

Understanding Factors that Influence Diabetes in Adult Female Pima Indians

Aditya Pendyala

2023-11-09

Introduction:

This dataset “diabetes” has data from a study conducted by The National Institute of Diabetes and Digestive and Kidney Diseases, on 768 adult female Pima Indians living near Phoenix. The purpose of the study was to investigate factors related to diabetes. The dataset has 9 variables:

1. pregnant: number of times pregnant
2. glucose: plasma glucose concentration at 2 hours in an oral glucose tolerance test
3. diastolic: diastolic blood pressure (mm Hg)
4. triceps: triceps skin fold thickness (mm)
5. insulin: 2-Hour serum insulin (mu U/ml)
6. bmi: body mass index (weight in kg/(height in meters squared))
7. diabetes: diabetes pedigree function
8. age: age (years)
9. test: test whether the patient shows signs of diabetes (coded 0 if negative, 1 if positive)

Purpose:

The goal of this project is to try and find variables that influence DPF(Diabetes Pedigree Function) in adult female Pima Indians. Furthermore, to formulate a prediction model, that predicts the DPF(Diabetes Pedigree Function) in an adult female Pima Indian.

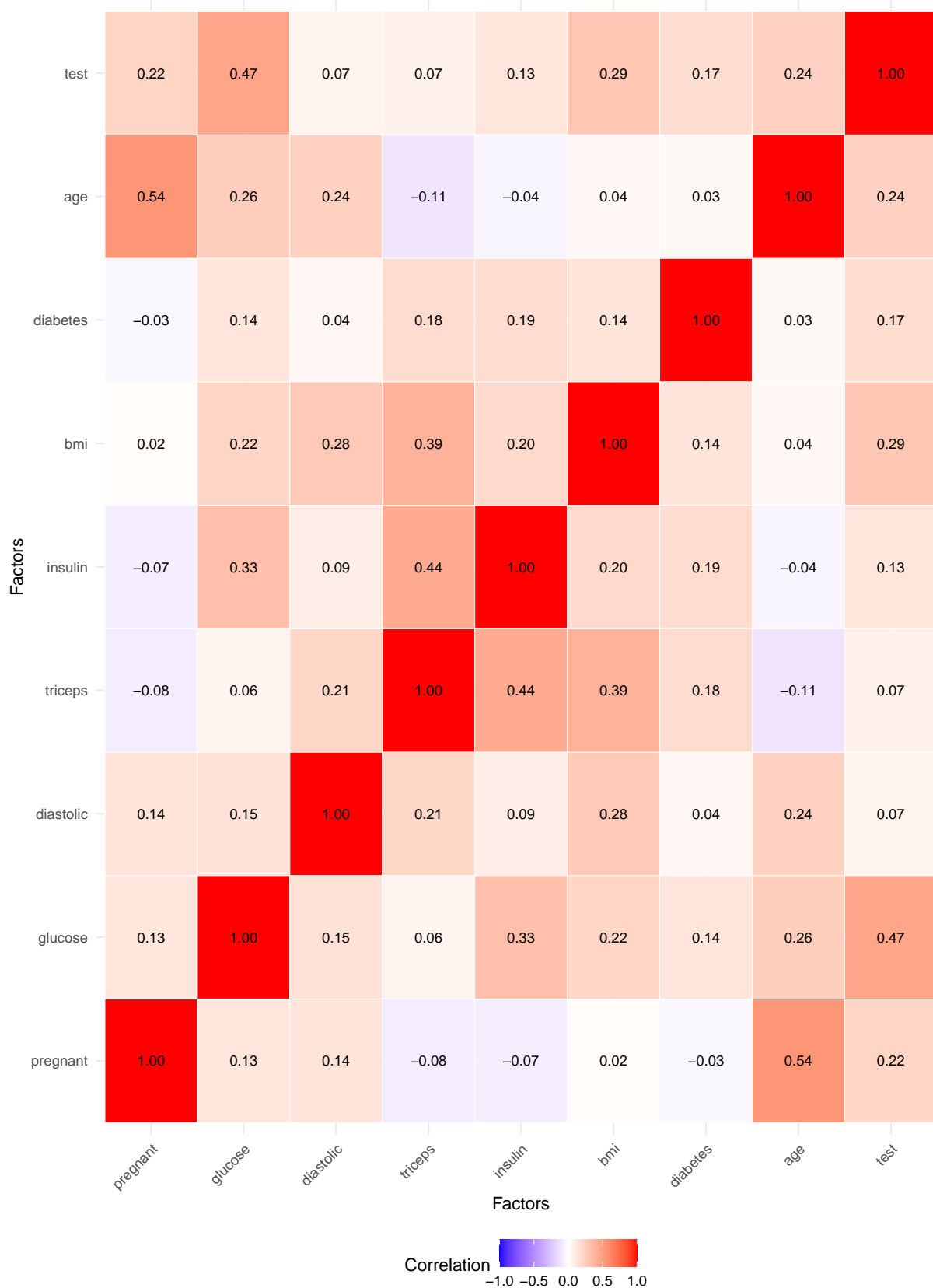
Exploring the dataset:

