# Abstract:

*"This project aimed to develop an accurate predictive model for diabetes using various machine learning algorithms. Diabetes is a chronic disease affecting people of all ages, and early detection is crucial for effective prevention and management. The dataset utilized in this study encompassed important medical features, including age, BMI, blood pressure, glucose levels, insulin levels, Diabetes Pedigree Function, Skin Thickness, Pregnancies. Multiple algorithms, such as logistic regression, decision trees, random forest etc were employed and evaluated to determine the most accurate predictor of diabetes. The outcomes revealed that '**Decision Tree** 'has the highest accuracy. This research underscores the potential of machine learning in accurately predicting diabetes, enabling proactive measures and personalized interventions. The models developed in this project provide valuable tools for healthcare professionals and individuals to assess diabetes risk and make informed decisions for disease prevention and management. Further research and refinement of the models hold promise for enhancing diabetes prediction and advancing healthcare outcomes"*

# Table of Contents:

# Introduction:

Diabetes is a chronic disease that affects individuals of all ages and is a significant public health concern. Early detection is vital for effective prevention and management. Machine learning algorithms have emerged as promising tools for developing accurate predictive models using medical data. This project aims to explore various machine learning algorithms and build a robust predictive model for identifying individuals at risk of diabetes.

The dataset used in this study includes essential medical features such as age, BMI, blood pressure, glucose levels, insulin levels, Diabetes Pedigree Function, Skin Thickness, and the number of pregnancies. These features are known to be influential in diabetes diagnosis and provide comprehensive insights into an individual's health status for assessing diabetes risk accurately.

Multiple machine learning algorithms, including logistic regression, decision trees, random forest, and others, were employed and thoroughly evaluated. The objective was to determine the algorithm that achieves the highest accuracy in predicting diabetes based on the dataset.

This research highlights the importance of ongoing research and model refinement in improving diabetes prediction and advancing healthcare outcomes. By exploring the potential of machine learning in diabetes detection, this project contributes to the broader understanding and utilization of data-driven approaches in healthcare.

## Existing Method:

Before developing the predictive model for diabetes in this project, several existing methods and approaches have been explored in the field of diabetes prediction. These existing methods include:

a.  **Medical Diagnostic Criteria:** Traditional medical diagnostic criteria are commonly used by healthcare professionals to identify individuals at risk of diabetes. These criteria involve assessing various factors such as fasting glucose levels, oral glucose tolerance tests, and HbA1c levels to diagnose diabetes or prediabetes. These criteria are based on established thresholds and guidelines provided by organizations like the American Diabetes Association (ADA) and the World Health Organization (WHO).

b.  **Risk Assessment Tools:** Various risk assessment tools have been developed to estimate an individual's risk of developing diabetes. These tools use statistical algorithms and predictive models based on factors such as age, BMI, family history, blood pressure, and glucose levels to calculate an individual's risk score. Examples of widely used risk assessment tools include the Finnish Diabetes Risk Score (FINDRISC) and the Diabetes Risk Calculator developed by the ADA.

c.  **Machine Learning Approaches:** With the advancements in machine learning techniques, researchers have applied these methods to diabetes prediction. Machine learning algorithms such as logistic regression, decision trees, support vector machines, random forest, and neural networks have been utilized to develop predictive models using large datasets containing various medical and lifestyle features. These models aim to accurately classify individuals into diabetes or non-diabetes categories based on their feature values.

d.  **Genetic Markers and Biomarkers:** Researchers have also explored the use of genetic markers and biomarkers in diabetes prediction. Genetic studies have identified certain genetic variants associated with increased susceptibility to diabetes. Biomarkers such as levels of specific hormones, cytokines, or inflammatory markers have also been investigated for their potential predictive value in diabetes.

e.  **Electronic Health Records (EHR):** Electronic health records contain a wealth of information about patients, including medical history, laboratory results, and demographic data. Researchers have utilized data mining and machine learning techniques to extract valuable insights from EHR data and develop predictive models for diabetes. These models aim to leverage the rich information available in EHRs to predict the likelihood of diabetes development.

