# A Computational Biology Framework

**A data analysis tool to support biomedical engineers in their research**

Tim Beishuizen, October 23rd 2018
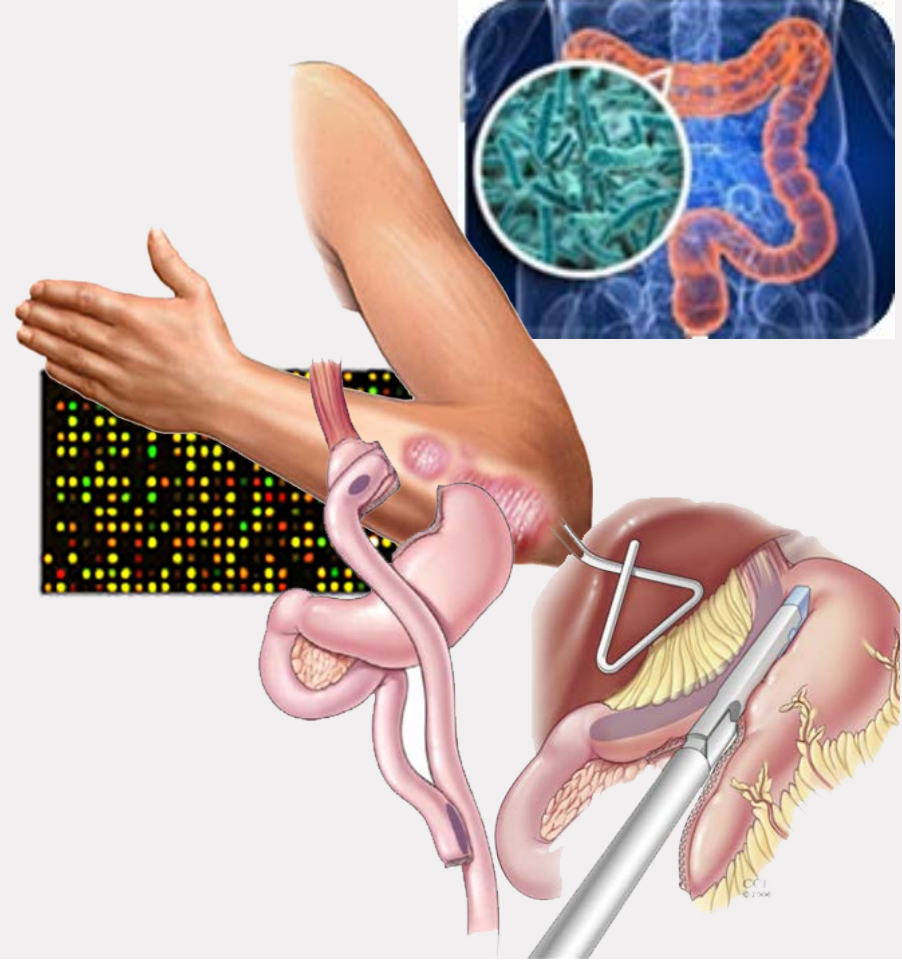
Supervisors:    Dr. Dragan Bošnački
                Dr. Ir. Joaquin Vanschoren
                Prof. Dr. Peter Hilbers

University of Technology Eindhoven, department Biomedical Engineering and department Computer Science and Engineering

# Contents

- Introduction
  - Research questions
- Project description
- Research lines
  - Feature selection
  - Missing value handling
  - Dataset exploration
- Framework

TU/e

# Introduction

- Biomedical data

- Processing

- Analysis

- Framework

TU/e

Main research question

## What aspects are of importance to be included in a framework for Biomedical Engineers for more efficient data analysis?

Datasets                 Preprocessing methods
Data issues              Analysis methods

# Sub-questions

- What feature selection methods show the best performance and as such should be added to the framework?

- What missing value handling methods show the best performance and as such should be added to the framework?

- Which initial analyses help find suitable preprocessing and data analysis algorithms?

TU/e

# Project description

- Statistics
- Machine learning

- Data matrix
  - Samples (s)
  - Features (f)
  - Output (y)
  - Classification

- Python

$$
\begin{array}{c}
\quad\quad f_1 \quad f_2 \quad f_3 \quad f_4 \\
\begin{array}{c} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{array}
\left(
\begin{array}{cccc}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & x_{23} & x_{24} \\
x_{31} & x_{32} & x_{33} & x_{34} \\
x_{41} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{array}
\right)
\left(
\begin{array}{c}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{array}
\right)
\end{array}
$$

TU/e

# Datasets

- **Microarray datasets**
  - Gene expression
  - Many features
  - Big datasets

- **Mass spectrometry datasets**
  - Protein mass
  - Many features
  - Differences between features

- Clinical datasets
  - Data out of a clinic
  - Missing values
  - Combining data issues

- Survey datasets
  - Questionnaire data
  - Missing values
  - Bias due to subjectivity

# Feature selection

# Feature selection

What feature selection methods show the best performance and as such should be added to the framework?

Feature selection methods
- Filter methods
- Wrapper methods
- Embedded methods

Datasets
- Microarray
- Mass spectrometry

TU/e

# Filter methods

- Separate features

TU/e

# Filter methods

- Separate features

- Evaluation function

- Selection

- Fast and simple

$$
\begin{array}{c}
\quad\ f_1 \quad f_2 \quad f_3 \quad f_4 \\
\begin{array}{c} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{array}
\begin{pmatrix}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & x_{23} & x_{24} \\
x_{31} & x_{32} & x_{33} & x_{34} \\
x_{41} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{pmatrix}
\begin{pmatrix}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{pmatrix}
\end{array}
$$

TU/e

# Wrapper methods

- Feature sets

TU/e

# Wrapper methods

- Feature sets

- Evaluation function

- Methods
  - Forward selection
  - Backward selection
  - Combinations

- Feature dependencies

- Longer computation time

$$
\begin{array}{cccc}
f_1 & f_2 & f_3 & f_4
\end{array}
$$

$$
\begin{array}{c}
s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5
\end{array}
\begin{pmatrix}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & x_{23} & x_{24} \\
x_{31} & x_{32} & x_{33} & x_{34} \\
x_{41} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{pmatrix}
\begin{pmatrix}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{pmatrix}
$$

