

STATS 4M03: Multivariate Analysis

Final Project

Diabetes Dataset

Submitted to

Dr. Eman M.S. Alamer

Department of Mathematics and Statistics

McMaster University

Hamilton, Ontario, Canada L8S 4K1

November 22, 2024

Reported by

Safi Khan (400402095)

Tonu Xu (400370837)

Xinyi Chen (400326045)

Rayyan Kazim (400406943)

Zesen Chen (400326984)

1 Introduction

1.1 Abstract

The study of diabetes is vital in understanding the progression of the disease and identifying key predictors. Throughout this paper, we perform data analysis on the diabetes dataset, primarily focusing on leveraging various supervised learning analysis methods that we learned in STATS 4M03\6M03: Multivariate Analysis. By using a variety of methods, our goal is to predict the onset of diabetes from detailed medical diagnostic measurements based on several contributing health factors — to uncover patterns and relationships between various clinical and lifestyle factors. Through this analysis, we hope to emphasize actionable insights for clinical decision-making and provide preventive strategies.

1.2 The Data

We will study the diabetes dataset from Rahman (2024), which contains 768 rows x 9 columns. Each column represents various health diagnostic metrics for predicting diabetes and each row corresponds to a unique patient record. Table 1 showcases a description for each of the columns in the dataset, and we will be using R as our main computing software (R Core Team, 2024).

Column	Description Of Column
Pregnancies	Integer: Number of times the patient has been pregnant.
Glucose	Integer: Plasma glucose concentration (mg/dL) after a 2-hour oral glucose tolerance test.
BloodPressure	Integer: Diastolic blood pressure (mm Hg).
SkinThickness	Integer: Triceps skinfold thickness (mm).
Insulin	Integer: 2-hour serum insulin (mu U/ml).
BMI	Float: Body mass index, defined as $\text{weight in kg}/(\text{height in m})^2$.
DiabetesPedigreeFunction	Float: A score indicating genetic predisposition to diabetes based on family history.
Age	Integer: Age of the patient (in years).
Outcome	Binary: Target variable where 1 indicates diabetes, and 0 indicates no diabetes.

Table 1: Column Descriptions of the Diabetes Dataset

1.2.1 Exploratory Data Analysis (EDA)

The true label of this dataset is Outcome, whereas the other variables are considered the predictors. Since our dataset is not that large, we determined that we should use all of the predictors as they all show reasonably high correlation with each other.

Figure 1 illustrates that there are a lot more individuals in the dataset who do not have diabetes. Figure 2 showcases that the ages listed in this dataset follow a right-skewed distribution, where majority of the individuals are aged 20-30.

Figure 3 illustrates the relatively high correlation between SkinThickness, BMI, and Insulin. We also note that Glucose is reasonably correlated with Insulin, BMI, and Age. Furthermore, Age is highly correlated with Pregnancies.

We cannot apply factor analysis on our dataset as it is not normally distributed. This can be confirmed by the Shapiro-Wilk test and normal QQ plot in Figure 4 (R Core Team, 2024).

