# Data Mining
## (Mining Knowledge from Data)

## K-Nearest Neighbors

Magda Friedjungová, Marcel Jiřina

# Models

- Deductive or <u>inductive</u>

- Methods/tasks of learning inductive models:

  - clasification
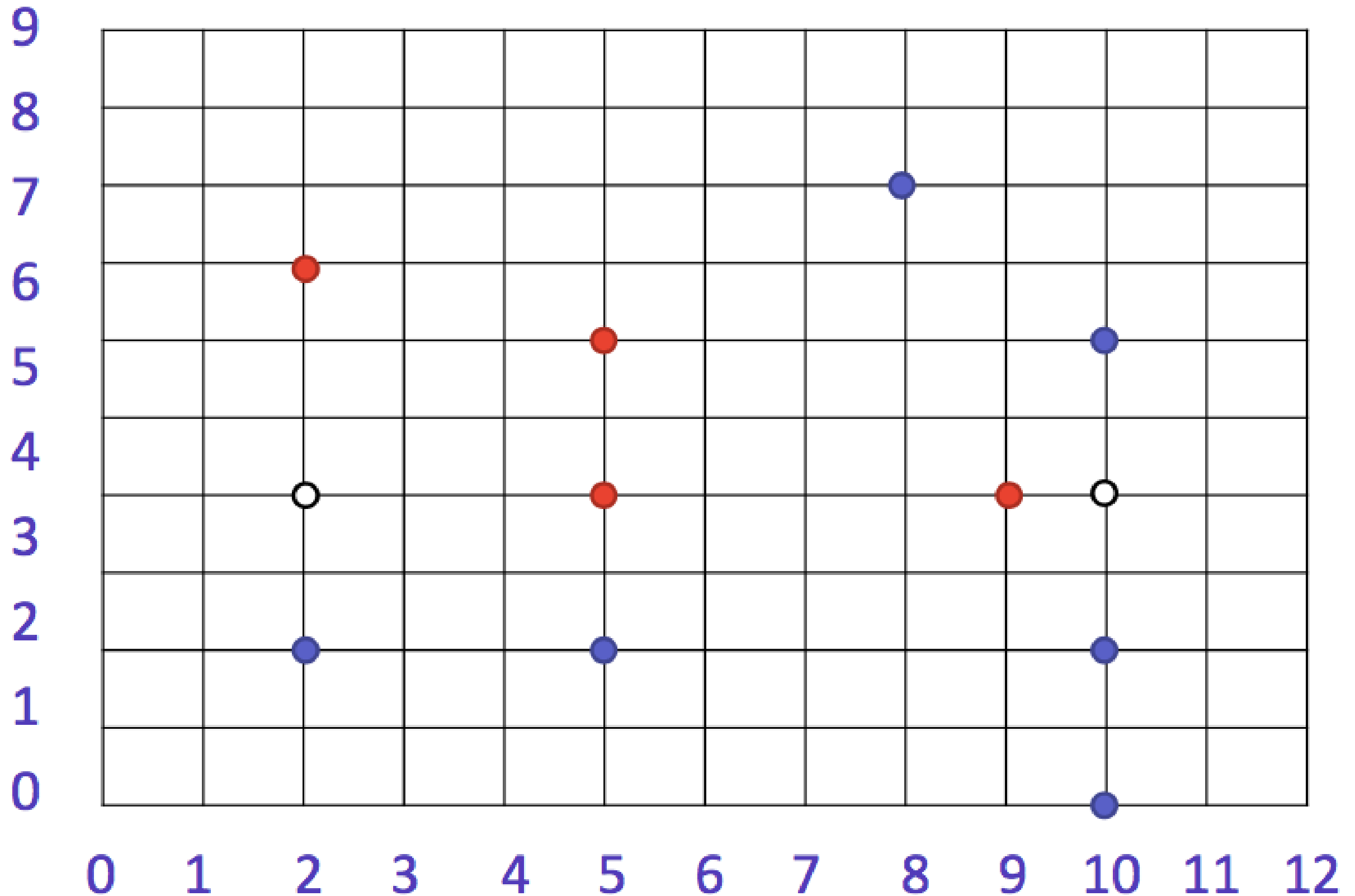
  - regression

  - prediction

  - clustering

# Overview of methods

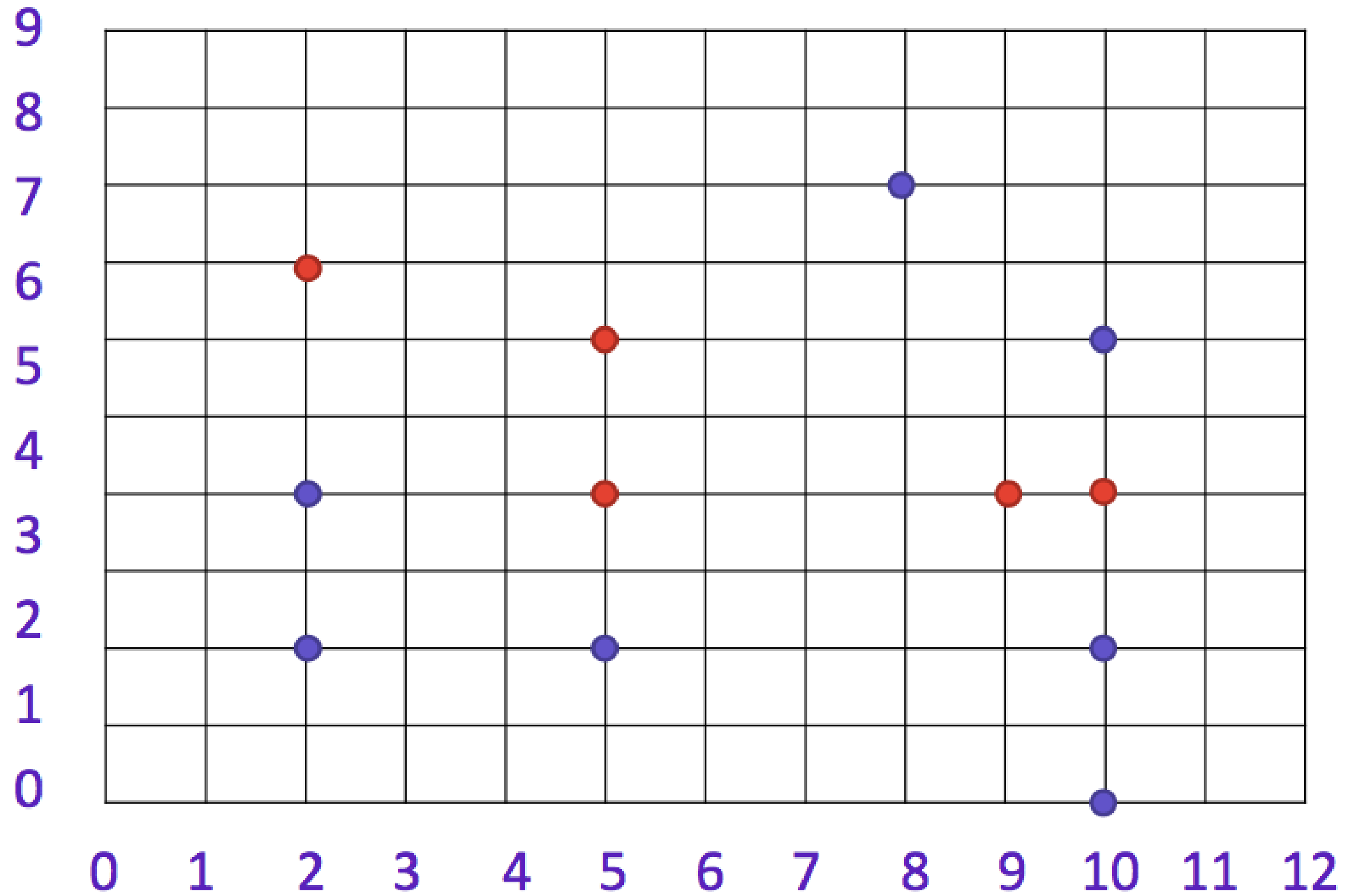| Task | Algorithms |
| --- | --- |
| **Classification** | K-Nearest Neighbors, Linear Separation, Decision Trees, Bayes Classifier, Neural Networks |
| **Regression, forecasting** | K-Nearest Neighbors, Linear Regression, Regression Trees, Neural Networks |
| **Clustering** | K-means, Hierarchical Clustering |
| **Rules, associations etc.** | Association Rules, K-means |

# Creating and using model

- Two stages:

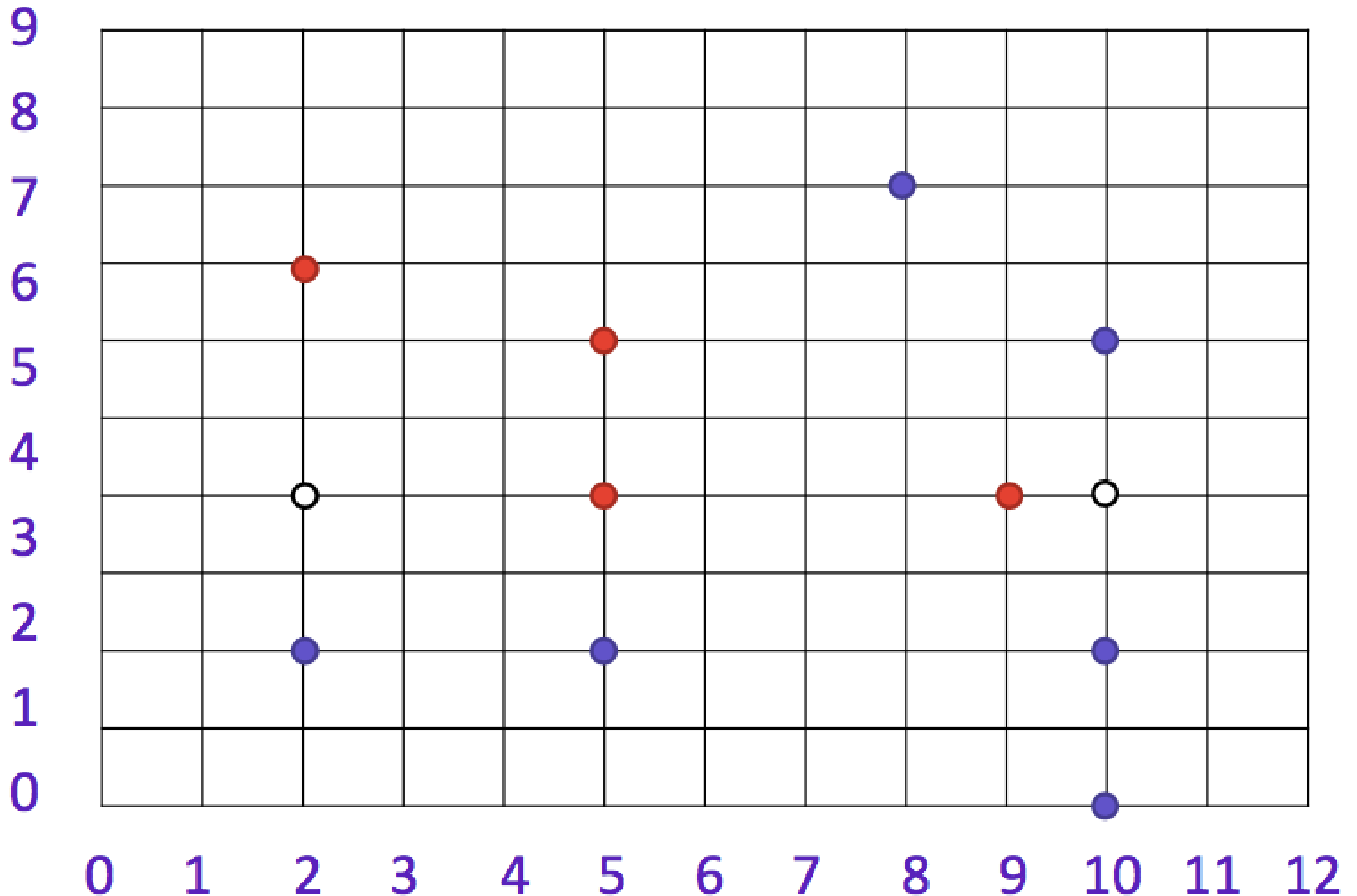    - Stage of learning, training

    - Stage of use, equipping

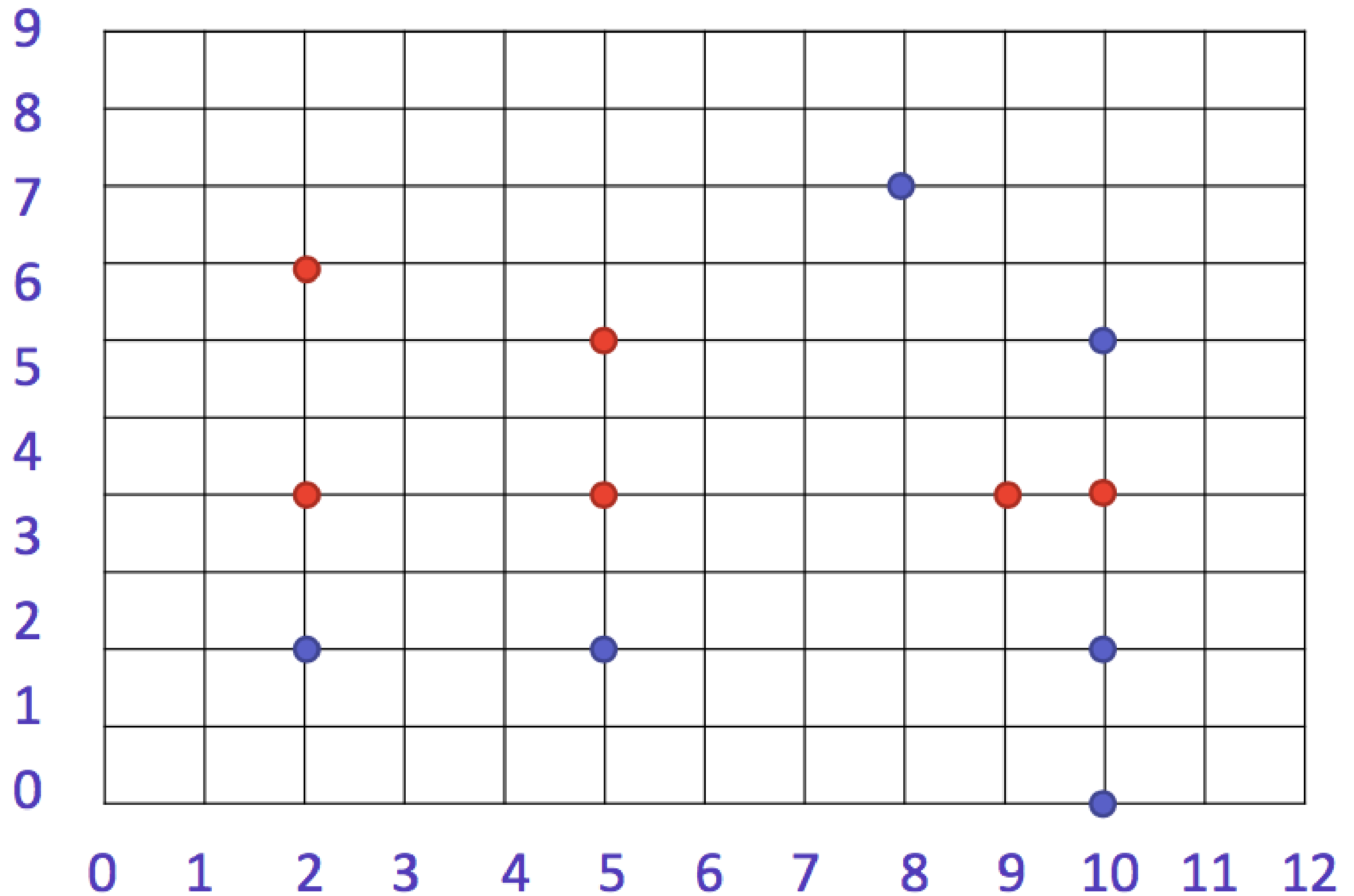Determine