

Predicting Heart Disease

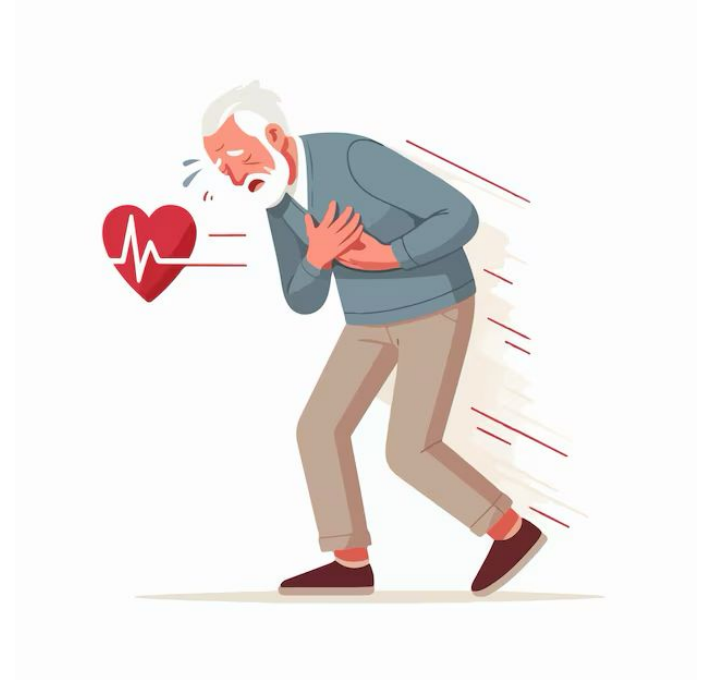
Implementations of Logistic Regression Models

Fennom Schalkwijk – 14619148
Kelt Paehlig – 14634716
Babet Wijsman –

Scientific Data Analysis – 5062SCDA6Y
Group Assignment – Project Group 9
17-12-2024

Motivation

- One of the leading causes of death (World Health Organization, 2021)
- Understanding risk factors to help research



