



# DIABETES PREDICATION

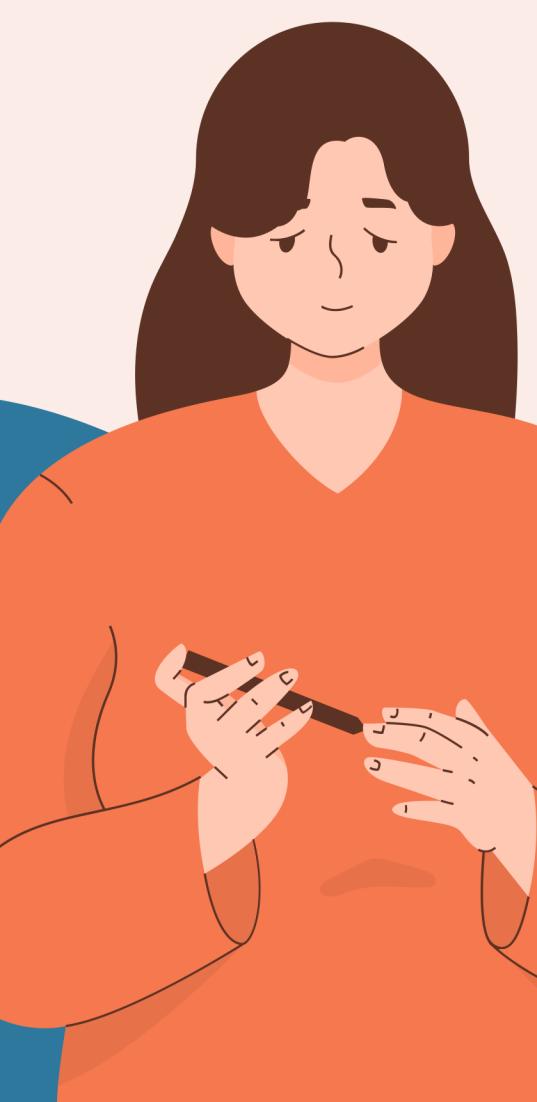
Name - Shweta Wadhwa  
Roll No. - 2210990847  
G-13

# PROBLEM STATEMENT

"Develop a machine learning model to predict the likelihood of diabetes in individuals based on various health indicators such as glucose level, blood pressure, BMI, and others. The model aims to assist healthcare professionals in early identification and intervention for individuals at risk of developing diabetes, thereby improving preventive healthcare strategies.

# KEY FEATURES

## Data Visualization



In my project, I visualize the dataset using plots like boxplots, pairplots, histograms, and heatmaps. These visuals aid in understanding feature distribution, relationships, and correlations efficiently.

## Model Evaluation

I assess machine learning models for diabetes prediction through techniques like train-test splitting, hyperparameter tuning, and performance metrics including accuracy, confusion matrix, and classification report

## Predictive Modeling:

My project aims to create a predictive model for diabetes using health indicators. I employ machine learning algorithms like kNN for classification tasks, distinguishing individuals as diabetic or non-diabetic based on their health data.

# DESIRED OUTCOMES

**EARLY DETECTION** - PROJECT'S PREDICTIVE MODEL FOR DIABETES AIDS IN EARLY DETECTION, FACILITATING TIMELY INTERVENTION AND REDUCING THE RISK OF COMPLICATIONS.

**PERSONALIZED HEALTHCARE** - YOUR ACCURATE PREDICTIVE MODEL ENABLES TAILORED INTERVENTIONS AND TREATMENTS, IMPROVING HEALTHCARE STRATEGIES AND PATIENT OUTCOMES.

**PUBLIC HEALTH IMPACT** - IMPLEMENTING YOUR PREDICTIVE MODEL IN HEALTHCARE SETTINGS INFORMS TARGETED PREVENTIVE STRATEGIES, REDUCING THE OVERALL BURDEN OF DIABETES IN COMMUNITIES BY IDENTIFYING HIGH-RISK INDIVIDUALS.



