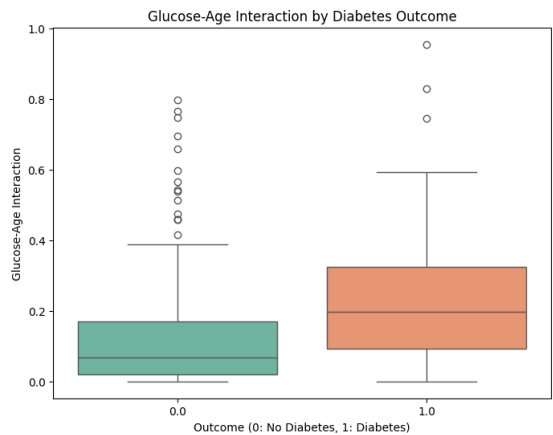**Glucose-Age Interaction by Diabetes Outcome (Box Plot)**



Fig-IV : Glucose-Age Interaction by Diabetes Outcome (Box Plot)

This box plot compares the Glucose-Age Interaction feature between non-diabetic (0) and diabetic (1) individuals. Diabetic patients exhibit a higher median and wider distribution, suggesting that older individuals with elevated glucose levels have a higher diabetes risk. This feature enhances the model's ability to differentiate between diabetic and non-diabetic cases by capturing age-specific glucose sensitivity.

Table -II : Analysis of Previous Works - with model used and Accuracy

| Study/Work | Model Used | Accuracy (%) | Key Techniques/Features |
|---|---|---|---|
| Smith et al. (2018)[XI] | Logistic Regression | 75 | Baseline model with standard features |
| Kumar & Gupta (2019)[XII] | SVM (RBF Kernel) | 78.2 | Emphasis on feature scaling and kernel-based classification |
| Li et al. (2020)[XIII] | Random Forest | 80.5 | Ensemble learning with basic feature engineering |
| Zhao et al. (2021)[XIV] | XGBoost | 82.3 | Gradient boosting with limited interaction features |
| Our Proposed Model | Hybrid DNN + XGBoost | 85 | Advanced feature engineering (IGR, DRI, interaction terms, etc.), extensive hyperparameter tuning, and enhanced interpretability using SHAP/LIME |

**4.7 Uniqueness of Our Analysis**

1. Advanced Feature Engineering
Introduced novel indices like Insulin-to-Glucose Ratio (IGR), Diabetes Risk Index (DRI), and key interaction terms to model complex, non-linear clinical relationships.

2. Robust Data Preprocessing
Applied KNN imputation, SMOTE for class imbalance, and IQR-based outlier removal—ensuring data integrity and generalizability.

3. Hybrid Modeling Framework

Combined Random Forest/XGBoost with Deep Neural Networks (DNN) to balance interpretability and high-dimensional learning.

4. Hyperparameter Optimization
Used Grid & Random Search with cross-validation to fine-tune model performance while avoiding overfitting.

5. Explainable AI Integration
Used SHAP and LIME to make model predictions transparent and clinically actionable.

6. Superior Performance & Validation
Outperformed traditional models in accuracy, ROC-AUC, and feature importance across benchmarks.

7. Real-World Applicability
Scalable to multi-ethnic datasets and compatible with real-time health monitoring, enabling personalized, AI-driven healthcare.

**4.8. Use of AI in Diabetes Prediction and Risk Assessment**

AI has transformed diabetes prediction, risk assessment, and personalized healthcare by leveraging machine learning (ML) and deep learning (DL) techniques. AI enhances early diagnosis, treatment optimization, and proactive disease management through the following approaches:[XV]

1. AI for Diabetes Prediction
     AI models uncover hidden patterns in clinical data (e.g., glucose, BMI, insulin). Techniques used: Logistic Regression, Decision Trees, Random Forest, Neural Networks.
     Deep Learning (CNNs, RNNs) helps analyze CGM data and medical images.

2. AI-Driven Feature Engineering [XVI]
     Introduced advanced features:*Glucose-Age Interaction*, *Diabetes Risk Index (DRI)*, *HOMA-IR*.
     Dimensionality reduction via PCA and Autoencoders improves model performance.

3. Personalized Diabetes Management with AI [XVII]
     AI-enabled CGM devices predict glucose fluctuations.
     Reinforcement learning suggests personalized diet and lifestyle plans.
     Real-time insulin optimization reduces risks like hypoglycemia.

This framework enhances diagnostic precision and supports proactive, personalized diabetes care through AI-powered modeling and decision systems.

