

End-to-End Predictive Predictive Healthcare Healthcare System

This presentation explores the development of a machine learning project aimed at predicting diabetes using the PIMA Diabetes dataset. This project utilizes various machine learning techniques to analyze and understand the data, with the ultimate goal of improving predictive accuracy for diabetes detection and empowering healthcare through data-driven insights.

Name – Naman Dixit Batch – 50
Sap ID – 500125539 Sem – 3



Project Overview

The project aims to develop a predictive model capable of identifying individuals at risk of developing diabetes. The PIMA Diabetes Dataset provides valuable information about patients, including their glucose levels, blood pressure, BMI, and other relevant attributes.

Objective

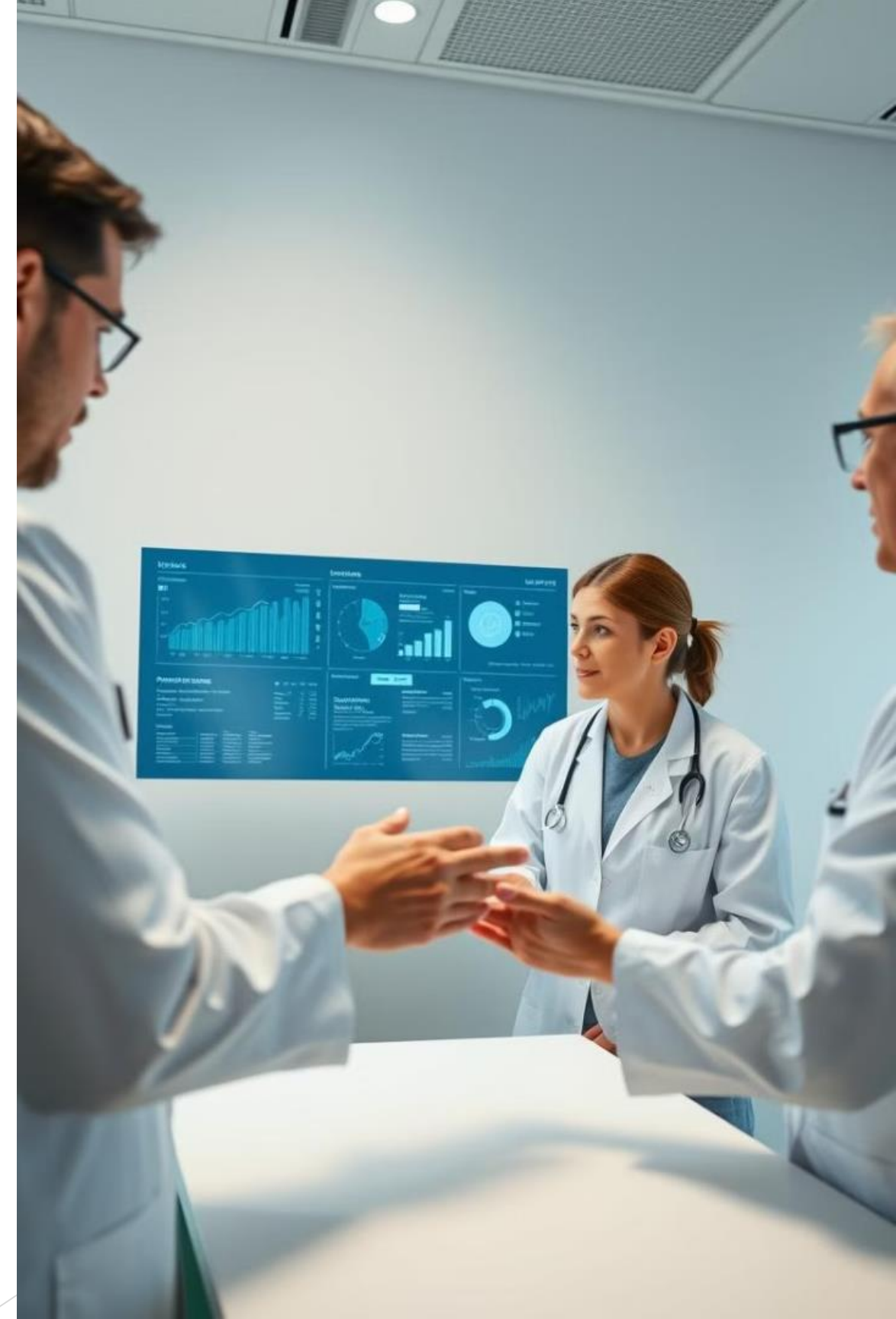
Build a model to predict diabetes using the PIMA Diabetes Dataset.