the class by the 1-NN method

# Result

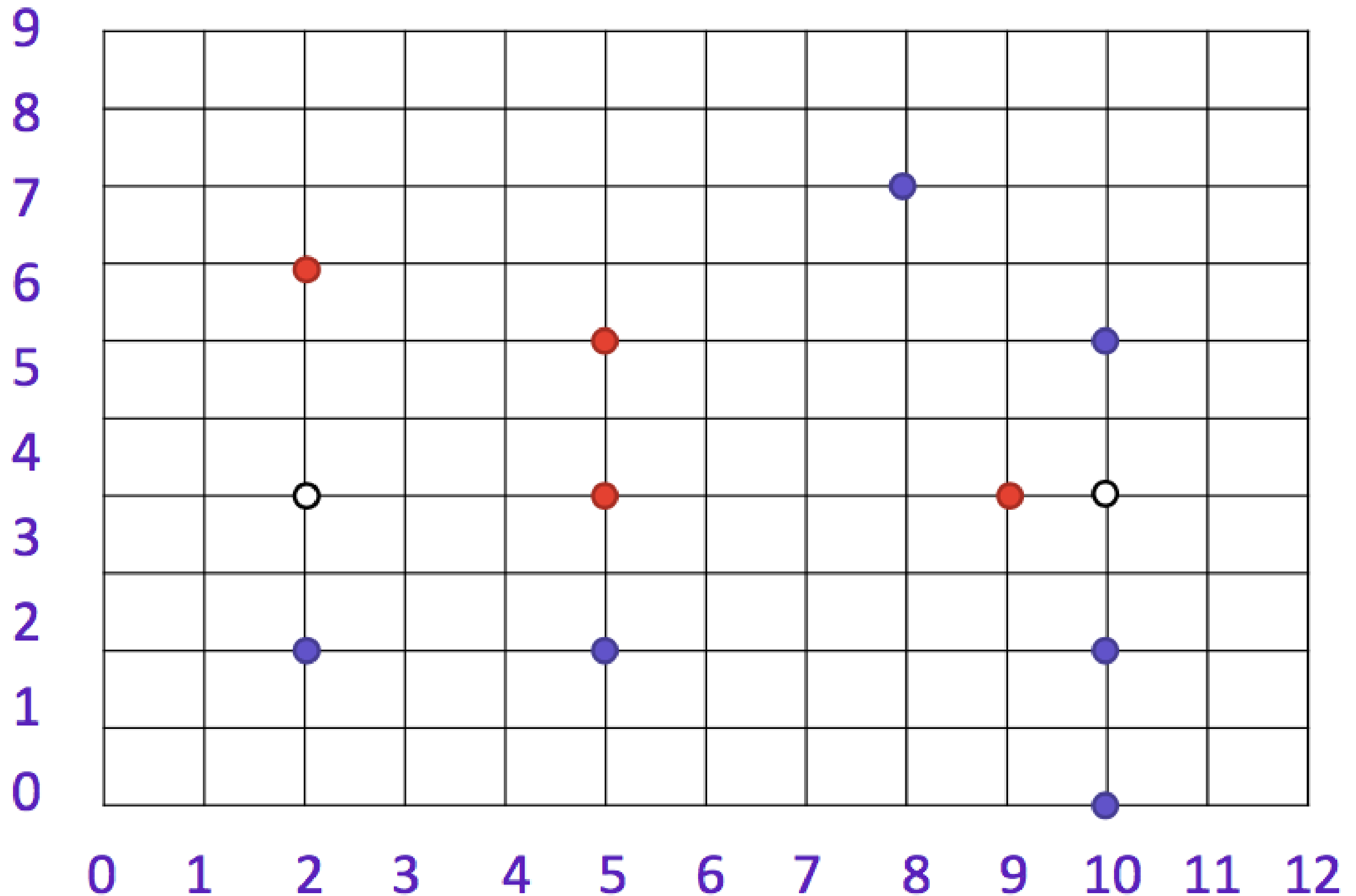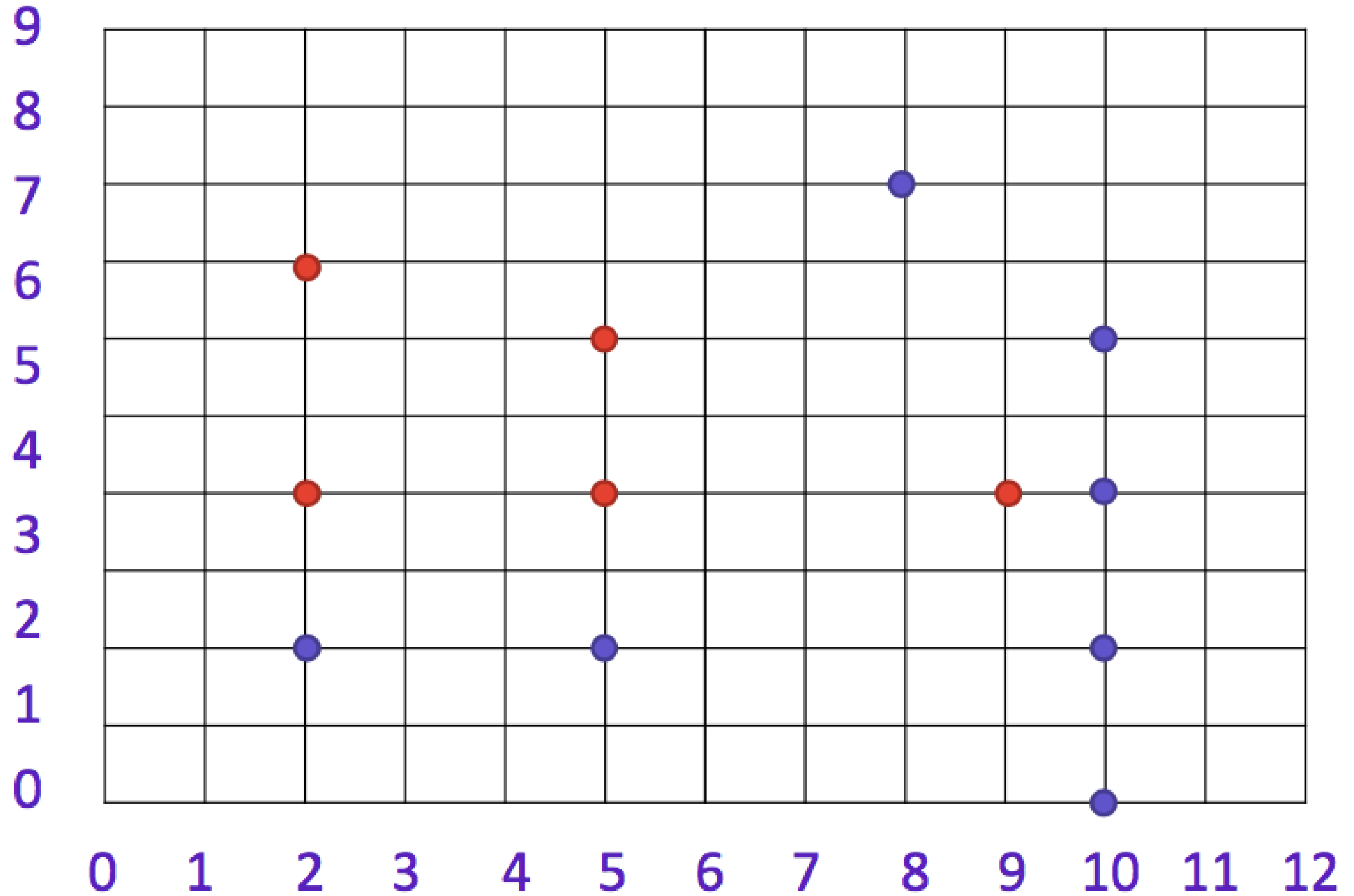# Determine the class by the 2-NN method

# Result

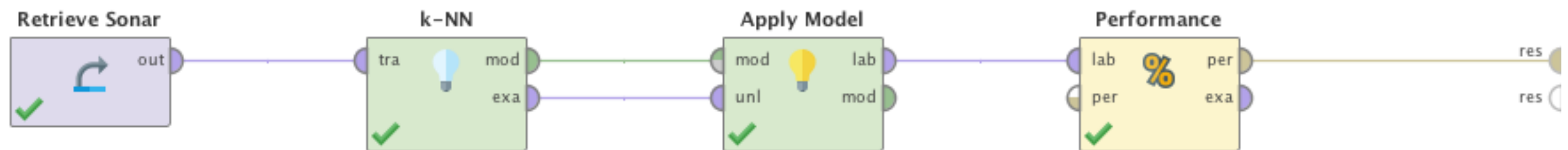# Determine the class by the 3-NN method

# Result

# K-NN on Sonar Data 1/3

- Load the Sonar dataset.

