# Statistical Machine Learning Task

## Predicting Prostate Cancer Status : A Machine Learning Approach

Agbatan Fiacre Luc KOUDERIN

African Institute for Mathematical Sciences (AIMS - Rwanda)
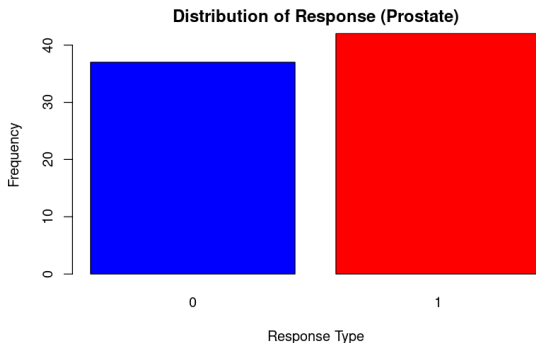
December 16, 2024

# Outline

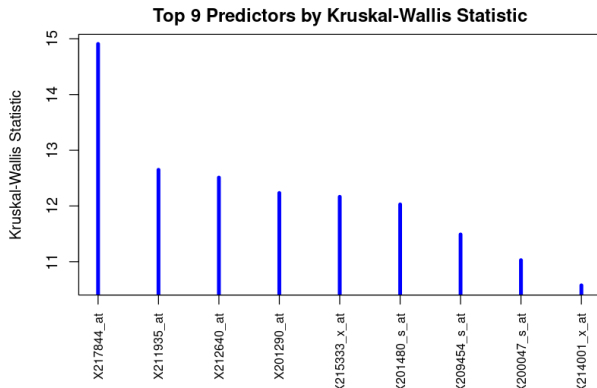- Data description and manipulation

- Models fitting and comparison

# Data description and manipulation

- 501 variables (one response variable, 500 covariates) with 79 observations
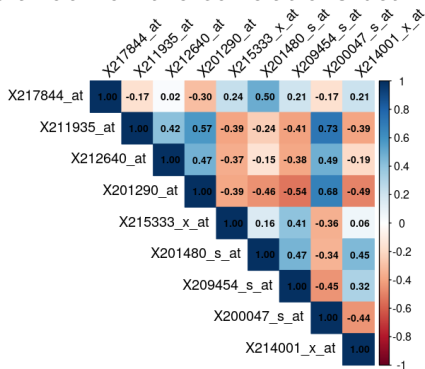- Distribution of the response variable named Y :



**Distribution of Response (Prostate)**

# Data description and manipulation

We are in a case of curve of dimensionality, let's see the best predictors.



Top 9 Predictors by Kruskal-Wallis Statistic

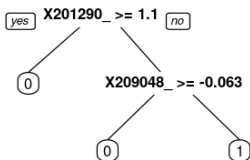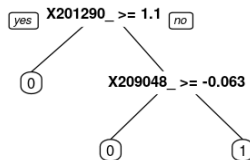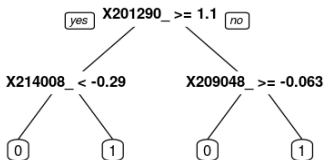# Data description and manipulation

- Have a look on the correlations between the top best predictors



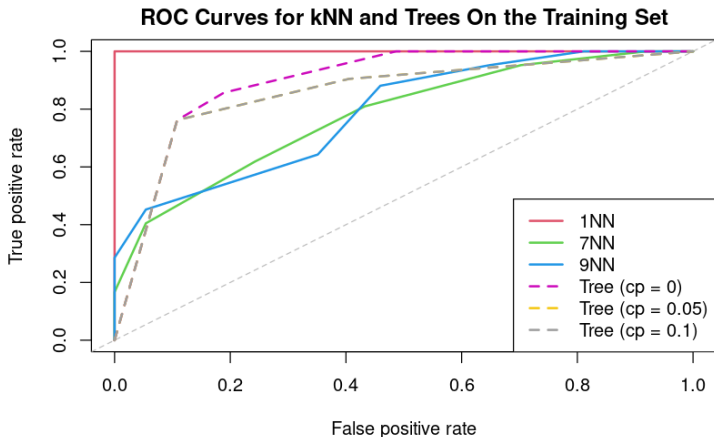- By doing the eigen decomposition of the correlation matrix we get the ratio = 17.02

# Models fitting and comparison

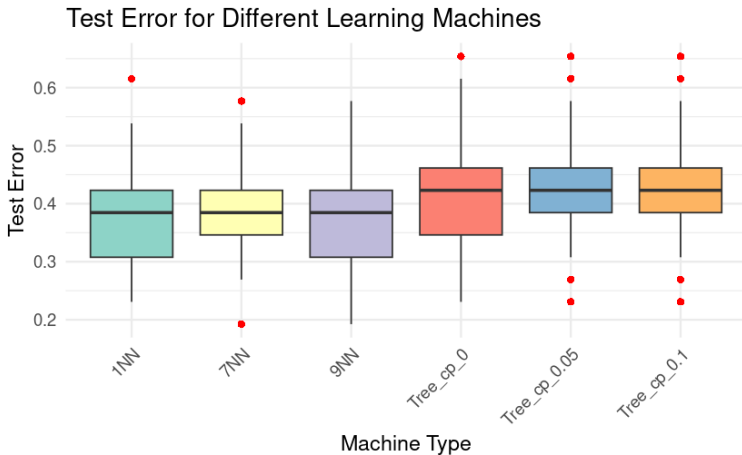- Let's start with a tree with differents values for the complexity parameter.

# Models fitting and comparison

- Let's add to the previous model a kNN model with differents values of k.



ROC Curves for kNN and Trees On the Training Set

# Models fitting and comparison

- Let's do predictions on the test set.



Test Error for Different Learning Machines

# Overall Conclusion

- A krustal Wallis Test can be done to have the best predictors variables.
- For this dataset, the kNN machine perform better than the Trees.
- Thers no signficant difference in the predictive merits of 1NN, 7NN and 9NN but in pratice, we will choose 9NN because its the least complex model between this three models.