Dataset

Utilizes the PIMA Diabetes Dataset, containing patient data on glucose levels, blood pressure, BMI, etc.

Goal

Enhance the accuracy of predicting diabetes for improved patient care.



Libraries and Tools

The project leverages various Python libraries to facilitate data manipulation, visualization, and model building. These libraries provide essential tools for data exploration, preprocessing, model training, and evaluation.



Python

The primary programming language used for data analysis and model development.



NumPy

A fundamental library for numerical computing, providing high-performance arrays and mathematical functions.



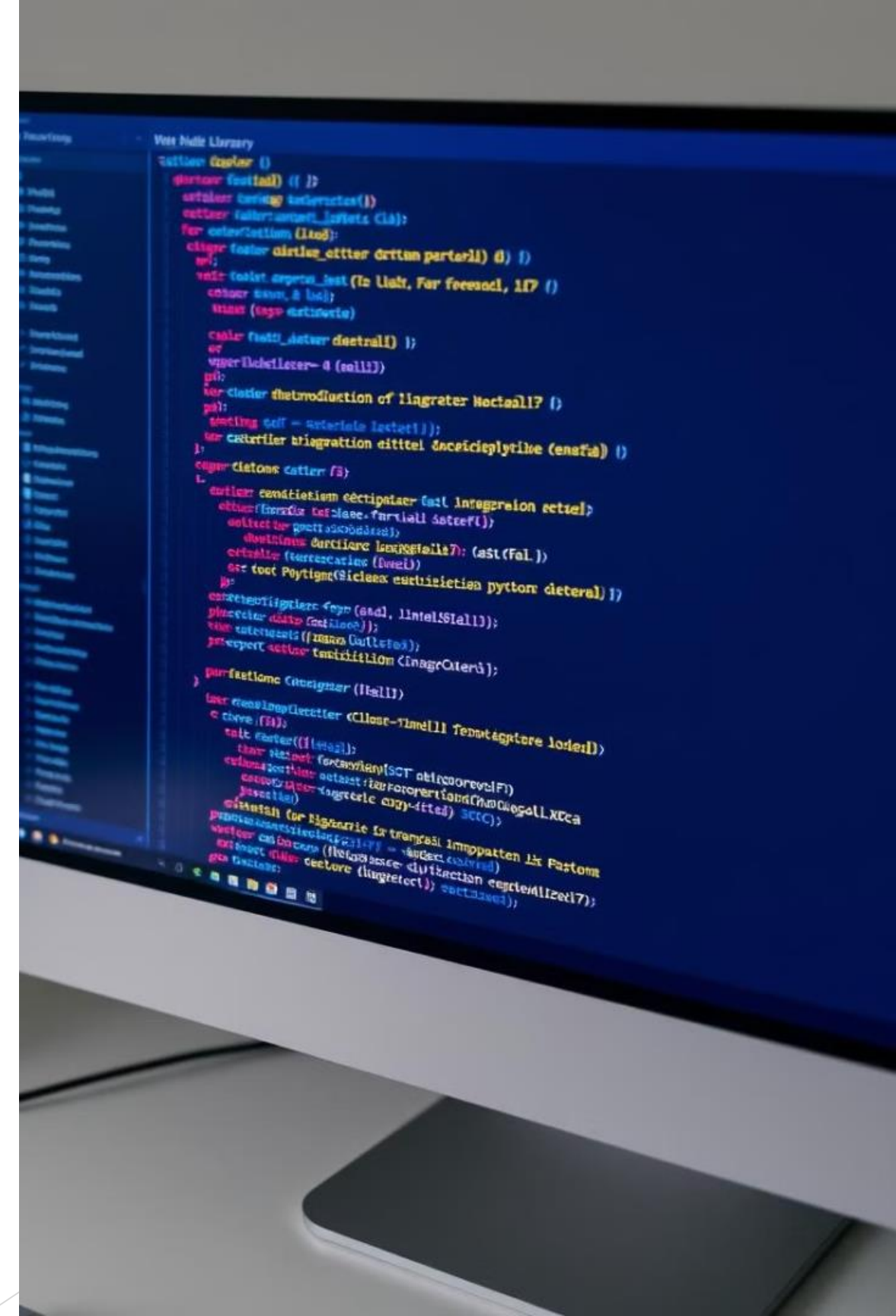
Pandas

A powerful library for data manipulation and analysis, providing data structures and functions for efficient data handling.



