



BEUTH HOCHSCHULE FÜR TECHNIK BERLIN  
University of Applied Sciences



# Getting Started with RapidMiner

Download tool and datasets:

<http://rapidminer.com/>

*Agathe Merceron*

*Beuth University of Applied Sciences*

*Berlin, Germany*





# Outline

Design and result perspectives in RapidMiner.

The Point Dataset Point1.csv:

Reading the data

Visual exploration to understand the data and guide the work

Clustering and the Loop operator: searching for the right k for k-means.





# Design and Result Perspectives

Organize your space to store processes, files and results.

Point1:

8 points defined by their x-, y-coordinates.





## Reading the Point1-Dataset

**Warning Read Wizard:** RapidMiner guesses the type of attributes looking at the first 100 by default.

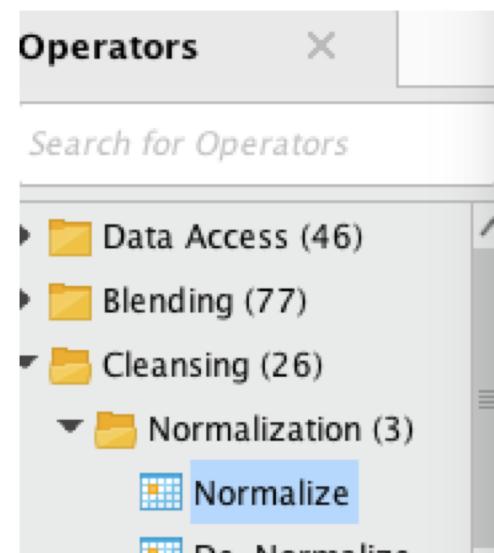
Metadata:

- Check the order of magnitudes of attributes
- Handy for classical transformations such as:

$$\frac{x - \text{average}}{\text{std deviation}}$$

$$\frac{x - \text{min}}{\text{range}}$$

- Transformation can be done inside RM:





# Exploring the Point1-Dataset

Plots:

- Histogram / color
- Box plots (quartile)
- Scatter
- and much more...

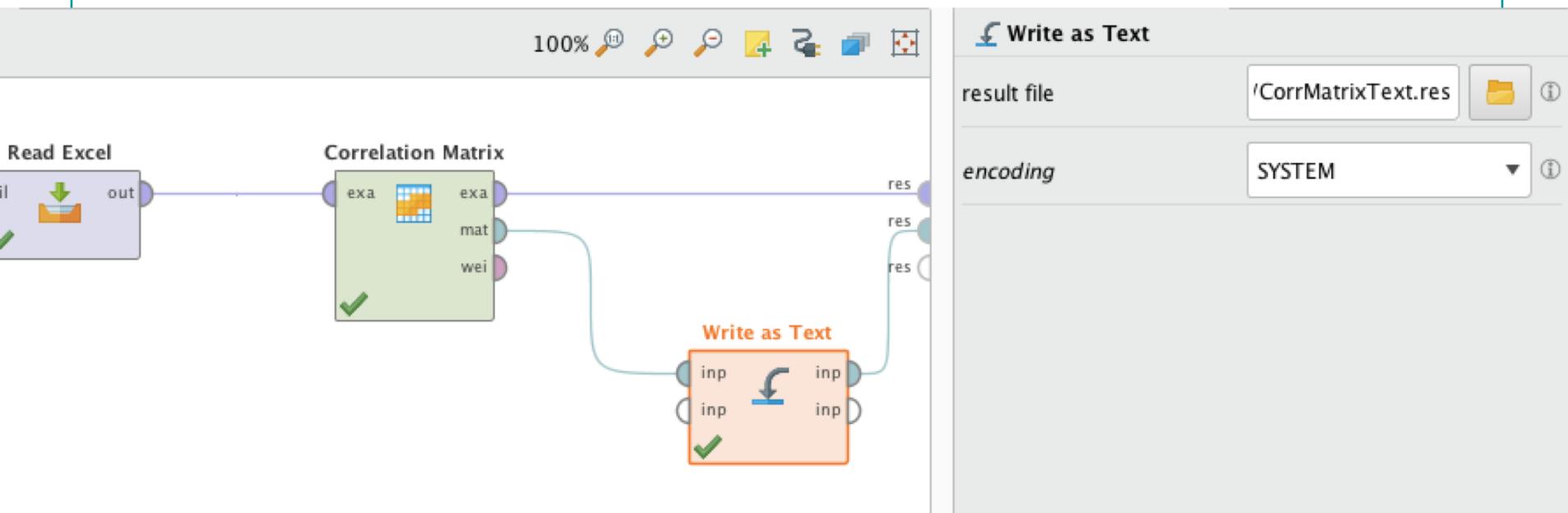
Correlation Operator:

- **Handy:** store the matrix with the Write as text Operator





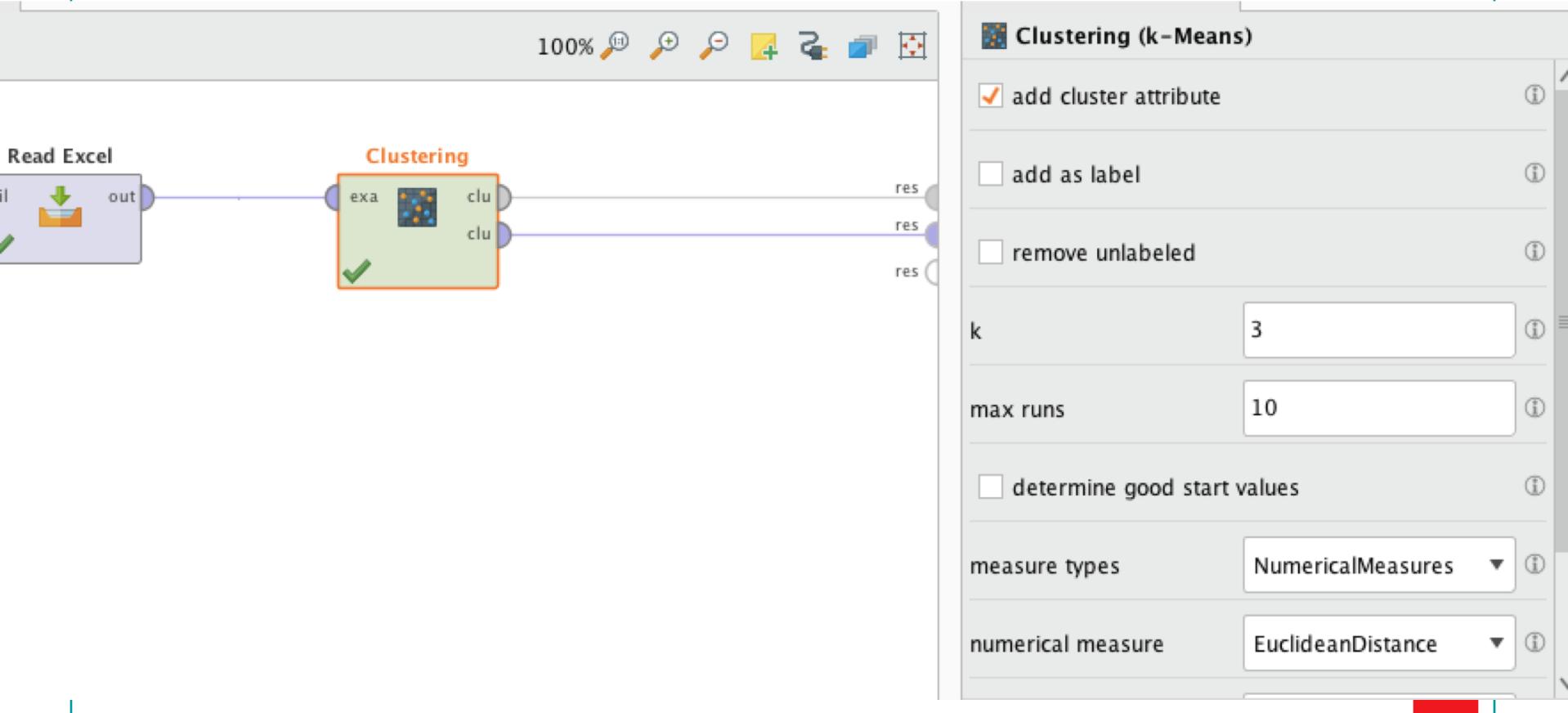
# Exploring the Point1-Dataset





# Looking for groups of points with Clustering

## K-means / Scatter Plot with Color Column cluster





# Looking for groups of points with Clustering

Loop Operator to discover the right k:

- **Warning:** skip the first port per inside the loop operator to get the results
- **Handy:** Log Operator to plot average within distance cluster or Davies Bouldin against k. Edit the parameters you want to plot, k and avg\_within\_distance for instance.





ceron)  
on)  
terLoop (aga  
ted)

Process X

Process >

Process

Read Excel

Loop Parameters

error handling

fail on error

synchronize

Select Parameters: configure operator

Select Parameters: **configure operator**  
Configure this operator by means of a Wizard.

