# Project Brief

**PROJECT NAME:** Diabetes Prediction Project – Classification

**BY:** Mukesh Kumar
**DATE:** 24 June 2024

| | |
|---|---|
| **PLAN** | <ul><li>Import data into the workspace.</li><li>Open the data and analyse the dependent and independent variables. Here the dependent variable is **labeled** and **categorical.**</li><li>Select the correct model to solve the problem. In this case we use classification type model to predict the outcome.</li><li>The model should have minimum false negative and false positive prediction outcomes. False positive could lead to unnecessary medical interventions and stress on the other hand false negative is more concerning, this can lead to undiagnosed cases and delayed or missed treatments, potentially causing serious health consequences.</li></ul> |
| **ANALYSE** | <ul><li>This is a diabetes dataset, there are 9 variables in this dataset with 8 independent and one dependent variables. The dependent variable is categorical where 0 indicates the person does not have diabetes and 1 indicates the has diabetes.</li><li>The data is slightly **imbalance** with 65% – 0 values and 35% – 1 values</li><li>Address the irregularities like null values, zeros, duplicates, outliers.</li><li>Understanding the relationships between variables in our dataset.</li><li>Standardise and encode the continuous and categorical variables.</li><li>Make ready the data to fit in the model.</li></ul> |
| **CONSTRUCT** | <ul><li>Choose the algorithm. We are using **Logistic regression, Support vector machine classifier, Random Forest Classifier and Decision Tree Classifier** algorithms and then select the most accurate one.</li><li>Train our model on the data.</li><li>Test the model on the test subset</li><li>Evaluate the prediction algorithm. Check the accuracy score, cross validation score and confusion matrix.</li></ul> |
| **EXECUTE** | <ul><li>After evaluating all the algorithms we found that Logistic Regression, Support Vector Machine Classifier and Random Forest Classifier have **accuracy 80%.** But **Random Forest Classifier have less false negatives and false positives** than the other two.</li><li>Now we can share the prediction model with the stakeholder to test the model and receive the feedback.</li></ul> |