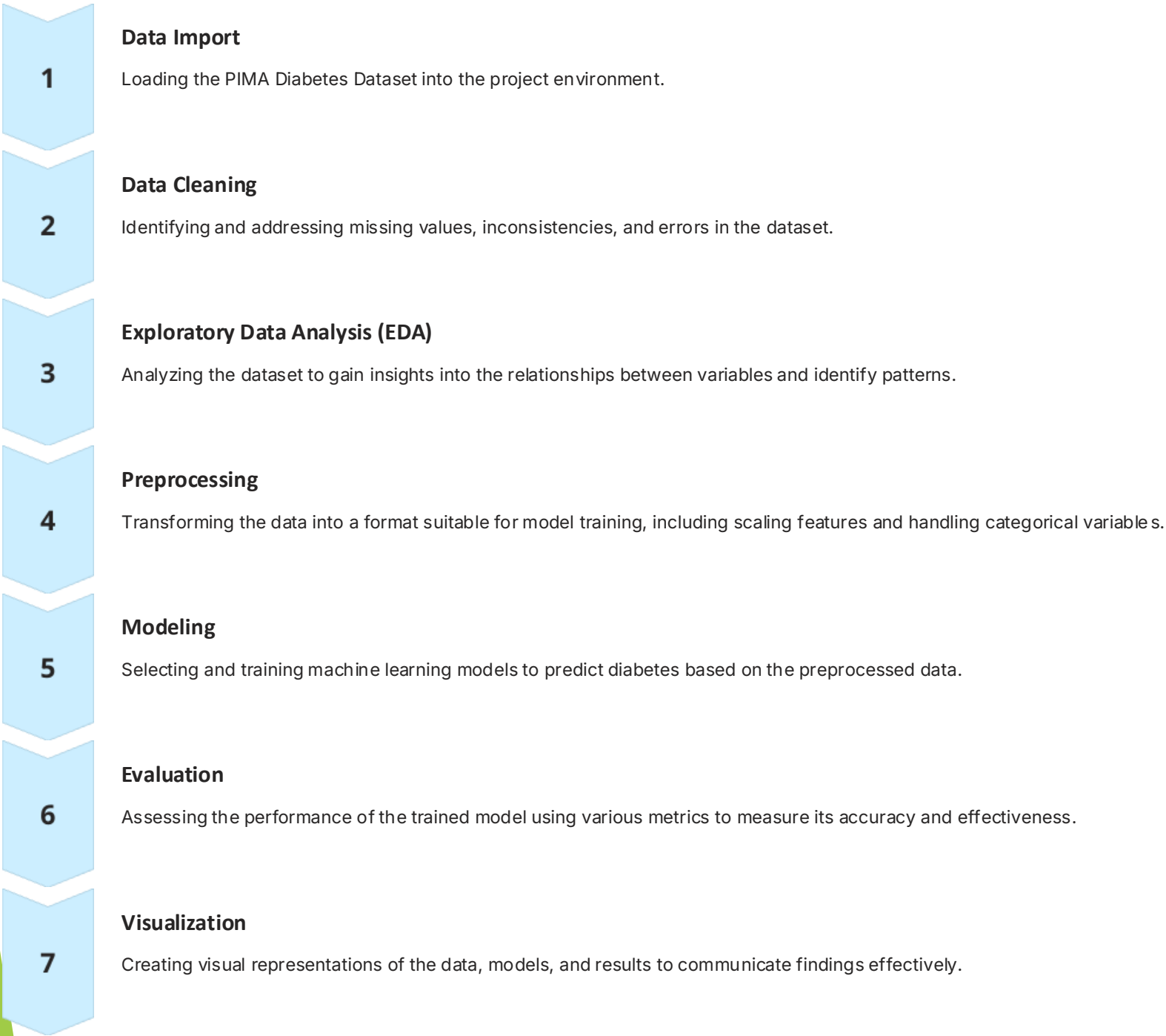
Scikit-learn

A comprehensive machine learning library offering various algorithms for classification, regression, clustering, and more.



Workflow Overview

The project workflow involves a series of steps, each contributing to the development of a robust and accurate predictive model.



Dataset Exploration

The PIMA Diabetes Dataset provides valuable insights into the characteristics of patients and their potential for developing diabetes.

Samples	768
Features	8
Missing Values	0



1 Key Statistics

The dataset comprises 768 samples and 8 features, with no missing values.

2 Data Distribution

The distribution of glucose levels, as shown in the histogram, reveals potential patterns and outliers.

Data Preprocessing