0.0

TU/e

# Embedded methods

- Machine learning algorithms

TU/e

# Embedded methods

- Machine learning algorithms

- Weights

- Selection

- Low computation time

- Some feature dependencies

$$
\begin{array}{c}
\begin{array}{cccc} f_1 & f_2 & f_3 & f_4 \end{array} \\
\begin{array}{c} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{array}
\left(
\begin{array}{cccc}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & x_{23} & x_{24} \\
x_{31} & x_{32} & x_{33} & x_{34} \\
x_{41} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{array}
\right)
\left(
\begin{array}{c}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{array}
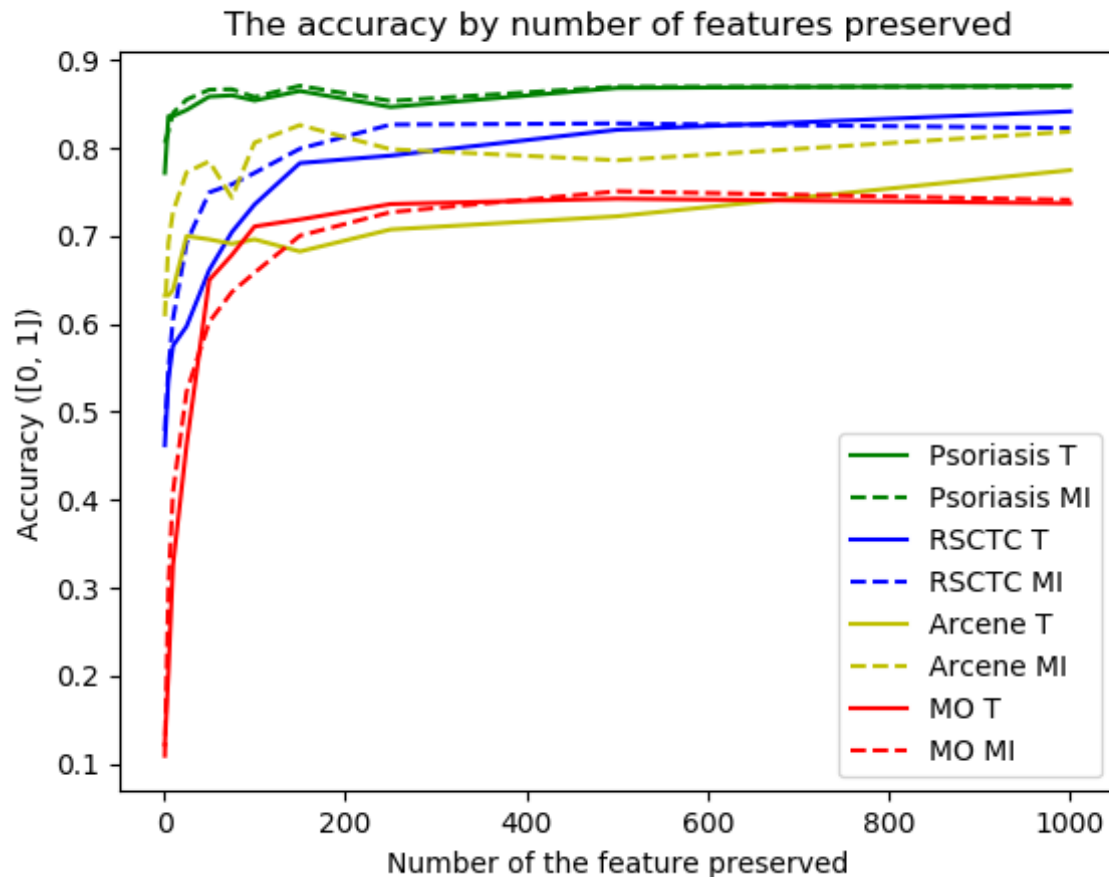\right)
\end{array}
$$

Alg:

TU/e

# Experiment 1

Filter methods

- Evaluation functions
  - T-test/ANOVA
  - Mutual Information
- Feature selection
  - Minimum 1
  - Maximum 1000
- Quality testing
  - Logistic regression accuracy ([0, 1])

TU/e

# Experiment 1

- Evaluation functions
- 200 features



The accuracy by number of features preserved
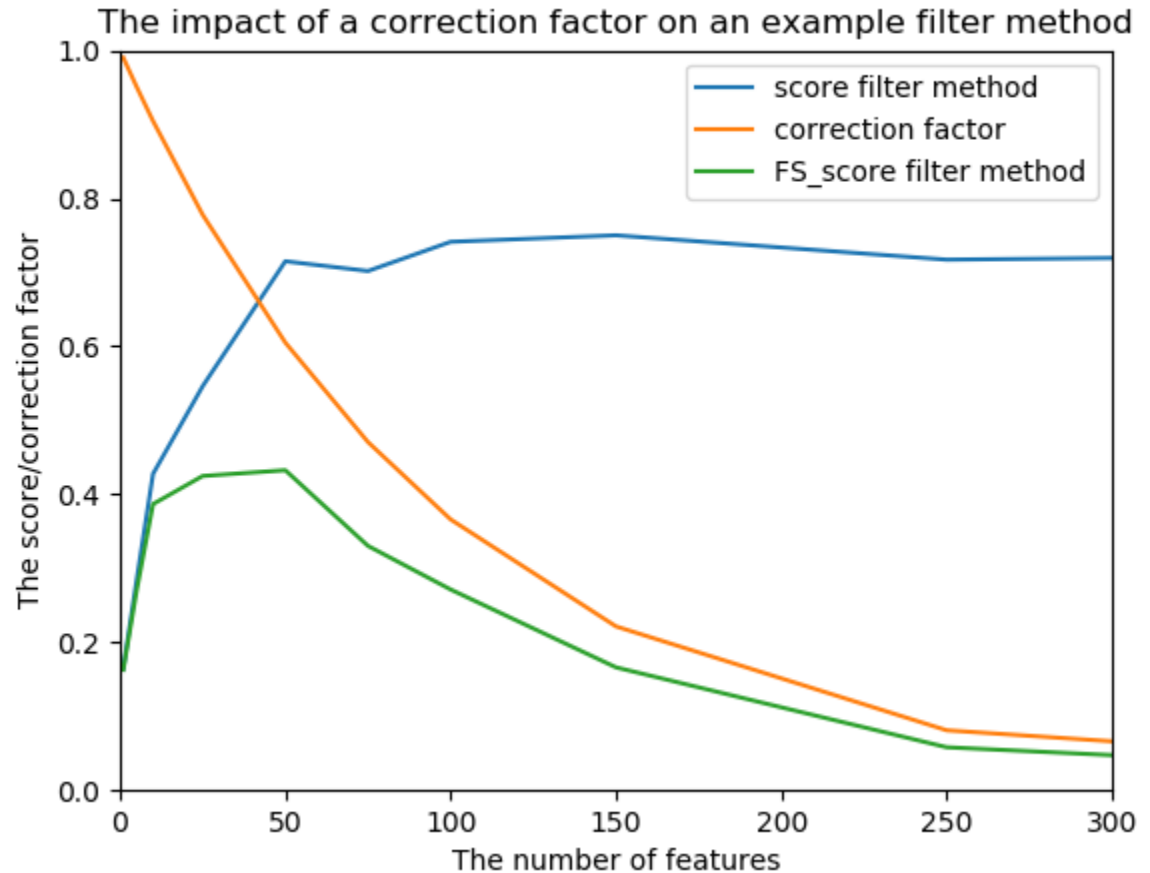
A Computational Biology Framework

