

Instructor: Prof. Dr. Katharina Jahn

Introduction to Focus Areas Data Science WS22/23 Project 1

Group 7 Maïke, Jule, Carlos, Abhinav

Goals

1. Exploratory data analysis
2. Classification
3. Performance metrics





Heart Disease (Cleveland)

Instances: 303

Attributes: 14 out of 76

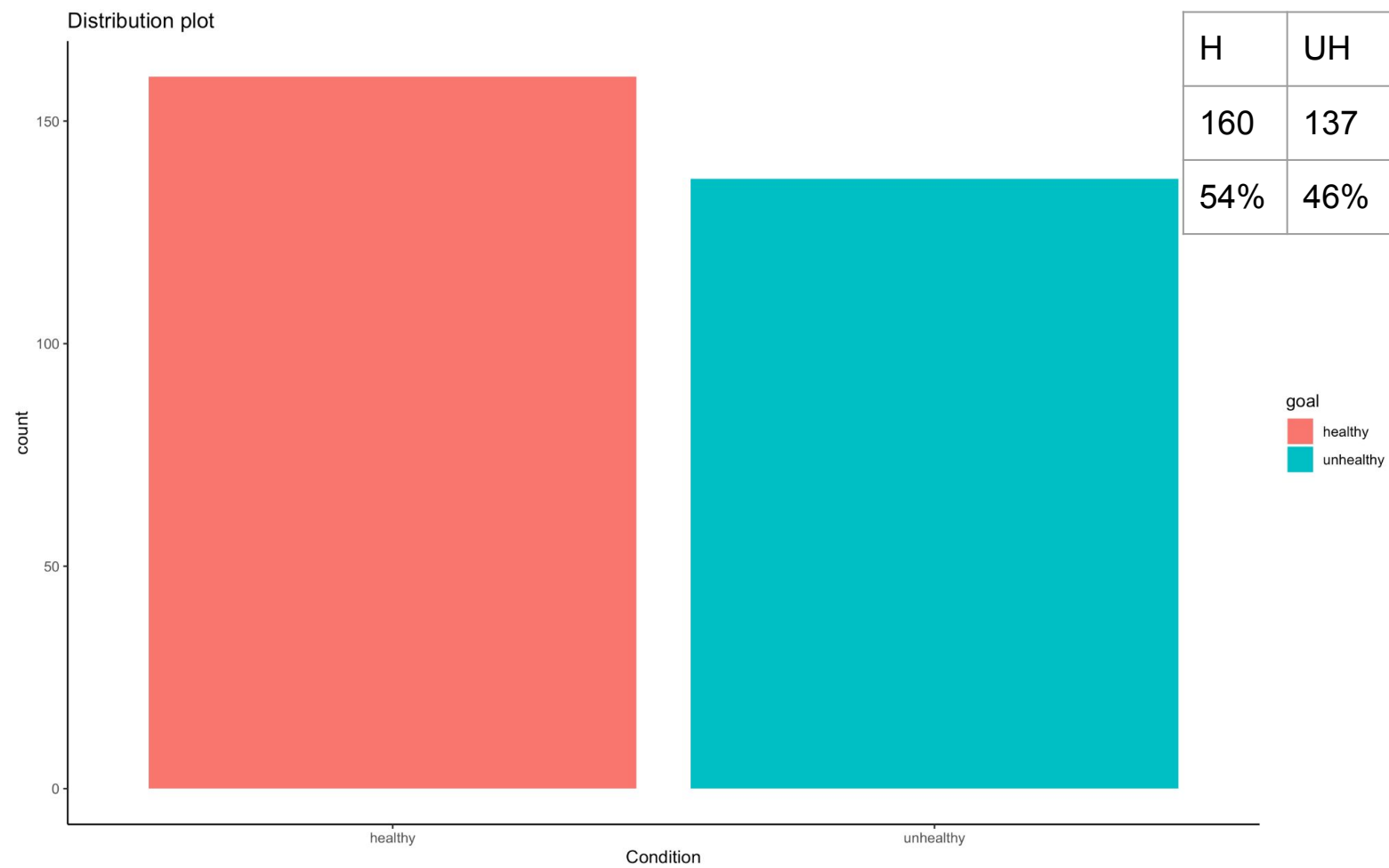
Type: categorical, numeric, integer

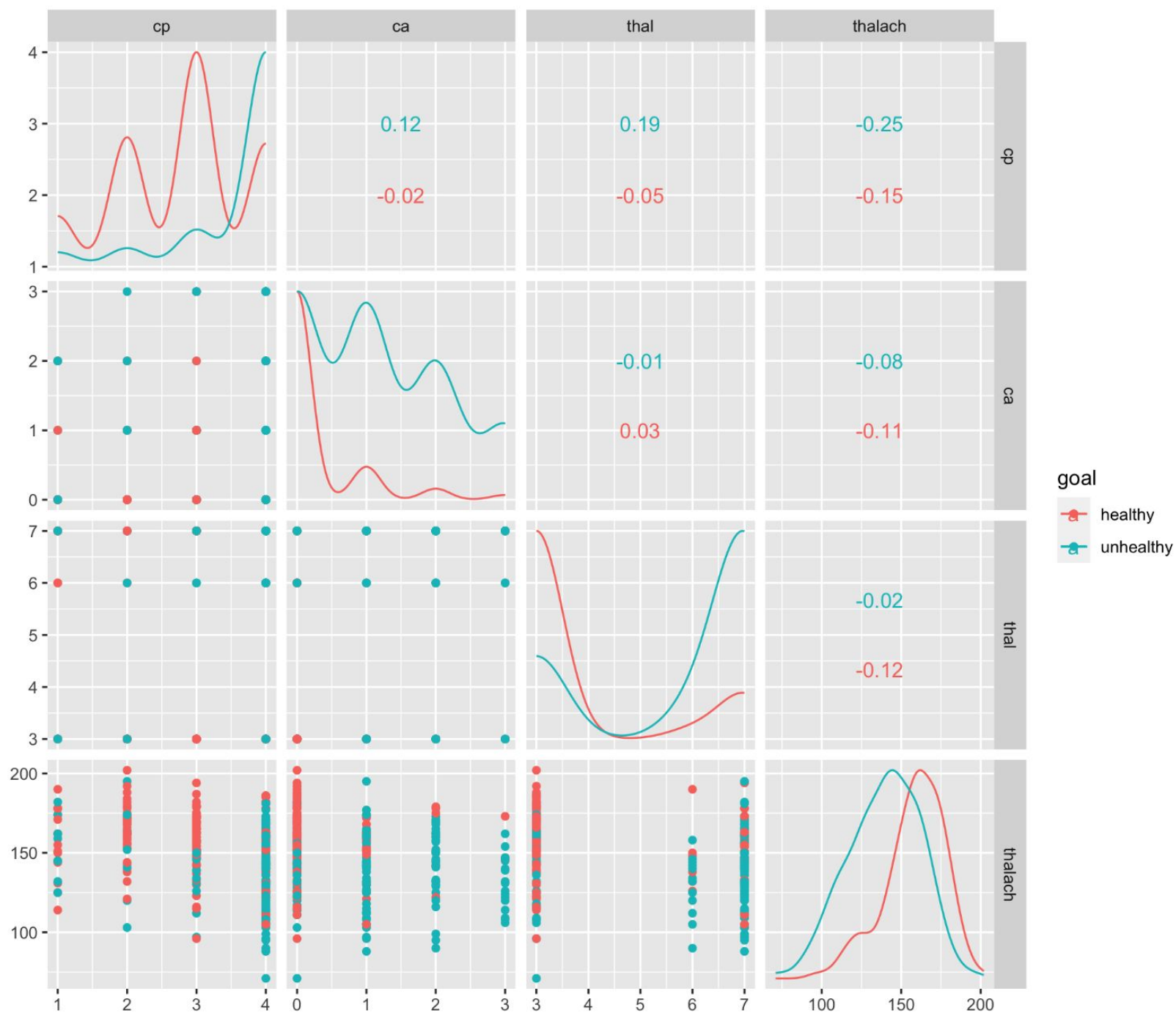
Data Type: continuous, discrete

Exploratory Data Analysis

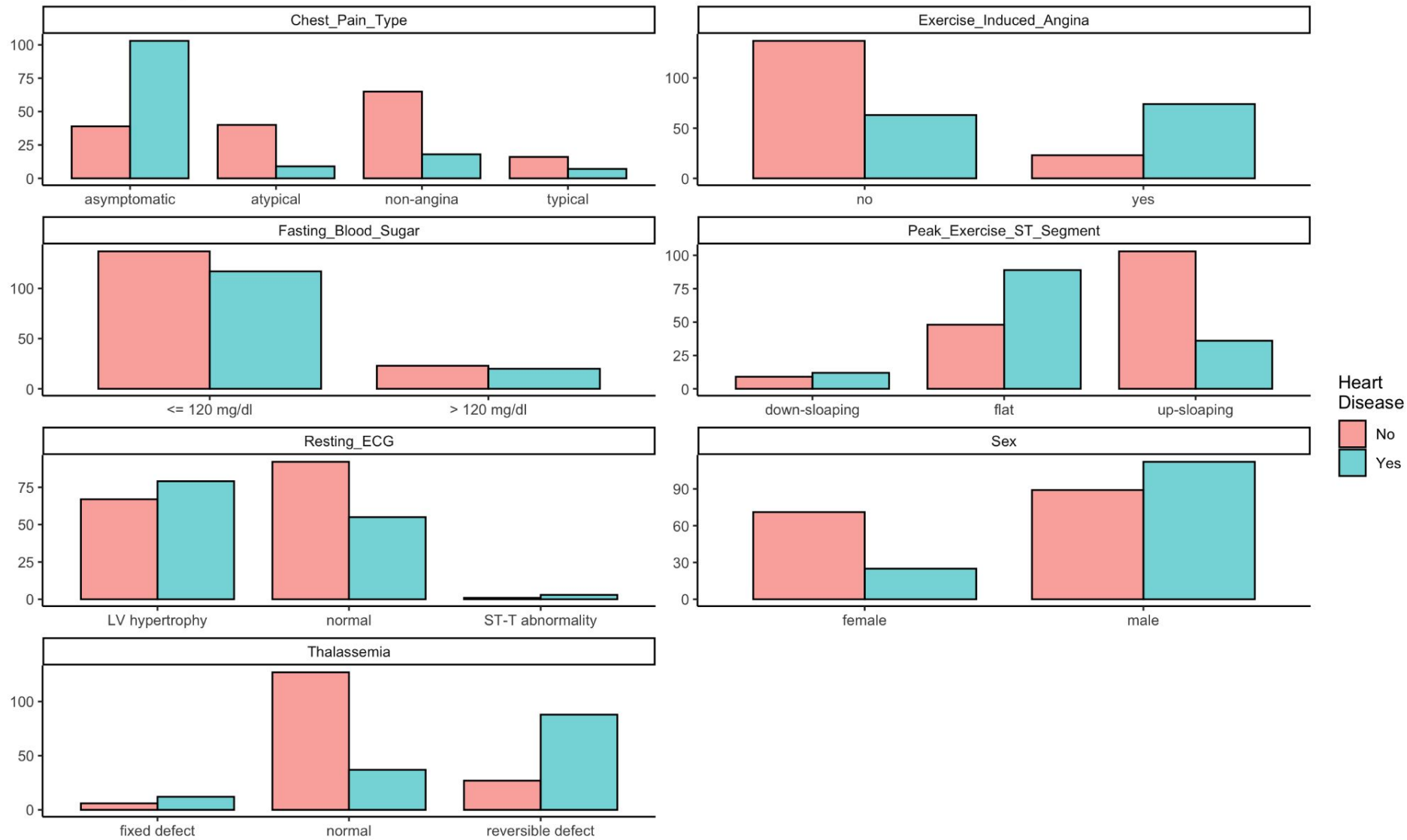
1. Barplot
2. Pairplot
3. Boxplot
4. Heatmap