Data preprocessing is essential to prepare the dataset for model training and ensure accurate predictions. This step involves various techniques to address data quality and format issues.

Handling Missing Values

Addressing any missing data points by imputing or removing them, based on the characteristics of the dataset.

Scaling Features

Standardizing the range of features to ensure that all variables contribute equally to the model's predictions.

Balancing the Dataset

Addressing class imbalance by oversampling the minority class using techniques like SMOTE, to improve model performance.

Dimensionality Reduction

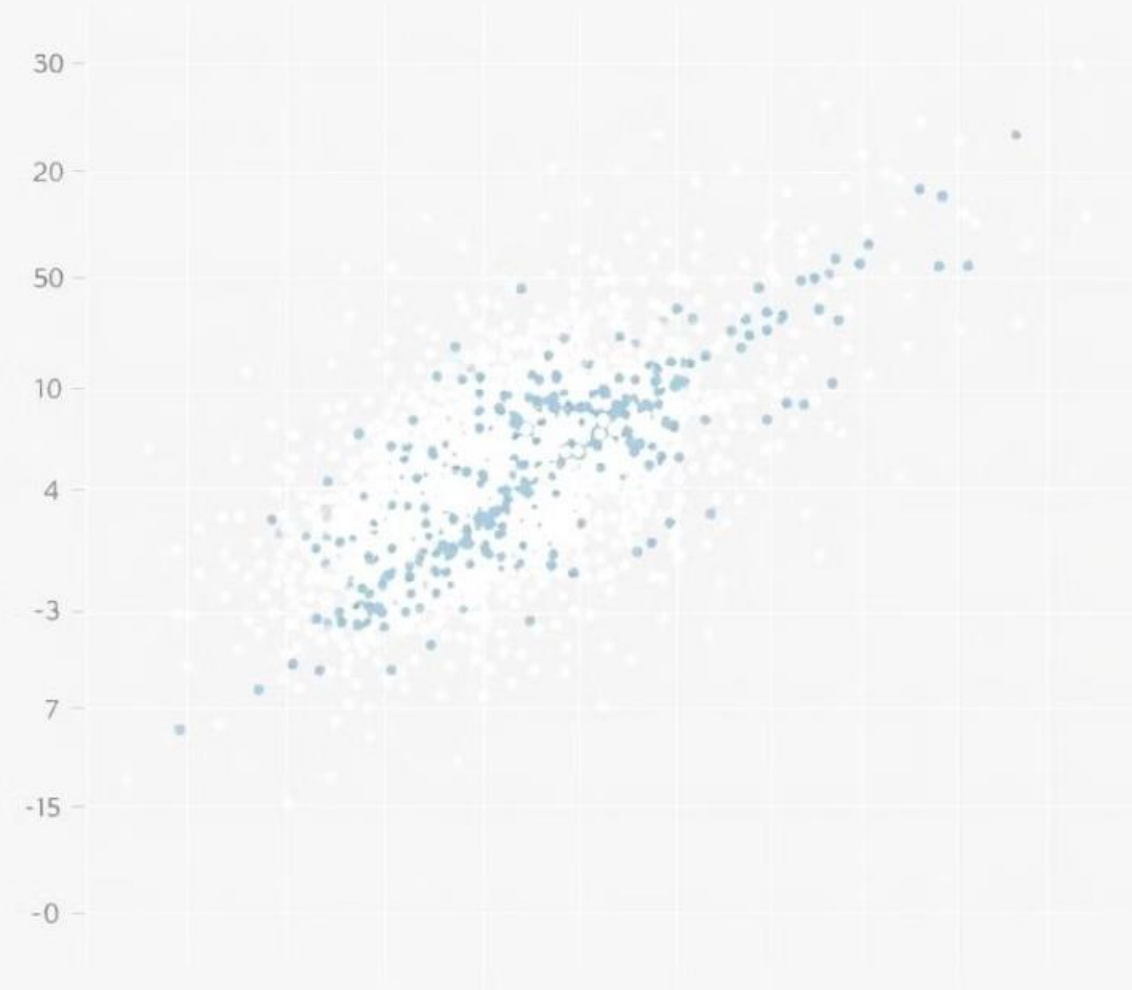
PCA is a technique used to reduce the dimensionality of the data by identifying the principal components, which capture the most variance in the dataset.

