



University of Minho
School of Engineering



Machine Learning and Decision-Making

ADI @ LEI/3º, MiEI/4º - 2º Semestre
Filipe Gonçalves, Inês Alves, Cesar Analide

Part VIII – April 2022

Contents

2

Model Validation

Feature Selection

Hands On

- Model Validation Techniques
- Feature Selection
- Hands On

Validation Techniques

3

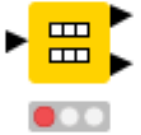
Model Validation

Feature Selection

Hands On

Cross Validation





Hold-out Validation

4

Model Validation

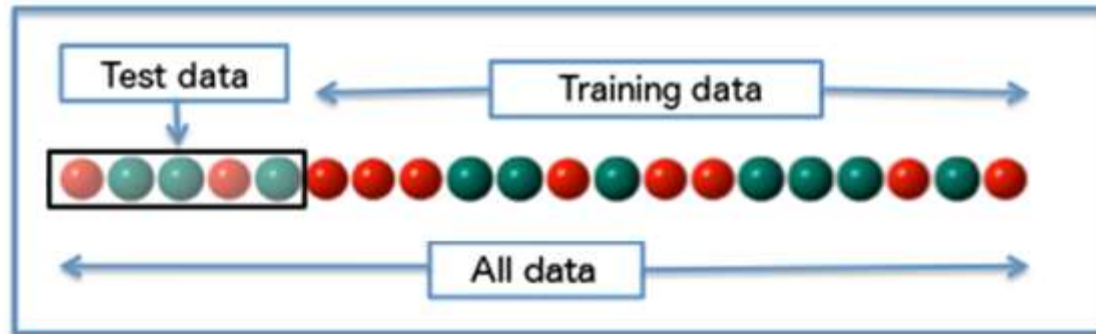
Feature Selection

Hands On

But before going into Cross Validation, a **model validation technique**, do you know you have already been using another model validation technique?

It is known as **Hold-Out Validation**!

In essence, it means we **validate the model on unseen data**, i.e., we use a “partitioning method” to split the learning and the testing data once. This means we **hold-out a subset of data** for testing (80/20; 75/25; 65/35...)!





Cross Validation

5

Model Validation

Feature Selection

Hands On

Cross validation is another model validation technique!

The goal is to have an accurate metric of how the model will perform in practice.

In essence, it consists in dividing the dataset into k folds. In each run of the model, $k-1$ folds are used for training and 1 fold (the remaining) is used as test. Keep repeating the process until all folds have been used for testing.

The final error metric is based on the mean value of all error metrics.