Dataset

- 300.000+ rows
- ~20 variables
- Heart disease as dependent variable

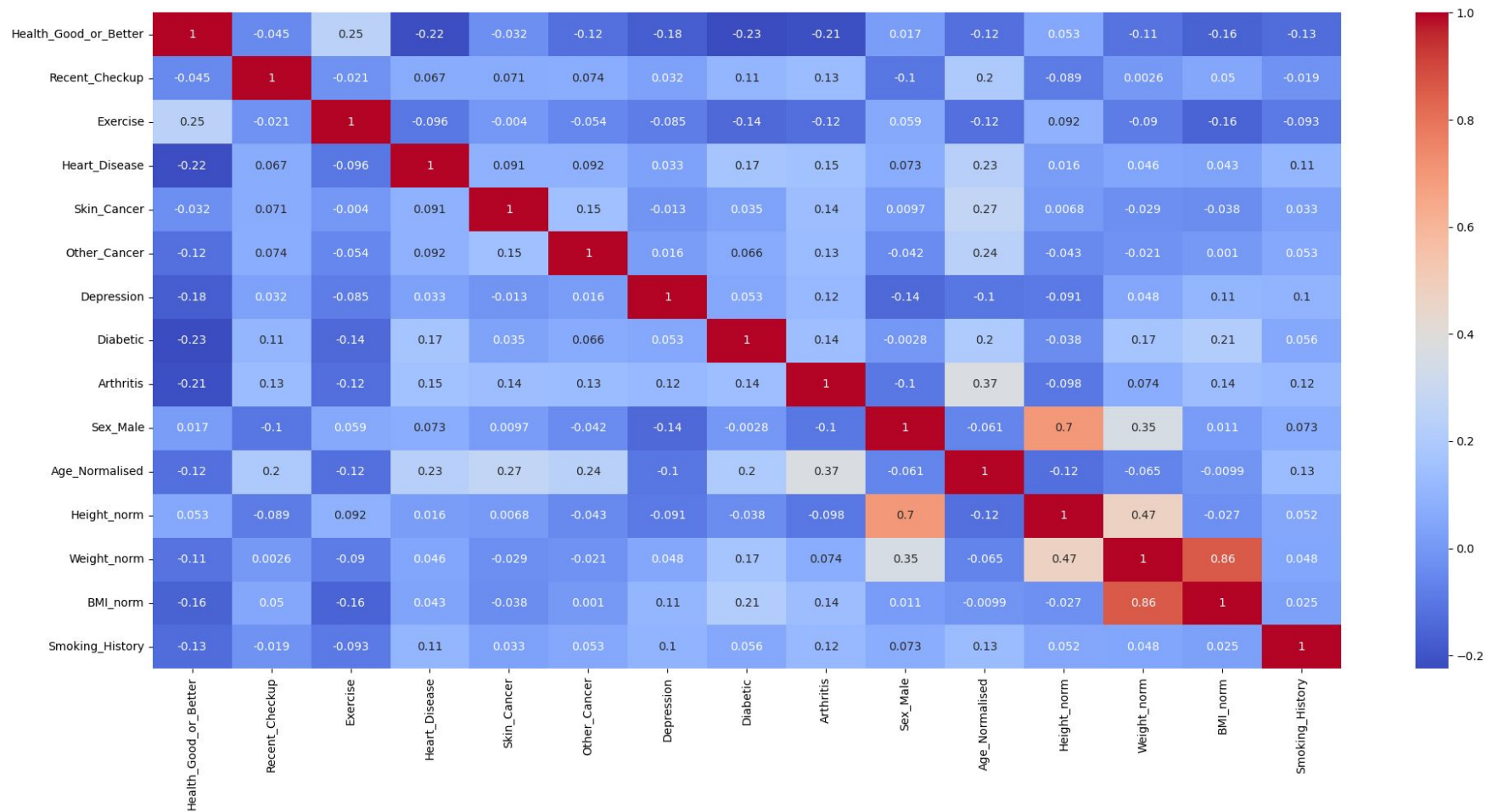


**U.S. CENTERS FOR DISEASE
CONTROL AND PREVENTION**

Dataset

Limitations

- Processing:
 - Loss of variability through grouping certain categories as binary
- Biases:
 - USA based dataset
 - Relatively uneven distribution of age among participants

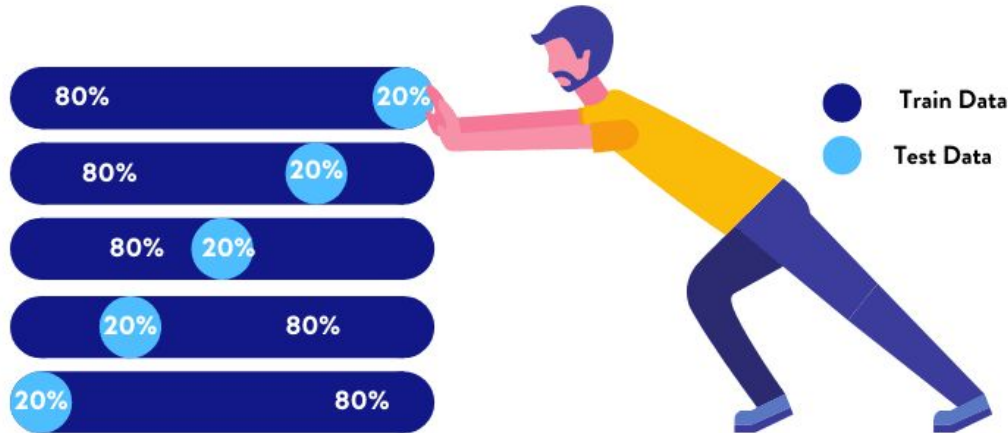


Research Questions and Hypothesis

- **RQ1:** Which features are the most predictive of heart disease when implementing a logistic regression model?
- **RQ2:** Does a logistic regression model based on PCA-generated components predict with higher accuracy than a logistic regression model fit on the 'clean' data?
 - **H²⁻⁰:** A logistic regression model using PCA-generated components does not achieve significantly higher accuracy compared to a logistic regression model using the original 'clean' data.
 - **RQ2^{sub-1}:** What implications do these components have on the interpretation of the predictive capacity of the features?
 - **RQ2^{sub-2}:** Does the usage of PCA components introduce any bias into the model?
- **RQ3:** Can we extract how many variables – and which – contribute to a high model accuracy?

Methods

- Logistic Regression
- Principal Component Analysis
- K-Fold Cross-Validation
 - Comparing PCA to Fully Fitted Model
 - Comparing Models With Select Variables



Methods: Coefficient Predictability

- RQ1: Which features are the most predictive of heart disease when implementing a logistic regression model?
 - Logistic regression predicts a binary outcome: heart disease (yes/no).

