

## **Dataset Link:**

The dataset used in the project can be found on Kaggle:

<https://www.kaggle.com/datasets/uciml/pima-indians-diabetes-database>

## **Description of Project:**

The code you provided is a Jupyter notebook that uses logistic regression to predict the likelihood of diabetes in patients. The goal of the project is to develop a machine learning model that can accurately predict whether a patient is likely to have diabetes based on several medical factors.

The notebook uses a dataset of patients that includes information such as age, sex, BMI, blood pressure, and glucose level. The dataset is preprocessed and cleaned before being used to train and test the logistic regression model.

The notebook walks through several steps of the machine learning process, including data exploration, data preprocessing, feature selection, model training and evaluation, and hyperparameter tuning. The goal is to develop a model that can accurately predict the likelihood of diabetes in new patients.

## **A brief explanation of the outputs of the Credit risk modeling Project**

1. Preprocessed dataset: The notebook outputs a clean, preprocessed version of the dataset that has been transformed and cleaned to prepare it for the machine learning algorithms.
2. Exploratory data analysis (EDA) outputs: The notebook includes several visualizations and summary statistics that help to explore the dataset and understand the relationships between variables. This includes histograms, scatter plots, box plots, and correlation matrices.
3. Model performance metrics: The notebook outputs the performance metrics of the logistic regression model such as accuracy, precision, recall, and F1 score. It may also include the confusion matrix that shows the number of true positive, true negative, false positive, and false negative classifications.
4. Predictions: The notebook outputs the predicted probabilities and the binary classification (diabetes or no diabetes) for each patient in the testing dataset.
5. Visualization of the model results: The notebook includes visualizations such as ROC (receiver operating characteristic) curves, precision-recall curves, and calibration plots to assess the model's performance.

6. Model parameters: The notebook includes the model coefficients or weights learned by the logistic regression algorithm, which can help to understand which features were most important in predicting diabetes.

### **Description of Output:**

The output of the notebook is a set of insights and recommendations for predicting the likelihood of diabetes in patients. The notebook provides visualizations and insights into the data, as well as code for training and evaluating the logistic regression model.

In addition to the model itself, the notebook also provides evaluation metrics for the model, such as accuracy, precision, and recall, as well as visualizations of the model performance over time. The notebook also provides feature importance rankings, which can help identify which factors are most important in predicting the likelihood of diabetes.

Based on the analysis, the notebook provides several recommendations for improving the accuracy of diabetes predictions, such as using different machine learning algorithms, improving feature selection methods, and gathering more data to better train the model.

### **Instructions on How to Run the Code/Project/File:**

1. To run the code in the notebook, you will need to have Jupyter Notebook installed on your computer. Once you have installed the Jupyter Notebook, you can download the notebook from the Kaggle website and open it in the Jupyter Notebook.
2. Before running the code, you will need to make sure that you have downloaded the necessary data files and saved them in the correct directory. The notebook provides instructions on how to download the data files and where to save them.
3. Once you have downloaded the data files and opened the notebook in Jupyter Notebook, you can run each cell of the notebook by clicking on the cell and then clicking the "Run" button in the toolbar or by using the keyboard shortcut "Shift + Enter".
4. It is recommended that you run the code cells in order, as some cells depend on the output of earlier cells.

**Note: Make sure to update the file paths in the code cells to match the location of the downloaded dataset and kernel files on your local machine.**