$$E = \frac{1}{k} \sum_{i=1}^k E_i$$



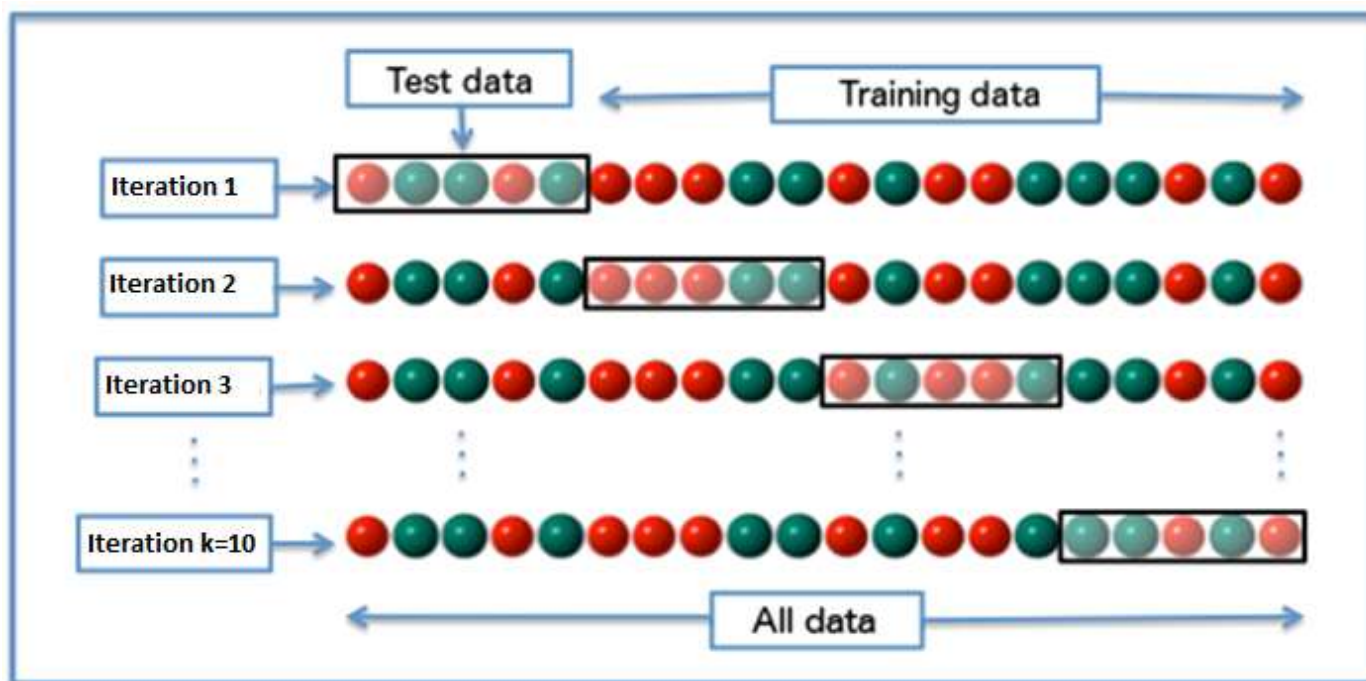
k-fold Cross Validation

6

Model Validation

Feature Selection

Hands On



Usually, $k=10$.



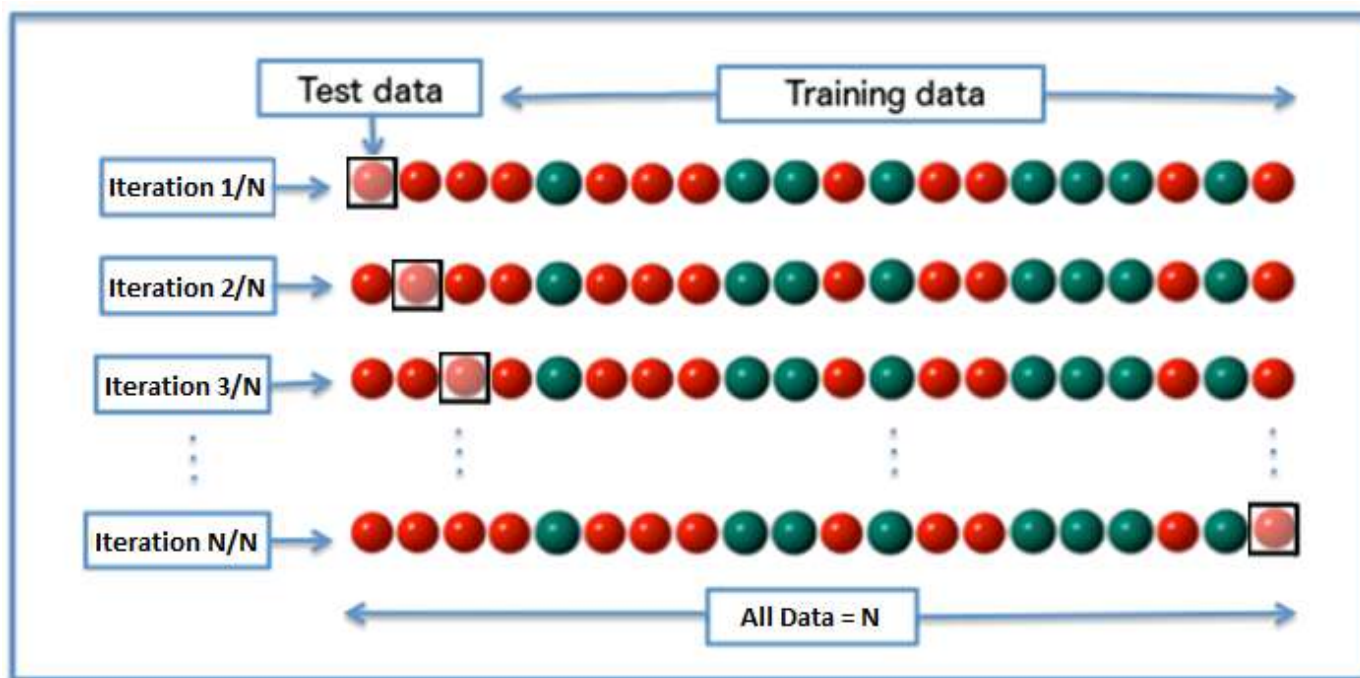
Leave-one-out Cross Validation ($k=N$)

7

Model Validation

Feature Selection

Hands On



The special case of having $k=N$. Expensive...
But a good approach when we have a **small dataset**.



Cross Validation - How many folds?

8

Model Validation

Feature Selection

Hands On

Well, ...

A **greater number of folds** will lead to a **better error estimate** of the model, a **lower bias** and **less overfitting**! However, it comes with a **higher computational cost**!

If we have a **large dataset**, a **smaller k** may be enough since we will have a larger amount of data for training. If we have a **small dataset**, we may want to **use leave-one-out cross validation** to maximize the amount of data for training...