# SUPERVISED LEARNING

SUPERVISED LEARNING IS THE TYPE OF MACHINE LEARNING IN WHICH MACHINES ARE TRAINED USING WELL "LABELLED" TRAINING DATA, AND ON BASIS OF THAT DATA, MACHINES PREDICT THE OUTPUT. THE LABELLED DATA MEANS SOME INPUT DATA IS ALREADY TAGGED WITH THE CORRECT OUTPUT

SUPERVISED LEARNING IS A PROCESS OF PROVIDING INPUT DATA AS WELL AS CORRECT OUTPUT DATA TO THE MACHINE LEARNING MODEL. THE AIM OF A SUPERVISED LEARNING ALGORITHM IS TO FIND A MAPPING FUNCTION TO MAP THE INPUT VARIABLE(X) WITH THE OUTPUT VARIABLE(Y)

# **SUPERVISED LEARNING CAN BE FURTHER DIVIDED INTO TWO TYPES -**

## **REGRESSION -**

- LINEAR REGRESSION
- REGRESSION TREES
- NON-LINEAR REGRESSION
- BAYESIAN LINEAR REGRESSION
- POLYNOMIAL REGRESSION

## **CLASSIFICATION -**

- RANDOM FOREST
- DECISION TREES
- LOGISTIC REGRESSION
- SUPPORT VECTOR MACHINES
- SUPERVISED LEARNING
- REGRESSION
- CLASSIFICATION
- SUPERVISED

**SUPERVISED  
LEARNING**

**REGRESSION**

**CLASSIFICATION**

# k-Nearest Neighbors (kNN)

The kNN algorithm is a simple and effective algorithm for classification tasks, where the model classifies new data points based on the majority class of their  $k$  nearest neighbors in the feature space. It's a versatile algorithm suitable for various types of data and can be easily implemented for classification tasks like predicting diabetes status based on health indicators.



# PIP

PIP IS A PACKAGE MANAGER FOR PYTHON PACKAGES, OR MODULES. IT IS THE STANDARD TOOL FOR INSTALLING PYTHON PACKAGES AND THEIR DEPENDENCIES IN A SECURE MANNER.

PIP IS AN ACRONYM OF "PIP INSTALL PACKAGES".

# LIBRARIES USED

- PANDAS
- NUMPY
- MATPLOTLIB
- SEABORN



# PANDAS

PANDAS IS A PYTHON LIBRARY USED FOR DATA MANIPULATION AND ANALYSIS. HERE ARE SOME OF THE THINGS YOU CAN DO WITH PANDAS:

- \* LOAD DATA FROM A VARIETY OF SOURCES, INCLUDING CSV FILES, EXCEL FILES, SQL DATABASES, AND HTML TABLES.
- \* CLEAN AND MANIPULATE DATA, INCLUDING DEALING WITH MISSING VALUES, OUTLIERS, AND DUPLICATE DATA.
- \* PERFORM DATA ANALYSIS, INCLUDING CALCULATING STATISTICS, GROUPING DATA, AND CREATING VISUALIZATIONS.
- \* MERGE AND JOIN DATAFRAMES.
- \* SAVE DATA TO A VARIETY OF FORMATS, INCLUDING CSV FILES, EXCEL FILES, AND SQL DATABASES.

# NUMPY

NUMPY (NUMERICAL PYTHON) IS A PYTHON LIBRARY USED FOR WORKING WITH ARRAYS. IT ALSO HAS FUNCTIONS FOR WORKING IN DOMAIN OF LINEAR ALGEBRA, FOURIER TRANSFORM, AND MATRICES. HERE ARE SOME OF THE FEATURES OF NUMPY:

- \* IT PROVIDES A HIGH-PERFORMANCE MULTIDIMENSIONAL ARRAY OBJECT
- \* IT HAS A LARGE COLLECTION OF MATHEMATICAL FUNCTIONS TO OPERATE ON THESE ARRAYS.
- \* IT ALLOWS YOU TO WORK WITH MATRICES AND VECTORS.
- \* IT IS EASY TO USE AND HAS A WELL-DOCUMENTED API.