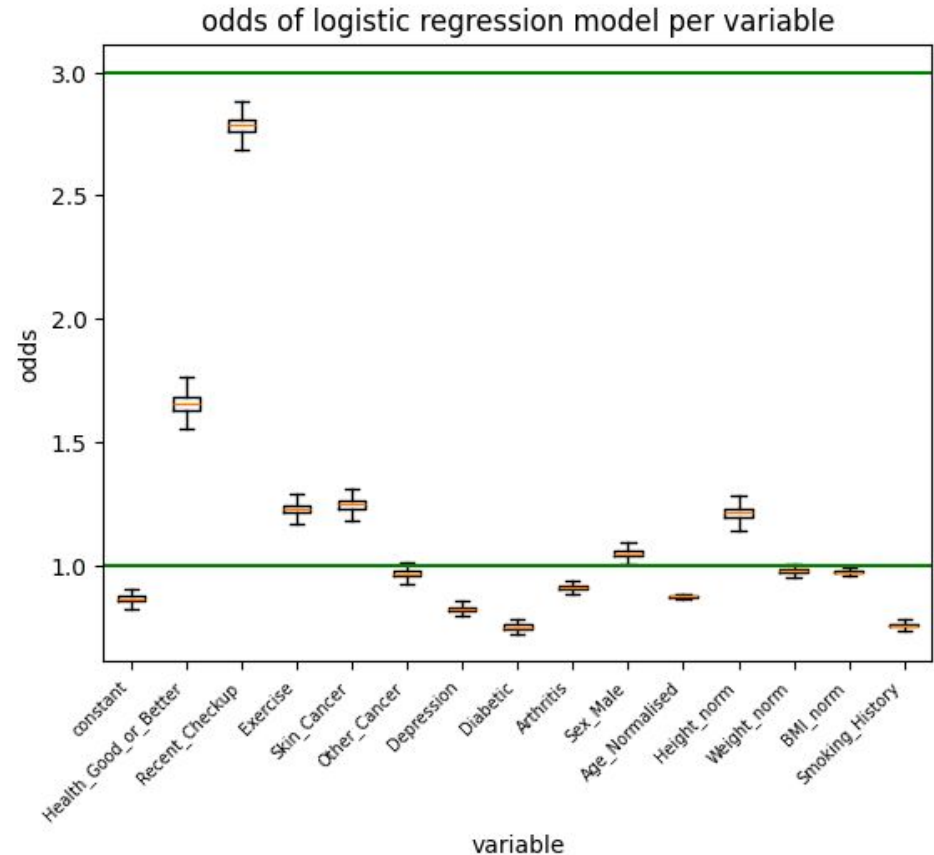
Use training set (75%) for regression
model

Test regression model → 0.75
accuracy

Calculate odds

Results: coefficient predictability

- < 1 : low risk
- 1 -3 moderate risk
- > 3 : high risk



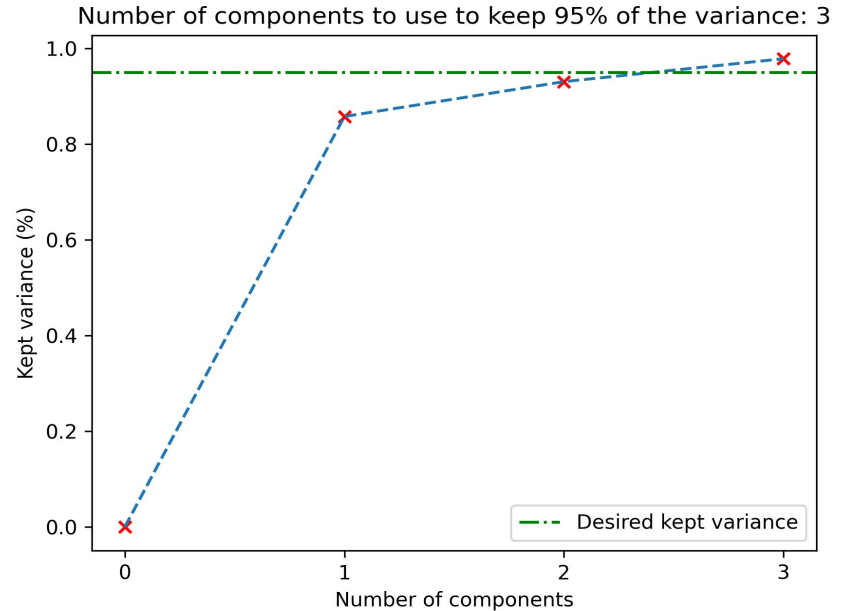
Limitations: coefficient predictability

- Interaction effects
- Multicollinearity
 - Weight, BMI, height
- Interpretation:
 - Recent checkups, exercise