In reality, **k depends on N!!**

Rule of thumb → **k=10!**

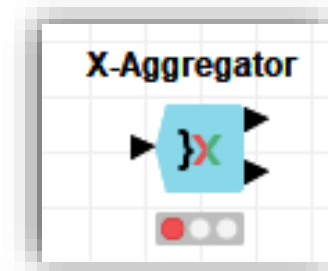
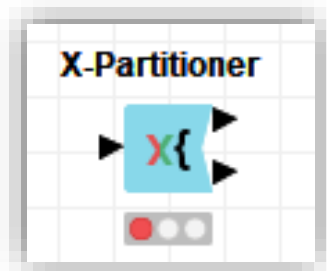
k-fold Cross Validation

9

Model Validation

Feature Selection

Hands On



Dialog - 0:80 - X-Partitioner

File

Standard settings | Flow Variables | Memory Policy

Number of validations: 10

Linear sampling: ☐

Random sampling: ☒

Stratified sampling: ☐

Class column: S quality

☐ Random seed: 0

Leave-one-out: ☐

OK Apply Cancel ?

Dialog - 0:81 - X-Aggregator (Aggregating)

File

Standard settings | Flow Variables | Memory Policy

Target column: S quality

Prediction column: S Prediction (quality)

☐ Add column with fold id

OK Apply Cancel ?

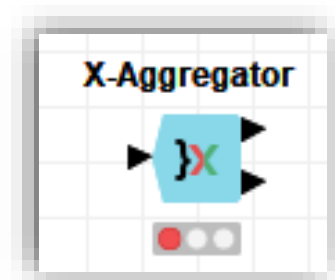
k-fold Cross Validation

10

Model Validation

Feature Selection

Hands On



▲ Error rates - 0:81 - X-Aggregator (Aggregating) — □ ×

File Hilite Navigation View

Table "default" - Rows: 10 Spec - Columns: 3 Properties Flow Variables

Row ID	D Error in %	I Size of Test Set	I Error Count
fold 0	38.281	128	49
fold 1	35.156	128	45
fold 2	44.531	128	57
fold 3	42.969	128	55
fold 4	42.969	128	55
fold 5	41.406	128	53
fold 6	41.406	128	53
fold 7	40.625	128	52
fold 8	41.406	128	53
fold 9	42.52	127	54

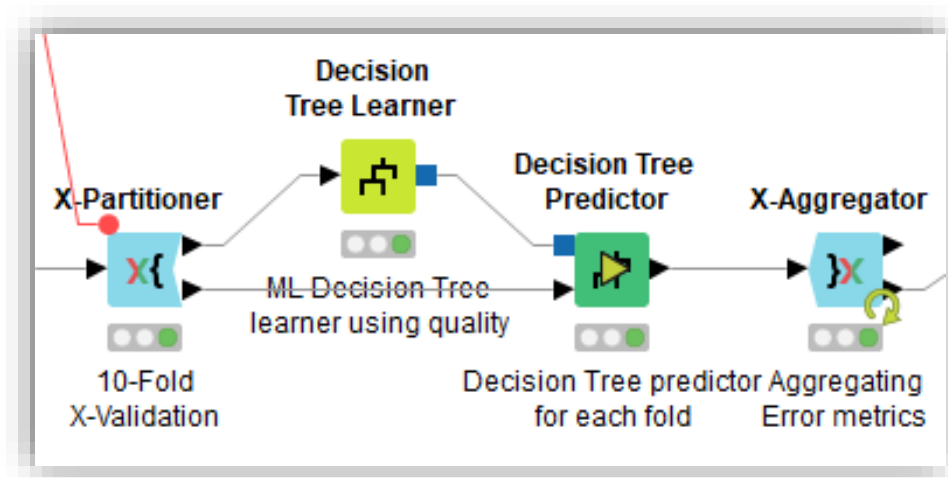
k-fold Cross Validation

11

Model Validation

Feature Selection

Hands On



Feature Selection

12

Model Validation

Feature Selection

Hands On

(dimensionality reduction)



Feature Selection

13

Model Validation

Feature Selection

Hands On

Feature Selection (or dimensionality reduction)

Rationale:

Which features should we use to create a predictive model? Select a sub-set of the most important features to reduce dimensionality.

The removal of unimportant features:

- May **affect significantly the performance of a model**
- **Reduces overfitting** (less opportunity to make decisions based on noise)
- **Improves accuracy**
- Helps **reducing the complexity** of a model (reduces training time)

Feature Selection

14

Model Validation

Feature Selection

Hands On

Feature Selection (or dimensionality reduction)

Rationale:

Which features should we use to create a predictive model? Select a sub-set of the most important features to reduce dimensionality.