- Learn the k-NN classifier.

- Apply the model on the training data.

- Visualize the output and confidence of the model.

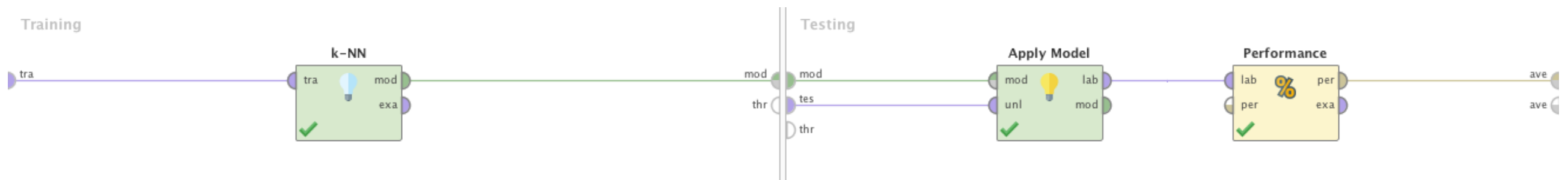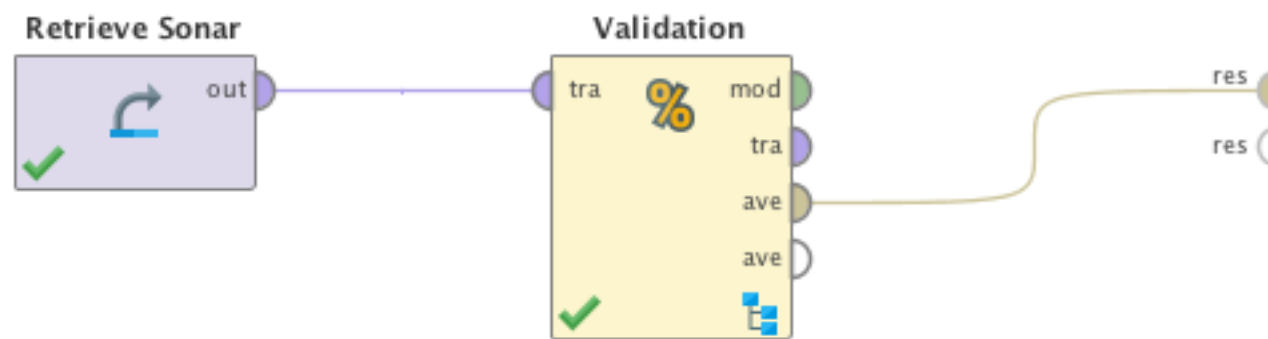- Calculate the performance on the training data.

# Sonar 1/3



- The accuracy of classification: 100%

accuracy: 100.00%

|  | true Rock | true Mine | class precision |
|---|---|---|---|
| pred. Rock | 97 | 0 | 100.00% |
| pred. Mine | 0 | 111 | 100.00% |
| class recall | 100.00% | 100.00% | |

# Sonar 2/3

- Calculate generalized performance (X-validation).

# Sonar 2/3

- The accuracy of x-validation: 82%

**accuracy: 82.14% +/- 8.95% (mikro: 82.21%)**

|  | true Rock | true Mine | class precision |
|---|---|---|---|
| pred. Rock | 75 | 15 | 83.33% |
| pred. Mine | 22 | 96 | 81.36% |
| class recall | 77.32% | 86.49% |  |

# Sonar 3/3

- Measure the classification accuracy for k={1, 2, 3, …, 20} and draw the plot.

- Explain the observation.

# Sonar 3/3

# Sonar 3/3

# Lazy Modeling

- What is lazy modeling?

  - Processing of training data take place just after receiving tha sample to be classified.

  - Response to the request by a combination of queries on training data.

  - Discard the calculated response and all intermediate results.

    What does this mean in practice?

# Measure "execution time" of "k-NN" and "apply"

- Start from the previous schema:

  - Remove the x-validation.