# FS_score

- Combination
  - Prediction score
  - Features preserved

- Correction factor

$$FS\_score = score * \beta^{\#features}$$
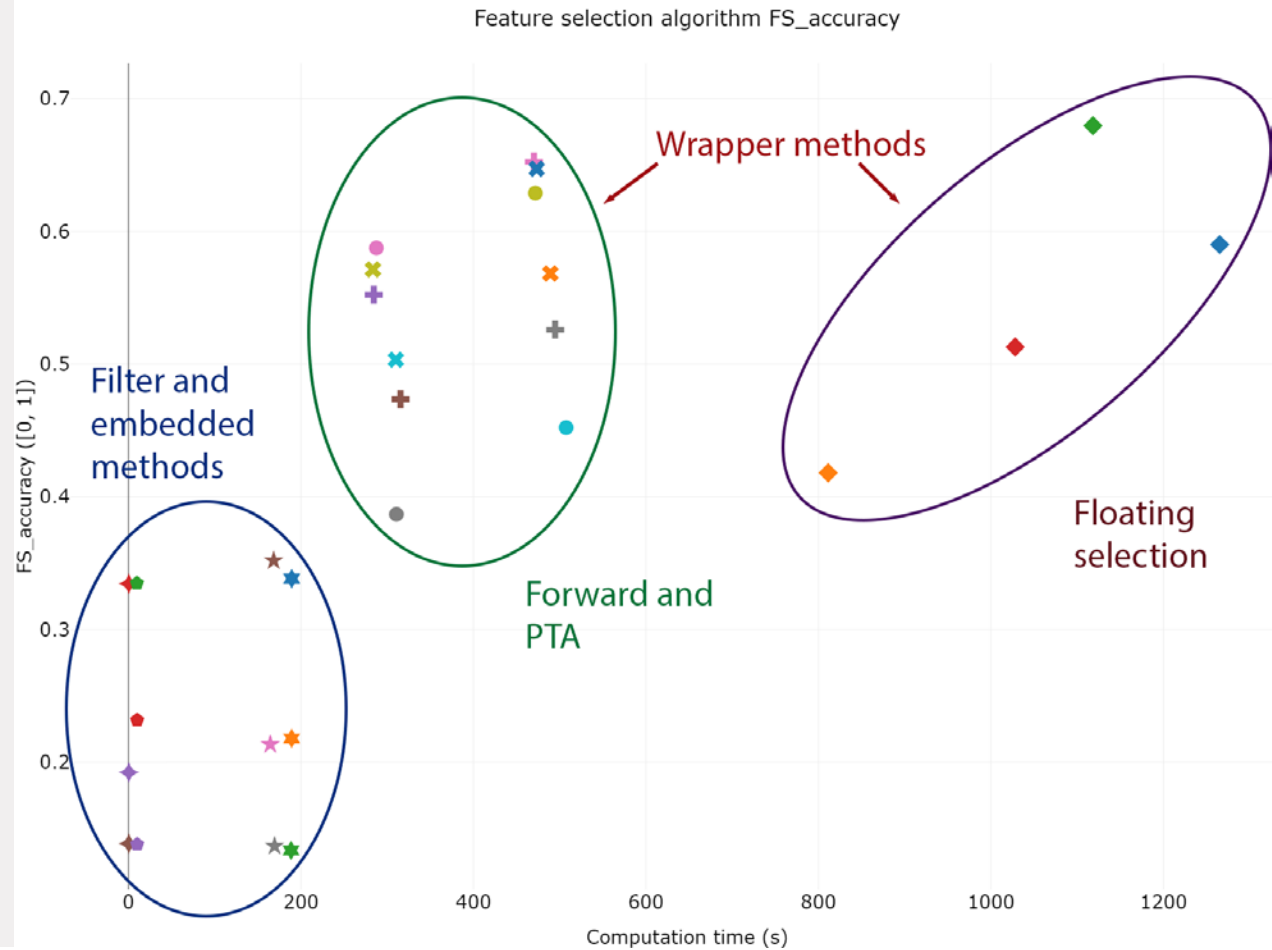
  - $\beta$=0.99

- Optimum

- Prediction scores



The impact of a correction factor on an example filter method

TU/e

# Experiment 2

Feature selection methods quality

- Quality testing
  - FS_score
  - Logistic regression accuracy ([0, 1])

TU/e

# Experiment 2

- Filter and embedded methods faster

- Wrapper methods better



Feature selection algorithm FS_accuracy

# Automated Machine learning

- Search for best analysis pipeline

- Preprocessing algorithms

- Machine learning algorithms

- Hyperparameters

A Computational Biology Framework

# Genetic Algorithms

- Solutions

- Evolving solutions
  - Mutations
  - Crossovers

- Selecting solutions

TU/e

# Tree-based Pipeline Optimization tool (TPOT)

- Genetic algorithms

- Searches for and stores best pipeline

- Issues with many features

TU/e

# Experiment 3

TPOT for many features

- FS_score
  - Accuracy ([0, 1])
- 2 Additions
  - Bias towards feature selection
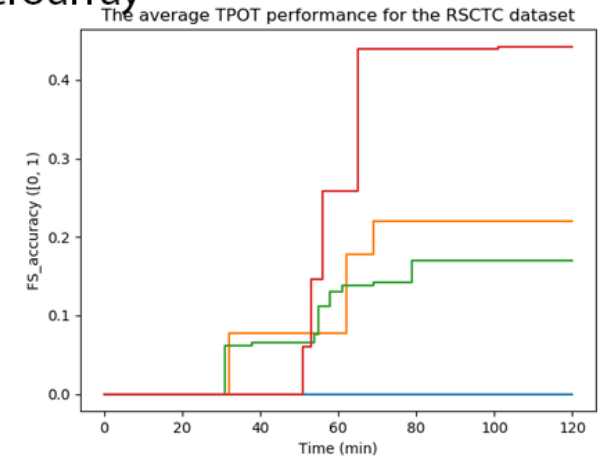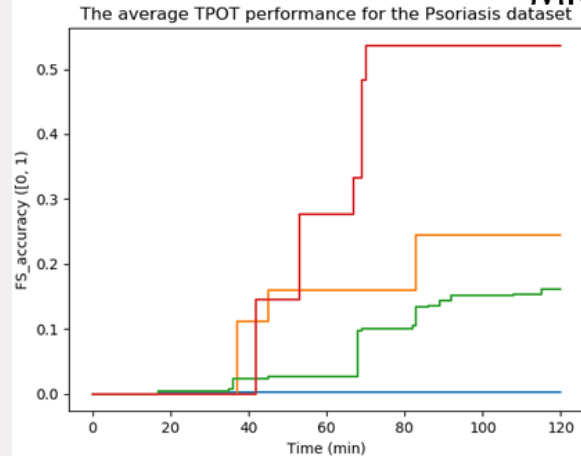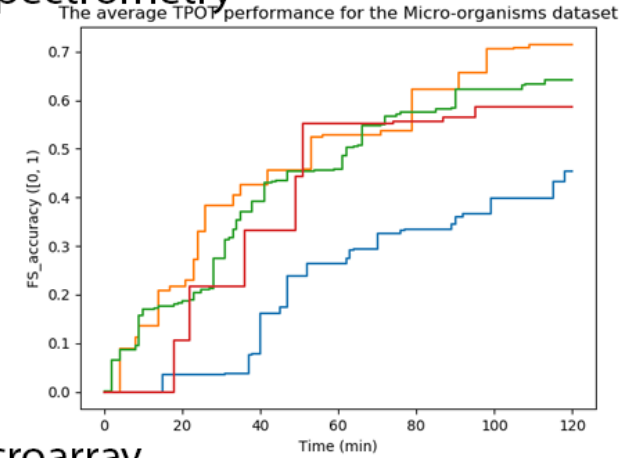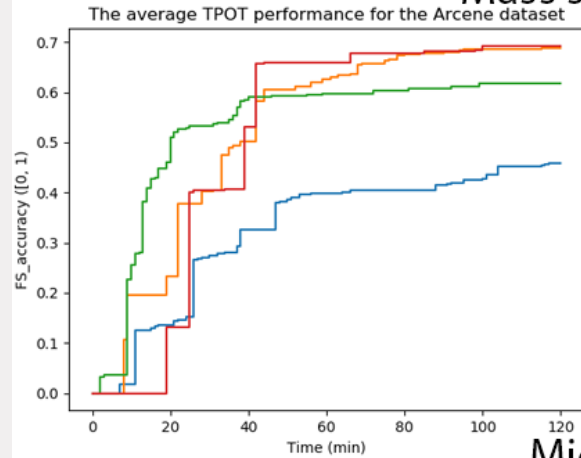  - New set feature selection algorithms

TU/e

# Experiment 3

- Improved

- Microarray both best

- Mass spectrometry debatable

## Legend
- **Blue** — No addition
- **Orange** — New algorithms
- **Green** — Feature selection bias
- **Red** — Both additions



Mass spectrometry

The average TPOT performance for the Arcene dataset

The average TPOT performance for the Micro-organisms dataset

Microarray

The average TPOT performance for the Psoriasis dataset

The average TPOT performance for the RSCTC dataset

TU/e

# Missing Value Handling

A Computational Biology Framework

# Missing Value Handling

What missing value handling methods show the best performance and as such should be added to the framework?