The removal of unimportant features:

- May **affect significantly the performance of a model**
- **Reduces overfitting** (less opportunity to make decisions based on noise)
- **Improves accuracy**
- Helps **reducing the complexity** of a model (reduces training time)

What can we remove:

- **Redundant features** (duplicate)
- **Irrelevant and unneeded features** (non-useful)

Feature Selection Methods:

- **Filter methods**
- **Wrapper methods**
- **Embedded methods**

Feature Selection

15

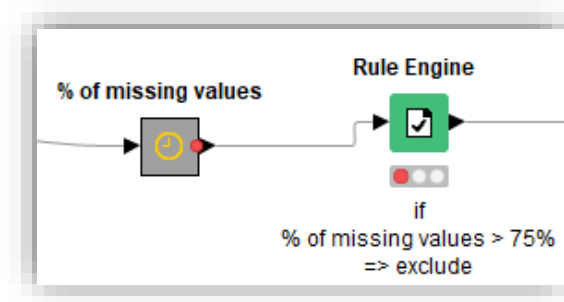
Model Validation

Feature Selection

Hands On

Filter Methods:

- i. Remove a feature if the **percentage of missing values** is **higher than** a threshold



- ii. Use the chi-square test to measure the **degree of dependency between a feature** and the **target class**
 - For each feature calculate X^2
 - Normalize X^2 and sort in descending order
 - Select n features with the highest importance (or those that are above the threshold)

Feature Selection

16

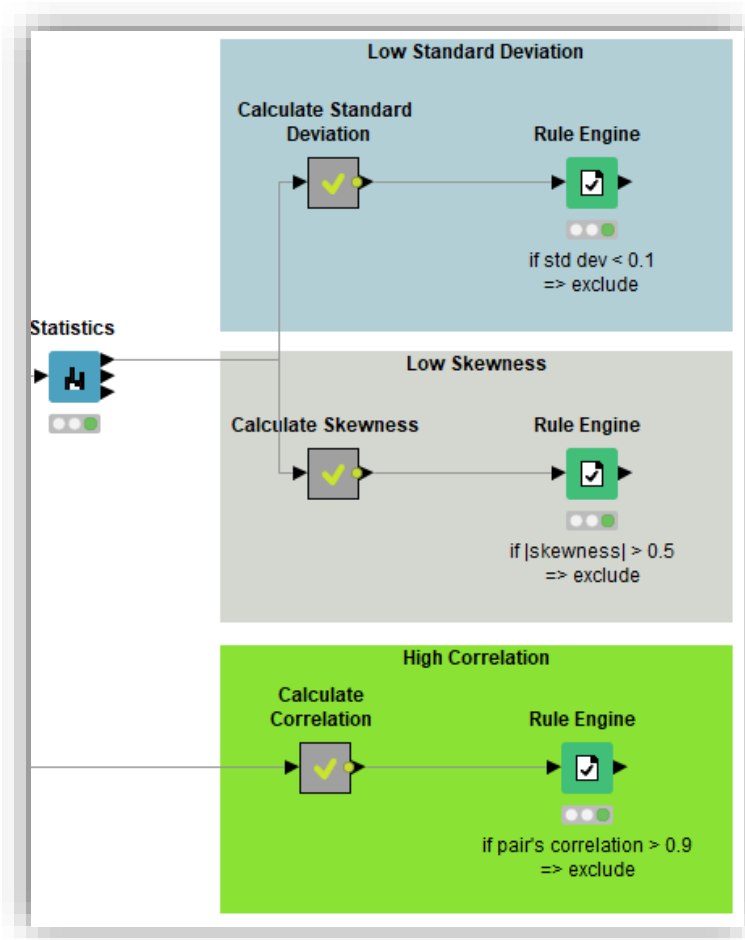
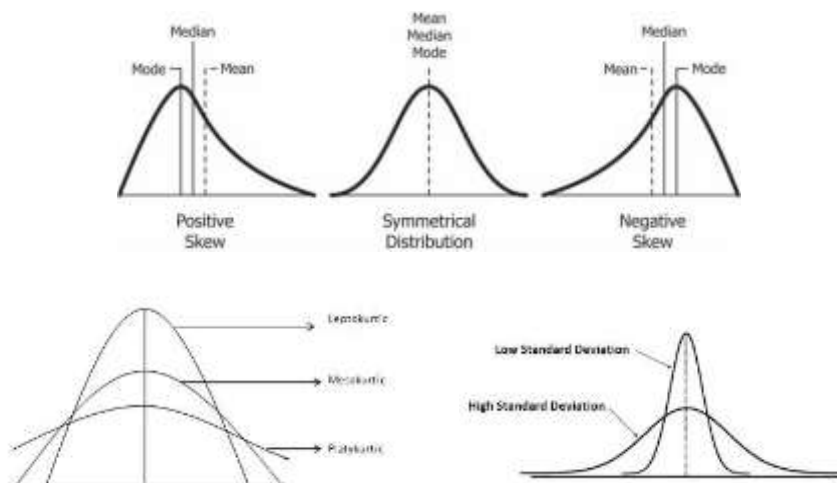
Model Validation