Operators

- Clustering (k-Means)
- Performance (Cluster Distance Performance)
- Log (Log)
- Multiply (Multiply)

Parameters

- add\_cluster\_attribute
- add\_as\_label
- remove\_unlabeled
- max\_runs
- determine\_good\_start\_values
- measure\_types
- mixed\_measure
- nominal\_measure

Selected Parameters

- Clustering.k

advanced parameters

page compatibility (7.2.002)

Grid/Range

Min	Max	Steps	Scale
2	5	6	linear

Value List

- 2
- 3
- 4
- 5

Loop Parameters

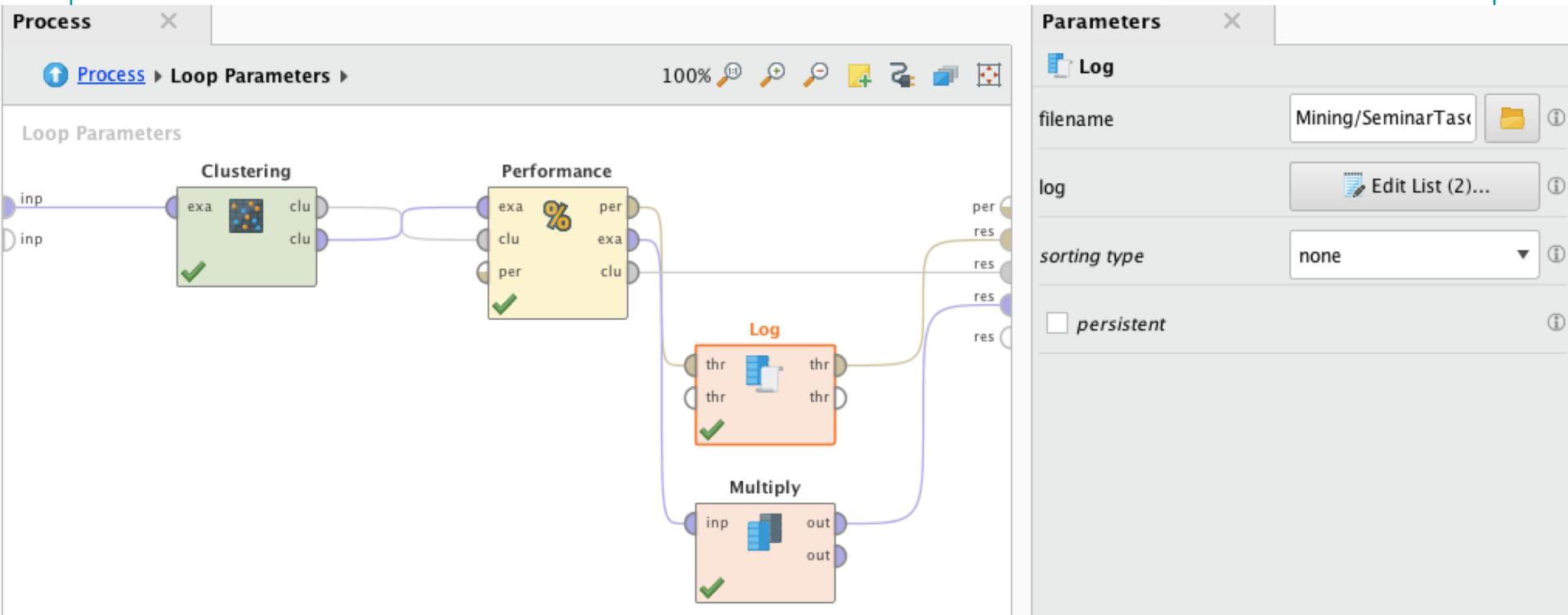
voidMiner Studio Core

Create, Settings, Grid, Search, Tune, Optimal, ...

This operator iterates over its subprocess for all the parameter combinations. The parameter settings can be set by the wizard provided in the configuration dialog.