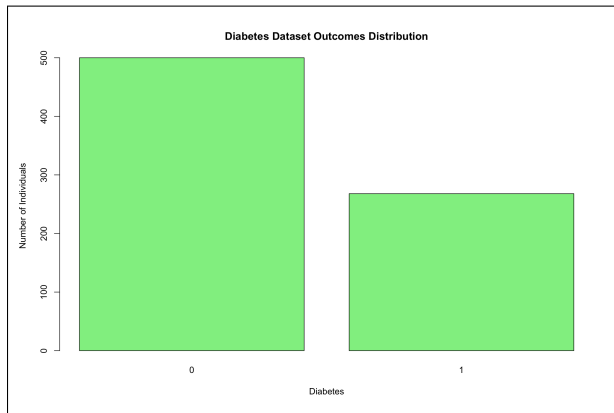


Figure 1: Distribution of Outcomes

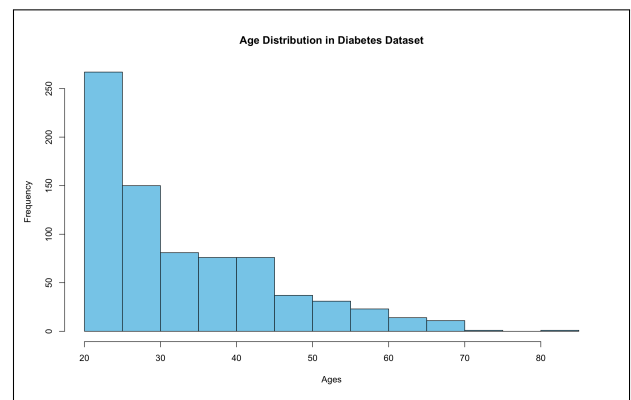


Figure 2: Age Distribution of the Dataset

1.2.2 Data Preparation

We prepared the data by scaling all the columns in the dataset except for column nine which is the response variable, Outcome, the true label of this dataset. We will split the data as follows: 75% for training to allow each of the models to learn effectively, and 25% will be reserved for testing to evaluate each of the models' performance on unseen data. Additionally, we will retain all variables in the model due to the relatively moderate to high association between each of the predictor variables.

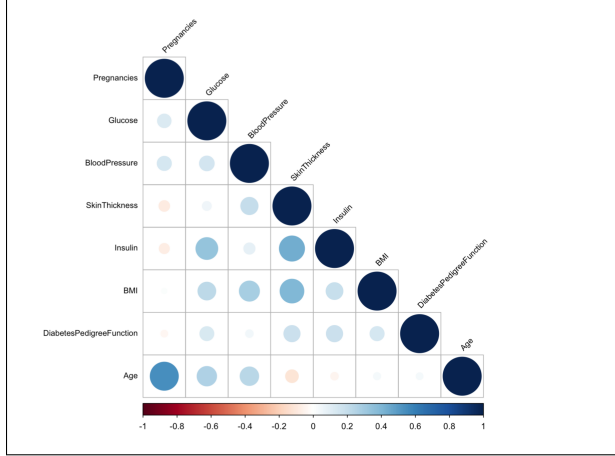


Figure 3: Correlation of Variables

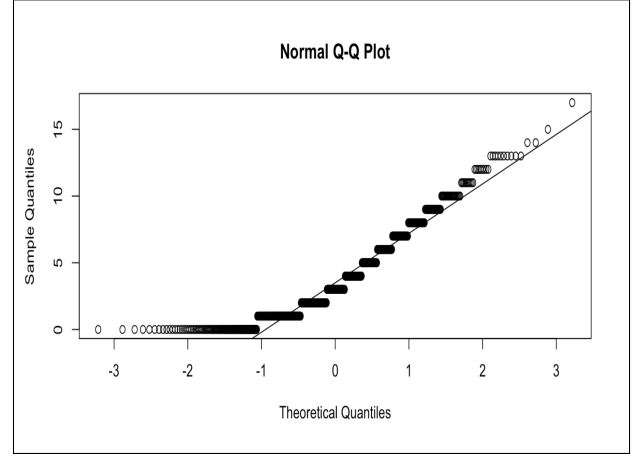


Figure 4: Normal Q-Q Plot

2 Methodology

2.1 k-Nearest Neighbours

We executed the k-Nearest Neighbours algorithm (k-NN), a non-parametric, supervised learning classifier that uses proximity to make classifications about the grouping of a dataset. (Peterson, 2009; Singh et al., 2016). Specifically, the algorithm takes an unlabelled observation and assigns it to the class that has the most labelled observations within its neighbourhood. (Alamer, 2024a). Additionally, we note that the optimal k-value will result in the best classification rate, based on the labelled points that are being treated as unlabelled by the algorithm (Alamer, 2024a).

2.2 Random Forest Classifiers

We also utilized the Random Forest Classifiers algorithm on our dataset (Zhou, 2012). Random Forest is a bootstrapping sampling method that combines the results of multiple decision trees to draw on a conclusion (Singh et al., 2016). The algorithm is an extension of the bagging algorithm which creates uncorrelated decision trees. For each tree, a random sample of \mathcal{M} predictors is taken at each decision tree split for each of the M trees, with the goal of improved predictive model accuracy (Alamer, 2024b).