Feature Selection

Hands On

Filter Methods:

- iii. Remove feature if **low standard deviation**
- iv. Remove feature if data are **highly skewed**
- v. Remove features that are **highly correlated** between each other



Feature Selection

17

Model Validation

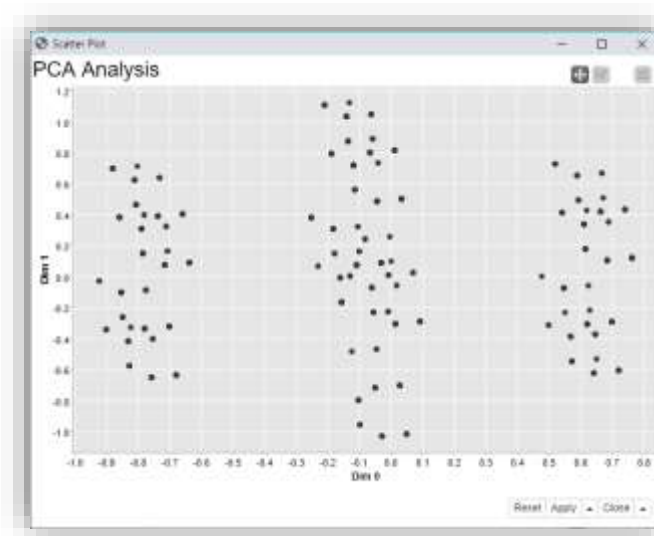
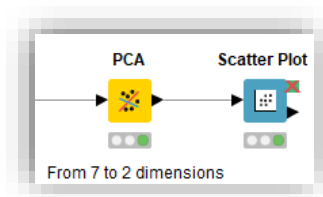
Feature Selection

Hands On

Filter Methods:

vi. **Principal Component Analysis (PCA)**

A technique to reduce the dimension of the feature space. The goal is to **reduce the number of features without losing too much information**. A popular application of PCA is for **visualizing higher dimensional data**.



Feature Selection

18

Model Validation

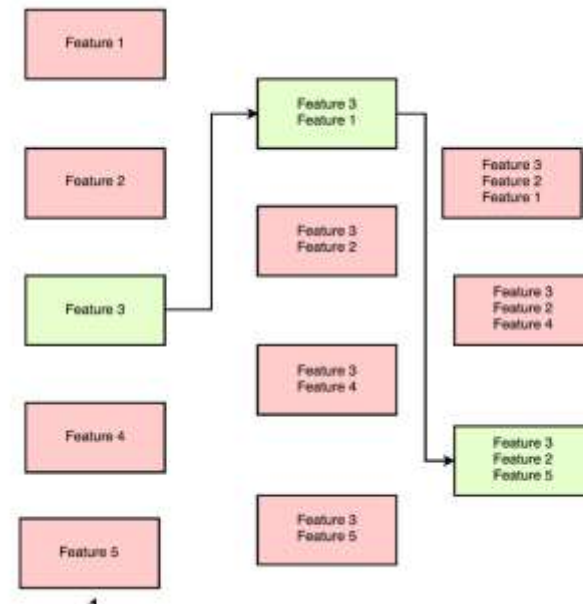
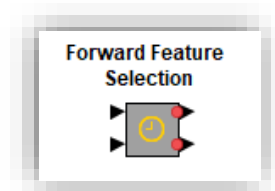
Feature Selection

Hands On

Wrapper Methods:

Use a **ML algorithm** to select the most important features! Select a set of features as a search problem, prepare different combinations, evaluate and compare them! Measure the “usefulness” of features based on the classifier performance.

Sequential Forward Selection



Feature Selection

19

Model Validation

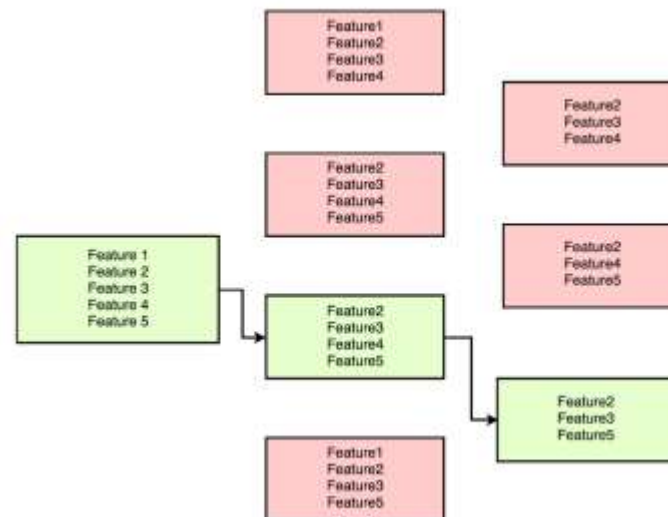
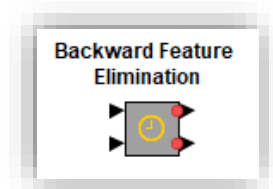
Feature Selection

Hands On

Wrapper Methods:

Use a **ML algorithm** to select the most important features! Select a set of features as a search problem, prepare different combinations, evaluate and compare them! Measure the “usefulness” of features based on the classifier performance.

