# Introduction to Machine Learning with Orange

# 1. Getting Started with Orange

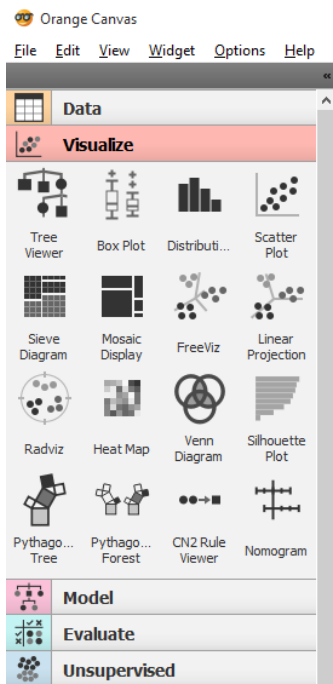Orange ([https://orangedatamining.com](https://orangedatamining.com)) is a machine learning and data visualization tool with interactive data analysis workflows which provides a large variety of tools for different types of data analytics such as data modelling, prediction, clustering, classification, neural networks, etc.
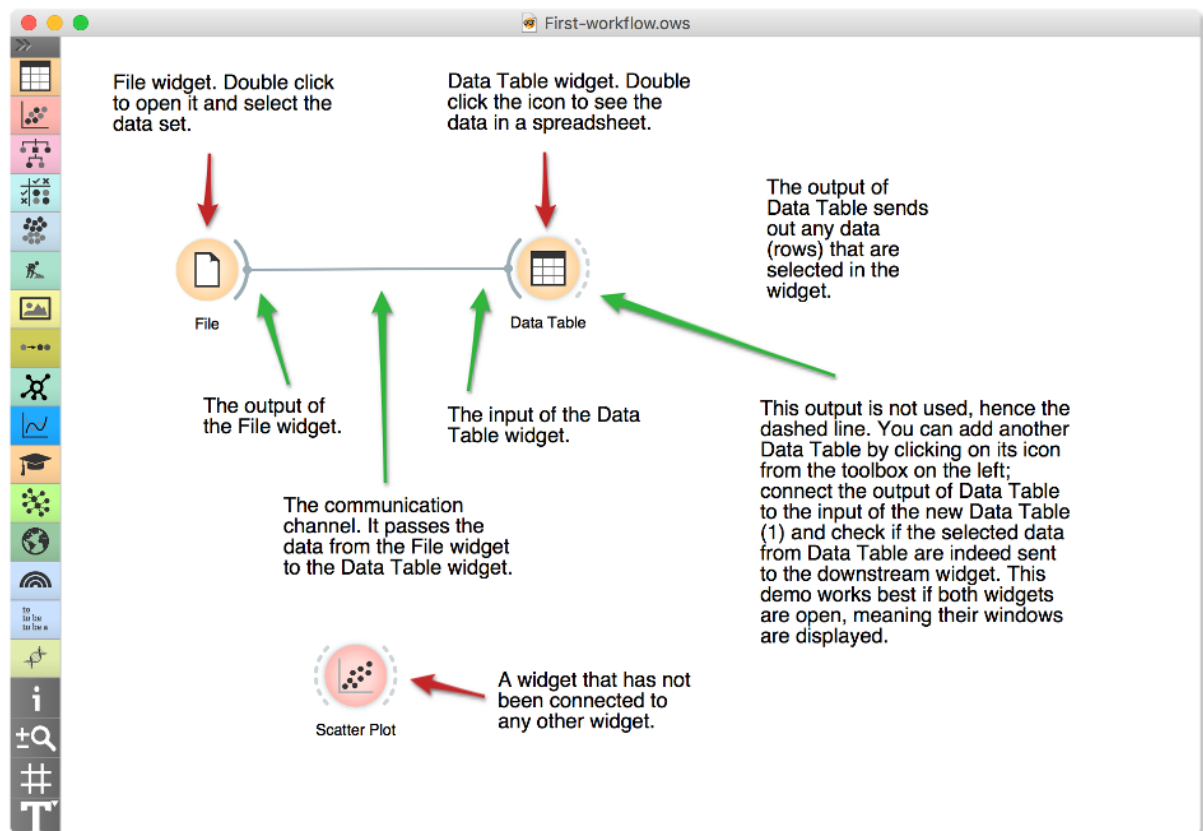
After starting up Orange from the desktop/start menu, you will see the welcome screen when you open Orange the first time. Click **New** to Start a new a recent data analysis workflow.



Orange workflows consist of components that read, process, and visualize data. We refer to these components as **Widgets.** The widgets are organized under different categories/tabs.



We can place the widgets on a drawing board called the "canvas" to design a workflow. Widgets communicate by sending information along their communication channel. Output from one widget can be used as input to another.
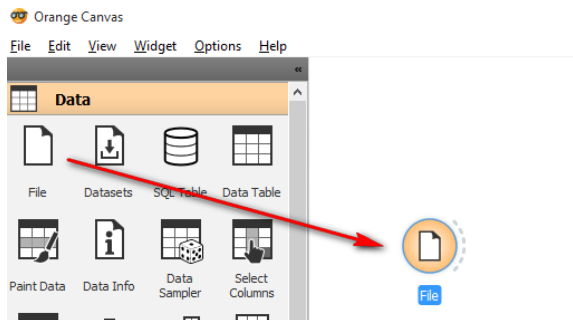
First-workflow.ows

File widget. Double click to open it and select the data set.

Data Table widget. Double click the icon to see the data in a spreadsheet.

The output of Data Table sends out any data (rows) that are selected in the widget.

File

Data Table

The output of the File widget.

The input of the Data Table widget.

The communication channel. It passes the data from the File widget to the Data Table widget.

This output is not used, hence the dashed line. You can add another Data Table by clicking on its icon from the toolbox on the left; connect the output of Data Table to the input of the new Data Table (1) and check if the selected data from Data Table are indeed sent to the downstream widget. This demo works best if both widgets are open, meaning their windows are displayed.