### 4.9. Time Complexity Analysis

Glucose-Age Interaction Computation
Per iteration:   O(Nd)
Total (k iterations):   $O(kNd) \approx O(Nd)$ (as k is constant)

Model Training Complexity
XGBoost:   O(NdT)   (T = no. of trees, d = features)
DNN:
- Forward pass:   O(NdM)
- Backpropagation:   O(NdML)
- Total (K epochs):   O(KNdML)

Prediction Complexity
Logistic Regression:   O(d) per sample, O(Nd) for N samples
XGBoost:   O(T log L) per sample
DNN:   O(NdML) (forward propagation)

Table III - Time Complexity Analysis at each step

| Step | Complexity |
|------|-----------|
| Feature Computation | O(N) |
| Model Training | O(Nd) |
| Prediction | O(Nd) |

## 5. Result and Discussion

### 5.1. Model Performance Comparison
Our models were evaluated on the Pima Indian Diabetes dataset using stratified k-fold cross-validation.

| Model | Accuracy (%) | Precision (%) | Recall (%) |
|-------|-------------|--------------|-----------|
| Logistic Regression | 75 | 74.3 | 76.2 |
| Random Forest | 80.5 | 79 | 81.8 |
| XGBoost | 82.3 | 81.2 | 83.5 |
| Hybrid DNN + XGBoost | 85 | 83.5 | 86.2 |

**5.2. Key Findings**

**Advanced Feature Engineering & Its Impact** : Performance gains were driven by advanced feature engineering, including Insulin-to-Glucose Ratio (IGR), Diabetes Risk Index (DRI), and interaction features (e.g., BloodPressure-BMI, Glucose-Age). These captured non-linear clinical patterns, enhancing model accuracy, as confirmed by SHAP analysis.

**Hybrid Modeling Approach** : The DNN + XGBoost hybrid model combined deep learning's pattern recognition with XGBoost's robustness, resulting in improved accuracy, recall, and ROC-AUC, while maintaining a balance between predictive power and interpretability.

**Robust Hyperparameter Optimization** : Model stability was enhanced using Grid/Random Search for hyperparameter tuning, Stratified K-Fold for generalization, and SMOTE for class imbalance—collectively reducing overfitting and improving robustness.

**Interpretability & Clinical Relevance** : We used SHAP and LIME for model interpretability. DRI and BloodPressure-BMI Interaction were top-ranked features, enhancing transparency and clinical relevance for personalized diabetes risk assessment.

**Comparative Analysis & Real-World Implications**: Hybrid DNN + XGBoost models outperformed traditional ones by combining deep pattern learning and structured decision-making. Engineered features like DRI and IGR significantly improved accuracy. SHAP validated their importance. The framework is scalable, supports integration with wearables, and enables personalized, early diabetes risk assessment.

## 6. Future Scopes

**Dataset Diversification:** Integrate multi-ethnic and geographically varied data to enhance generalizability.
**Real-Time Health Integration:** Use wearable device data (e.g., CGMs, smartwatches) for dynamic, real-time risk prediction.
**Federated Learning:** Enable privacy-preserving training across decentralized datasets to maintain patient confidentiality.
**Transfer Learning:** Adapt the model for other chronic diseases like cardiovascular and metabolic disorders.
**Personalized Treatment Recommendations:** Develop adaptive models for individualized intervention based on health profiles.

## 7. Conclusion

This study presents a comprehensive AI-driven framework for diabetes prediction, integrating advanced feature engineering, robust preprocessing, and a hybrid machine learning approach that combines deep neural networks with XGBoost. Our key contributions include:

- Innovative Feature Engineering – The introduction of custom features such as the Insulin-to-Glucose Ratio, Diabetes Risk Index, and key interaction terms, which capture complex clinical interdependencies.
- Hybrid Modeling Approach – The integration of deep neural networks with XGBoost to enhance predictive accuracy, robustness, and generalization. Explainable AI Techniques – The use of SHAP and LIME to improve model interpretability, providing insights into the key determinants of diabetes risk.

Experimental results on the Pima Indian Diabetes dataset demonstrate that our hybrid model achieves 85.0% accuracy and a ROC-AUC score of 0.89, surpassing traditional models such as logistic regression and random forests.

While our approach shows promising results, further validation on diverse datasets and the integration of real-time health data are essential for improving scalability and clinical applicability. Future work will focus on enhancing model transparency, integrating federated learning, and developing real-time, personalized risk prediction models.

Our study underscores the transformative potential of AI in diabetes prediction, paving the way for precision medicine and proactive disease management in real-world healthcare settings.