Results

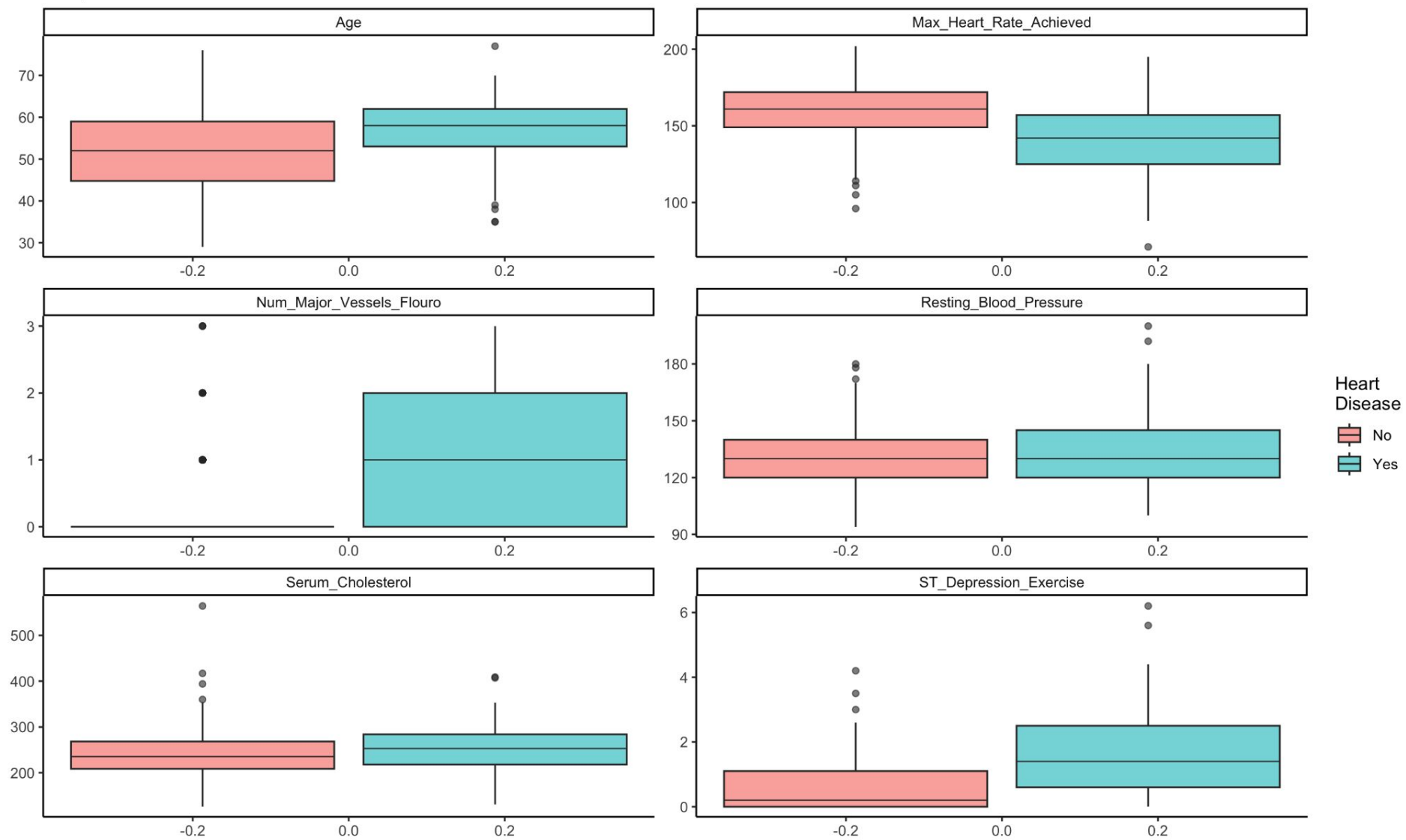




Effect of Categorical Variables

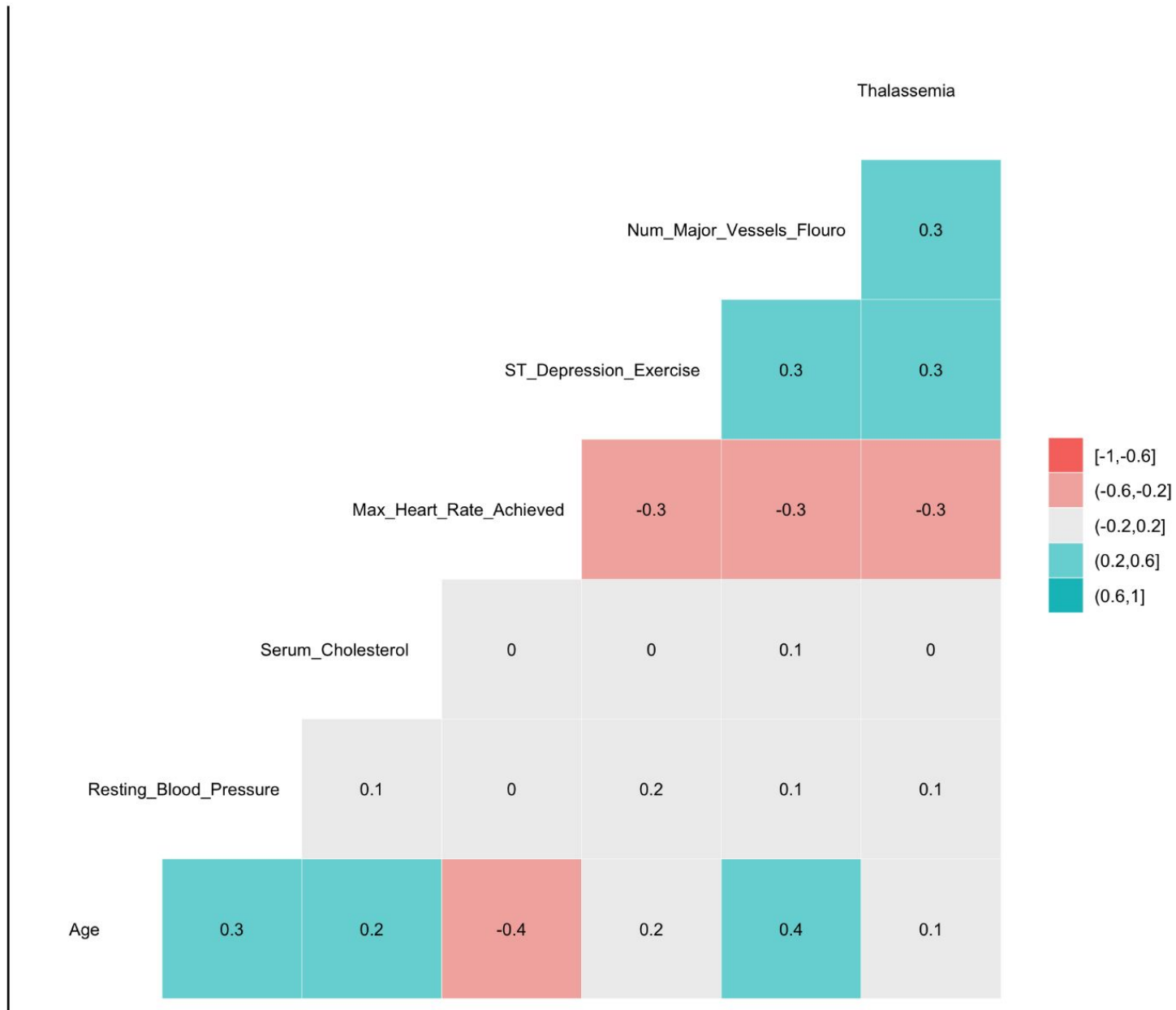


Boxplots for Numeric Variables



Heat Map

Pearson correlation



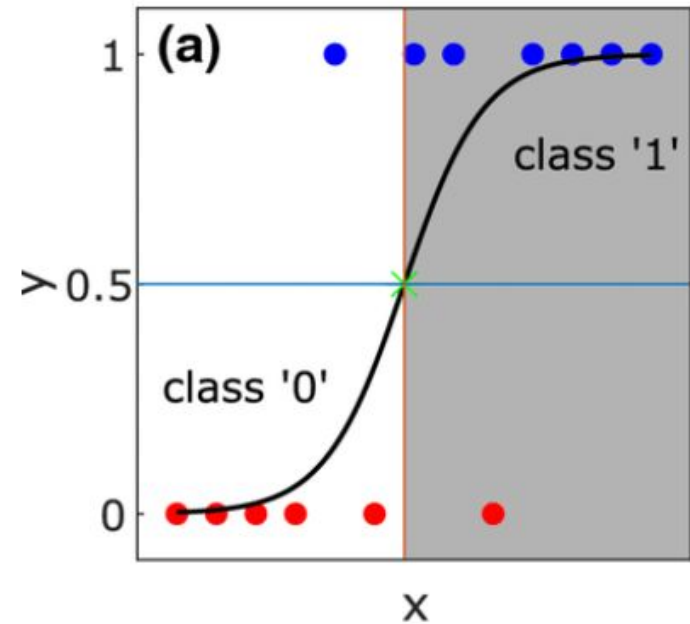
Classifiers

1. Logistic regression
2. Boosted Logistic regression
3. Random forest
4. K-nearest neighbour



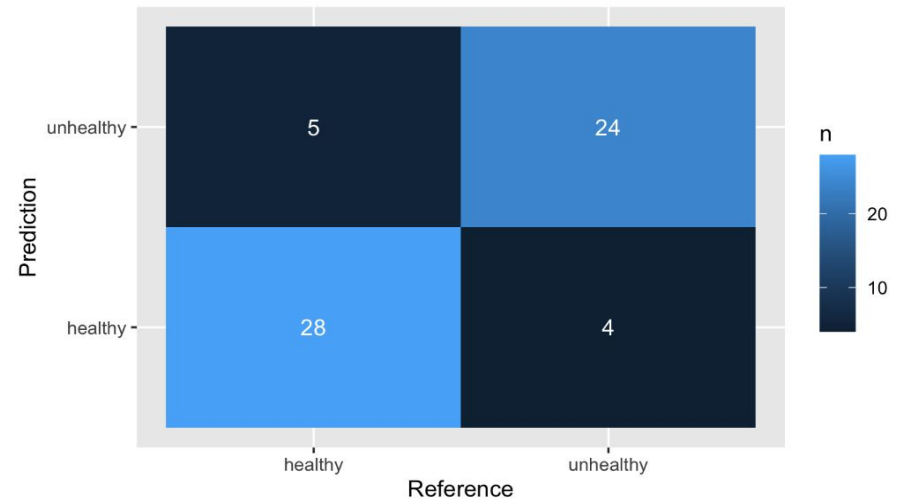
Logistic regression

- A logistic regression model is used