Scatter Plot

A widget that has not been connected to any other widget.

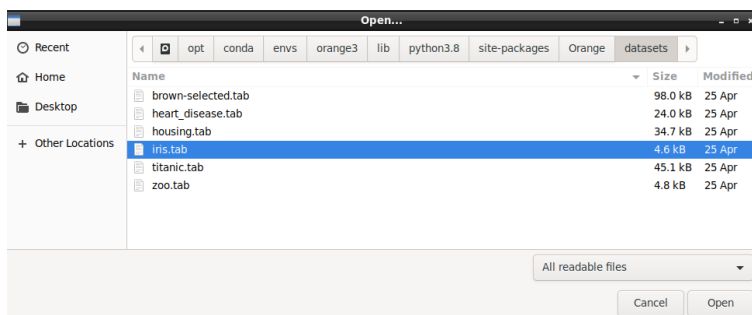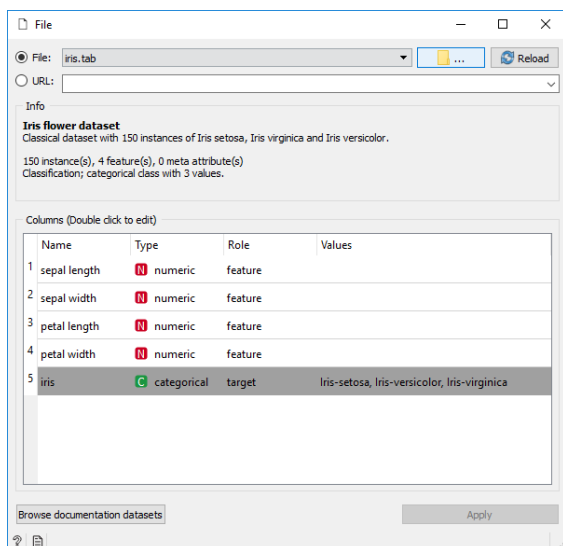# 2. Data Exploration and Visualization

## Preparation

In this step, we will set up a data analysis workflow. The famous Iris data set[1] is used for demonstration purpose.

Drag the **File** widget on **Data** tab to the canvas.



Double click the *File* widget in the canvas. Click **…** to browse the sample dataset and select the iris dataset (**iris.tab**).





---

[1] The Iris data set, introduced in 1936 by Ronald Fisher, a British statistician and biologist, consists of 50 samples from each of three species of Iris (Iris setosa, Iris virginica and Iris versicolor) with 4 measured features from each sample: the length and the width of the sepals and petals, in centimeters.

About the column (variable) types:

- **Categorical**: This refers to a variable with discrete values, for example eye color, nationality, ticket class (first, second, third), and so on. In Orange, these would be marked with a green C.
- **Numeric** or **continuous** variable: This refers to a variable with numbers as values, for example cholesterol level, heart rate, age, and so on. In Orange, these would be marked with a red N.

About the column (variable) roles:

- **Feature**: In statistics they would be called **independent variables**. This all refers to descriptions of data samples. If you have patients in rows, that variables are in columns and they describe these patients (i.e. with name, date of birth, cholesterol level, heart rate, and so on).
- **Target**: Define the target variable or **dependent variable**. This is the variable you are trying to predict, for example the survival of Titanic passengers or the price of a house in Boston.

Close the widget after the file is selected.

Click on *File widget* in the data widgets menu. Add a **Data Table** widget to the canvas. Connect the two widgets by dragging a line from the *output channel* of file to the *input channel* of Data Table.
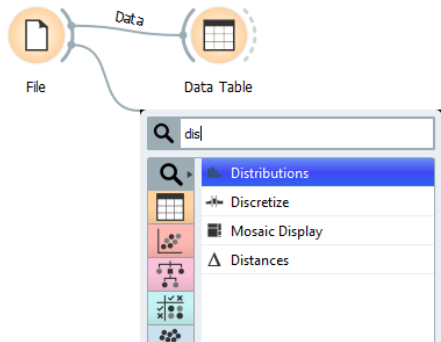


Fig. 1: File to Data Table Canvas Configuration

Double-click the **Data table** widget in the canvas to show the Iris data set. The dataset has 150 instances and 4 features.
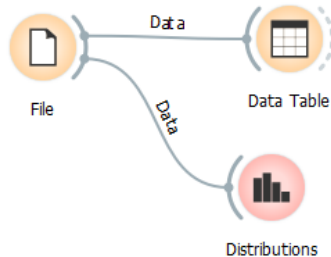
**Data Table**  — □ ✕

### Info

150 instances (no missing values)
4 features (no missing values)
Discrete class with 3 values (no missing values)
No meta attributes

### Variables

☑ Show variable labels (if present)
☐ Visualize numeric values
☑ Color by instance classes

### Selection

☑ Select full rows

| | iris | sepal length | sepal width | petal length | petal width |
|---|---|---|---|---|---|
| 1 | Iris-setosa | 5.1 | 3.5 | 1.4 | 0.2 |
| 2 | Iris-setosa | 4.9 | 3.0 | 1.4 | 0.2 |
| 3 | Iris-setosa | 4.7 | 3.2 | 1.3 | 0.2 |
| 4 | Iris-setosa | 4.6 | 3.1 | 1.5 | 0.2 |
| 5 | Iris-setosa | 5.0 | 3.6 | 1.4 | 0.2 |
| 6 | Iris-setosa | 5.4 | 3.9 | 1.7 | 0.4 |
| 7 | Iris-setosa | 4.6 | 3.4 | 1.4 | 0.3 |
| 8 | Iris-setosa | 5.0 | 3.4 | 1.5 | 0.2 |
| 9 | Iris-setosa | 4.4 | 2.9 | 1.4 | 0.2 |
| 10 | Iris-setosa | 4.9 | 3.1 | 1.5 | 0.1 |
| 11 | Iris-setosa | 5.4 | 3.7 | 1.5 | 0.2 |
| 12 | Iris-setosa | 4.8 | 3.4 | 1.6 | 0.2 |
| 13 | Iris-setosa | 4.8 | 3.0 | 1.4 | 0.1 |

## Distribution

Create a **Distributions** widget and link to the **File** widget as follows. From **File** widget, drag from the right side of the widget and type the name of the widget and select the appropriate widget.