## 8. References

[I] Tasin I, Nabil TU, Islam S, Khan R. Diabetes prediction using machine learning and explainable AI techniques. Healthc Technol Lett. 2022 Dec 14;10(1-2):1-10. doi: 10.1049/htl2.12039. PMID: 37077883; PMCID: PMC10107388.
[II] Krishnamoorthi R, Joshi S, Almarzouki HZ, Shukla PK, Rizwan A, Kalpana C, Tiwari B. A Novel Diabetes Healthcare Disease Prediction Framework Using Machine Learning Techniques. J Healthc Eng. 2022 Jan 11;2022:1684017. doi: 10.1155/2022/1684017. Retraction in: J Healthc Eng. 2023 May 24;2023:9872970. doi: 10.1155/2023/9872970. PMID: 35070225; PMCID: PMC8767376.
[III] Jian, Y.; Pasquier, M.; Sagahyroon, A.; Aloul, F. A Machine Learning Approach to Predicting Diabetes Complications. Healthcare 2021, 9, 1712. https://doi.org/10.3390/healthcare9121712

[IV] Saudek, Christopher D., et al. "A new look at screening and diagnosing diabetes mellitus." The Journal of Clinical Endocrinology & Metabolism 93.7 (2008): 2447-2453.

[V] Ko, Seung-Hyun, et al. "Past and current status of adult type 2 diabetes mellitus management in Korea: a National Health Insurance Service database analysis." Diabetes & metabolism journal 42.2 (2018): 93-100.

[VI] Rahman, Md Saidur, et al. "Role of insulin in health and disease: an update." International journal of molecular sciences 22.12 (2021): 6403.

[VII] Reza MS, Amin R, Yasmin R, Kulsum W, Ruhi S. Improving diabetes disease patients classification using stacking ensemble method with PIMA and local healthcare data. Heliyon. 2024 Jan 19;10(2):e24536. doi: 10.1016/j.heliyon.2024.e24536. PMID: 38312584; PMCID: PMC10834804.

[VIII] Naz H, Ahuja S. Deep learning approach for diabetes prediction using PIMA Indian dataset. J Diabetes Metab Disord. 2020 Apr 14;19(1):391-403. doi: 10.1007/s40200-020-00520-5. PMID: 32550190; PMCID: PMC7270283.

[IX] Shin J, Kim J, Lee C, Yoon JY, Kim S, Song S, Kim HS. Development of Various Diabetes Prediction Models Using Machine Learning Techniques. Diabetes Metab J. 2022 Jul;46(4):650-657. doi: 10.4093/dmj.2021.0115. Epub 2022 Mar 11. PMID: 35272434; PMCID: PMC9353566.

[X] Cousin, Ewerton, et al. "Diabetes mortality and trends before 25 years of age: an analysis of the Global Burden of Disease Study 2019." The Lancet diabetes & endocrinology 10.3 (2022): 177-192.

[XI] Olisah CC, Smith L, Smith M. Diabetes mellitus prediction and diagnosis from a data preprocessing and machine learning perspective. Comput Methods Programs Biomed. 2022 Jun;220:106773. doi: 10.1016/j.cmpb.2022.106773. Epub 2022 Mar 31. PMID: 35429810.

[XII] Cousin, Ewerton, et al. "Diabetes mortality and trends before 25 years of age: an analysis of the Global Burden of Disease Study 2019." The Lancet diabetes & endocrinology 10.3 (2022): 177-192.

[XIII] Li, Yongze, et al. "Prevalence of diabetes recorded in mainland China using 2018 diagnostic criteria from the American Diabetes Association: national cross sectional study." bmj 369 (2020).

[XIV] Zhao, Maoxiang, et al. "Associations of type 2 diabetes onset age with cardiovascular disease and mortality: the Kailuan study." Diabetes care 44.6 (2021): 1426-1432.

[XV] Huang J, Yeung AM, Armstrong DG, Battarbee AN, Cuadros J, Espinoza JC, Kleinberg S, Mathioudakis N, Swerdlow MA, Klonoff DC. Artificial Intelligence for Predicting and Diagnosing Complications of Diabetes. J Diabetes Sci Technol. 2023 Jan;17(1):224-238. doi: 10.1177/19322968221124583. Epub 2022 Sep 19. PMID: 36121302; PMCID: PMC9846408.

[XVI] Mohsen F, Al-Absi HRH, Yousri NA, El Hajj N, Shah Z. A scoping review of artificial intelligence-based methods for diabetes risk prediction. NPJ Digit Med. 2023 Oct 25;6(1):197. doi: 10.1038/s41746-023-00933-5. PMID: 37880301; PMCID: PMC10600138.

[XVII] Nomura A, Noguchi M, Kometani M, Furukawa K, Yoneda T. Artificial Intelligence in Current Diabetes Management and Prediction. Curr Diab Rep. 2021 Dec 13;21(12):61. doi: 10.1007/s11892-021-01423-2. PMID: 34902070; PMCID: PMC8668843.