2.3 Binary Logistic Regression

We performed binary logistic regression on our dataset to explore the relationship and significance between our predictor variables and the binary response variable, Outcome (Faraway, 2016). We identify the most significant predictors using backward elimination, a step-wise technique employed in regression to minimize the risk of overfitting. This means that predictors whose p -value is greater than $\alpha = 0.05$ level of significance will be removed and we will run the algorithm again, repeating this process until we end up with only statistically significant predictors in our model.

2.4 Boosting

Boosting is one of the ensemble methods which combines multiple weak learners, typically decision trees, to create a strong predictive model (Chen, 2015; Friedman, 2001). For this analysis, we used XGBoost (Extreme Gradient Boosting), which is efficient and well optimized for large datasets, to predict the presence of diabetes based on our predictor variables.

Model Parameters; Learning rate (`eta`): 0.1, Maximum tree depth (`max_depth`): 6, Evaluation metric: Area Under the ROC Curve (AUC), and Number of boosting rounds (`nrounds`): 100

3 Discussion

3.1 k-Nearest Neighbours

After tuning the algorithm using a 5-fold cross-validation to get the optimal parameters for the k-nearest neighbours algorithm (k-NN), we found that the best value for k is 3, which can be observed in Figure 5. This result indicates that if 2 out of the $k = 3$ neighbours have a positive response for diabetes, the individual will be assigned a positive response for diabetes, i.e. Outcome = 1. Conversely, if 2 out of the $k = 3$ neighbours have a non-positive response for diabetes, the individual will be assigned a non-positive response, i.e. Outcome = 0. After inputting $k = 3$ and running the k-NN algorithm, we found that the MCR given by k-NN classification is 0.2916667, i.e. $\sim 71\%$ of the individuals predicted to have a positive or non-positive response for diabetes were correctly classified.

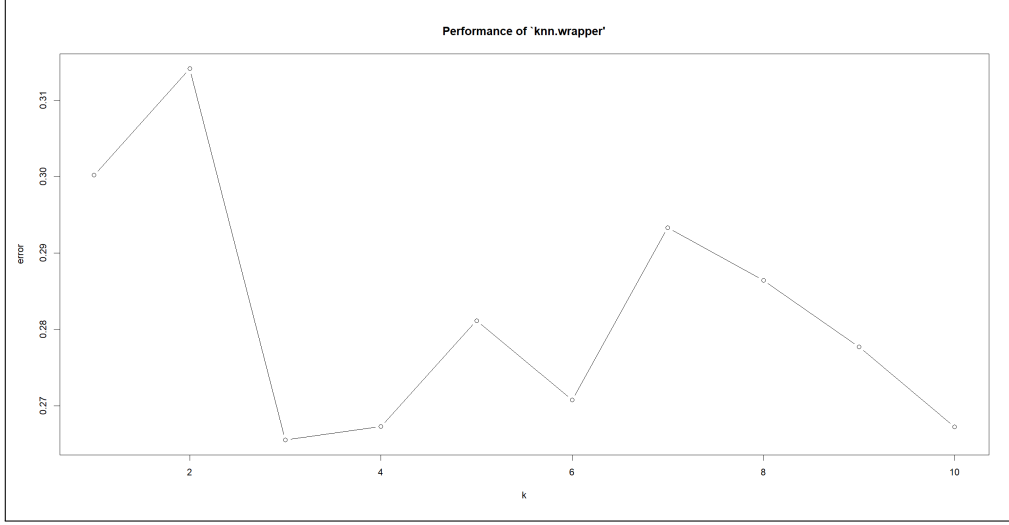


Figure 5: Output Plot From k-NN Classification

3.2 Random Forest Classifiers

After tuning the algorithm using a 5-fold cross-validation to get the optimal parameters for the random forest algorithm, we found that the best value for \mathcal{M} is 4 and the best value of M is 200.

After inputting $\mathcal{M} = 4$ and $M = 200$ into the random forest algorithm, we found that the MCR given by random forest classification is 0.28125, i.e. $\sim 72\%$ of the individuals predicted to have a positive or non-positive response for diabetes were correctly classified. Consequently, Figure 6 highlights that Glucose and BMI are the most important predictor variables, and excluding them from the model will lead to worse accuracy in predicting the Outcome. Conversely, SkinThickness and BloodPressure are the least important predictors and excluding them will not greatly affect the accuracy of the model.