In orange, you may extend a workflow by dragging a line from the output channel of a widget and select the widget that you want to add. You may type the name of the widget for filtering.



Double click the **Distribution** widget.



Set **Split by** to **iris** and **Fitted distribution** to **Normal**.

- Explore the distribution of the different attributes (e.g. sepal length) across the three different species of Iris.
- Change the bin width and smoothing parameters and observe the output.

## Box Plot

Box plots shows the distribution of attribute values and is useful for identifying outliers. Modify the workflow as follows.



Double click **Box Plot**. You may select an attribute (e.g. sepal length), and analyse the distribution of this attribute across different subgroups (in this example, iris species).



For continuous attributes, the box plots shows

- The mean (blue vertical line)
- The blue highlighted area is the entire standard deviation of the mean.
- The median (yellow vertical line)
- The thin blue line represents the area between the first (25%) and the third (75%) quantile
- The thin dotted line represents the entire range of values (from the lowest to the highest value in the data set for the selected parameter).

## Scatter Plot

Under **Visualize** tab in the widget menu, insert a **Scatter plot**.



Set up the following workflow.



Double click the scatter plot widget in the canvas. Notice that the three species of iris are not clearly separated from each others.

Click *Find Information Projections* in the top left corner. Click **Start**.

The rank of the result indicates how good the attribute pair can separate the species in the scatter plot. The pair "Petal length and petal width" fits the best in this case.



Double click the first entry and close the popup window.

From **File** menu, click **Save**.



Save the model as "**iris visualization.ows"** on your desktop.

# 3. Decision Tree Classifier

Decision tree is one of the oldest, but still popular, machine learning methods. We will construct a decision tree to predict the species of iris.

Select **File→New** to start a new canvas.
Let us load *iris* data set, build a tree (widget *Tree*) and visualize it in a *Tree Viewer*.





Double click and open the **Tree viewer**.

- We can read the tree from top to bottom. In this example, the column *petal length* best separates the iris variety *setosa* from the others, and in the next step, *petal width* then almost perfectly separates the remaining two varieties.
- Trees place the most useful feature at the root based on how well they distinguish between classes. It then splits both subsets further, again by their most useful features, and keeps doing so until it reaches subsets in which all data belongs to the same class (leaf nodes in strong blue or red) or until it runs out of data instances to split or out of useful features (the two leaf nodes in white).

Create the workflow as shown in the diagram below.

- Select the *Iris* data set in the **File** widget. In the *Scatter Plot*, we first find the best visualization of this data set, that is, the one that best separates the instances from different classes.
- Then we connect the *Tree Viewer* to the *Scatter Plot*. Data instances (particular irises) from the selected node in the *Tree Viewer* are shown in the *Scatter Plot*.





Click the **Tree** widget and review the hyper-parameters of the decision tree. Change the **"Min. number of instances in leaves"** to 1 and observe the change in the decision tree.
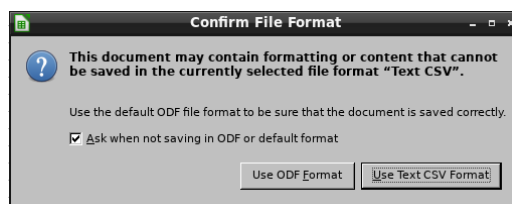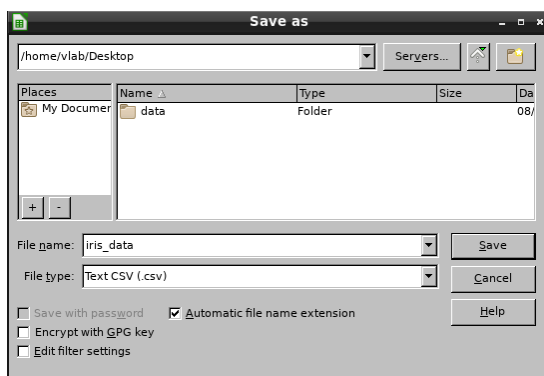
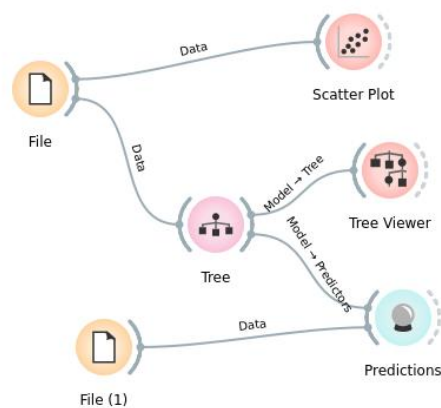Start **LibreOffice Calc** from Start Menu (If you are using MS Windows, you use may use **MS Excel**).



Create a spreadsheet as follows. Save the file as **data.csv** on your desktop. Click **Use Text CSV Format** (if you are using the OpenOffice calc).
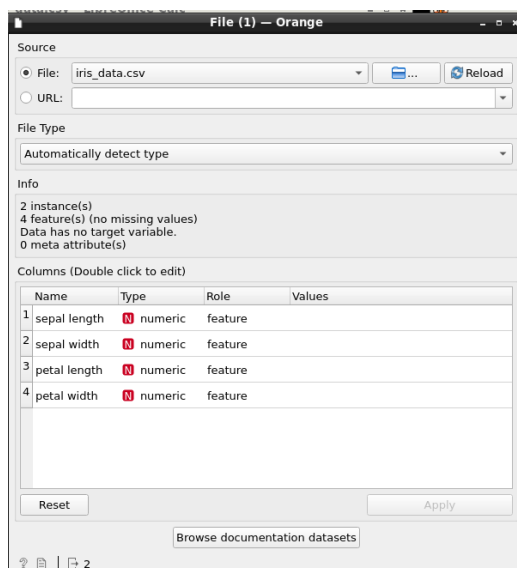
Create the workflow as shown in the figure below.