Missing value handling methods
- List Deletion
- Single imputation
- Multiple imputation

Datasets
- Clinical
- Survey

## List deletion

- Available Case Analysis (ACA)
- Complete Case Analysis (CCA)
  - Weighted Case Analysis (WCA)

$$
\begin{array}{c}
 & f_1 \quad f_2 \quad f_3 \quad f_4 \\
\begin{array}{c} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{array}
\left(
\begin{array}{cccc}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & \ldots & x_{24} \\
x_{31} & x_{32} & \ldots & \ldots \\
\ldots & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{array}
\right)
\left(
\begin{array}{c}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{array}
\right)
\end{array}
$$

TU/e

# Single imputation

A Computational Biology Framework

# Single imputation

- Mean imputation
  - Median
  - Mode
- Hot deck imputation

$$
\begin{array}{c}
\phantom{x} \\
s_1 \\
s_2 \\
s_3 \\
s_4 \\
s_5
\end{array}
\begin{pmatrix}
f_1 & f_2 & f_3 & f_4 \\
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & x_{13} & x_{24} \\
x_{31} & x_{32} & x_{53} & x_{24} \\
x_{31} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{pmatrix}
\begin{pmatrix}
y_1 \\
y_2 \\
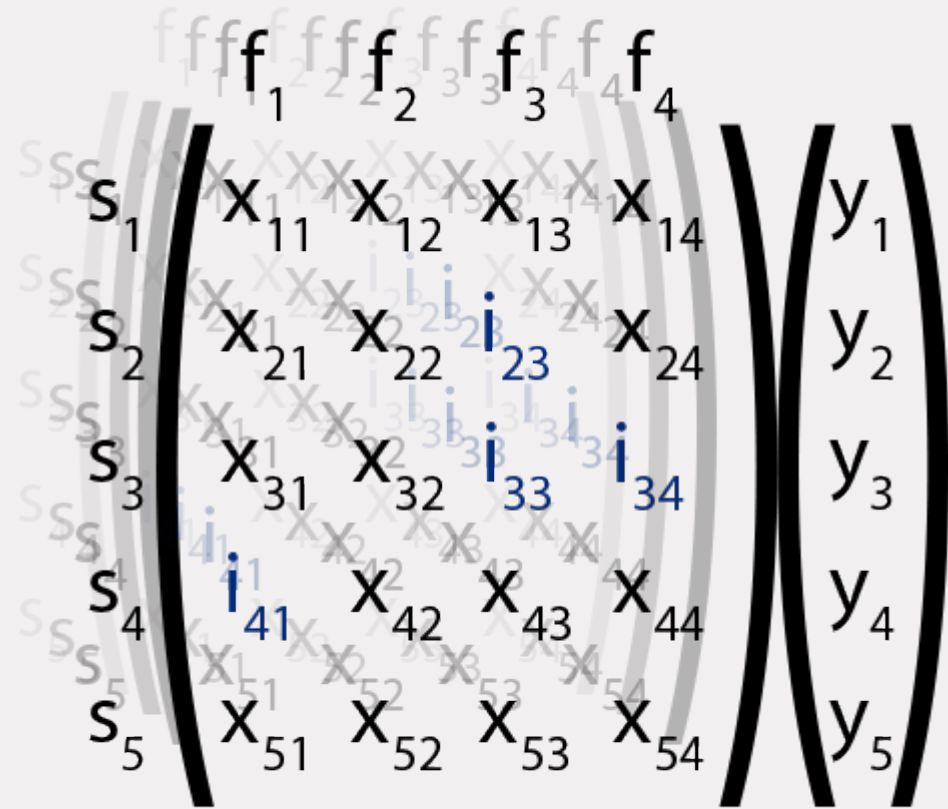y_3 \\
y_4 \\
y_5
\end{pmatrix}
$$

TU/e

## Single imputation

- Mean imputation
  - Median
  - Mode
- Hot deck imputation

- Advanced imputation
  - Regression imputation
  - Nearest neighbour imputation (kNN)

$$
\begin{array}{c c c c c}
 & f_1 & f_2 & f_3 & f_4 \\
s_1 & x_{11} & x_{12} & x_{13} & x_{14} \\
s_2 & x_{21} & x_{22} & \dots & x_{24} \\
s_3 & x_{31} & x_{32} & \dots & \dots \\
s_4 & \dots & x_{42} & x_{43} & x_{44} \\
s_5 & x_{51} & x_{52} & x_{53} & x_{54}
\end{array}
\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{pmatrix}
$$

TU/e

# Multiple Imputation

- Multiple datasets
- Multiple Imputation Chained Equations (MICE)

$$
\begin{array}{c}
 & f_1 \quad f_2 \quad f_3 \quad f_4 \\
\begin{array}{c} s_1 \\ s_2 \\ s_3 \\ s_4 \\ s_5 \end{array}
\begin{pmatrix}
x_{11} & x_{12} & x_{13} & x_{14} \\
x_{21} & x_{22} & i_{23} & x_{24} \\
x_{31} & x_{32} & i_{33} & i_{34} \\
i_{41} & x_{42} & x_{43} & x_{44} \\
x_{51} & x_{52} & x_{53} & x_{54}
\end{pmatrix}
\begin{pmatrix}
y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5
\end{pmatrix}
\end{array}
$$

TU/e

# Experiment 1

Feature distributions
- Mean
- Standard deviation
- T-test similarity
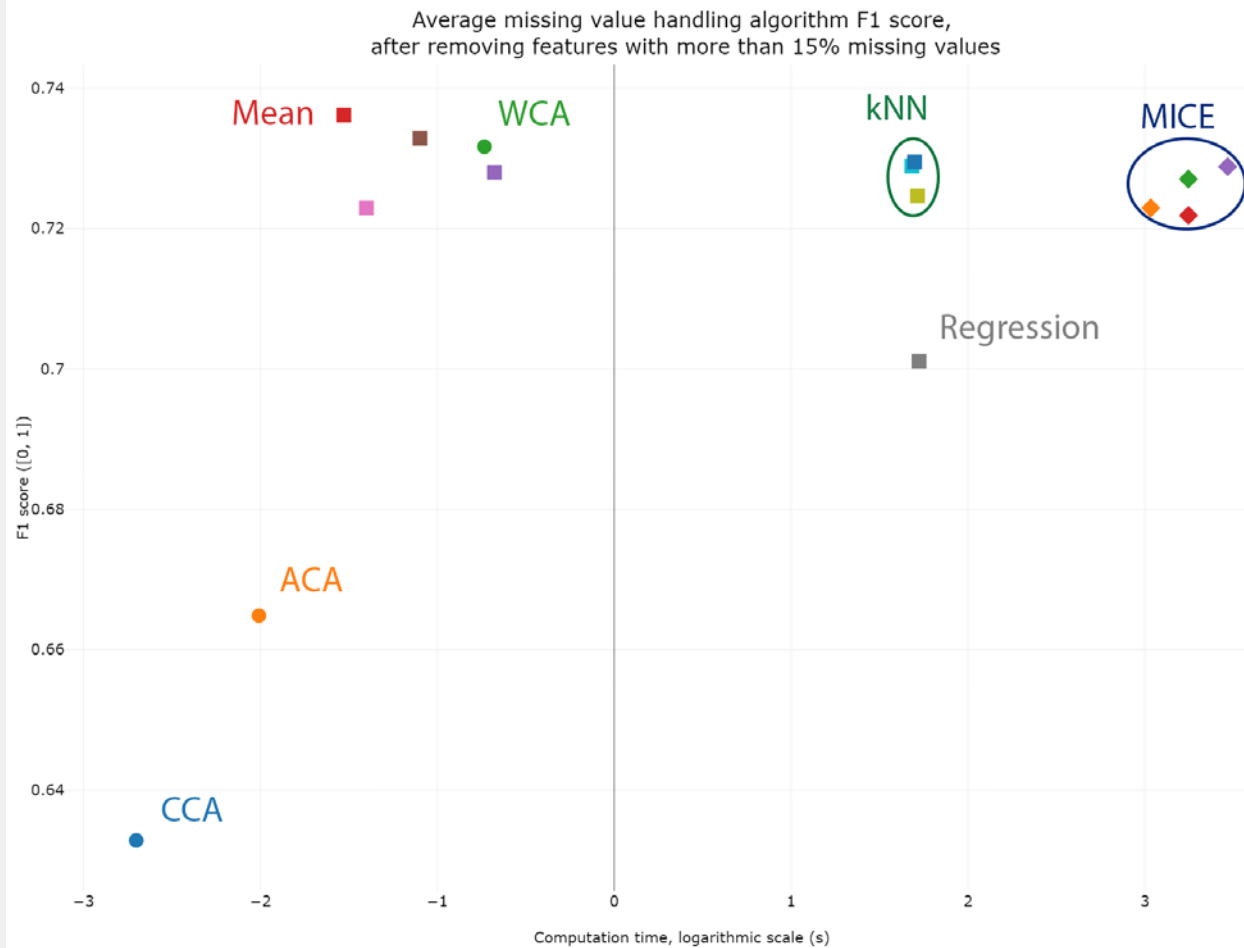
**Results**
- > 15% missing values removed