3.3 Binary Logistic Regression

From (Table 2), we see that the coefficients for BloodPressure, SkinThickness, Insulin, and Age have high p-values, indicating that at $\alpha = 0.05$ level of significance, they are not significant in predicting the Outcome. At this step, the MCR given by the full binary logistic regression model is 0.3697917.

Table 3 showcases the output from the reduced Binary Logistic Regression Model. We observe that the p -values for all the predictors in the model are less than $\alpha = 0.05$, indicating that they're statistically significant in predicting the response variable, Outcome. Furthermore, the MCR has

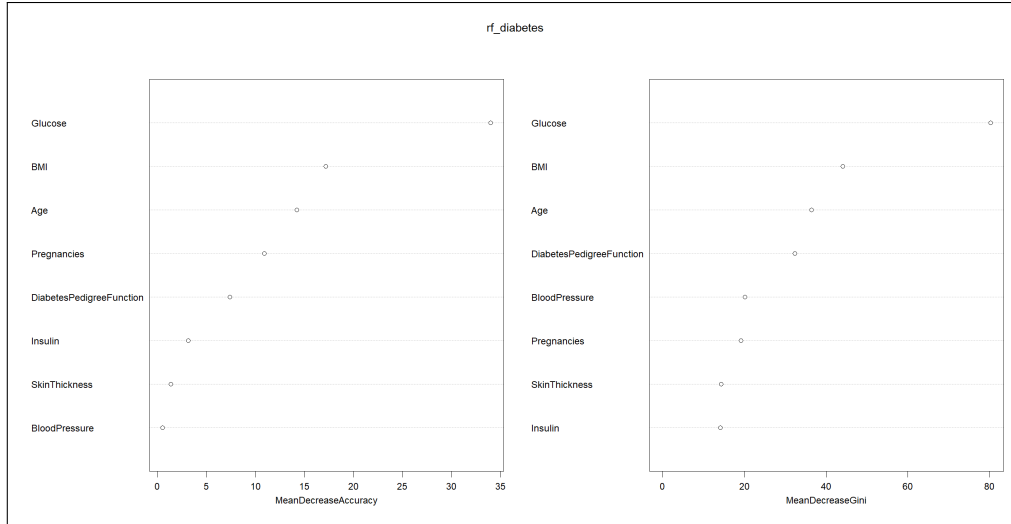


Figure 6: MeanDecreaseAccuracy & MeanDecreaseGini Plot from Random Forest

Coefficients	Estimate	Std. Error	Z-Value	Pr(> Z)
(Intercept)	-0.863576	0.111504	-7.745	9.57E-15
Pregnancies	0.411598	0.128257	3.209	1.33E-03
Glucose	1.003291	0.133428	7.519	5.51E-14
BloodPressure	-0.15522	0.122145	-1.271	0.2038
SkinThickness	-0.008376	0.123519	-0.068	0.94594
Insulin	-0.160385	0.120504	-1.331	1.83E-01
BMI	0.699525	0.135793	5.151	2.59E-07
DiabetesPedigreeFunction	0.295516	0.114215	2.587	0.00967
Age	0.212147	0.127142	1.669	0.0952

Table 2: Binary Logistic Regression Full Model Output

Coefficients:	Estimate	Std. Error	z value	Pr(> z)	Odds Ratio
(Intercept)	-8.118918	0.730397	-11.116	<2e-16	—
Pregnancies	0.154186	0.032151	4.796	1.62E-06	1.167
Glucose	0.030679	0.003781	8.113	4.92e-16	1.031
BMI	0.080086	0.016048	4.990	6.03e-07	1.083
DiabetesPedigreeFunction	0.863621	0.336274	2.568	0.0102	2.372

Table 3: Binary Logistic Regression Reduced Model Output with Odds Ratio

decreased to 0.2395833, i.e. $\sim 76\%$ of the individuals predicted to have a positive or non-positive response for diabetes were correctly classified, showcasing an improvement in the performance of the model.