- Model improvement by error minimization



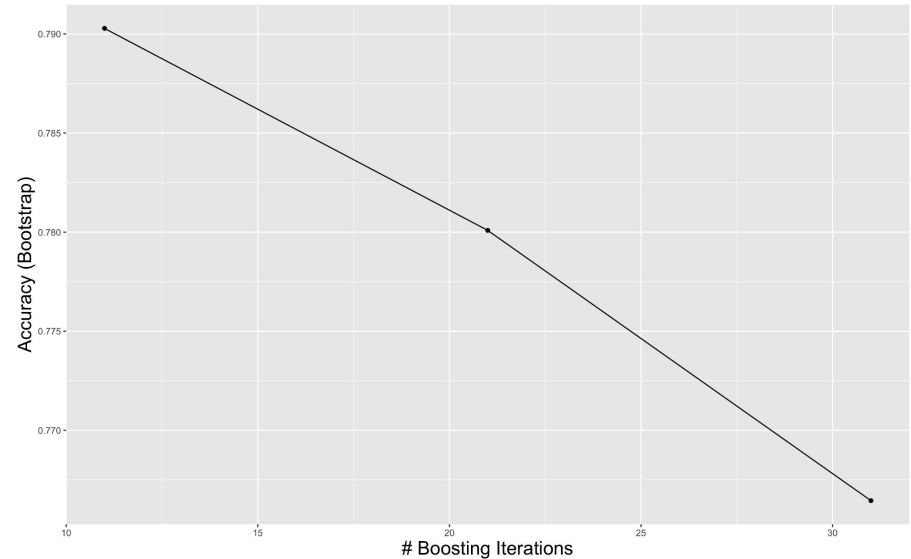
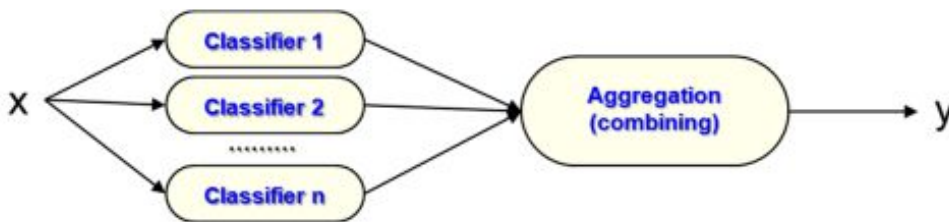
2

Logistic model



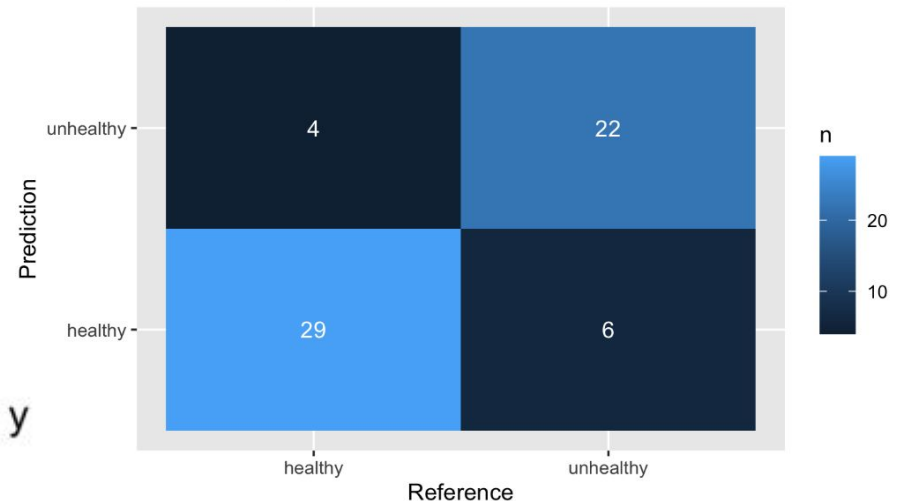
Boosted Logistic regression

- First classification based on each feature separately (decision stumps)
- Combining of the classifications by applying different weights and a logistic regression



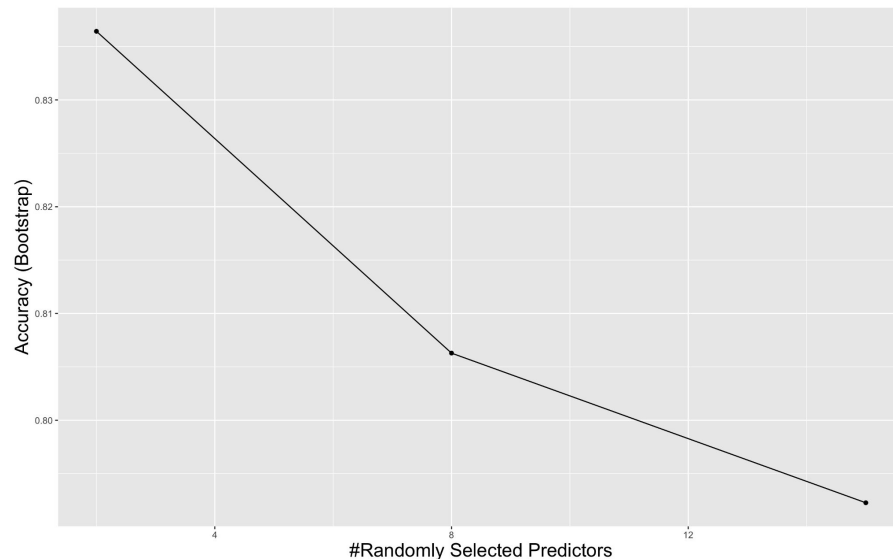
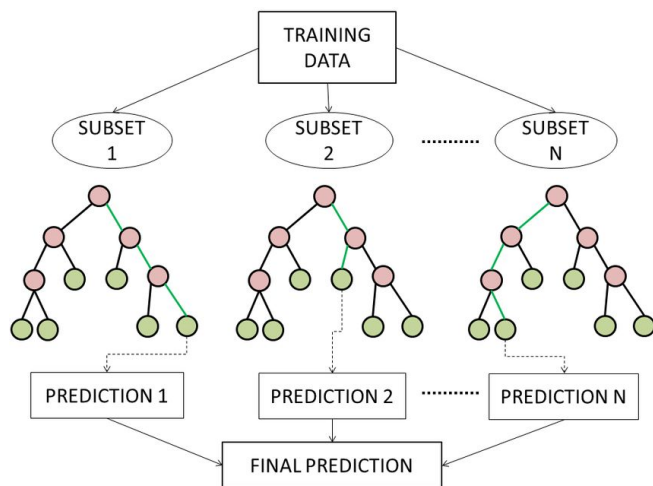
3