Add a **Predictions** widget and another **File** widget to the canvas and connect the workflow as follows. There should be two **File** Widgets in the workflow.

- The **File** widget one on the top is used to train a decision tree classifier.
- The **File** widget (File 1) at the bottom is used to provide the data for prediction. Double click it and select **iris_data.csv** that you have created.



Click the **prediction** widget on the canvas to view the predicted class for the various animals in the given dataset.
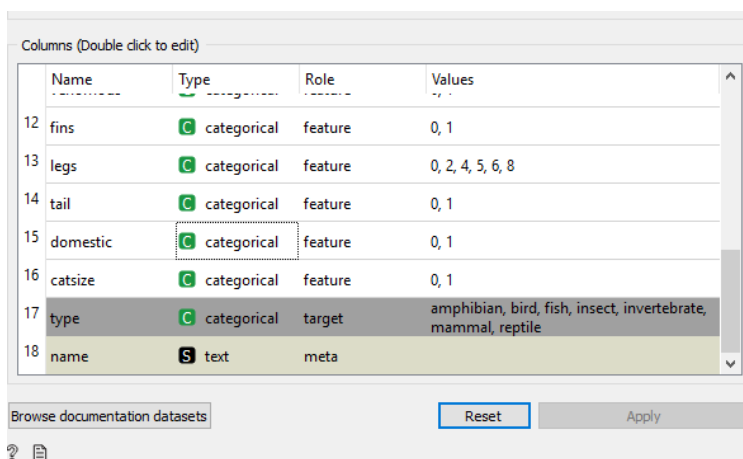
## 4. The Zoo Dataset

In this section, we will explore the zoo dataset and develop a decision tree classification model to predict the type of animal (e.g. Mammal, Bird, Reptile, Fish, Amphibian, Insect and Invertebrate) given various features of the animals.

Start a new workflow (**File -> New**). Add a **File** widget.

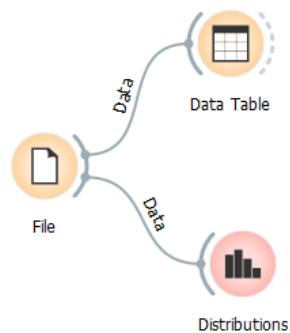Click **…** to browse the sample dataset. Choose the ***zoo data set (zoo.tab)*** as the data source.



Scroll down the column windows to view the "type" column which is a "target" variable. This indicates that we are going to build a machine-learning model to predict the type of animal given the various "feature" columns.

Remark:

- **string** or **text**: variables with text as value, for example the address, the ID of a product, or the name of the person.
- **meta variable**: variables that adds additional information on a data sample, but it is not used in computation (e.g. classification and clustering)

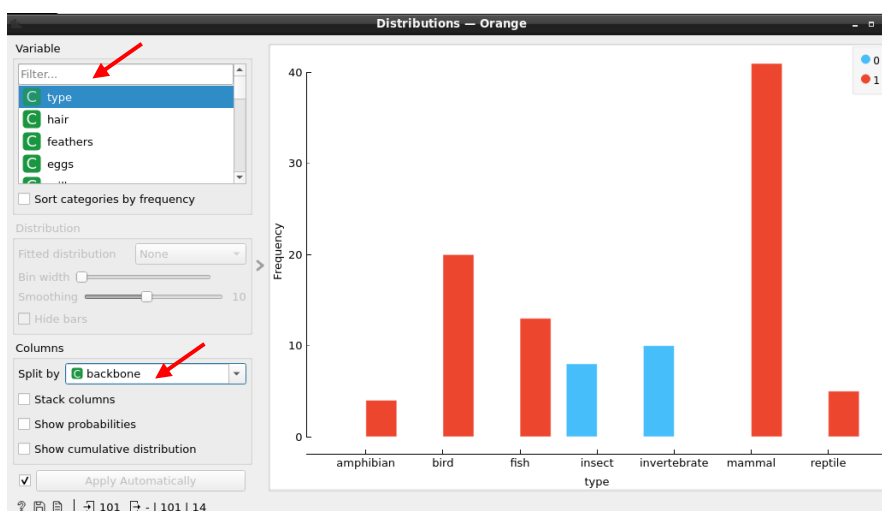Create the following workflow in the canvas.



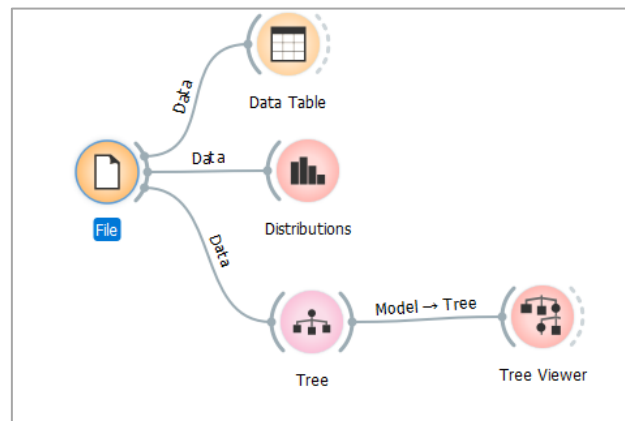Double click the **Data Table** widget to view the data.



Close the **Data Table**.

