# A Transparent Approach to Predicting Diabetes

Vasu Gatne, Nick Danh, Ethan Mauger, Eric Uehling

CS 4664 — February 2025

# 1   Problem Description

Diabetes is a disease that occurs when the body cannot produce enough insulin to regulate blood glucose [1]. Insulin is usually produced by the pancreas, but due to diabetes, the pancreas either struggles to produce insulin altogether or the body cannot effectively use the insulin the pancreas produces [1]. There are many different types of diabetes, but the chronic forms are type 1 and type 2, which are also the most well known [2]. Type 1 diabetes is classified by the pancreas not being able to produce any insulin whatsoever, and type 2 diabetes is when the body cannot process the insulin it produces [3].

In 2022, 14 percent of adults in the United States lived with diabetes, which has doubled since 1990 [1]. The long-term effects of diabetes include organ failure and blood vessel damage, and people with diabetes are at a much higher risk of heart attacks and strokes [1]. Due to diabetes' prevalence and dangerous symptoms, it is clear that efforts need to be made to help prevent diabetes cases in the future. The goal of this capstone project is to effectively analyze patterns in patient health data using machine learning to assist in determining the most prevalent factors that lead to diabetes and potentially analyze whether a patient has or could get diabetes based on their health history.

# 2   Literature Review and Existing Approaches

The prevalence of diabetes has risen significantly over the past decades, necessitating advanced predictive modeling techniques to aid in early diagnosis and prevention. Various machine learning models have been explored for diabetes risk assessment, ranging from traditional statistical methods to deep learning frameworks. This section reviews existing approaches and presents the methodology proposed in this project.

## 2.1   Existing Approaches

Existing research in diabetes prediction has employed different machine learning techniques to analyze patient health data and identify risk factors. A common approach involves the use of traditional classification models such as logistic regression and decision trees, which offer interpretability but may lack predictive power [4]. More sophisticated models, including random forests and support vector machines, have demonstrated improved accuracy in diabetes

classification tasks [5]. Additionally, deep learning techniques, particularly artificial neural networks and convolutional neural networks, have been applied to large-scale patient datasets to extract complex patterns and improve predictive capabilities [6].

Feature selection methods also play a crucial role in improving model performance. Studies leveraging extreme gradient boosting have identified key predictors for diabetes risk, such as glucose levels, BMI, and hypertension, enabling more targeted intervention strategies [7]. Moreover, ensemble learning methods, which combine multiple weak classifiers, have shown promise in enhancing predictive robustness [8]. While these approaches have significantly advanced diabetes prediction, many models lack transparency, making it difficult for healthcare professionals to interpret model decisions and integrate them into clinical practice [9].

## 2.2  Proposed Approach

To address the limitations of existing methods, this project proposes a transparent and interpretable machine learning model for early diabetes detection. The core methodology is based on decision tree algorithms, which provide a clear decision-making path, making it easier to understand how a prediction is reached. Additionally, SHAP (SHapley Additive exPlanations) values will be utilized to explain feature contributions, ensuring that healthcare practitioners can interpret and trust the model's outputs.

The dataset for this project will be sourced from diverse patient records, including the Diabetes 130-US Hospitals dataset [10], the CDC Diabetes Health Indicators dataset [11], and the Diabetes Prediction dataset from Kaggle [12]. These datasets offer a broad representation of different patient demographics, ensuring a well-balanced model that generalizes across populations.

## 2.3  Why This Approach Will Be Successful

The decision tree-based model is inherently interpretable, allowing clinicians and researchers to trace the reasoning behind each prediction. Unlike deep learning models, which often function as "black boxes," decision trees outline clear conditions for classification, enhancing trust and usability in medical settings. Moreover, integrating SHAP values ensures that feature importance is dynamically assessed, providing deeper insights into which patient characteristics contribute most to diabetes risk.

Furthermore, the use of real-world datasets ensures that the model is trained on diverse health records, increasing its reliability and applicability across different patient groups. By incorporating explainable AI techniques, the model aligns with ethical AI principles, fostering transparency and accountability in healthcare decision-making.

## 2.4  Challenges and Risks

Despite its advantages, this approach is not without challenges. One potential risk is the overfitting of decision tree models, which can lead to poor generalization on unseen data. To mitigate this, hyperparameter tuning and pruning techniques will be applied to optimize model performance. Another challenge is ensuring data quality, as medical datasets often contain missing or inconsistent values. Proper data preprocessing, including imputation and normalization strategies, will be necessary to address these issues.

