## Context:

This dataset is originally from the National Institute of Diabetes and Digestive and Kidney Diseases. The objective of the dataset is to diagnostically predict whether or not a patient has diabetes, based on certain diagnostic measurements included in the dataset. Several constraints were placed on the selection of these instances from a larger database. In particular, all patients here are females at least 21 years old of Pima Indian heritage.

## Problem Statement:

Build a model to accurately predict whether the patients in the dataset have diabetes or not?

## Dataset Description:

The dataset consists of several medical predictor variables and one target variable, Outcome. Predictor variables includes the number of pregnancies the patient has had, their BMI, insulin level, age, and so on.

**Pregnancies**: Number of times pregnant

**Glucose**: Plasma glucose concentration a 2 hours in an oral glucose tolerance test

**BloodPressure**: Diastolic blood pressure (mm Hg)

**SkinThickness**: Triceps skin fold thickness (mm)

**Insulin**: 2-Hour serum insulin (mu U/ml)

**BMI**: Body mass index (weight in kg/(height in m)^2)

**DiabetesPedigreeFunction**: Diabetes pedigree function

**Age**: Age (years)

**Outcome**: Class variable (0 or 1) 268 of 768 are 1, the others are 0

## Approach:

Following pointers will be helpful to structure your findings.

1.     Perform descriptive analysis. It is very important to understand the variables and corresponding values. We need to think through - Can minimum value of below listed columns be zero (0)? On these columns, a value of zero does not make sense and thus indicates missing value.

- Glucose
- BloodPressure
- SkinThickness
- Insulin
- BMI

        How will you treat these values?

2.     Visually explore these variables, you may need to look for the distribution of these variables using histograms. Treat the missing values accordingly.

3.      We observe integer as well as float data-type of variables in this dataset. Create a count (frequency) plot describing the data types and the count of variables.

4.     Check the balance of the data by plotting the count of outcomes by their value. Describe your findings and plan future course of actions.

5.     Create scatter charts between the pair of variables to understand the relationships. Describe your findings.

6.     Perform correlation analysis. Visually explore it using a heat map.

*(Note: Do not focus on visualization aspects when working with SAS)*

7.     Devise strategies for model building. It is important to decide the right validation framework. Express your thought process. Would Cross validation be useful in this scenario?

*(Note: if you are working with SAS, ignore this question and perform stratified sampling to partition the data. Create strata of age for this.)*

8.     Apply an appropriate classification algorithm to build a model. Compare various models with the results from KNN.

*(Note: if you are working with SAS, ignore this question. Apply logistic regression technique to build the model.)*

9.     Create a classification report by analysing sensitivity, specificity, AUC(ROC curve) etc. Please try to be as descriptive as possible to explain what values of these parameter you settled for? any why?

10. Create a dashboard in tableau by choosing appropriate chart types and metrics useful for the business. The dashboard must entail the following:

    a) Pie chart to describe the diabetic/non-diabetic population

    b) Scatter charts between relevant variables to analyse the relationships

    c) Histogram/frequency charts to analyse the distribution of the data

    d) Heatmap of correlation analysis among the relevant variables

    e) Create bins of Age values – 20-25, 25-30, 30-35 etc. and analyse different variables for these age brackets using a bubble chart.