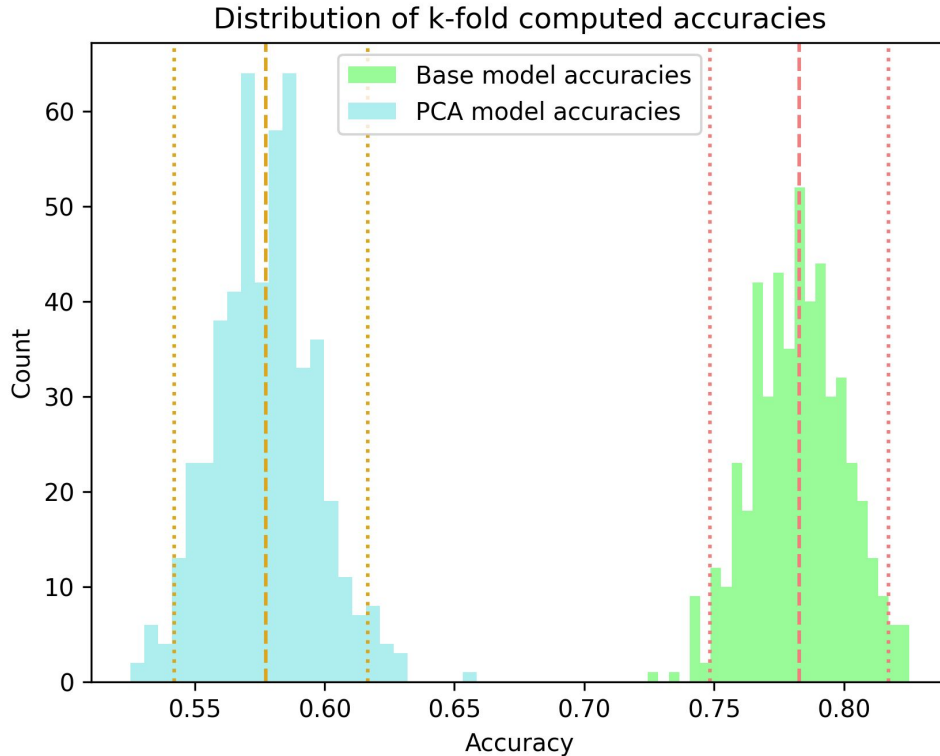
Principal Component Analysis

- 'Summary' of the dataset
- 3 Components for retainment of 95% of the variance

- PC 1: BMI, weight, height
- PC 2: Height, weight
- PC 3: Age



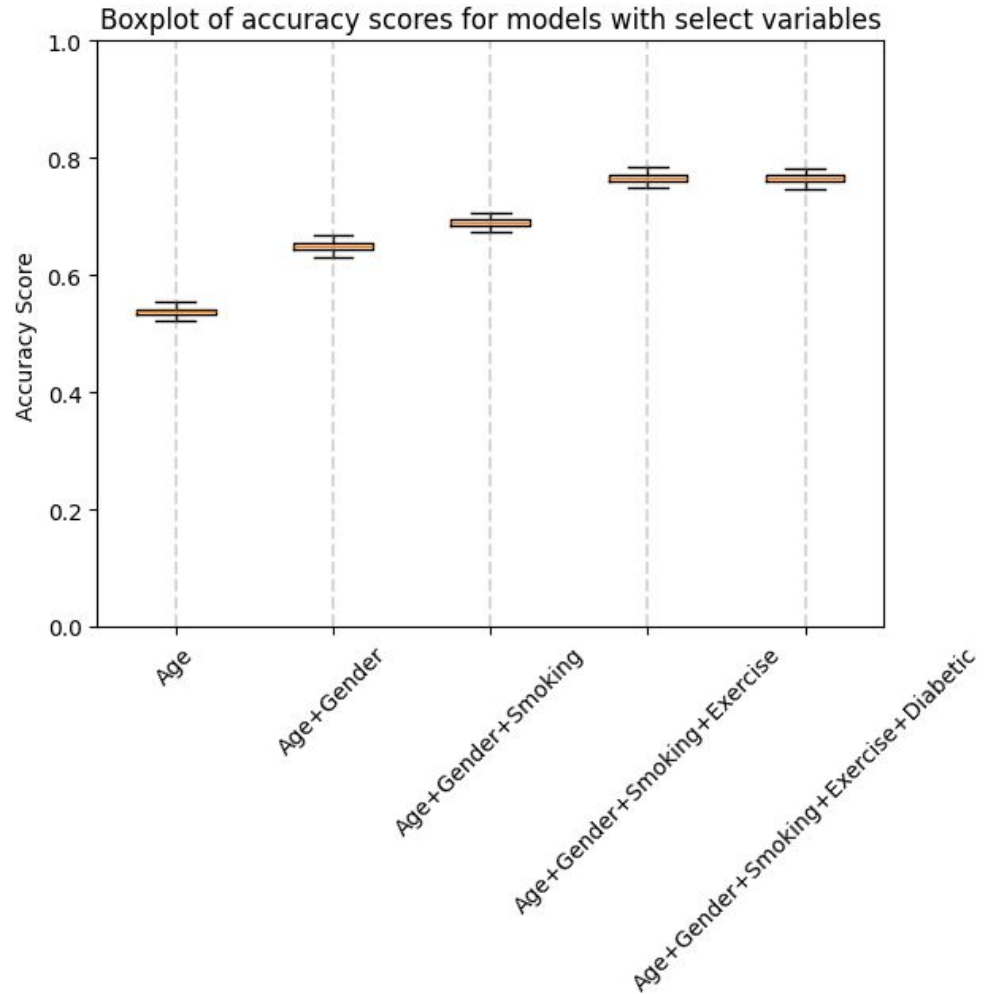
Results: Cross-Validation I



- We fail to reject our null-hypothesis
- Can we answer **RQ2**?
- Product of overfitting?
- AIC suggests this is not the case