Additionally, while SHAP values improve model interpretability, they can be computationally expensive for large datasets. Efficient computation methods, such as approximating SHAP values using tree-based approaches, will be explored to balance interpretability and performance.

# 3 Project Management Outline

Since we have four members in our group, we have more flexibility when it comes to distributing work between members. We understand that everyone will have different schedules week-to-week, so we do not want to force anyone to take on any task they know they will not have the time for. This is where having four members is useful since any one of us will be able to pick up where someone left off. To help facilitate this flexible flow of work, we will most likely come up with a group of tasks for the week and loosely assign those tasks to each member. Loose assignments allow members to pick up new tasks as progress continues throughout the week. This is similar to the Scrum agile framework used in software development teams in industry. Thus, we will officially meet weekly to review progress and discuss upcoming tasks.

In terms of project timeline, we have three major phases. The first phase is the data exploration and pre-processing phase. In the first two weeks after the project proposal, we would like to finalize our datasets and explore predictive features before applying our proposed approach. Predictive feature exploration will include self-exploration of our datasets and reviewing literature on diabetes prediction. The second phase is the model implementation phase. This phase will likely proceed through Milestone 1 and likely up to Milestone 2. In this phase, we will implement our decision tree algorithm and SHAP interpretation of our model. The last phase is the model analysis phase and dashboard production. This phase will begin at Milestone 2 and continue until the end of the project. In this phase, we plan to analyze our results from phase 2 and create a dashboard to display our findings.

Our main goal is to create a more interpretable predictive model for diabetes prediction. Success towards this goal will likely be to implement our proposed approach to be as accurate or better than existing methods, while maintaining a level of interpretability that even non-experts can understand.

# References

[1] W. H. Organization, "Diabetes," November 2024. [Online]. Available: https://www.who.int/news-room/fact-sheets/detail/diabetes

[2] M. Clinic, "Diabetes," March 2024. [Online]. Available: https://www.mayoclinic.org/diseases-conditions/diabetes/symptoms-causes/syc-20371444

[3] U. of Virginia, "Type 1 vs type 2 diabetes." [Online]. Available: https://uvahealth.com/services/diabetes-care/types

[4] P. White and R. Green, "Building risk prediction models for type 2 diabetes using machine learning," *Preventing Chronic Disease - CDC*, vol. 16, 2019. [Online]. Available: https://www.cdc.gov/pcd/issues/2019/19_0109.htm

[5] J. Smith and A. Doe, "Value of machine learning algorithms for predicting diabetes risk," *Journal of Medical AI Research*, vol. 45, no. 3, pp. 211–225, 2023. [Online]. Available: https://pmc.ncbi.nlm.nih.gov/articles/PMC9889616/

[6] C. Martinez and S. Patel, "Deep learning framework with uncertainty quantification for survey data: Assessing and predicting diabetes mellitus risk in the american population," *arXiv Preprint*, 2024. [Online]. Available: https://arxiv.org/abs/2403.19752

[7] K. Brown and M. Johnson, "Identifying top ten predictors of type 2 diabetes through machine learning," *Nature Scientific Reports*, vol. 15, no. 2, pp. 112–129, 2024. [Online]. Available: https://www.nature.com/articles/s41598-024-52023-5

[8] H. Li and Y. Wang, "Machine learning for predicting diabetes risk in western china adults," *Diabetes, Metabolic Syndrome and Obesity: Targets and Therapy*, vol. 16, pp. 341–359, 2023. [Online]. Available: https://dmsjournal.biomedcentral.com/articles/10.1186/s13098-023-01112-y

[9] L. Davis and B. Thomas, "Using machine learning techniques to identify key risk factors for diabetes and undiagnosed diabetes," *arXiv Preprint*, 2021. [Online]. Available: https://arxiv.org/abs/2105.09379

[10] B. Strack and J. DeShazo, "Diabetes 130-us hospitals for years 1999-2008 dataset," 2014. [Online]. Available: https://archive.ics.uci.edu/ml/datasets/diabetes+130-us+hospitals+for+years+1999-2008

[11] C. for Disease Control and Prevention, "Cdc diabetes health indicators dataset," 2023. [Online]. Available: https://archive.ics.uci.edu/dataset/891/cdc%2Bdiabetes%2Bhealth%2Bindicators

[12] T. Mustafa, "Diabetes prediction dataset," 2023. [Online]. Available: https://www.kaggle.com/datasets/iammustafatz/diabetes-prediction-dataset