TU/e

# Experiment 2

Missing value handling quality

- >15% missing values removed
- Quality output
  - Logistic regression F1-score ([0, 1])

TU/e

## Experiment 2

- List Deletion

- Mean Imputation

- Regression, kNN and MICE



Average missing value handling algorithm F1 score, after removing features with more than 15% missing values

TU/e

# Dataset exploration

# Dataset exploration

Which initial analyses help find suitable preprocessing and data analysis algorithms?

Exploration
- Dataset issues
- Meta-features
- Package *metalearn*

Datasets
- Clinical
- Mass spectrometry

# Analysis 1

Dataset exploration with package *metalearn*

- Issues <> meta-features

**Result**

- Issues recognizable
  - Multicollinearity
  - Feature irrelevance
- Needs clarification

TU/e

# Analysis 2

Dataset exploration with package *metalearn*

- *metalearn* addition
  - meta-features
  - Outlier detection
  - Plots

**Results**

- All issues explained
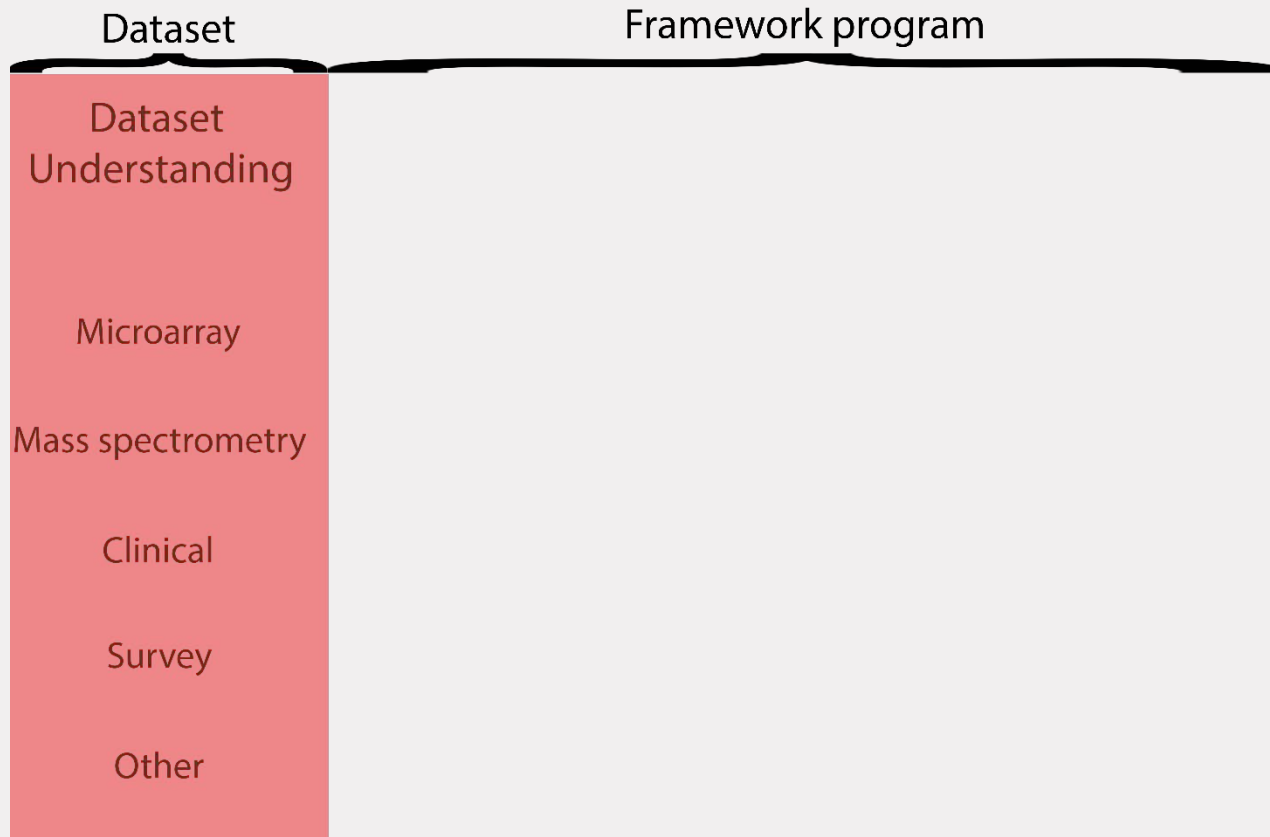- Better understanding