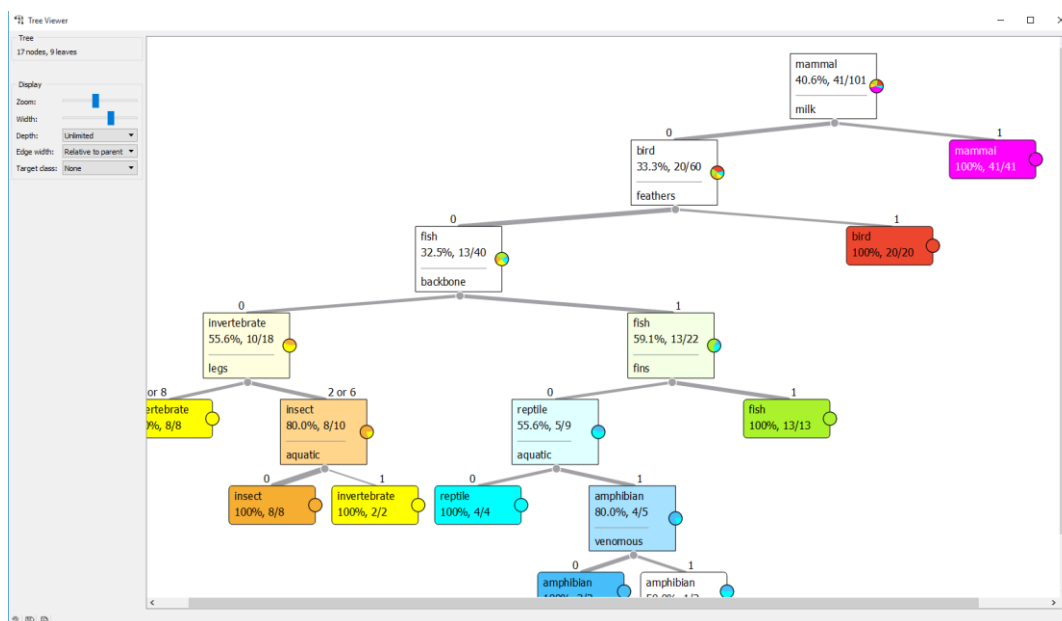Double click the **Distribution Table**. Select "type" under **variable** and "backbone" under **Split by**.

Close the **Distribution Table**.

Feed the data from the zoo dataset to the *Tree* widget (under the *Model category* from the *widget menus)*. Connect the output of the tree widget to a tree viewer widget (under the **visualize** category from the widget menu). The resulting configuration is shown below.



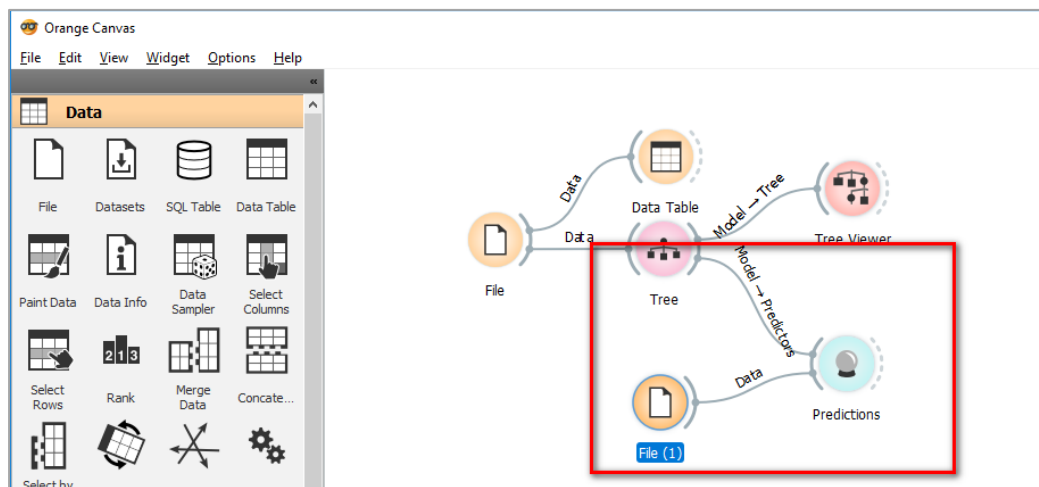Click the **Tree Viewer** to view the classification tree.



**Save** the model as "**zoo classification.ows**".

Use **OpenOffice Calc** or **Excel** to create a spreadsheet as follows.
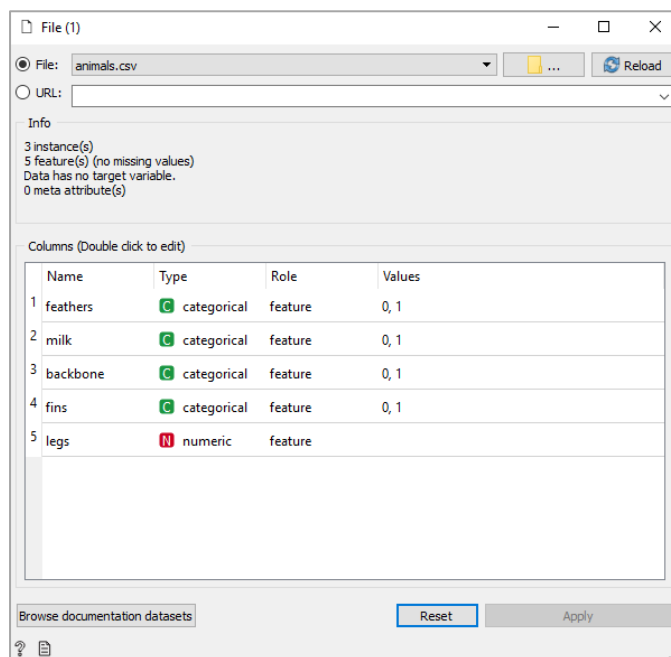
|   | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | feathers | milk | backbone | fins | legs |
| 2 | 1 | 0 | 1 | 0 | 2 |
| 3 | 0 | 0 | 1 | 1 | 0 |
| 4 | 0 | 1 | 1 | 0 | 4 |

Save the file as **animals.csv** on your desktop.