Backward Feature Elimination



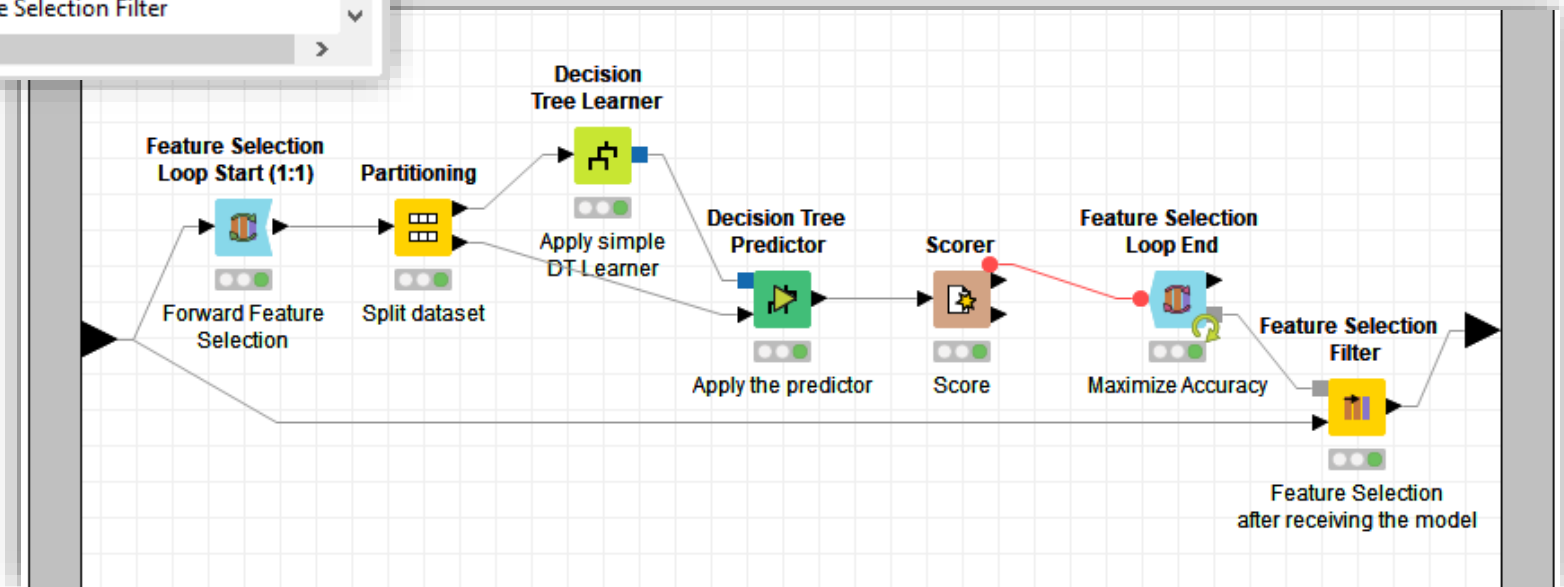
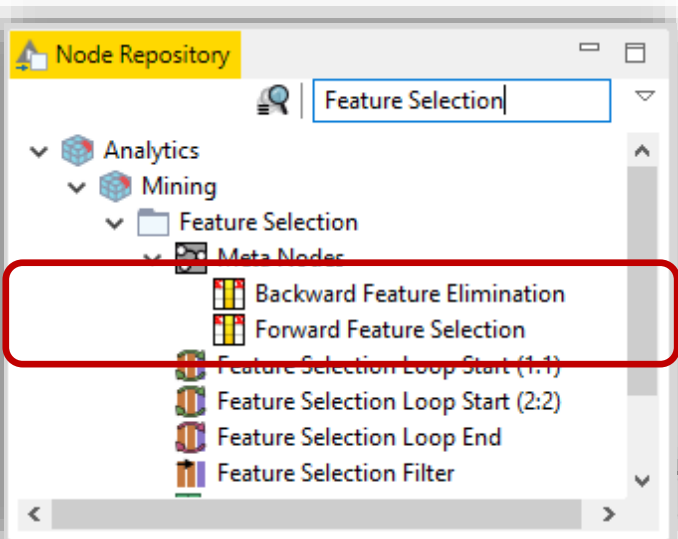
Feature Selection - Wrapper Methods

