# Supervised Learning on Diabetes Dataset

Siyi Mo

# Project purpose

The purpose of the project is to use supervised learning techniques to build a machine learning model that can predict whether a patient has diabetes or not, based on certain diagnostic measurements. The project involves three main parts: exploratory data analysis, preprocessing and feature engineering, and training a machine learning model.
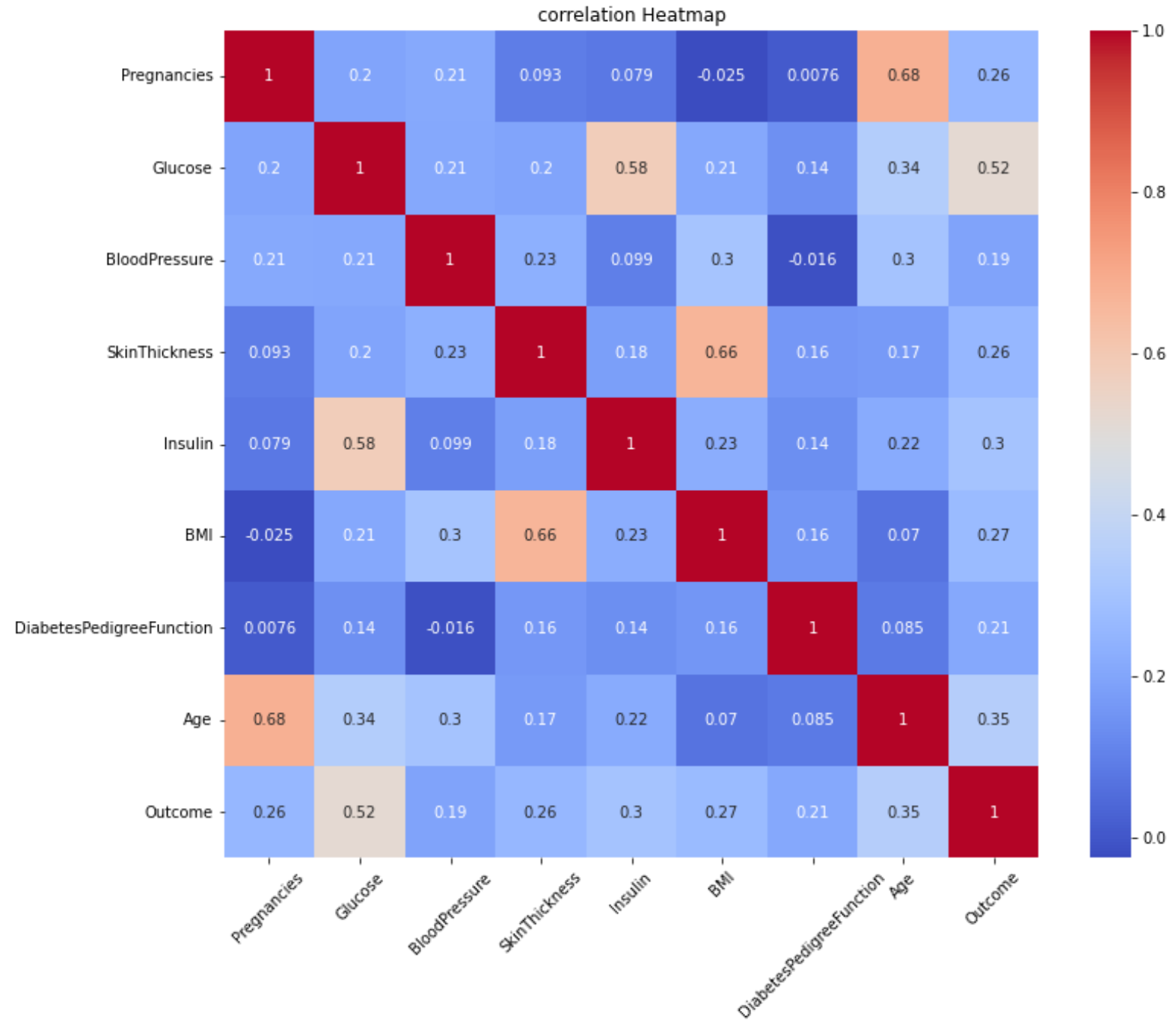
# EDA Results

Glucose has the highest correlation to the target variables among all the predictive variables.

Age and insulin also have relatively high level correlation to the target variable.

Some of the predictive variables are not mutually exclusive instead there are some correlation between the variables

People with diabetes have higher mean value in each metrics. ( pregnancies, glucoses, blood pressure, skin thickness, insulin, BMI and diabetes pedigree. )



correlation Heatmap

# Model Performance

|  | Random forest | Logistic regression |
|---|---|---|
| Accuracy score | 0.77 | 0.71 |
| Precision score | 0.9 | 073 |
| Recall score | 0.46 | 0.41 |
| F1 score | 0.61 | 0.52 |
| Roc auc score | 0.71 | 0.66 |

# Summary

The project is used random forest and logistic regression to predict diabetes. The process of the project consists of data cleaning (define missing values and remove outliers), data preprocessing (data scaling and transforming), modelling and model evaluation. Random forest performed better than logistic regression in this case based on the metrics scores returned by the evaluation metrics.