TU/e

# Framework cBioF

# Framework cBioF

What aspects are of importance to be included in a framework for Biomedical Engineers for more efficient data analysis?
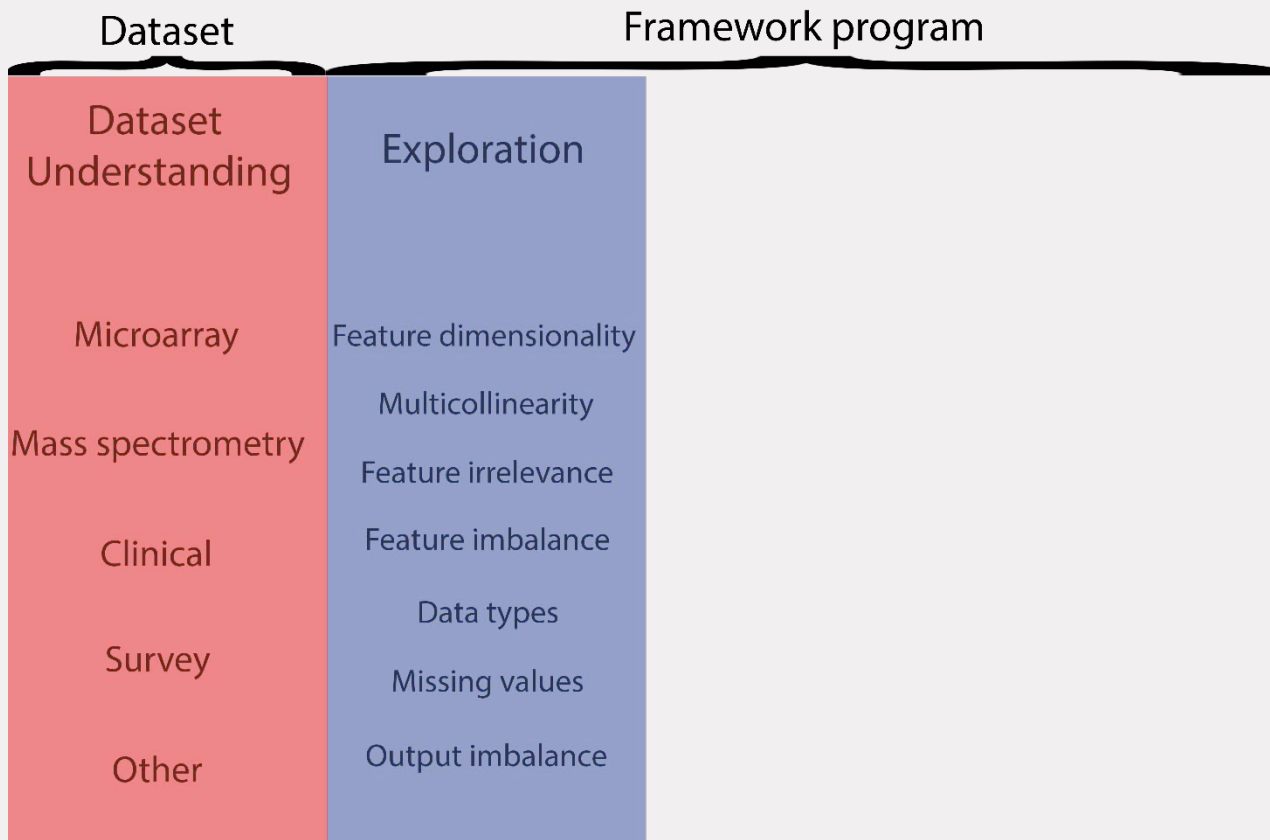
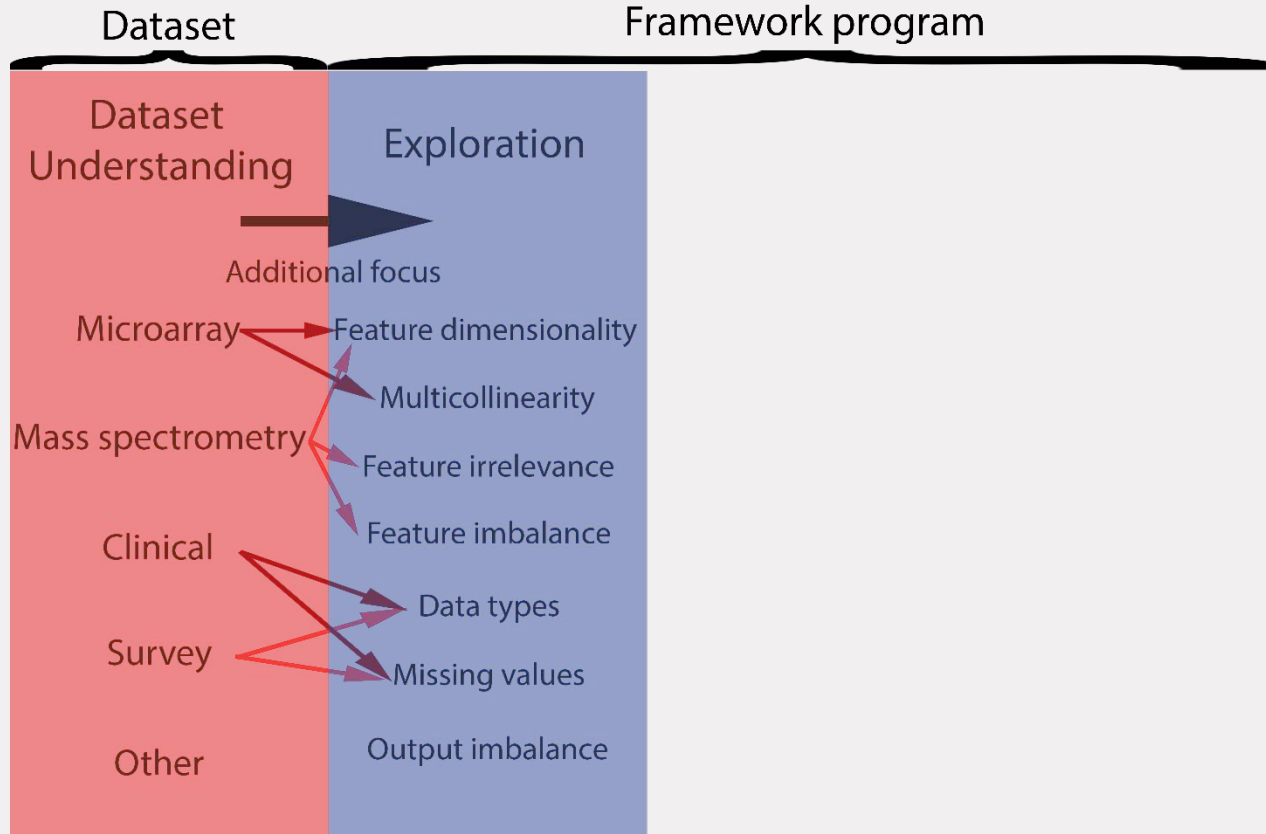Phases
- Datasets
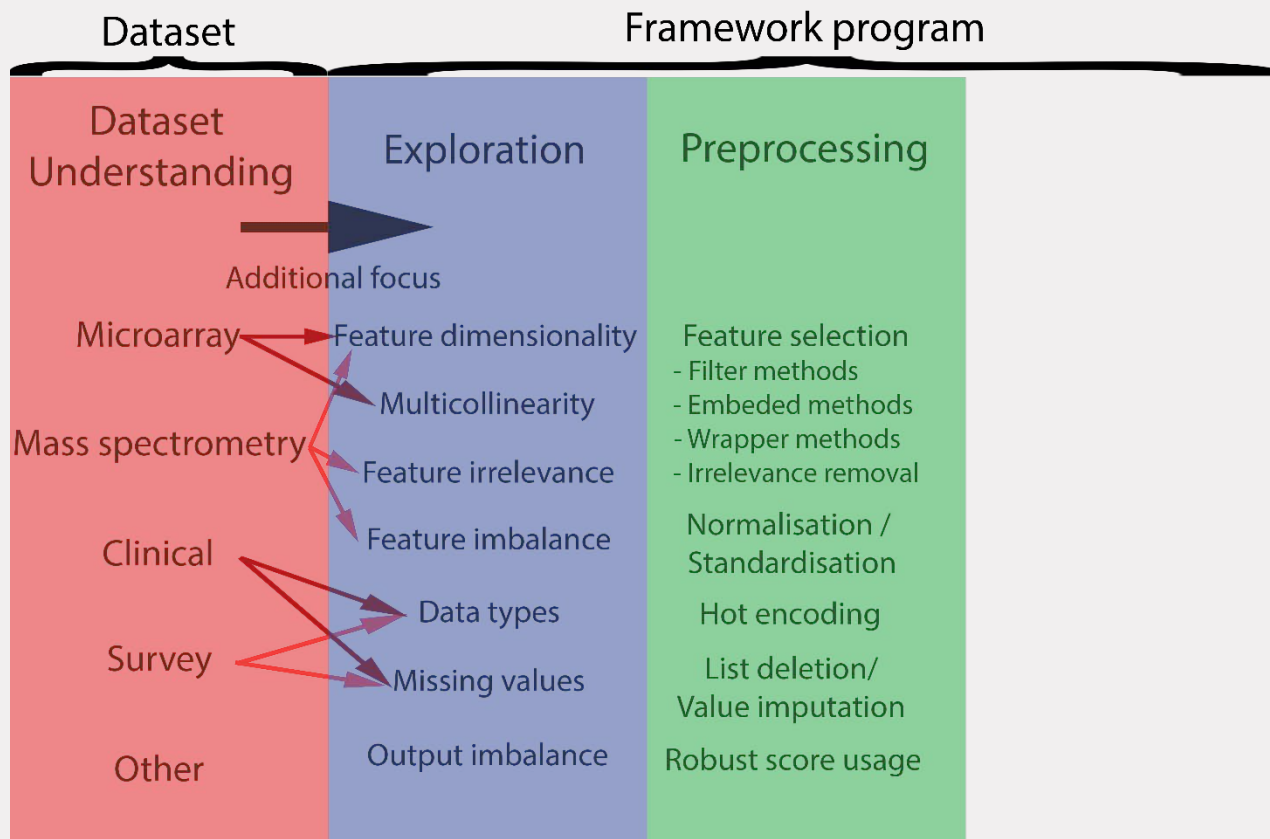- Exploration
- Preprocessing
- Analysis