20

Model Validation

Feature Selection

Hands On



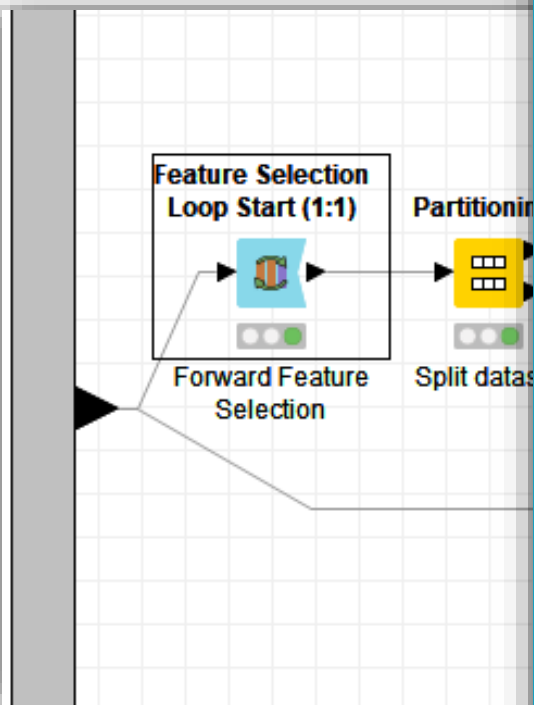
Feature Selection Loop Start

21

Model Validation

Feature Selection

Hands On



Dialog - 0:96:95 - Feature Selection Loop Start (1:1) (Forward Feature)

File

Options Flow Variables Memory Policy

The list on the left contains 'static' columns such as the target column.
The columns to choose from need to be in the list on the right.

☒ Manual Selection ☐ Wildcard/Regex Selection

Static Columns

Filter

S quality

☒ Enforce exclusion

Variable Columns ('Features')

Filter

D fixed acidity
D volatile acidity
D citric acid
D residual sugar
D chlorides
D free sulfur dioxide
D total sulfur dioxide
D pH

☐ Enforce inclusion

Feature selection strategy Forward Feature Selection

☐ Use threshold for number of features

Select threshold for number of features 20

OK Apply Cancel ?