Add a **Predictions** widget and another **File** widget to the canvas and connect the workflow as follows.



Double click the newly created **File** widget (File (1)) and select **animals.csv** that you have created.



Click the **prediction** widget on the canvas to view the predicted class for the various animals in the given dataset.

| | Tree | features | milk | backbones | fins | legs |
|---|---|---|---|---|---|---|
| 1 | bird | 1 | 0 | 1 | 0 | 2 |
| 2 | bird | 0 | 0 | 1 | 1 | 0 |
| 3 | mammal | 0 | 1 | 1 | 0 | 4 |

*Note: Your predictions may be different than the one shown in the diagram, which depend on the setting of the hyper-parameters in the Tree widgets, which may affect the rules in the decision tree model.*
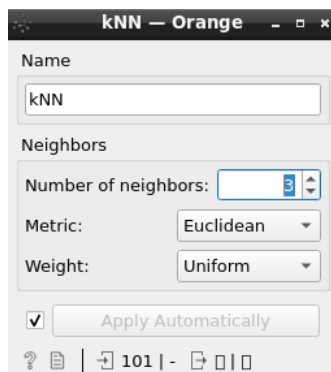
Use the **Tree Viewer** to understand how the predicted type of animal is computed.

Modify the workflow as follows

- Add a **KNN**, **SVM**, and **Logistic Regression** classifiers widgets. Connect the output of the dataset (the **File** widget) to the input of the classifiers. Connect the output of the classifier widgets to the prediction.
- Connect the output of Logistic Regression to a **Data Table**.



Double click the kNN widgets to review the hyper-parameters. Set the Number of **neighbors** to **3**.



Click the **prediction** widget to examine the predicted animal types.

Save your workflow as **zoo.ows**.

# 5. Evaluation of classification models

On April 15, 1912, during her maiden voyage, the widely considered "unsinkable" RMS Titanic sank after colliding with an iceberg. Unfortunately, there weren't enough lifeboats for everyone onboard, resulting in the death of 1502 out of 2224 passengers and crew.

In this section, we will analyse a simplified version of the Titanic Dataset provided by Orange. We will understand what sort of people were more likely to survive the Titanic sinking, develop a classification model to predict whether passengers will survive by building classification models and evaluate the model's accuracy.

Start a new workflow (**File -> New**).

## Data Exploration

Start a new workflow (**File -> New**). Create the following workflow.



Double click the Dataset widget. Use the **titanic dataset.**



Close the widget.

Double click **Data Table** to review the dataset.

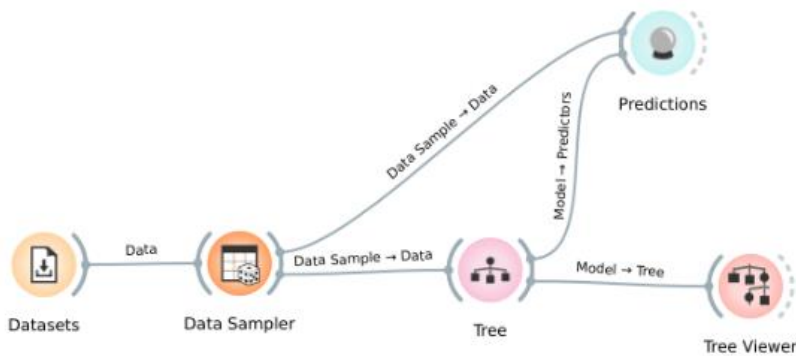| | survived | status | age | sex |
|---|---|---|---|---|
| 1 | yes | first | adult | male |
| 2 | yes | first | adult | male |
| 3 | yes | first | adult | male |
| 4 | yes | first | adult | male |
| 5 | yes | first | adult | male |
| 6 | yes | first | adult | male |
| 7 | yes | first | adult | male |
| 8 | yes | first | adult | male |
| 9 | yes | first | adult | male |
| 10 | yes | first | adult | male |
| 11 | yes | first | adult | male |
| 12 | yes | first | adult | male |
| 13 | yes | first | adult | male |

Double click the **Distribution** widget for exploration of the data. For the **Split by** option, Select **survived**
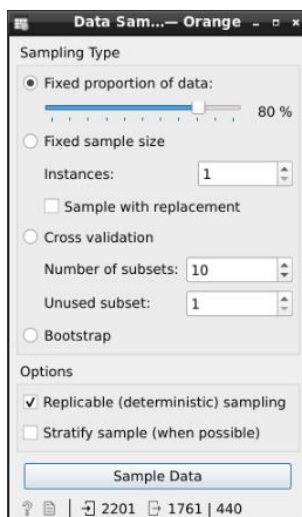


Select the different variable and explore which attribute will best split the data to predict the survival of a passenger.

## Training and evaluating the classification model

We will split the data into training set (80%) and test set (20%) and train a decision tree classifier on the training set.



For click the **Data Sampler**. Select Fixed proportion of data and select 80%. To ensure that that the result can be replicated, check the option **Replicable (deterministic) sampling**. If not, the data will be sampled randomly and you may get different results every time.



Double click the **Tree** widget. Check the option **Limit the maximal tree depth** to **1**. Uncheck all other options.

Click the **Tree Viewer** to view the decision tree.



Double click the **Predictions** widget. Check **Show performance scores** and **select (Average over classes)**.



What is the **Classification Accuracy (CA)** of the model on the training data?