Boosted Logistic model



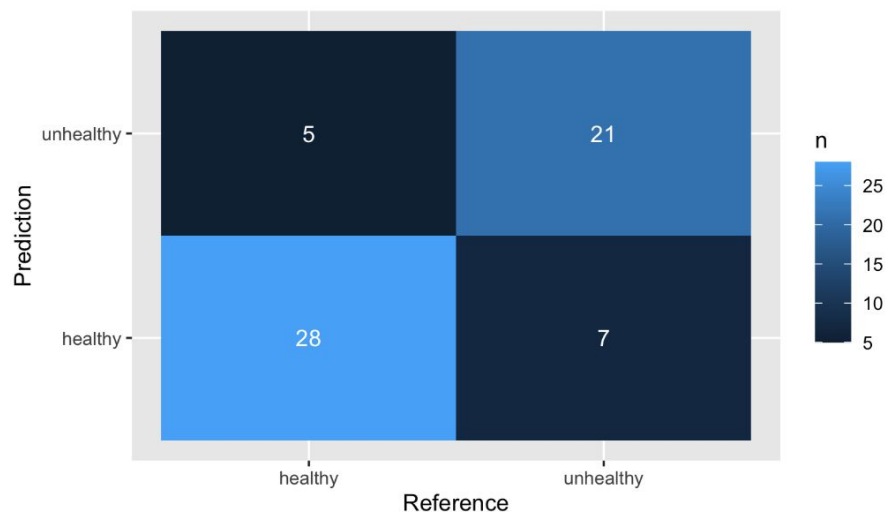
Random Forest

- Multiple Decision Trees are randomly created
- A data point is classified by each of them
- The class with the most 'votes' is the final classification



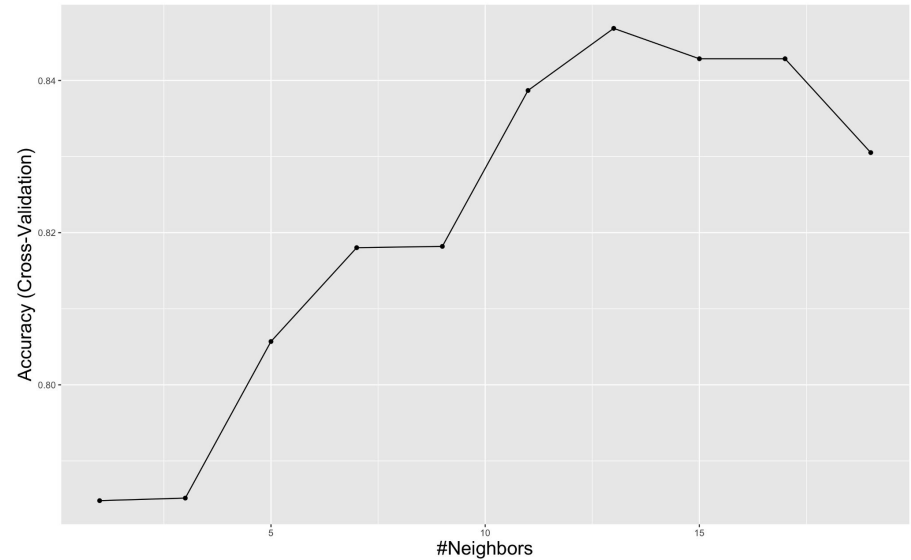
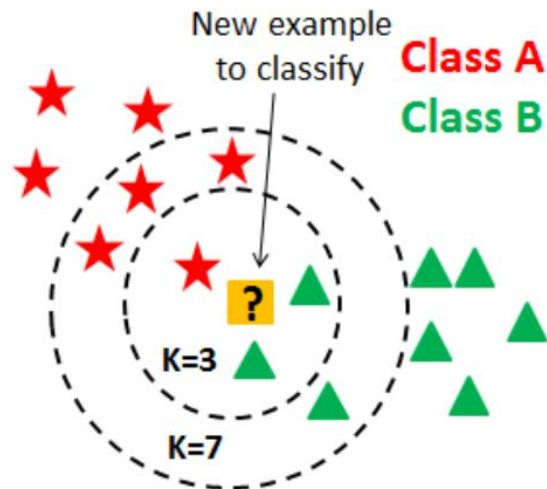
1

Random Forest model



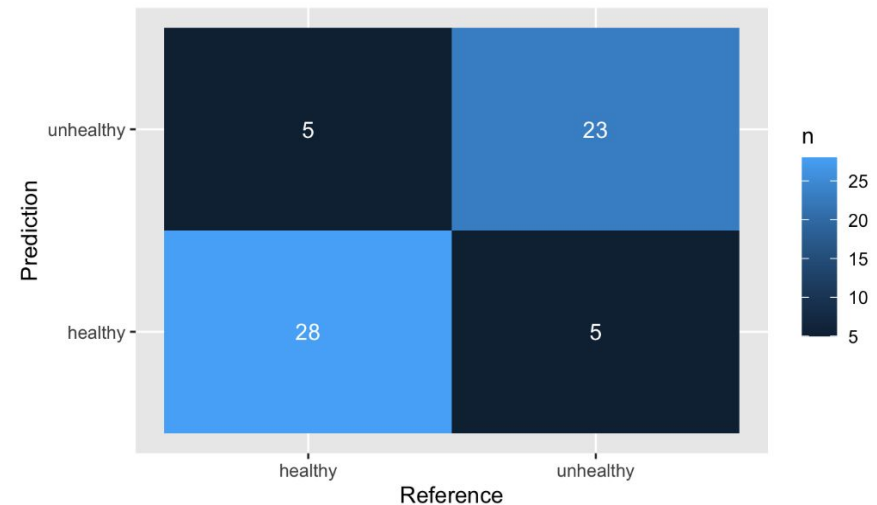
kNN

- A new datapoint is classified by looking at the k nearest neighbors
- Best tuned model $\rightarrow k = 13$