Furthermore, Table 3 highlights the Odds Ratios (OR) for the predictors in the model. Note that each OR is greater than 1, suggesting increased likelihoods of the presence of diabetes with respect to the predictors. Specifically, the OR for Pregnancies is 1.167, indicating that females with a history of pregnancies are 1.167 times more likely to have diabetes than those without. Most significantly, the OR for DiabetesPedigreeFunction is 2.372, meaning that females who indicate a genetic predisposition to diabetes based on family history are 2.372 times more likely to develop diabetes than those without a genetic predisposition. The OR for the other predictors can be interpreted similarly.

3.4 Boosting

Boosting achieved an accuracy of 73.44%, with an AUC of 0.802. The AUC is the area under the ROC curve (Figure 7). The ROC Curve is a plot of the model's Sensitivity (True Positive Rate) against the Specificity (True Negative Rate), illustrating how well the model distinguishes between classes across all possible thresholds (Fawcett, 2006). The AUC is the area under the ROC curve, summarizing the model's overall ability to discriminate between positive and negative classes. An AUC of 0.5 represents random guessing, while an AUC closer to 1.0 indicates excellent performance (Robin et al., 2011; Fawcett, 2006).

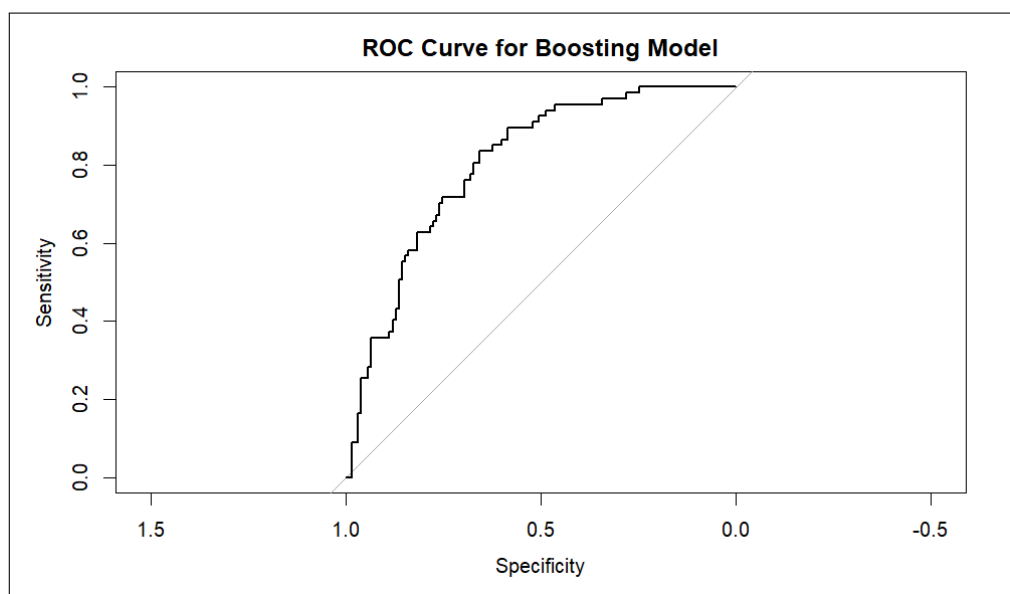


Figure 7: ROC Curve for Boosting Model

Boosting proved to be an effective method for diabetes classification through achieving an AUC of 0.802 and identifying Glucose, BMI, and Age as the most critical predictors in Figure 8.