We start by creating two data-sets for dependent variables and independent variable(DPF). We aim to find a prediction model for DPF which is the “diabetes” variable.

Here the “diabetes” variable is the value from the Diabetes Pedigree Function (DPF), which calculates diabetes likelihood depending on the subject’s age and his/her diabetic family history.

First step in analyzing the dataset is to find the correlation of the factors. By doing this we can get a better understanding of which variables can be the possible factors of predicting DPF.

Correlation Heatmap for Diabetes factors



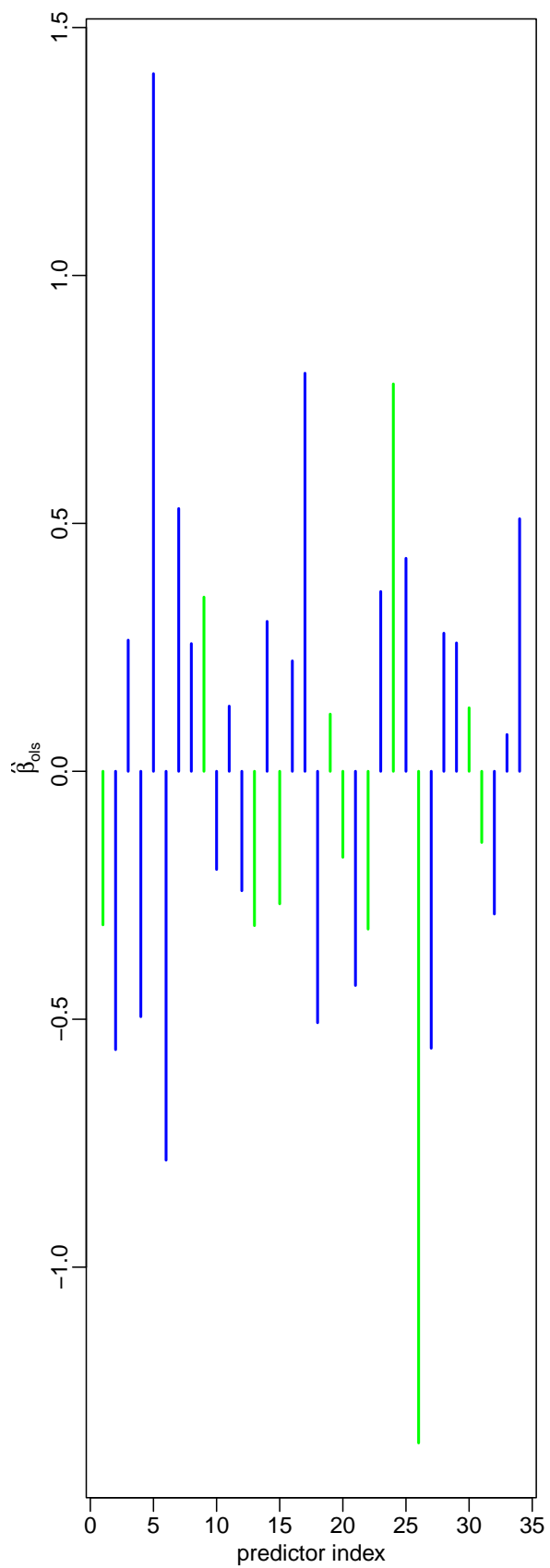
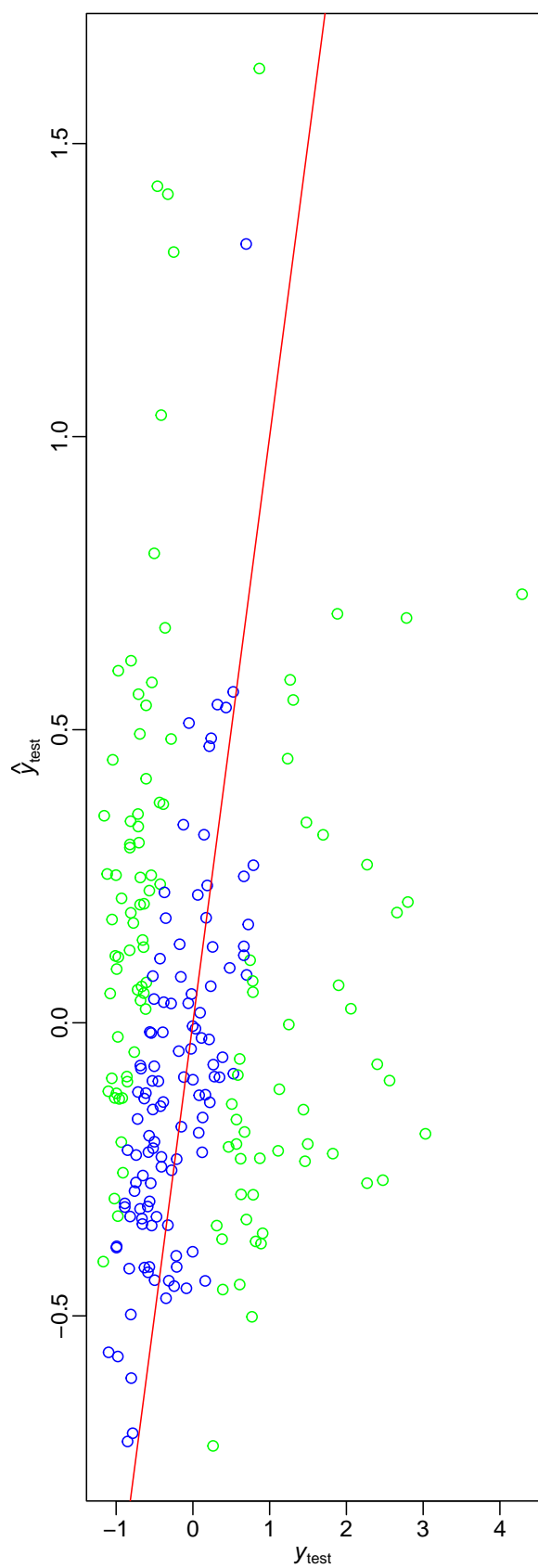
Looking at the heat map, we can see that 'diabetes' variable has some correlation with 'glucose', 'triceps', 'insulin', and 'bmi'. These are some variables that we can look out for while creating a prediction model for DPF.

Prediction and analysis:

To continue with predictions and analysis, we standardize the data before regression and create new variables which are products of each variable with one another. We also include squares of each variable.

After creating the new data, we continue by standardizing the data and continue by splitting the standardized data into testing and training data-sets(30-70 ratio).

First prediction model that we consider is, OLS. We perform OLS for DPF vs all the dependent variables. We now plot y-test and estimated-y-test. We are also plotting the predictor index for variables.



Looking at the 2nd plot, it's visible that many variables may not be truly associated with DPF.
So, to remove variables that aren't helpful for the plot, we perform a backward elimination procedure.