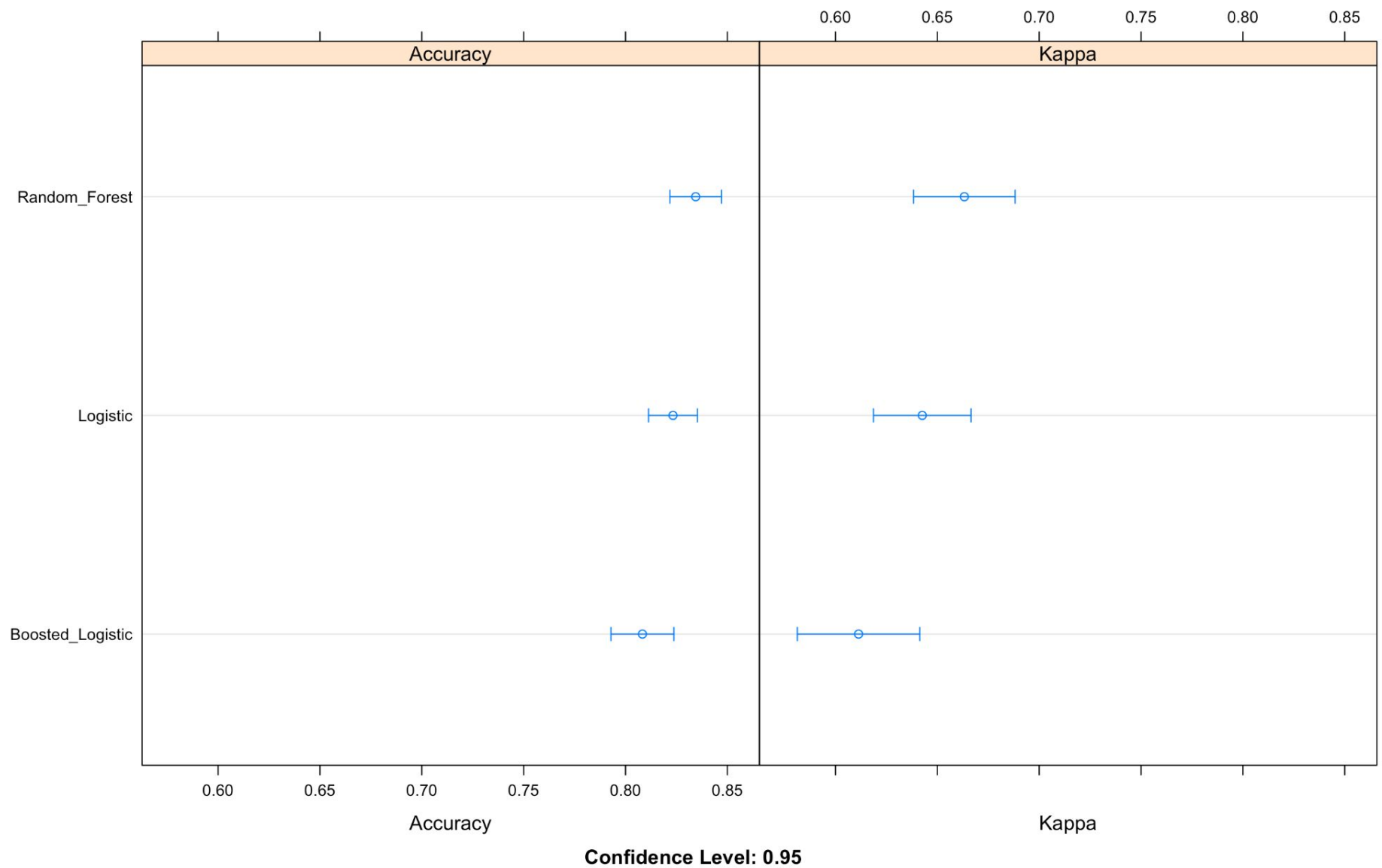
4

kNN model



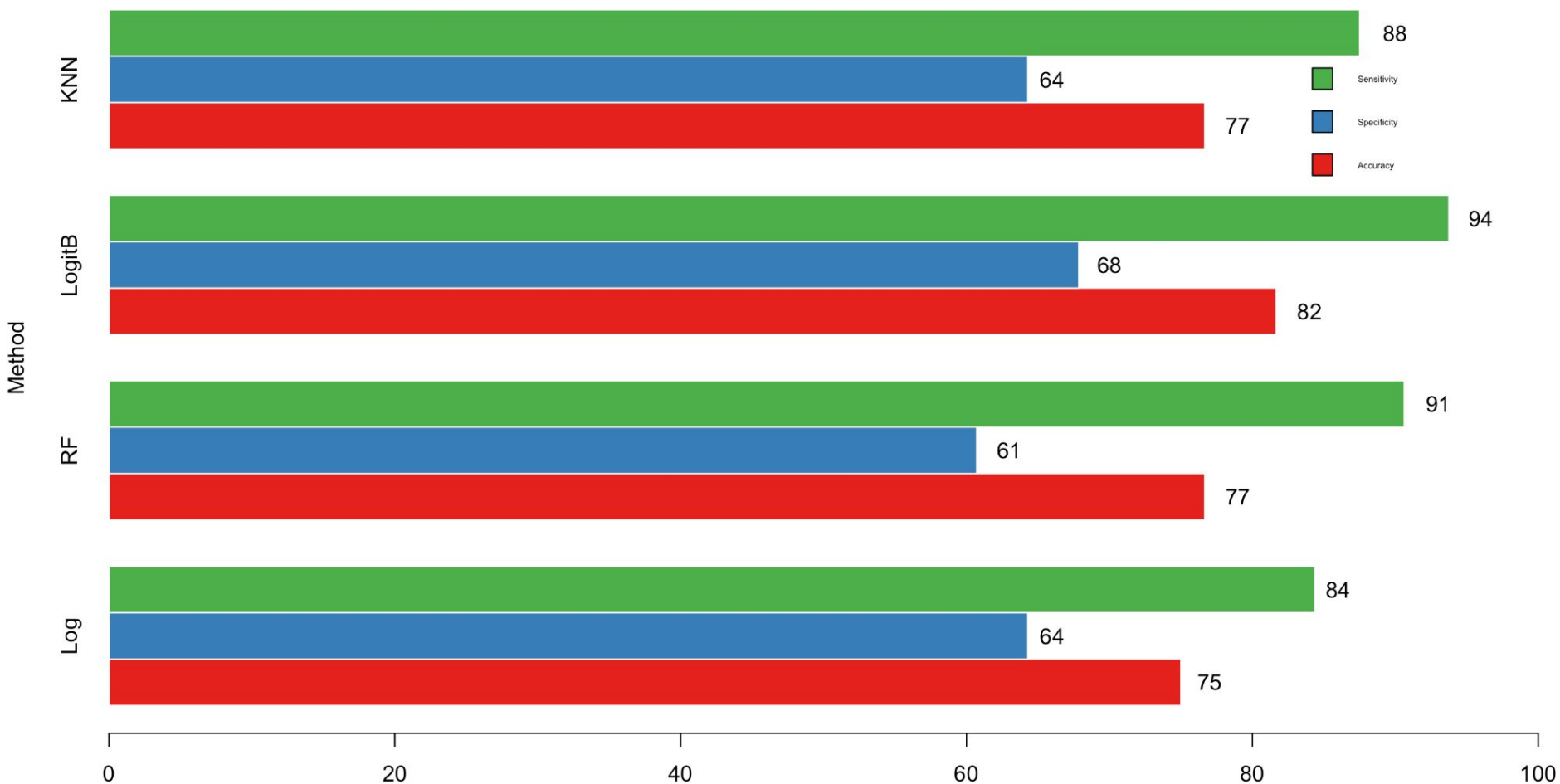
Performance Metrics

1. Accuracy dotplot
2. Comparison graph
3. PR curve
4. PRG curve
5. ROC curve
6. Calibration curve



The amount of agreement correct by the agreement expected by chance is Cohen's **Kappa**