1 Principal Component Analysis (PCA)

Reduces the dimensionality of the data by identifying the principal components, which capture the most variance in the dataset.

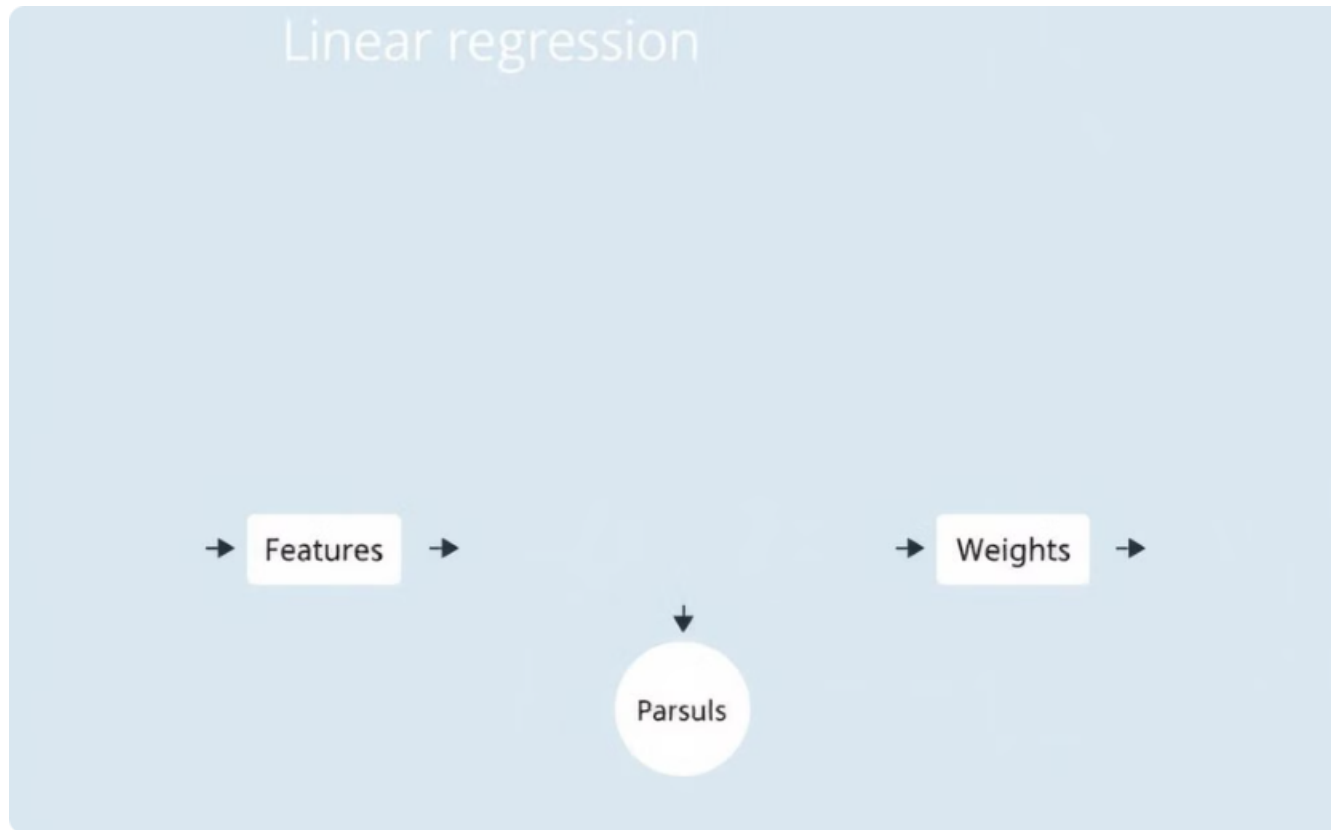
2 Visualization

The scatter plot illustrates the distribution of data points in the two-dimensional space defined by the first two principal components.