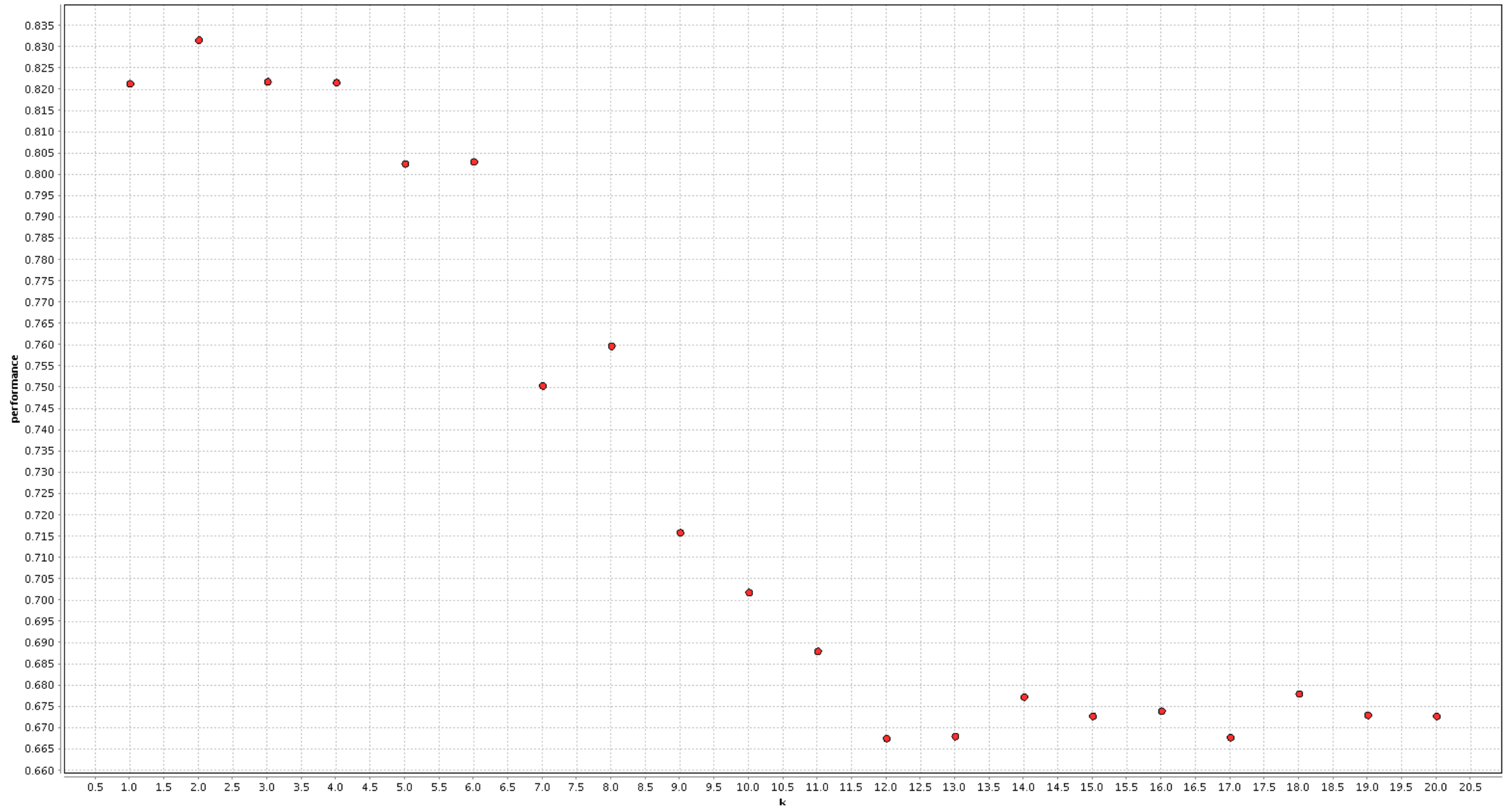  - Instead of the Sonar data use the Generate data. Target function set to "checkerboard classification", number of attributes to 2.

  - In the Loop set "Generate Data.number_examples" on 1…1000 and 10 steps.

# Measure "execution time" of "k-NN" and "apply"



**Loop Parameters**

inp — res — res

**Select Parameters: configure operator**

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

**Operators**
- Generate Data (Generate Data)
- k-NN (k-NN)
- Apply Model (Apply Model)
- Log (Log)

**Parameters**
- target_function
- number_of_attributes
- attributes_lower_bound
- attributes_upper_bound
- gaussian_standard_deviation
- largest_radius
- use_local_random_seed
- local_random_seed

**Selected Parameters**
- Generate Data.number_examples

**Grid/Range**

| Min | Max | Steps | Scale |
|-----|-----|-------|-------|
| 1.0 | 1000.0 | 10 | linear |

**Generate Data** — out

**k-NN** — tra, mod, exa

**Apply Model** — mod, lab, unl, mod

**Log** — thr, thr, thr, thr

...lue specifies the process value to log.
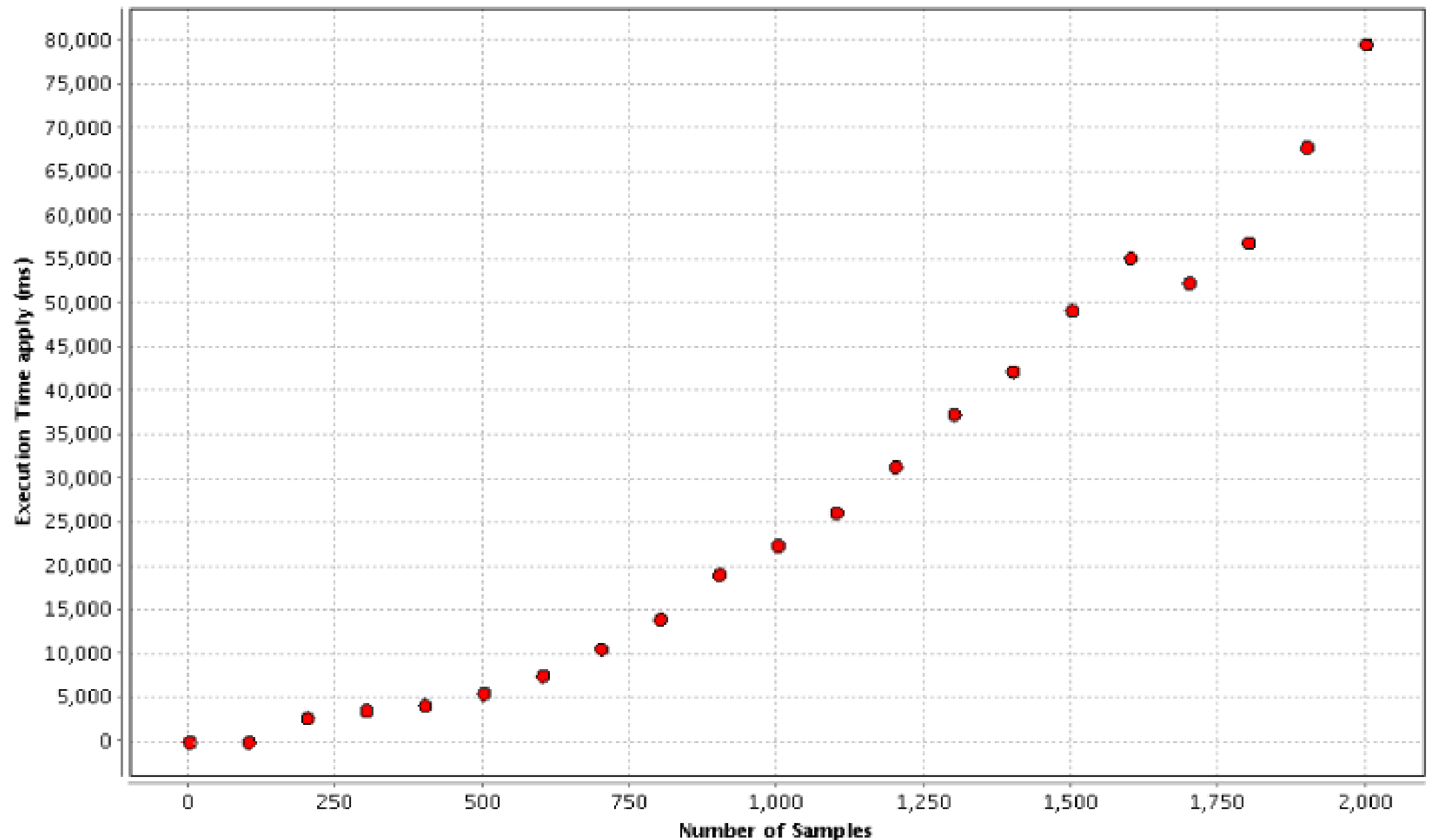
| column name | value | | |
|-------------|-------|---|---|
| number of samples | Generate Data | parameter | number_exa... |
| execution time k-nn (ms) | k-NN | value | execution-time |
| execution time apply (ms) | Apply Model | value | execution-time |