Performance Chart

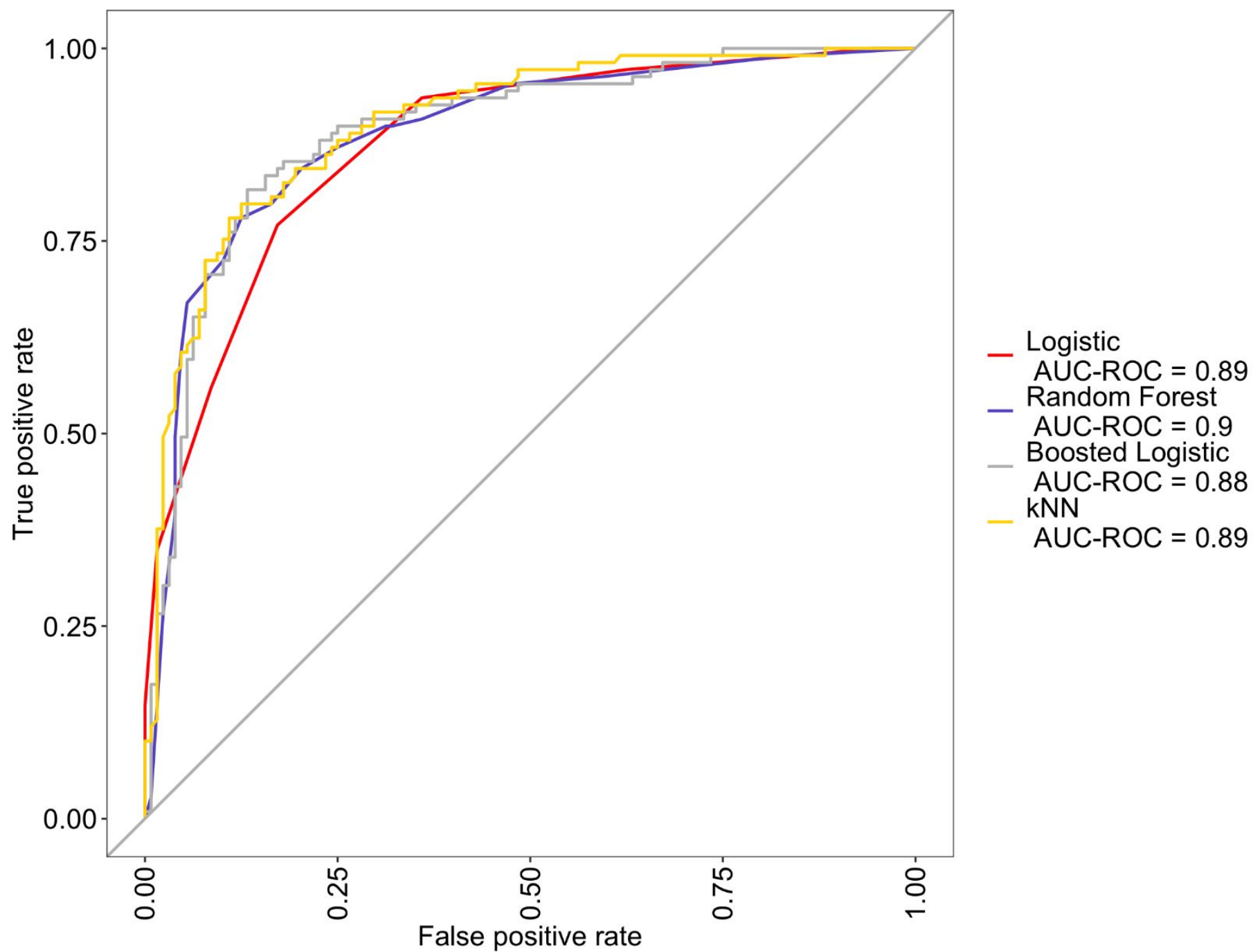


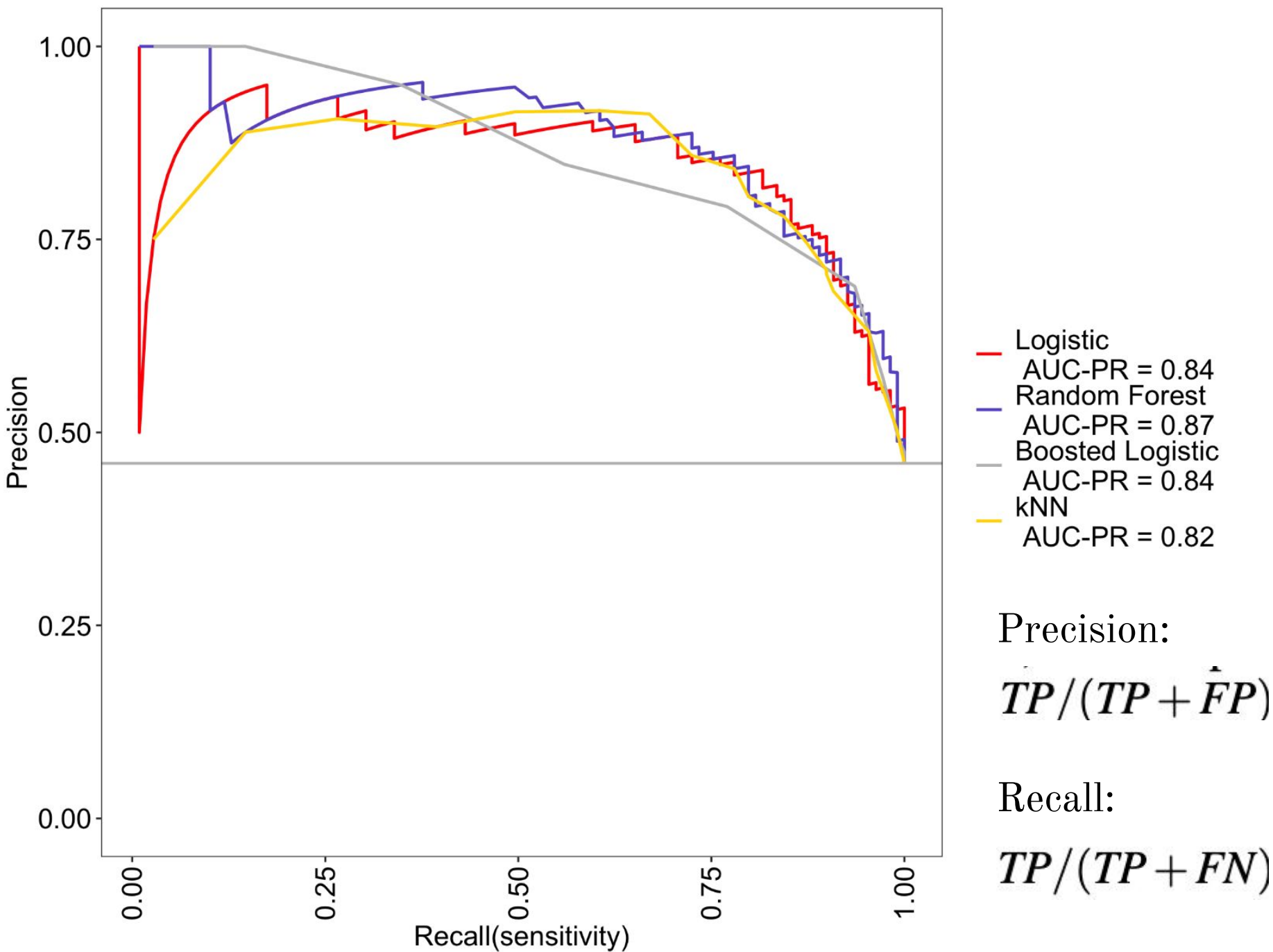
Logistic Optimal Informedness = 0.68370126146789

Random Forest Optimal Informedness = 0.673165137614679

Boosted Logistic Optimal Informedness = 0.598767201834862

kNN Optimal Informedness = 0.654816513761468

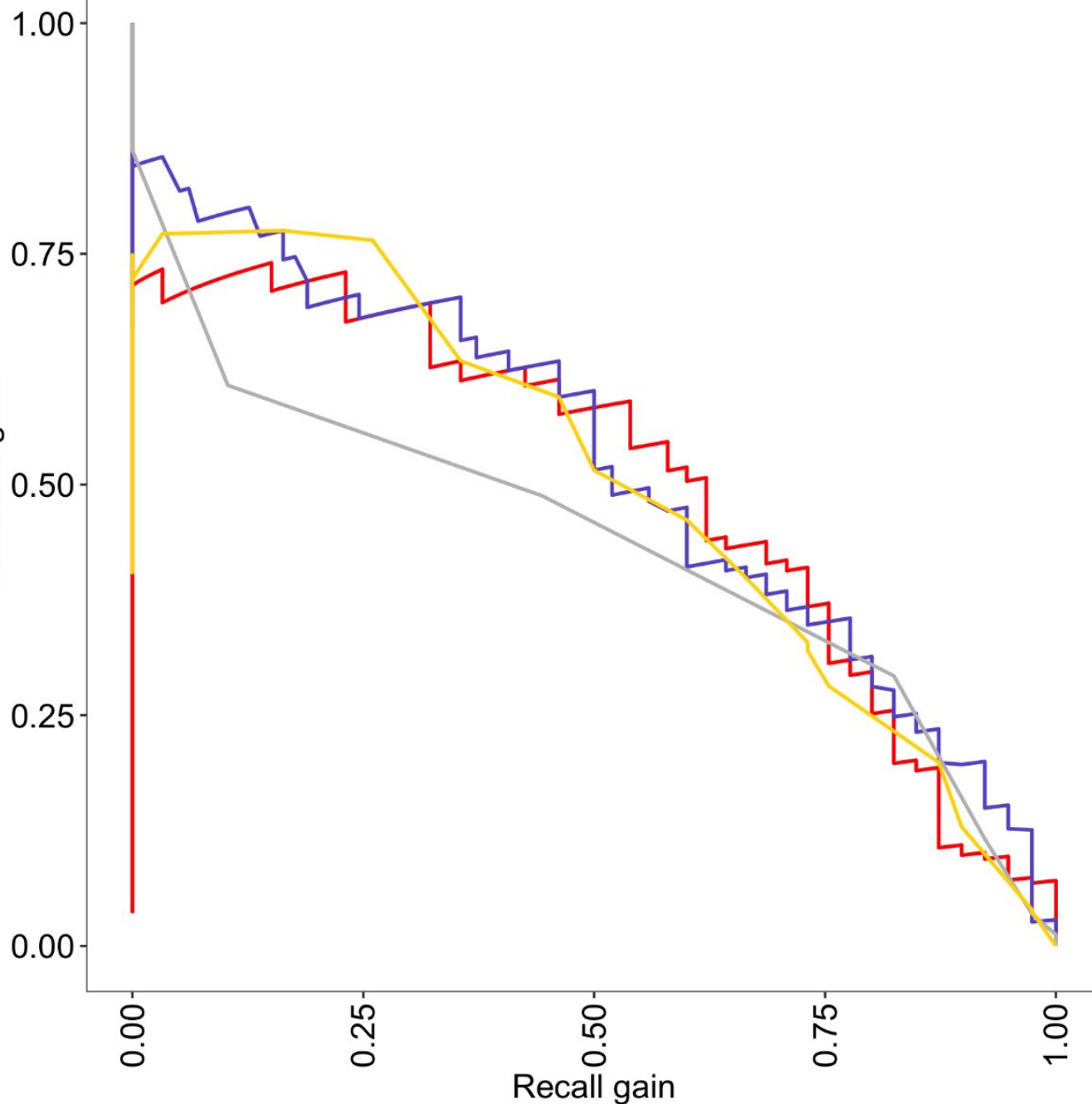




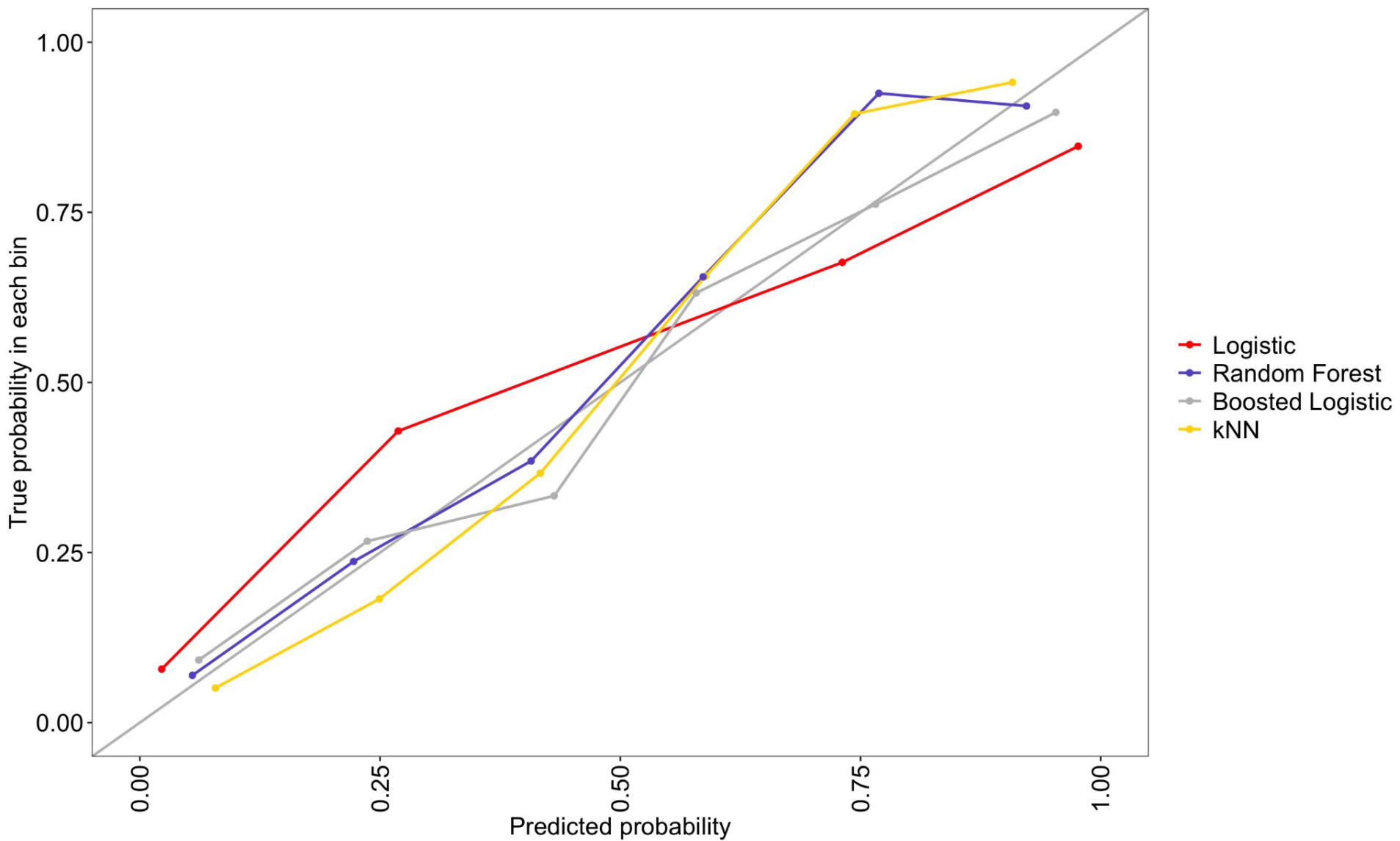
$$precG = \frac{prec - \pi}{(1 - \pi)prec} = 1 - \frac{\pi}{1 - \pi} \frac{FP}{TP}$$

$$recG = \frac{rec - \pi}{(1 - \pi)rec} = 1 - \frac{\pi}{1 - \pi} \frac{FN}{TP}$$