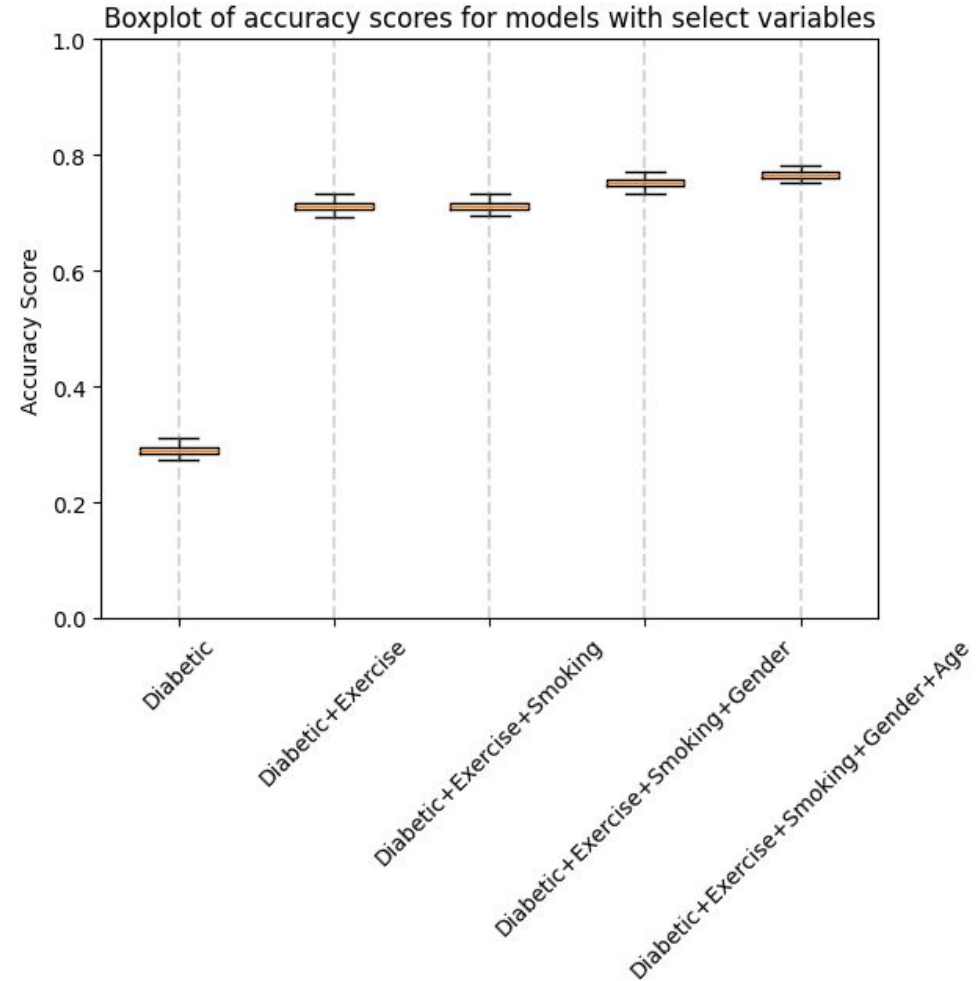
Results: Cross-Validation II

- **RQ3:** Can we extract how many variables – and which – contribute to a high model accuracy?
- Variable selection based on literature
 - Age
 - Gender
 - Smoking
 - Exercise
 - Diabetes
- Arbitrary order



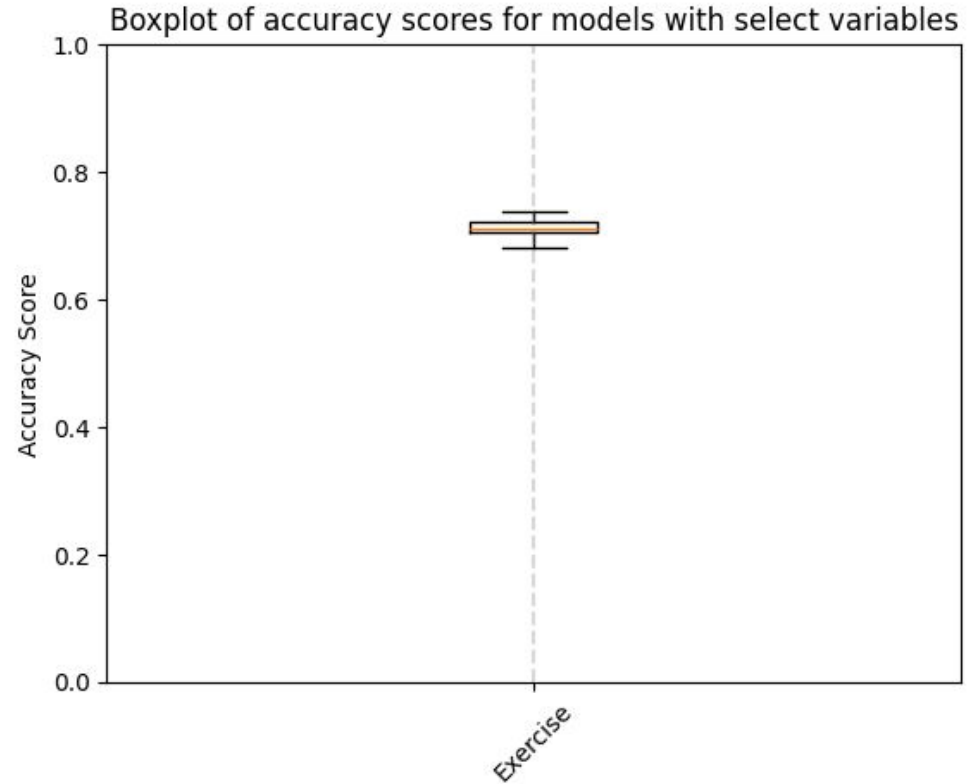
Results: Cross-Validation II

- Reversed order of variables



Results: Cross-Validation II

- Is exercise *that* predictive?
- It seems like it, but we cannot say for certain.



Conclusion and Limitations

- Some insights into variable predictability
- PCA did not yield better results

Conclusion and Limitations

- Potential overfit
- No analysis of potential collinearity
- Correlation \neq causation

Predicting Heart Disease

Implementations of Logistic Regression Models

Thank you for listening

Any questions?

Fennom Schalkwijk – 14619148
Kelt Paehlig – 14634716
Babet Wijsman –

Scientific Data Analysis – 5062SCDA6Y
Group Assignment – Project Group 9
17-12-2024