# Measure "execution time" of "k-NN" and "apply"



The execution time is for k-NN approximately constant.

# Measure "execution time" of "k-NN" and "apply"



The execution time for Apply is approximately quadratic.

# Normalization

- k-NN method uses the Euclidean metric:

$$\sqrt{\sum_{i=1}^{k}(x_i - y_i)^2}$$

- What if all dimensions ranges of values are 0 and 1, but only in one dimension the range is between 0 and 1000?
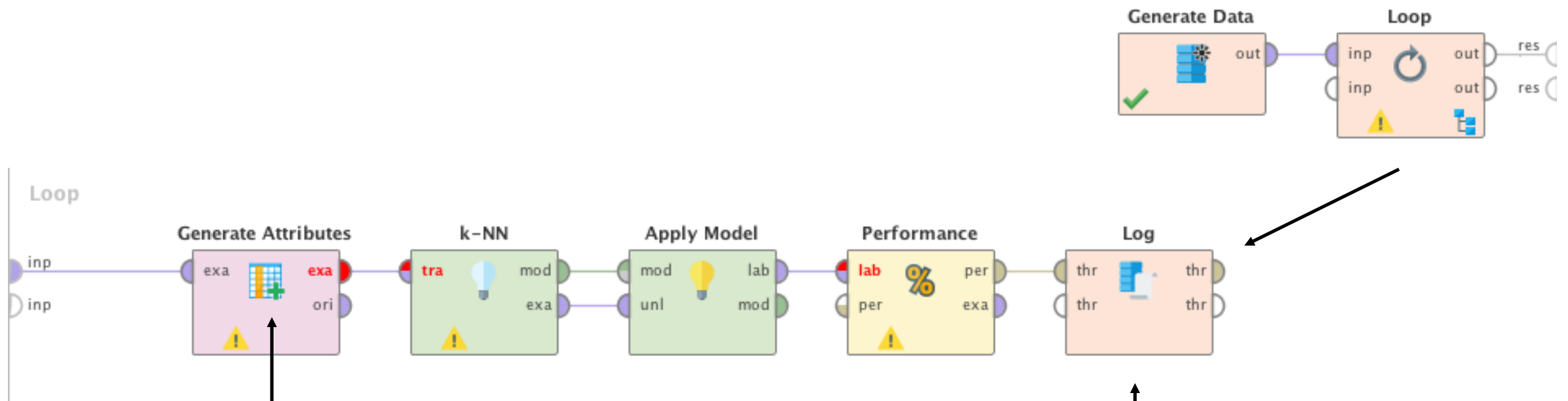
- How k-NN will behave?

# Normalization

- Use the schema from the previous task. As not all parameters can be changed in the "Loop Parameter" use "Loop" instead.

  - The size of a dimension can be changed by the block "Generate Attributes". Set it to:

    "att1*pow(1.5,eval(%{iteration}))"

  - Set the logging properly.

  - Place the data generator outside the Loop (we want to compare the same data, but modified other ways).