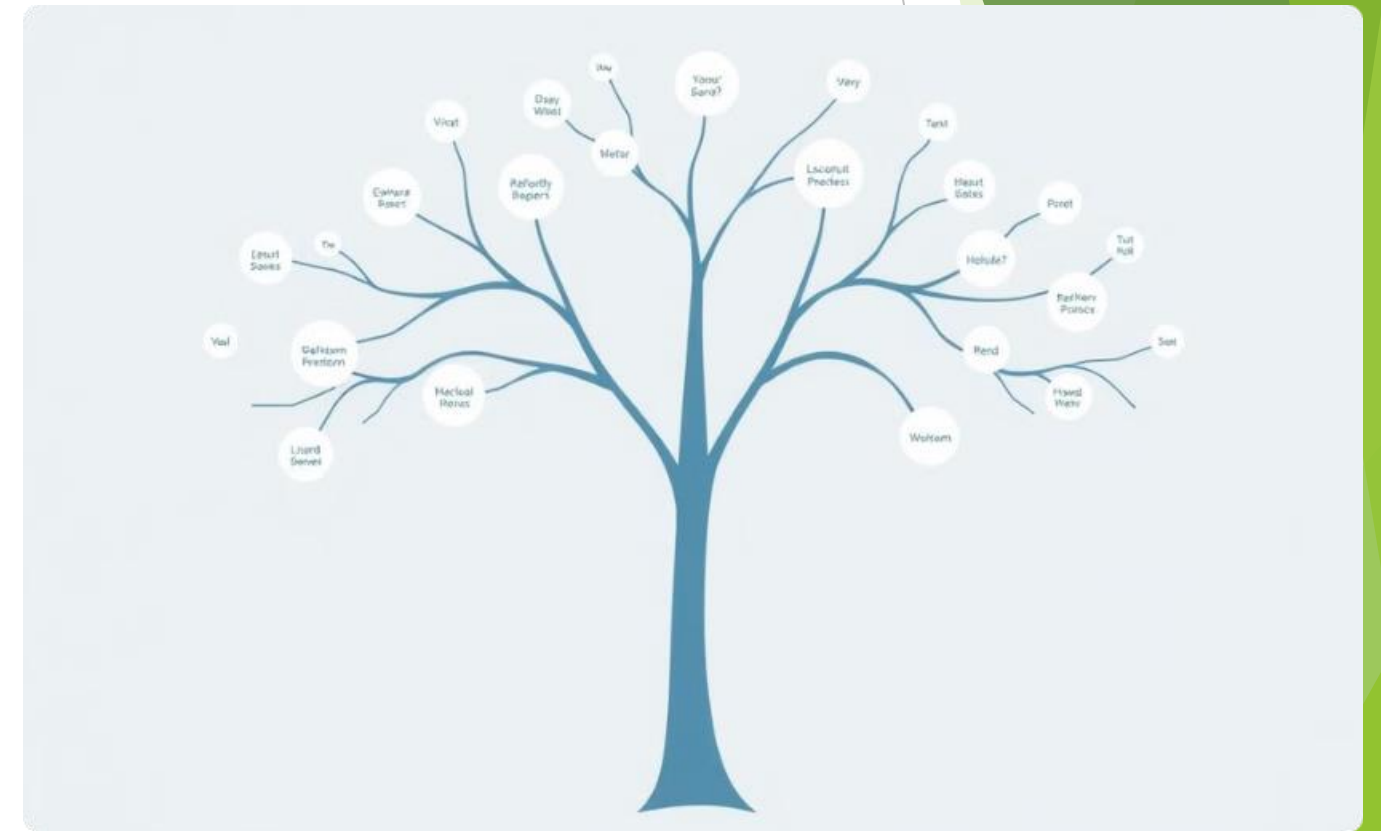
Model Building

The project employs various machine learning models to predict diabetes, each suited to specific aspects of the problem.



Linear Regression

A model used for predicting the risk score for diabetes, providing a continuous measure of risk.



Random Forest Classifier

A model used for classifying individuals as having or not having diabetes, providing a binary prediction.

Model Evaluation

The performance of the models is evaluated using various metrics to assess their accuracy, precision, and ability to identify true positive and negative cases.

Metric	Linear Regression	Random Forest
Mean Squared Error (MSE)	0.12	N/A
Accuracy	N/A	0.78
Precision	N/A	0.75
Recall	N/A	0.80
F1-Score	N/A	0.77

	Accuracy matex	Accuracy naxex	Recall matex	Recallc matrix	Contic natie
	9.1.6.2	54.1.60	52.35.2	57.8.83	52.41