Precision gain



- Logistic
AUC-PRG = 0.48
- Random Forest
AUC-PRG = 0.51
- Boosted Logistic
AUC-PRG = 0.44
- kNN
AUC-PRG = 0.45



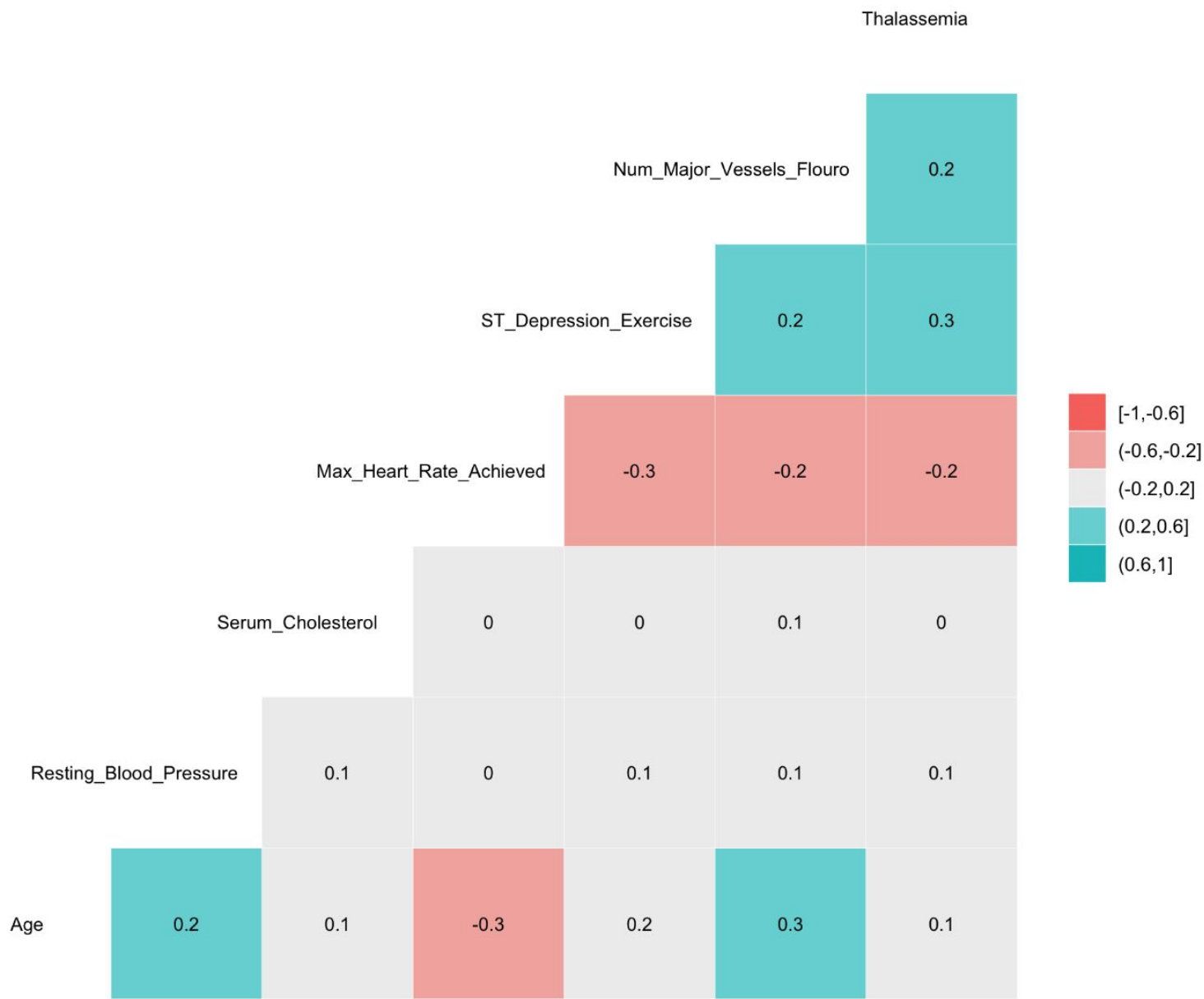
References

1. Pisula T. "An Ensemble Classifier-Based Scoring Model for Predicting Bankruptcy of Polish Companies in the Podkarpackie Voivodeship." Journal of Risk and Financial Management. 2020; 13(2):37. <https://doi.org/10.3390/jrfm13020037>.
2. Detrano, R et al. "International application of a new probability algorithm for the diagnosis of coronary artery disease." The American journal of cardiology vol. 64,5 (1989): 304-10. [https://doi.org/10.1016/0002-9149\(89\)90524-9](https://doi.org/10.1016/0002-9149(89)90524-9).
3. J. H. Gennari, P. Langley, and D. Fisher. 1989. "Models of incremental concept formation." Artif. Intell. 40, 1–3 (Sep. 1989), 11–61. [https://doi.org/10.1016/0004-3702\(89\)90046-5](https://doi.org/10.1016/0004-3702(89)90046-5)
4. Kibler, Dennis, David W. Aha, and Marc K. Albert. "Instance-based prediction of real-valued attributes." Computational Intelligence 5.2 (1989): 51-57.
5. R Core Team (2022). "R: A language and environment for statistical computing." R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
6. Illustrations
https://www.researchgate.net/figure/Data-classification-by-logistic-regression-aClassification-of-1D-data-showing-the-fitted_fig1_353913155
7. Random Forest sketch
https://www.researchgate.net/figure/Example-of-a-Random-Forest-workflow_fig2_342028855

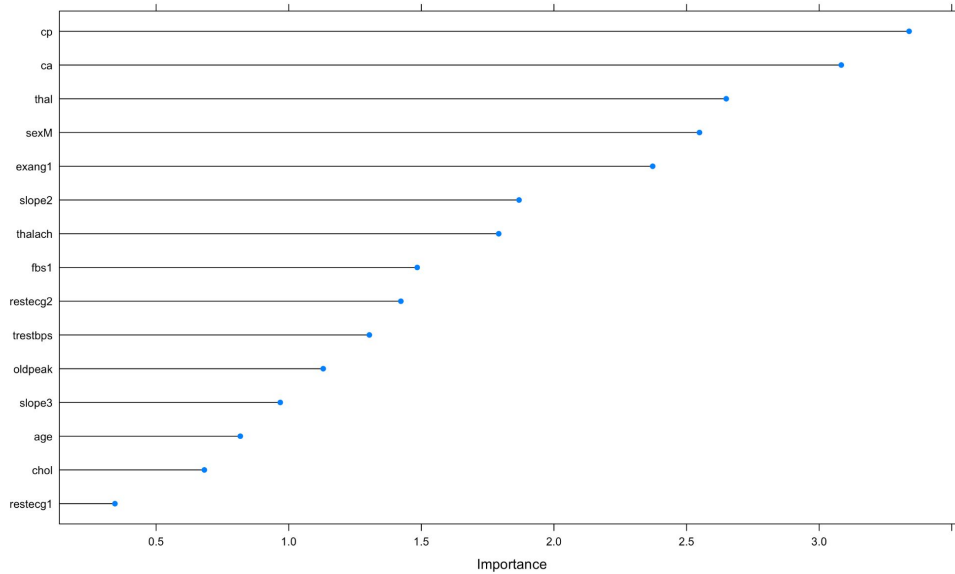
Dataset description can be found [here](#) , and data file [here](#).

Thanks for listening!
Questions ?