# cBioF layout



A Computational Biology Framework

TU/e

# cBioF layout

Dataset      Framework program

| Dataset Understanding | Exploration |
|---|---|
| Microarray | Feature dimensionality |
| Mass spectrometry | Multicollinearity |
| | Feature irrelevance |
| Clinical | Feature imbalance |
| | Data types |
| Survey | Missing values |
| Other | Output imbalance |

TU/e

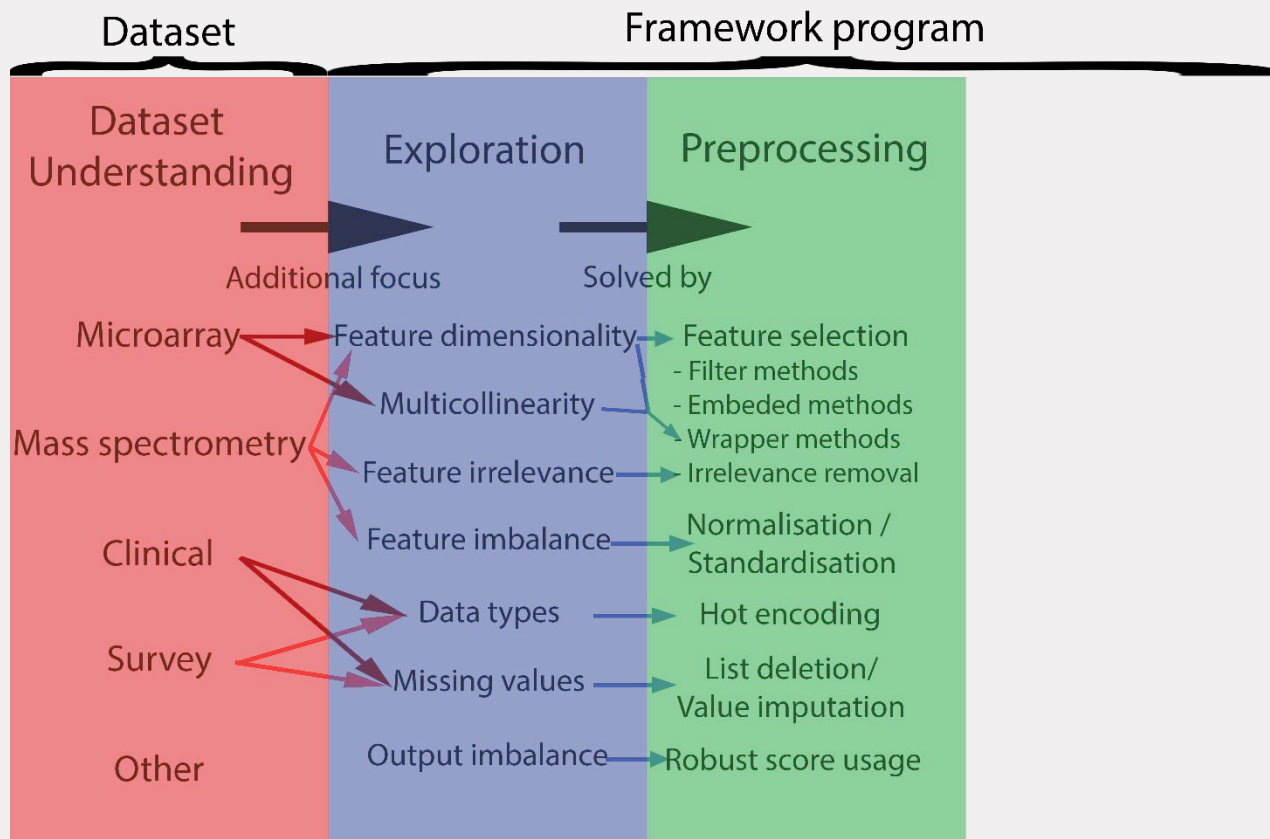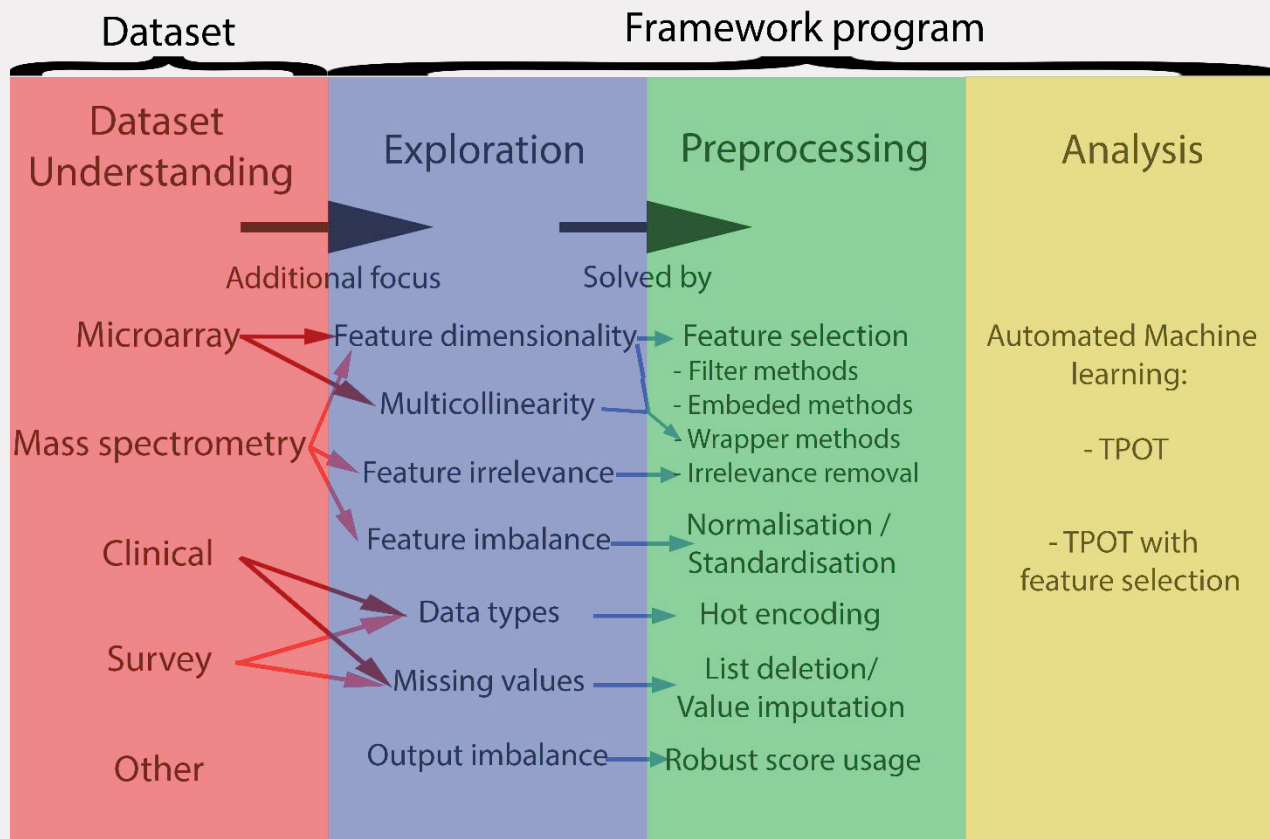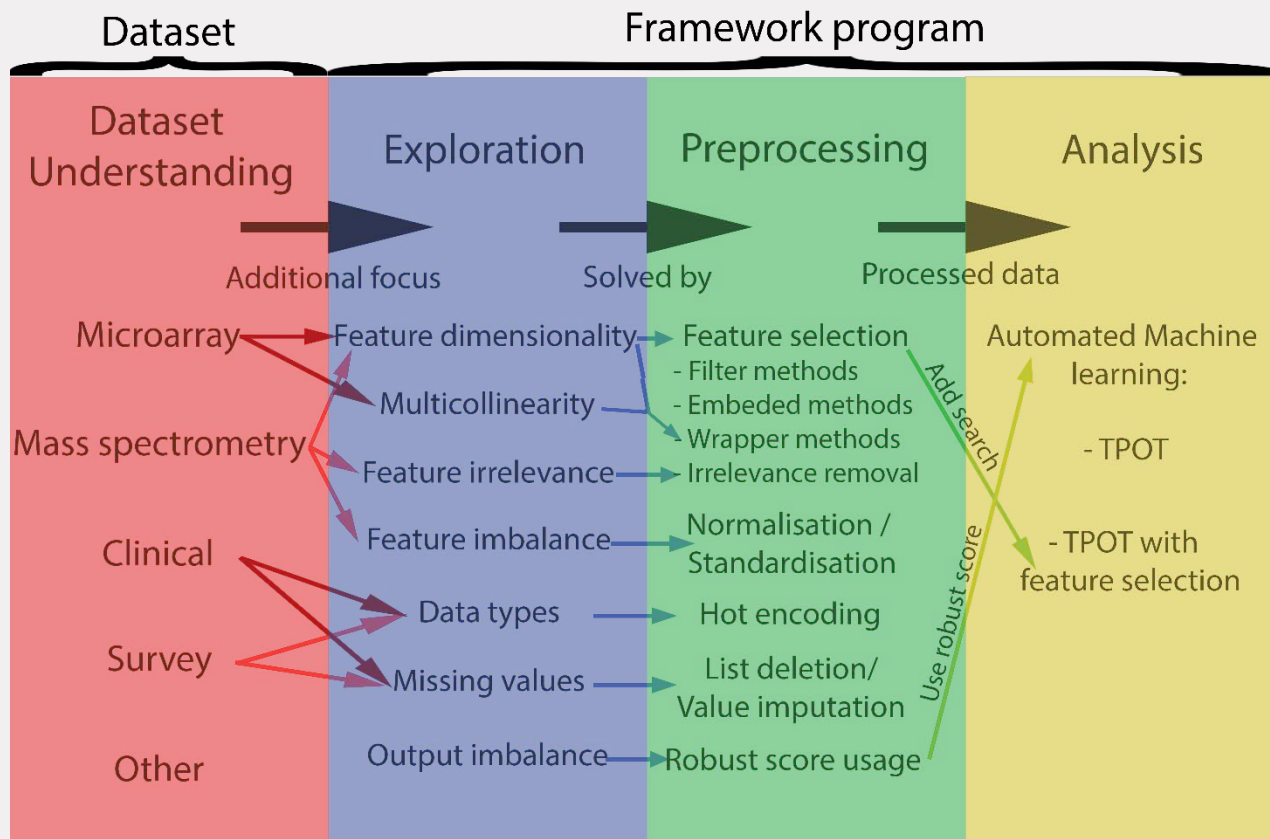# cBioF layout