	9.1.9.5	38.04.3	30,50.2	32.918	50.97
ronegy	9.8.1.4	31.93.9	46,69.5	57,510	41.17
elerage	9.7.2.1	38.55.5	38.48.9	51.3.38	20.02
eleract	9.0.8.6	36.33.5	58.45.5	35.9.89	94.8.3
elerract	9.7.1.2	29.3.15	59,53.3	367.39	55.3.3
vidersact	7.0.8.8	30.7.23	43.55.6	3.45.14	54.4.5
elersatt	7.1.7.3	59.5,78	59,14.3	5.88.60	52.4.2
elariate	9.1.3.2	56.3.12	33,43.1	386.24	74.1.0
elersaic	5.9.7.6	20.7.13	50.13.7	3.85.72	52.3.9
scorege	9.3.3.9	50.3.13	52.3,76	55.2.18	54.4.1
	7.0.9.8	30.3.34	54.13.1	3.85.38	55.3.1
olierage	9.0.2.4	38,7.25	50,59.5	3.71.32	31.5.3
elerratt	9.9.2.6	30.3.43	40.16.6	4.38.34	51.1.5
eleragt	9.7.8.3	25.53.3	55.3.13	59.5.18	46.5.3
eleragit	9.0.3.4	59.3.31	30,73.3	565.35	95.1.1
eleragit	9.4.391	30.5.23	30,55.5	363.16	37.5.0
eleragic	9.1.3.6	55.3.25	49.4.1.9	5.83.39	527.1
elersats	9.3.2.5	59.1.33	25.4.1.7	363.35	51.4.2
	9.0.4.8	38.9.54	30.50.5	30.9.86	37.8.5