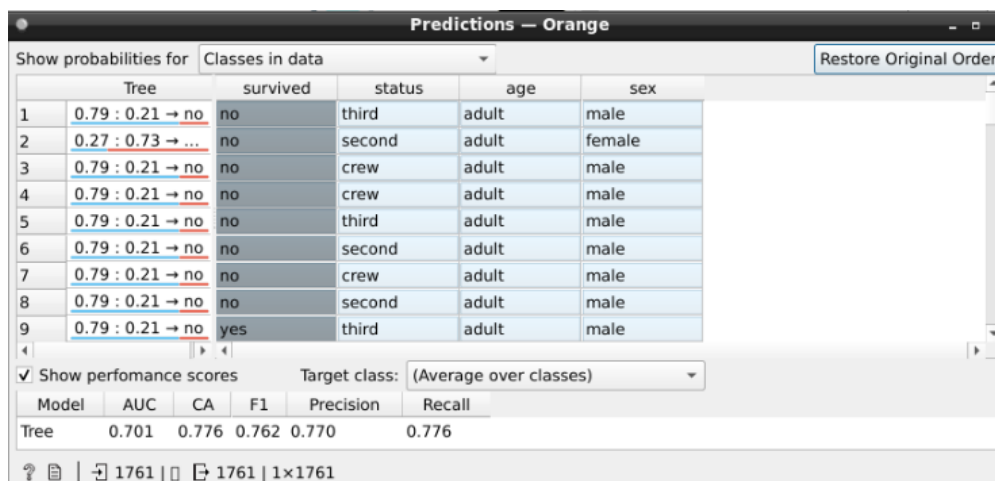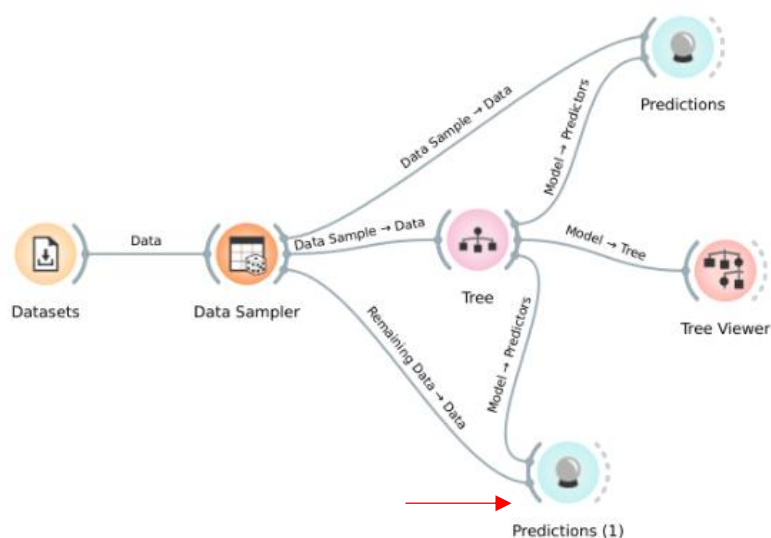Next, add another **Predictions** widget (**Predictions (1)** in the diagram below). To send the remaining 20% of the data to the second **Predictions** widget, double click the connection line and connect **Remaining Data** output of the Data Sampler widget to the Data input of the **Predictions (1)** widget.



Click the **Predictions(1)** widget to view the **classification accuracy** metrics on the **test set**.

**Exercise**

Change the **Limit the maximal tree depth** to 2,3, ... and observe the change in decision Tree. Observe the change in the classification accuracy on the training and test set.

# 6. Cross validation
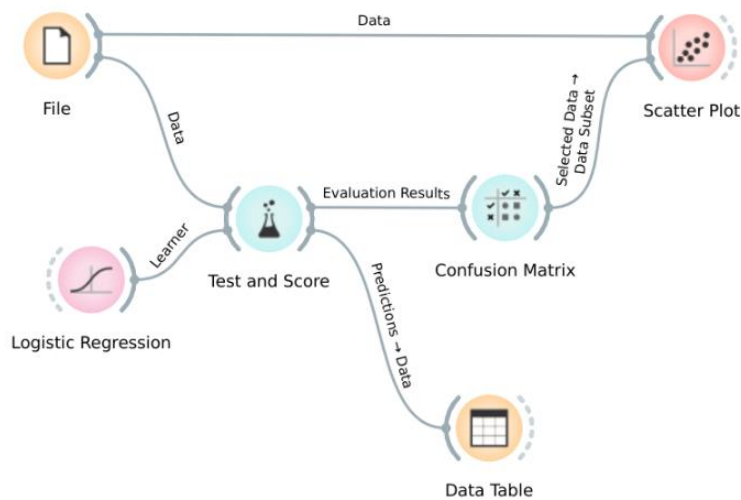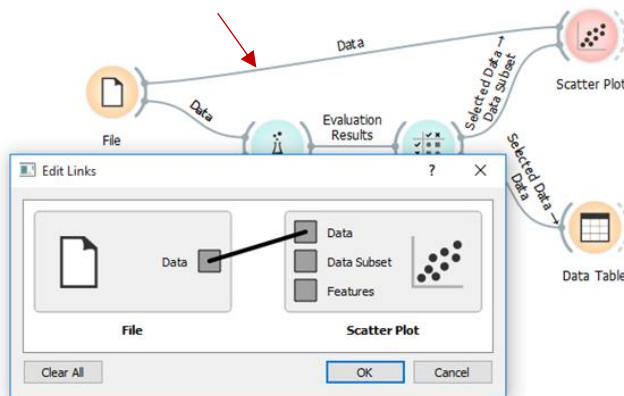
In this section, we will use logistic regression to classify the type of iris species. The performance of the model will be evaluated using cross validation and compared with other classification models.

Start a new workflow (**File -> New**).

Create a **data source** using the **Iris data set** and feed the data to a **scatter plot**. Add the **logistic regression widget** in the canvas and connect the **output channel** to the input channel of the **test & score widget (**available under the **evaluate category** in the menu**)**. Connect the output of **test & score widget** to the input channel of the **confusion matrix widget** (available under the **evaluate category** in the menu). The workflow configuration is illustrated below.



Double click the line between **File** and **Scatter Plot** and configure as follows.

Double click the line between **Confusion Matrix** and **Scatter Plot** and configure as follows.



Double click **Test & Score** to open the widget.