# cBioF layout

TU/e

# cBioF layout

# cBioF layout

# Conclusions

## Feature selection

- 200 features threshold
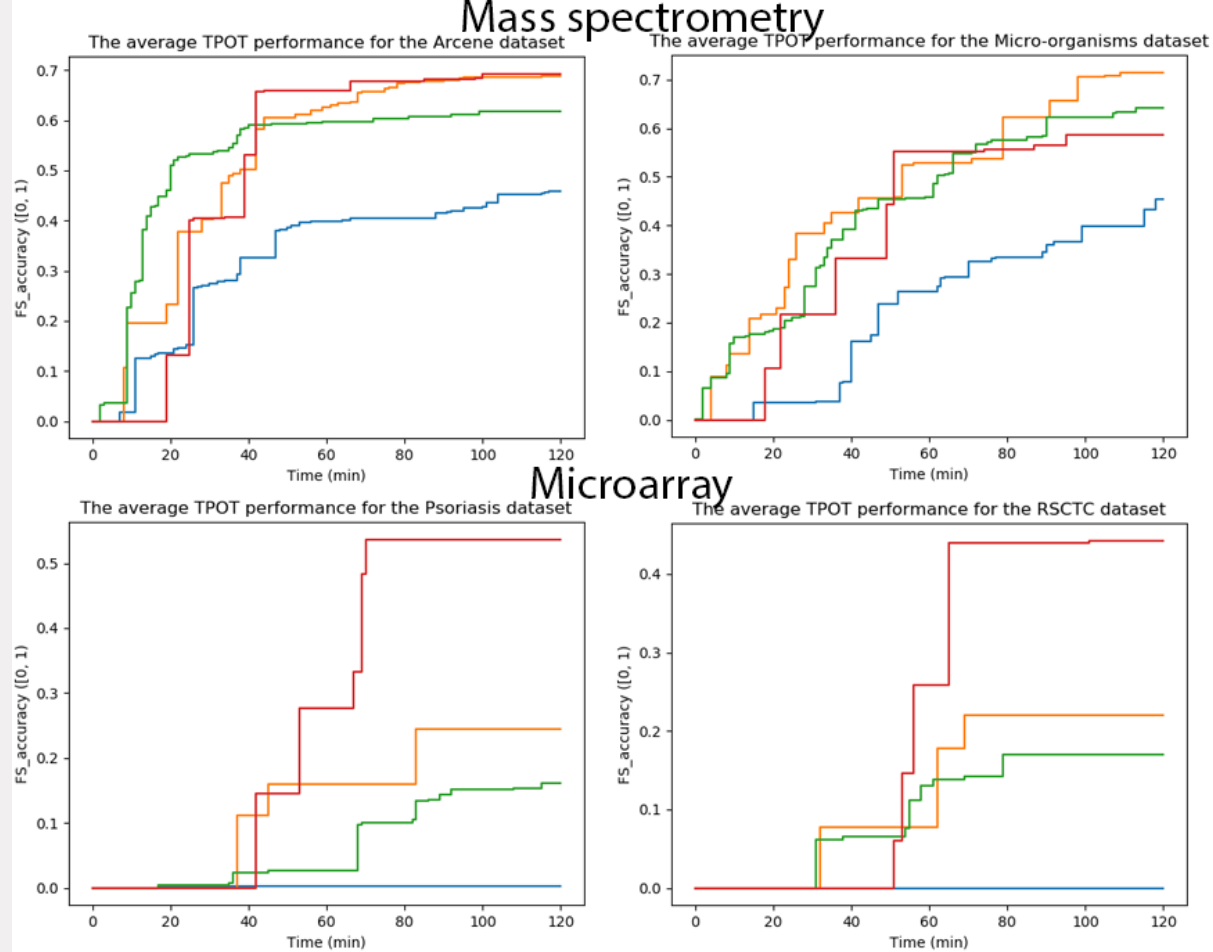- Wrapper methods perform best

## Missing value handling

- > 15% missing values
- Mean imputation

## Dataset exploration

- Possible with additions *metalearn*

TU/e

# Conclusion TPOT

- Problems many features

- FS_score

- Feature selection added:
  - Bias
  - Algorithms

- Better performance



Mass spectrometry

The average TPOT performance for the Arcene dataset

The average TPOT performance for the Micro-organisms dataset

Microarray

The average TPOT performance for the Psoriasis dataset

The average TPOT performance for the RSCTC dataset

TU/e

# Future work

## Feature selection

- 200 features threshold
- Combination filter and wrapper

## Missing value handling

- New datasets

## Framework

- Multicollinearity
- Other dataset issues

TU/e