Each of these existing methods has its advantages and limitations in terms of accuracy, scalability, and interpretability. The goal of this project is to

contribute to the existing methods by developing an accurate predictive model for diabetes using machine learning algorithms and evaluating its performance against these existing approaches.

# Proposed Method with Architecture:

The proposed method in this project involves implementing a machine learning architecture to develop an accurate predictive model for diabetes. This architecture consists of several components and steps aimed at effectively utilizing the dataset and optimizing the model's performance.

1. **Data Preprocessing:** The first step involves preprocessing the dataset to ensure its quality and suitability for training the model. This includes handling missing values, normalizing, or scaling features, and splitting the data into training and testing sets.

2. **Feature Selection:** In this step, relevant features are selected from the dataset based on their importance in diabetes prediction. Techniques such as correlation analysis or feature importance ranking are used to identify the most informative features of the selected model. (Decision tree)

3. **Model Selection:** Multiple machine learning algorithms, including logistic regression, decision trees, random forest, and others, are considered for building the predictive model. Each algorithm is evaluated based on its performance metrics, such as accuracy, precision, recall, and F1 score.

4. **Model Training:** The selected machine learning algorithm is trained on the pre-processed dataset using the training set. The model learns from the data and adjusts its internal parameters to optimize its performance in predicting diabetes.

5. **Model Evaluation:** The trained model is then evaluated using the testing set. Its performance is assessed using various evaluation metrics, including accuracy, precision, recall, and F1 score. This step helps determine the model's effectiveness and generalization ability.

6. **Model Deployment:** Once the model is deemed satisfactory, it can be deployed for real-world applications. Healthcare professionals and individuals can utilize the model to assess the risk of diabetes for individuals based on their relevant medical features.

**The proposed architecture** combines the strengths of different machine learning algorithms to create a predictive model that accurately identifies individuals at risk of diabetes. By leveraging the dataset and following the outlined steps, this architecture aims to optimize the model's performance and provide valuable tools for healthcare professionals and individuals in disease prevention and management.

## Methodology:

In this section we shall learn about the various classifiers used in machine learning to predict diabetes. We shall also explain our proposed methodology to improve the accuracy. Three different methods were used in this paper. The different methods used are defined below. The output is the accuracy metrics of the machine learning models. Then, the model can be used in prediction. The algorithms I used in my project are:

a. **Logistic regression**: Logistic regression is a widely used algorithm for binary classification tasks, making it suitable for

predicting diabetes. It models the relationship between the input features and the probability of the outcome class. By fitting a logistic function to the training data, it estimates the probabilities and predicts the presence or absence of diabetes.

b. **Random forest**: Random Forest is an ensemble learning algorithm that combines multiple decision trees to make predictions. It leverages the concept of bagging, where each tree is trained on a random subset of the data. By aggregating the predictions of the individual trees, random forest produces more robust and accurate predictions. It is known for its ability to handle high-dimensional datasets and reduce overfitting.

c. **Decision tree:** The decision tree algorithm is a powerful tool for predictive modelling. It builds a tree-like structure where each internal node represents a feature and each leaf node corresponds to a class label. By recursively splitting the data based on feature thresholds, it creates a set of rules for classifying instances. Decision trees can capture complex relationships between features and the target variable.

**Model Selection**: Different machine learning algorithms, including logistic regression, decision trees, random forest are considered for model selection. Each algorithm is evaluated based on its suitability for the dataset and its ability to accurately predict diabetes. Ultimately the Decision Tree algorithm was chosen as it gave the highest accuracy.

**Validation and Testing**: The final trained model is validated and tested using a separate unseen portion of the dataset. This ensures that the model can generalize well to new data and provides reliable predictions for diabetes.

**Results Analysis**: The results obtained from the model's evaluation and testing phases are analysed to determine the accuracy and effectiveness of the predictive model. This analysis may involve comparing the performance of different algorithms, examining the contribution of individual features, and identifying any limitations or areas for further improvement.