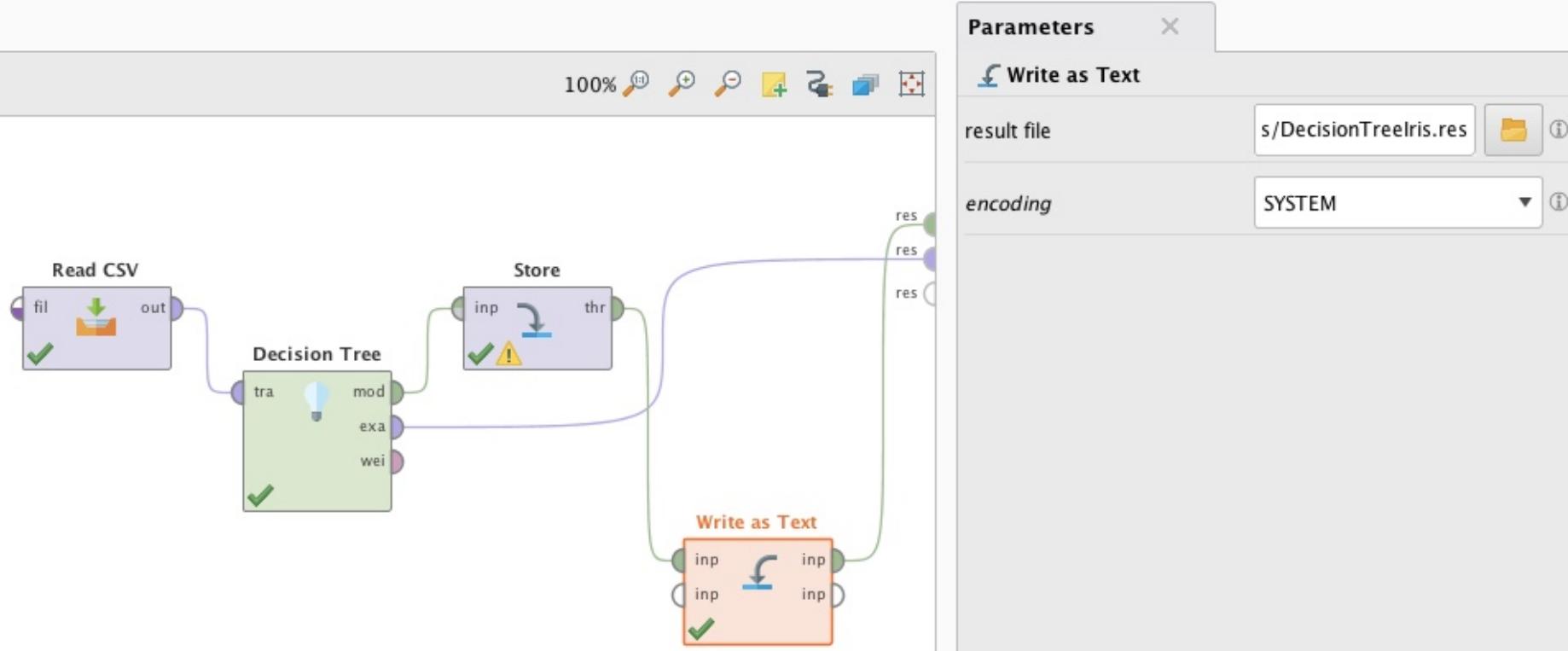
# Looking for groups of points with Clustering





# Classification with the iris dataset

Class should be of type Label





# Classification

Decision tree:

- **label** to set the attribute to predict

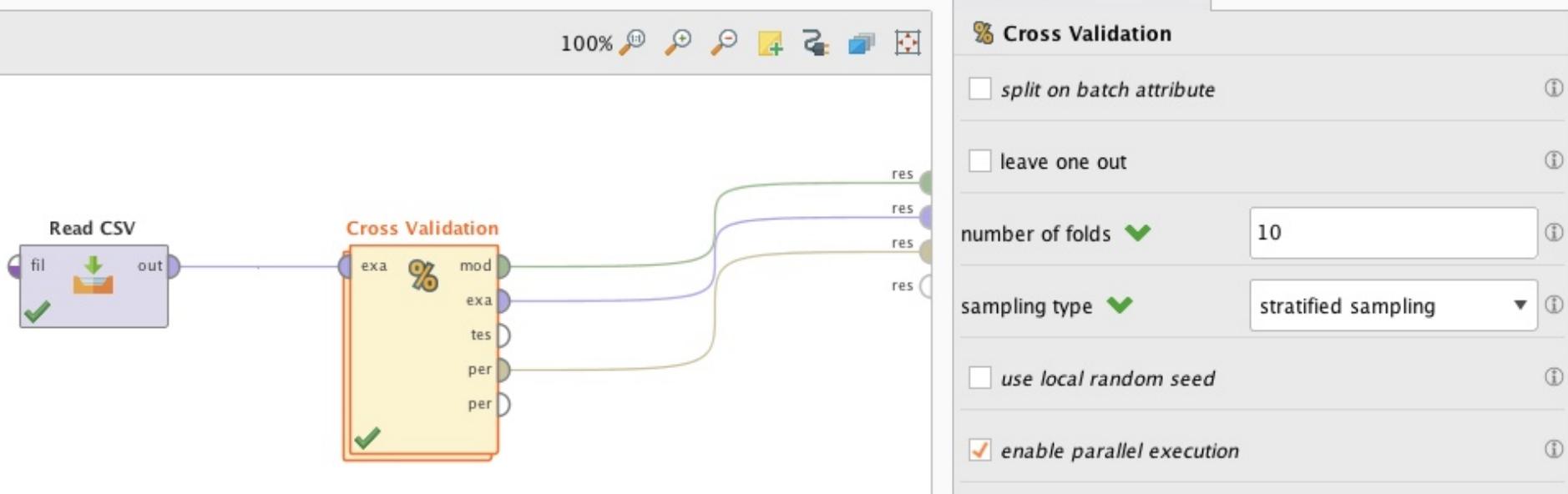
Validation:

- X Validation Operator
- **Handy:** Write as Text Operator the output of Performance Operator to check the result step by step.



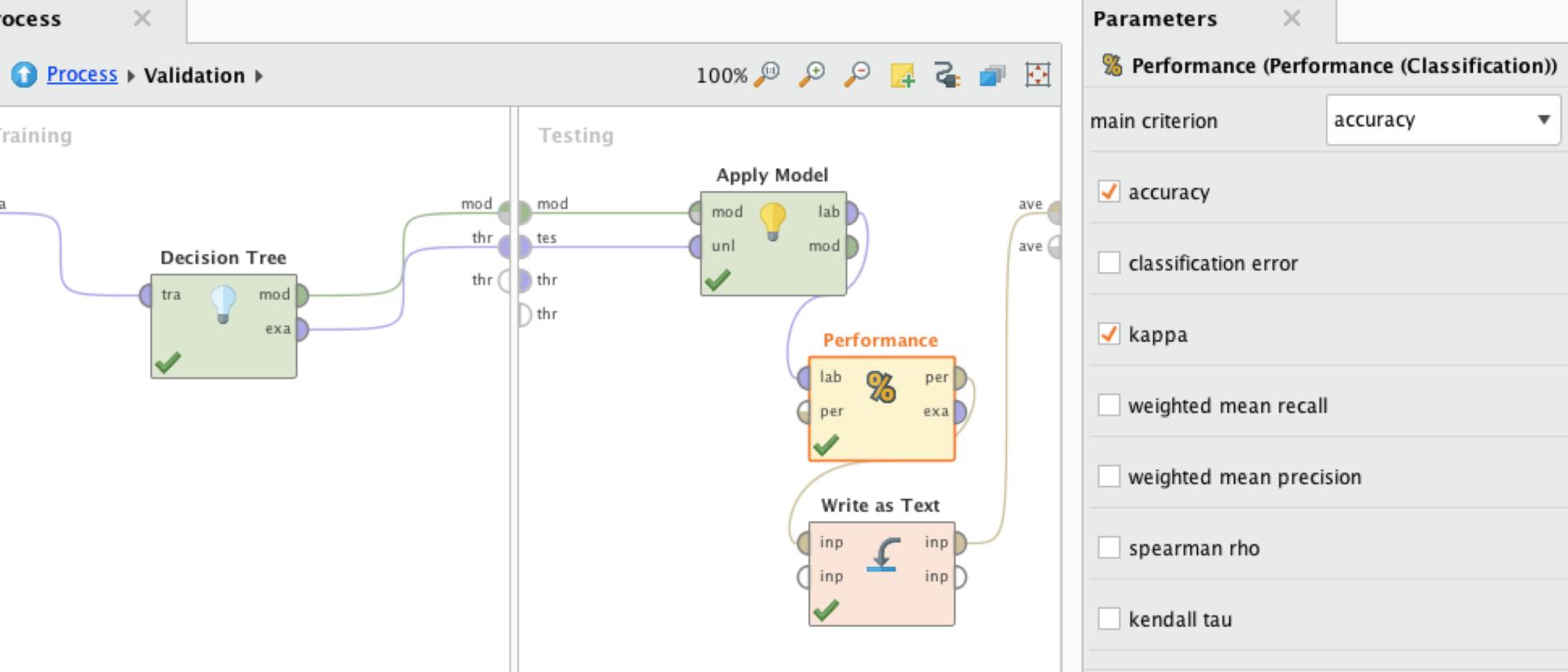


# Classification





# Classification





## References

Rapid Miner website: [rapidminer.com](http://rapidminer.com) with good tutorials.

- <https://community.rapidminer.com/>
- <https://docs.rapidminer.com/>

Many videos on youtube, see for example:

<https://www.youtube.com/watch?v=C8Ko3-2f-pA&list=PLssWC2d9JhOZLbQNZ80uOxLypglgWqbJA&index=16>





## References

Thank you for your attention!



Schloß Charlottenburg, Berlin