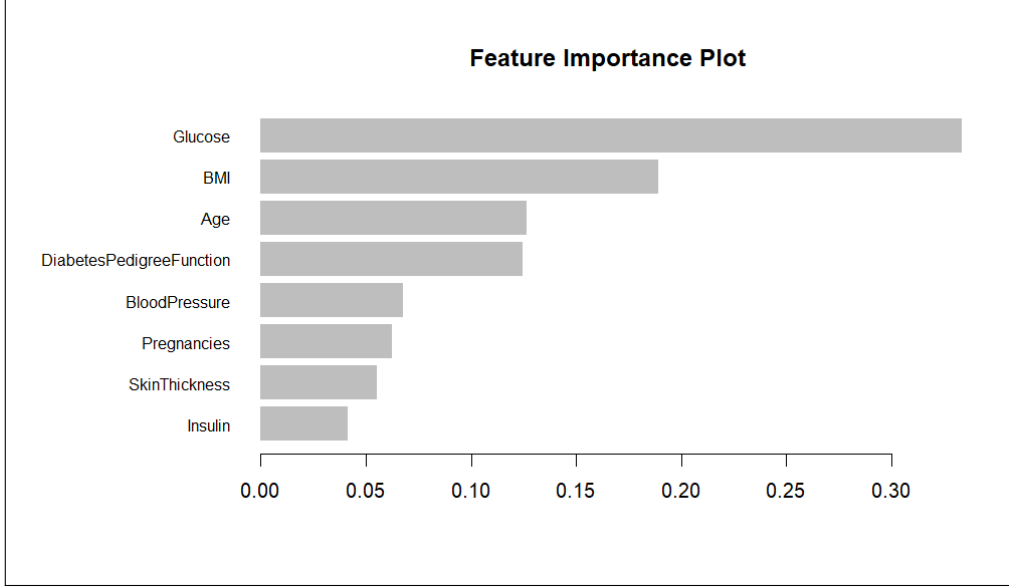


Figure 8: Feature Importance Plot

However, the moderate specificity given by the model (58.21%) suggests potential overemphasis on predicting positive cases, leading to false positives. Future work could improve model balance through hyperparameter tuning or additional data preprocessing techniques.

Summarizing the performance results from k-Nearest Neighbours classification, Binary Logistic Regression, Random Forest, and Boosting, we see that each of the models outputted an accuracy of 70.833%, 76.042%, 71.88%, and 73.44% respectively, indicating that Binary Logistic Regression was the best supervised learning analysis method in predicting the onset of diabetes.

4 Conclusion

Throughout this work, we presented various statistical supervised learning analysis methods that could be used to study the Diabetes Dataset from Rahman (2024). We effectively showcased the performance accuracy of the k-Nearest Neighbours algorithm, Binary Logistic Regression algorithm, and the ensemble methods: Random Forest and Boosting. Furthermore, the conclusions determined by each of the various methods align with the general scientific consensus on predictors for diabetes in females such as Glucose levels, BMI, Age, and whether an individuals family history has been shown to have a predisposition to diabetes or not.

Ultimately, although various supervised learning analysis methods have proven to be effective

in predicting the onset of diabetes based on several detailed medical diagnostic measurement predictor variables, we remain conscious of challenges associated with some of these methods, such as the large sample size required for logistic regression to output stable results, choosing ‘ k ’ in k-NN through computationally expensive techniques like cross-validation, and the high memory usage and computational cost associated with random forest (Singh et al., 2016).

5 Bibliography

References

- Alamer, E. (2024a). Lecture notes in multivariate analysis, lecture 14.
- Alamer, E. (2024b). Lecture notes in multivariate analysis, lecture 16.
- Chen, T. (2015). Xgboost: extreme gradient boosting. *R package version 0.4-2*, 1(4).
- Faraway, J. J. (2016). *Extending the linear model with R: generalized linear, mixed effects and nonparametric regression models*. Chapman and Hall/CRC.
- Fawcett, T. (2006). An introduction to roc analysis. *Pattern recognition letters*, 27(8):861–874.
- Friedman, J. H. (2001). Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232.
- Peterson, L. E. (2009). K-nearest neighbor. *Scholarpedia*, 4(2):1883.
- R Core Team (2024). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rahman, M. H. (2024). Kaggles predicting diabetes onset based on diagnostic measures.
- Robin, X., Turck, N., Hainard, A., Tiberti, N., Lisacek, F., Sanchez, J.-C., and Müller, M. (2011). proc: an open-source package for r and s+ to analyze and compare roc curves. *BMC Bioinformatics*, 12:77.
- Singh, A., Thakur, N., and Sharma, A. (2016). A review of supervised machine learning algorithms. In *2016 3rd International Conference on Computing for Sustainable Global Development (INDIACom)*, pages 1310–1315.
- Zhou, Z.-H. (2012). *Ensemble methods: foundations and algorithms*. CRC press.