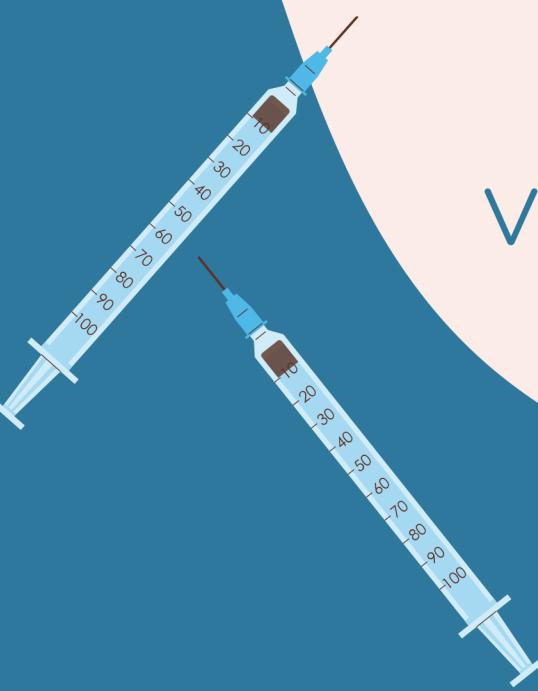
IT IS COMPATIBLE WITH OTHER PYTHON LIBRARIES, SUCH AS PANDAS AND MATPLOTLIB



# MATPLOTLIB

MATPLOTLIB IS A PYTHON LIBRARY THAT PROVIDES A COMPREHENSIVE SET OF TOOLS FOR CREATING STATIC, ANIMATED, AND INTERACTIVE VISUALIZATIONS IN PYTHON. IT IS A POPULAR CHOICE FOR DATA VISUALIZATION BECAUSE IT IS FLEXIBLE, POWERFUL, AND EASY TO USE.

MATPLOTLIB CAN BE USED TO CREATE A WIDE VARIETY OF PLOTS, INCLUDING HISTOGRAMS, SCATTER PLOTS, BAR CHARTS, LINE PLOTS, PIE CHARTS, BOX PLOTS, AND MORE. IT CAN ALSO BE USED TO CREATE MORE COMPLEX VISUALIZATIONS, SUCH AS 3D PLOTS, HEATMAPS, AND CONTOUR PLOTS.



# SEABORN

SEABORN IS A POWERFUL PYTHON VISUALIZATION LIBRARY THAT ENHANCES STATISTICAL GRAPHICS WITH ATTRACTIVE DEFAULT STYLES AND COLOR PALETTES.

IT PROVIDES BEAUTIFUL DEFAULT STYLES AND COLOR PALETTES TO MAKE STATISTICAL PLOTS MORE ATTRACTIVE.

INTEGRATED WITH MATPLOTLIB AND PANDAS, SEABORN PRIORITIZES VISUALIZATION TO AID IN DATA EXPLORATION AND UNDERSTANDING. ITS DATASET-ORIENTED APIs FACILITATE SWITCHING BETWEEN VARIOUS VISUAL REPRESENTATIONS FOR A COMPREHENSIVE DATASET ANALYSIS.

# CONCLUSION

IN CONCLUSION, THIS PROJECT HAS SUCCESSFULLY LEVERAGED MACHINE LEARNING TECHNIQUES TO DEVELOP A PREDICTIVE MODEL FOR DIABETES BASED ON HEALTH INDICATORS. THROUGH THOROUGH DATA VISUALIZATION AND EXPLORATION, KEY INSIGHTS INTO THE DISTRIBUTION, RELATIONSHIPS, AND CORRELATIONS AMONG VARIOUS FEATURES HAVE BEEN UNCOVERED. BY EMPLOYING THE K-NEAREST NEIGHBORS ALGORITHM, THE MODEL ACHIEVES ACCURATE CLASSIFICATION OF INDIVIDUALS INTO DIABETIC OR NON-DIABETIC CATEGORIES, FACILITATING EARLY DETECTION AND INTERVENTION. THE PERSONALIZED APPROACH ENABLED BY THE MODEL ALLOWS HEALTHCARE PROFESSIONALS TO TAILOR INTERVENTIONS BASED ON INDIVIDUAL RISK PROFILES, POTENTIALLY IMPROVING PATIENT OUTCOMES AND REDUCING THE OVERALL BURDEN OF DIABETES IN COMMUNITIES.

MOVING FORWARD, THE PREDICTIVE MODEL DEVELOPED IN THIS PROJECT HOLDS PROMISING IMPLICATIONS FOR PUBLIC HEALTH INITIATIVES AND PREVENTIVE STRATEGIES. BY IMPLEMENTING THE MODEL IN HEALTHCARE SETTINGS, TARGETED INTERVENTIONS CAN BE DEPLOYED MORE EFFECTIVELY, LEADING TO BETTER MANAGEMENT OF DIABETES AND ITS ASSOCIATED COMPLICATIONS. FURTHER RESEARCH COULD EXPLORE THE INTEGRATION OF ADDITIONAL FEATURES OR ADVANCED MACHINE LEARNING ALGORITHMS TO ENHANCE THE MODEL'S PREDICTIVE PERFORMANCE AND BROADEN ITS APPLICABILITY IN DIVERSE HEALTHCARE CONTEXTS.

**THANK  
YOU**