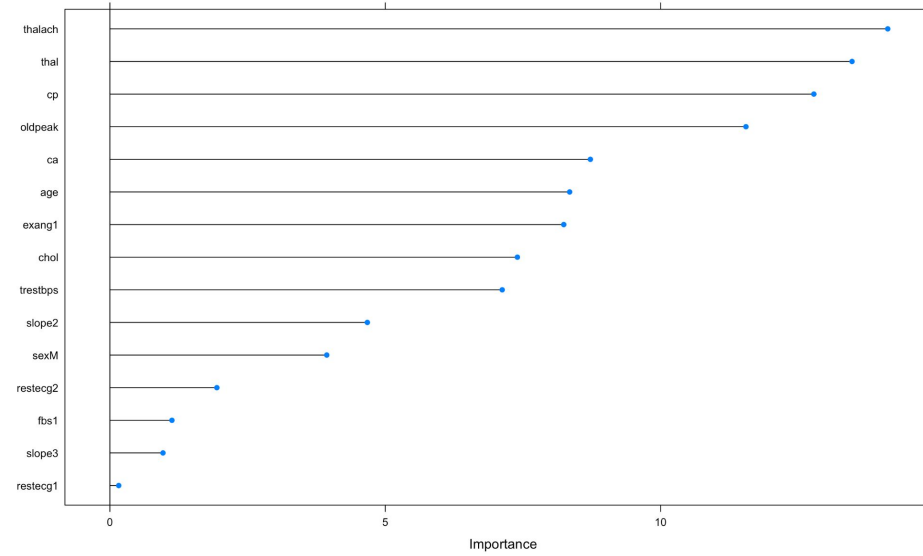
Heat Map
Kendall correlation



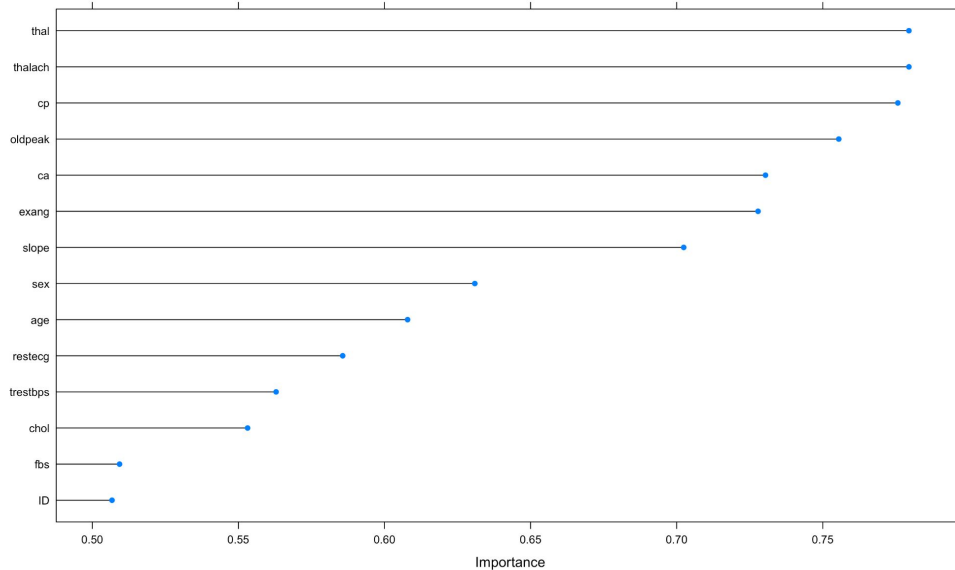
Logistic regression: features



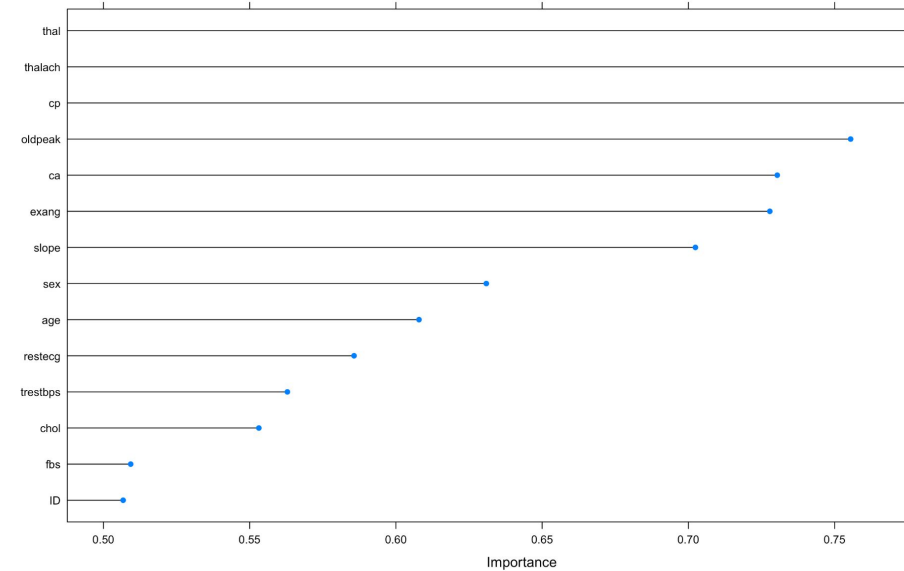
Random forest: features



Boosted Logistic regression: features

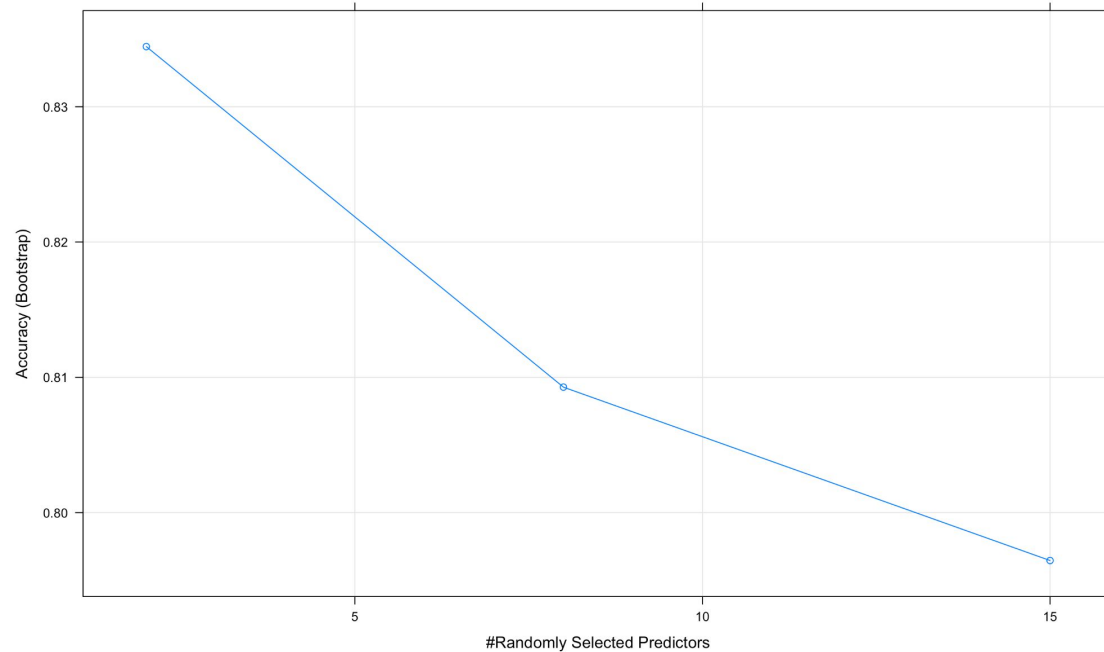


KNN: features



Top features: chest pain type (cp), number of major vessels (ca), maximum heart rate achieved (thalach), thallium scintigraphy (thal)

Random Forest

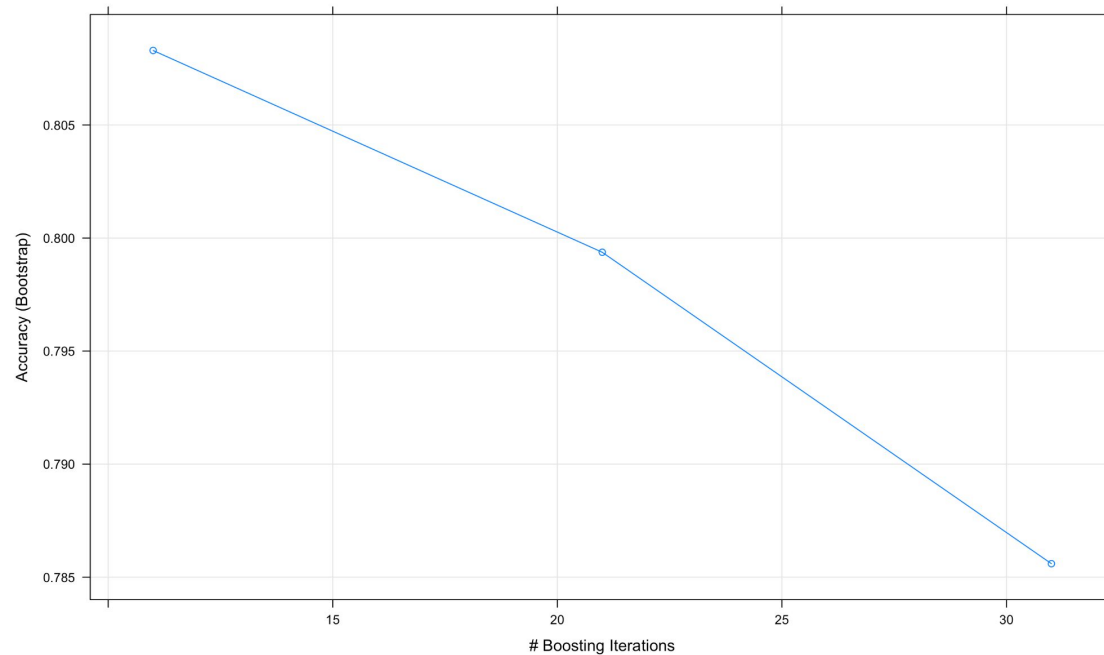


Accuracy \Rightarrow **Downward trend**

Possible reason: **Tuning parameters**

- *mtry* (by default) \rightarrow bagging
- Fast
- Good variance-bias tradeoff

Boosted Logistic



- *nIter* (by default) \rightarrow decision stump
- Speedy
- Weak