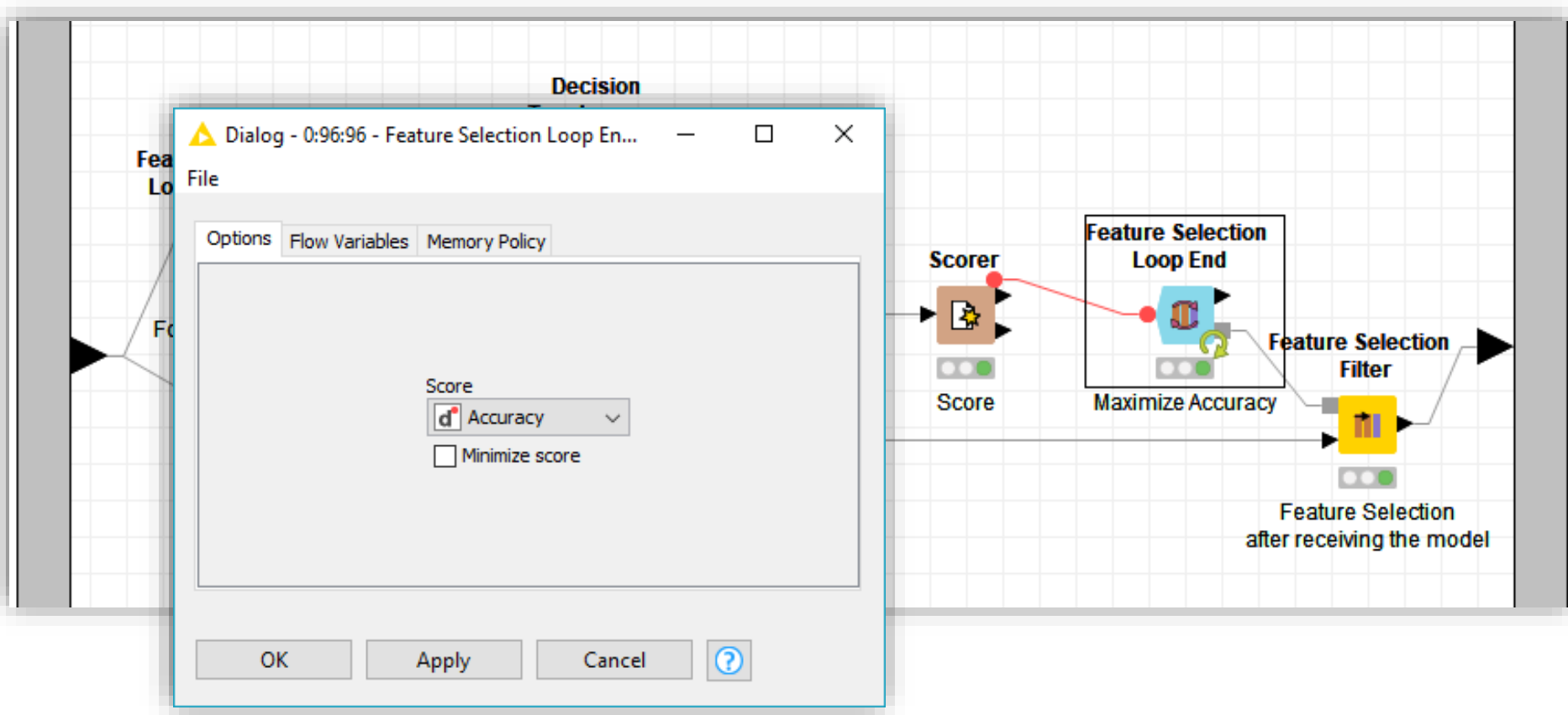
Feature Selection Loop End

22

Model Validation

Feature Selection

Hands On



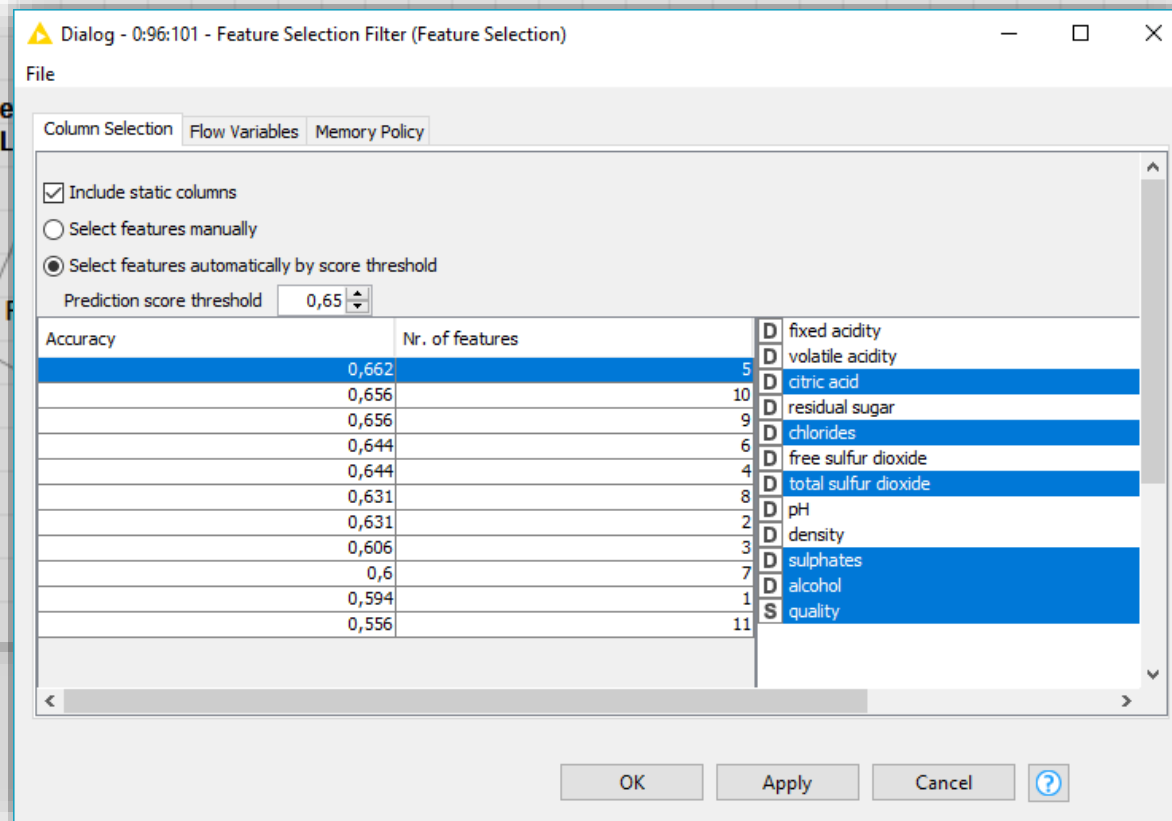
Feature Selection Filter

23

Model Validation

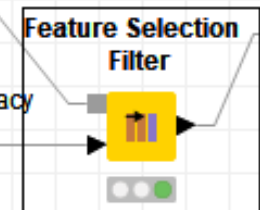
Feature Selection

Hands On



Feature Selection
Loop End

Maximize Accuracy



Feature Selection
after receiving the model

Hands On

24

Model Validation

Feature Selection

Hands On

HANDS ON