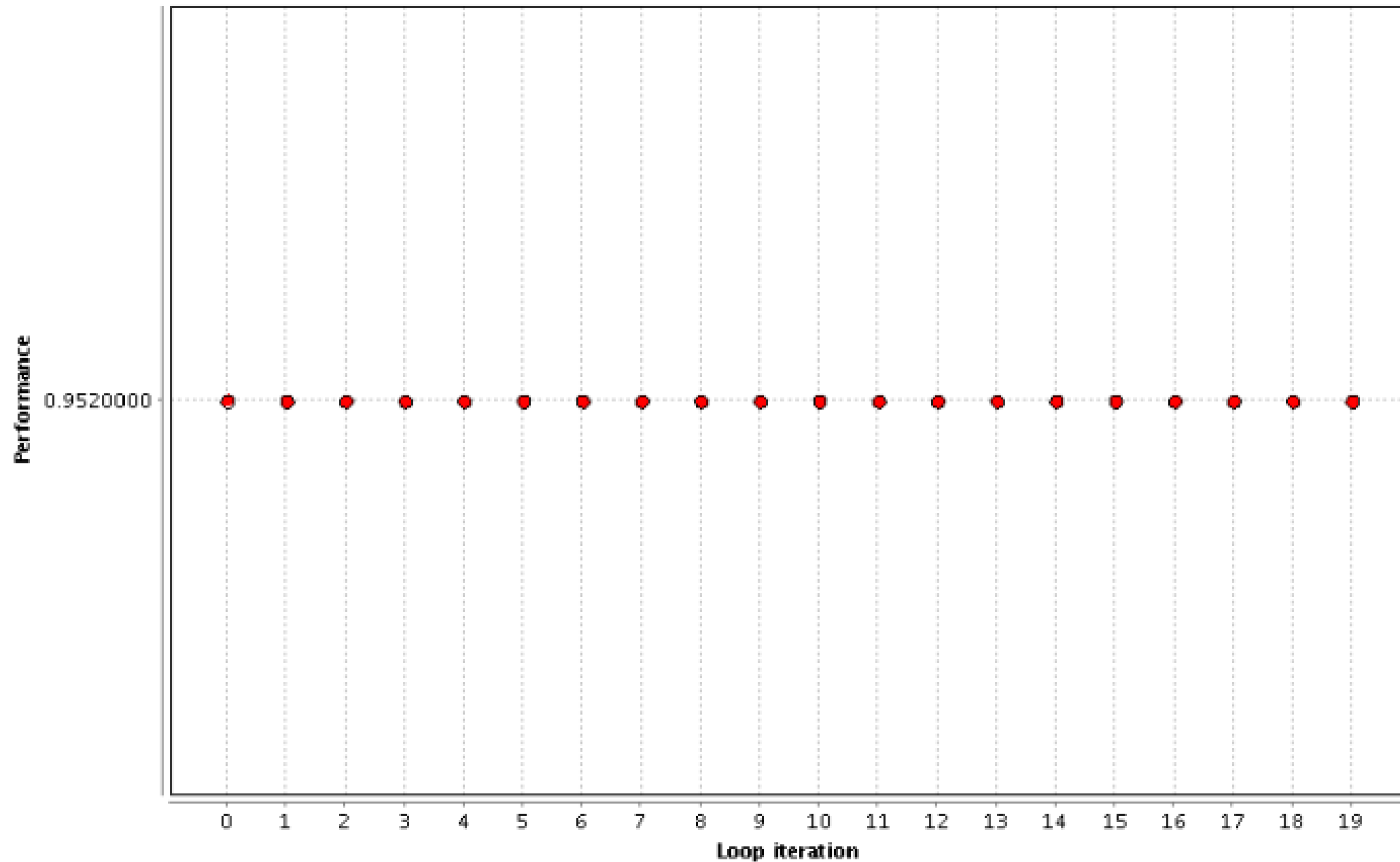# Normalization

# Normalization



- k-NN favors parameters of a large scale.

# Normalization

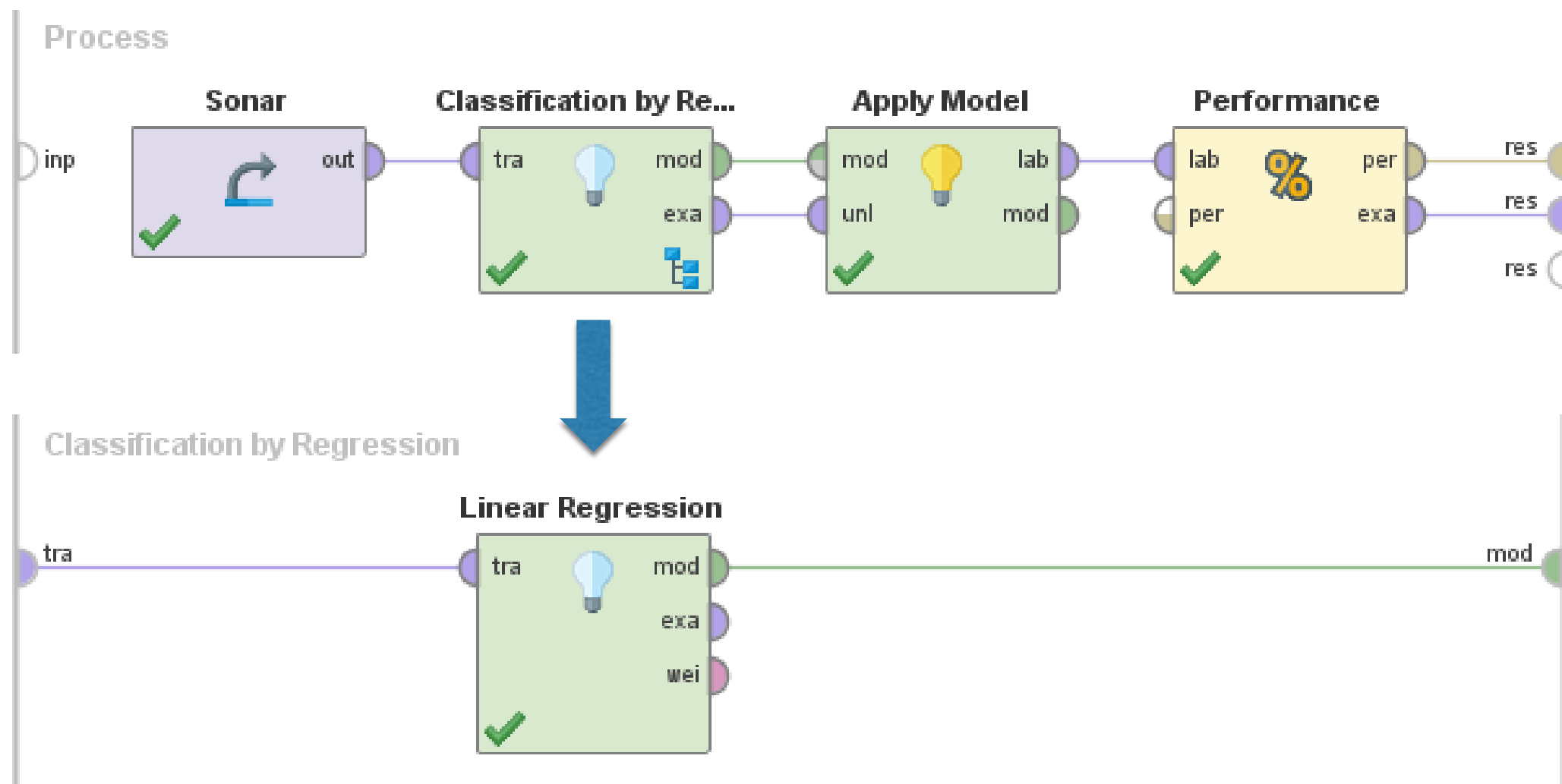- Repeat the experiment, but insert the "Normalize" before the "k-NN".

- What happens?

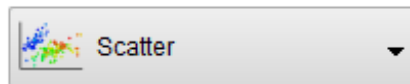# Normalization



- The accuracy is now 92,5% all the time.

# Linear Separation

- Operator "Classification by Regression" -> Insert „Linear Regression"

# Linear Separation