# Implementation:

| Library Name | Version |
|---|---|
| pandas | 2.0.3 |
| seaborn | 0.12.2 |
| matplotlib | 3.7.2 |
| numpy | 1.25.1 |

1. **Data Loading:** The code utilizes the **pandas** library to load the diabetes dataset, which is stored in a CSV file. The dataset contains various medical features related to diabetes, such as age, glucose levels, blood pressure, BMI, etc. By using the **read_csv()** function, the code reads the dataset into a pandas DataFrame. It then displays the first few rows of the dataset using the **head()** function to provide an initial glimpse of the data.

2. **Data Preprocessing:** After loading the dataset, the code performs data preprocessing steps to ensure its quality and suitability for model training. It starts by exploring the dataset using functions like **info()** and **describe()**, which provide information about the number of instances, data types, and summary statistics of the features. The code then checks for missing values by using the **isnull()** function, which returns a Boolean DataFrame indicating the presence of null values in each cell. To handle these missing values, the code applies mean imputation, filling the null values with the mean value of the corresponding feature using the **fillna()** function.

**3. Data Scaling:** To ensure that all features are on a similar scale and to prevent certain features from dominating the model training process, the code performs feature scaling using the StandardScaler from the scikit-learn library. This scaling process transforms each feature to have a mean of 0 and a standard deviation of 1, making the dataset more suitable for machine learning algorithms. The scaled data is then stored in a new DataFrame called **data_scalled.**

**4. Model Building:** The code proceeds to build three different machine learning models: **Decision Tree, Random Forest**, and **Logistic Regression**. Each model is built using the respective algorithms provided by **scikit-learn**. For each model, the code splits the pre-processed dataset into training and testing sets using the **train_test_split()** function, with 80% of the data used for training and 20% for testing. This division allows the models to be trained on a portion of the data and evaluated on unseen data to assess their performance.

**5. Model Evaluation:** Once the models are trained, the code evaluates their performance using various evaluation metrics. For each model, it makes predictions on the testing set using the **predict()** function and compares them to the true labels. The code then calculates metrics such as **accuracy, precision, recall, and F1 score** using functions from the metrics module of **scikit-learn**. These metrics provide insights into how well the models are performing in predicting the presence or absence of diabetes.

**6. Feature Importance:** In the case of the **Decision Tree** model, the code determines the importance of each feature in making predictions. It retrieves the feature importance using the **feature_importances_** attribute of the trained Decision Tree classifier. To visualize the importance of each feature, the code creates a bar graph using the **plot()** function from pandas, allowing easy interpretation of the relative significance of each feature.

**7. Model Saving:** The code saves the trained Decision Tree model using the **pickle** module. By using the **dump()** function, the model is serialized and stored in a file, making it possible to load and reuse the model in the future without needing to retrain it.

**8. Model Prediction:** To demonstrate the prediction capability of the saved Decision Tree model, the code loads the model using the **loads()** function from pickle. With the model loaded, the code makes a prediction on a sample input, using the **predict()** method, to showcase how the model can be used to predict the outcome of diabetes based on the given features.

**9. Final Checks:** Finally, the code performs additional checks on the original dataset to ensure that the model is working correctly and providing accurate predictions. It displays the first few rows of the dataset using the **head()** function and the last few rows using the **tail()** function. These checks allow verification that the dataset remains intact and that the trained model is capable of producing reliable predictions.

# Conclusion

Overall, this project contributes to the advancement of diabetes prediction using machine learning. The developed model has the potential to enhance diabetes risk assessment and improve healthcare outcomes. Further research and refinement of the models will continue to enhance the accuracy and effectiveness of diabetes prediction, leading to early intervention and personalized care for individuals at risk of diabetes.

*About me:*
*Name: Trisha Pal*
*3rd year CSE student at NITMAS*
*trishapal845@gmail.com*