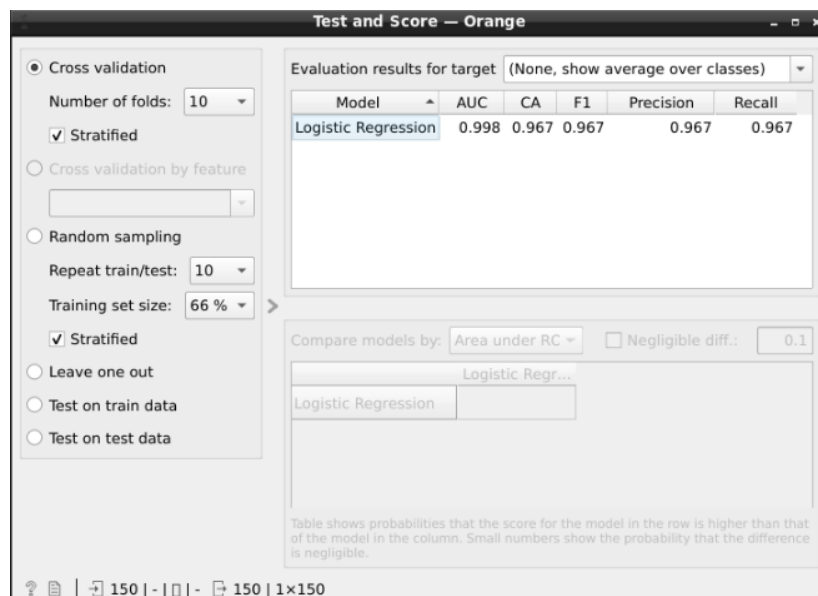Under Sampling, select **Cross validation** and choose **10**.  Check the **Stratified** option.

- A 10-fold cross validation is used to evaluation the predictive model. The dataset is first shuffled randomly and split into 10 groups. For each group (1/10 of the data), the data is used as the **validation set**. The remaining data is used as the **training set** to fit the model and evaluate it on the test set.
- In **stratified sampling**, samples for the training/validation/test set are selected in the same proportion as in the dataset.

The various performance metrics of logistics regression such as **classification accuracy (CA)**, **Precision** and **Recall** is shown in **Evaluation Results**.

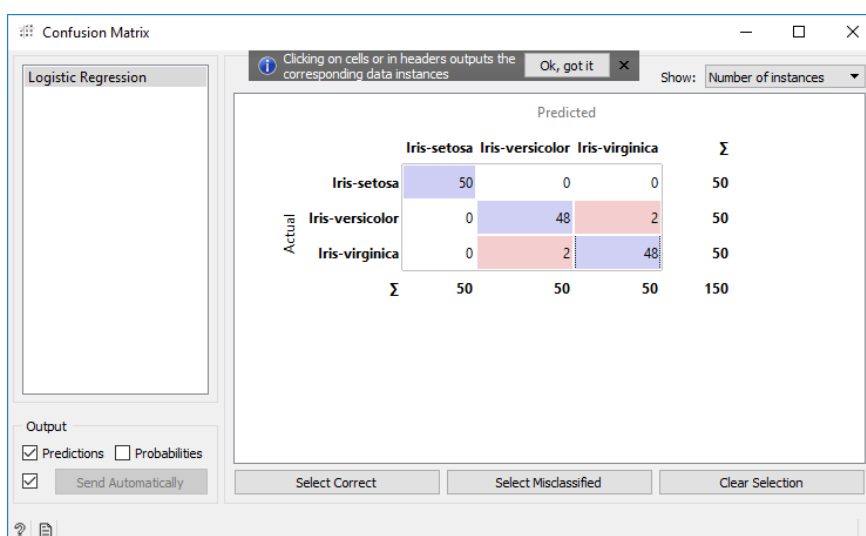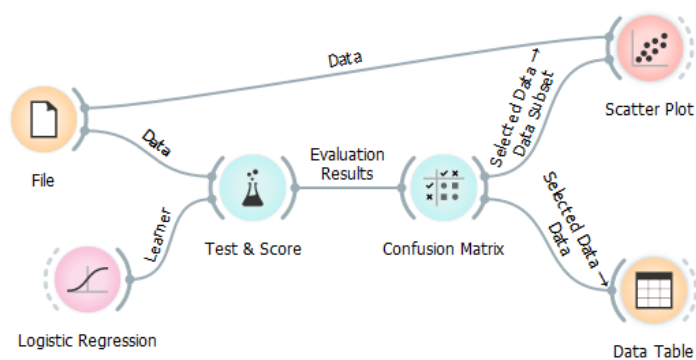Double click the **Data table** (with input from **Test and Score**).



You may view the **Fold** of each sample and their predicted probability for each class.

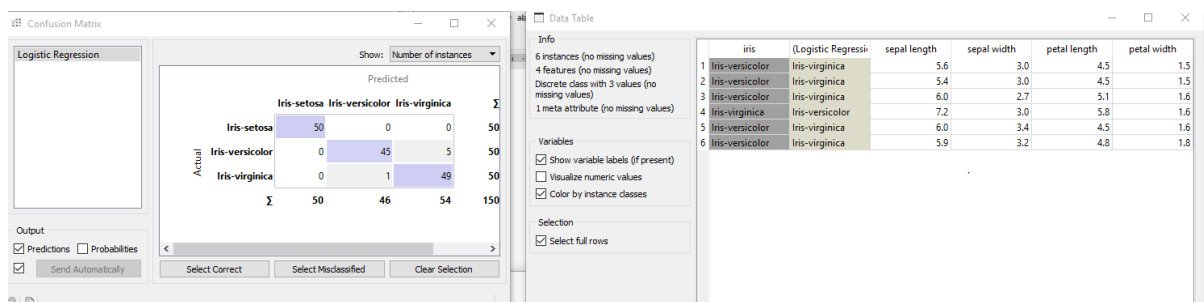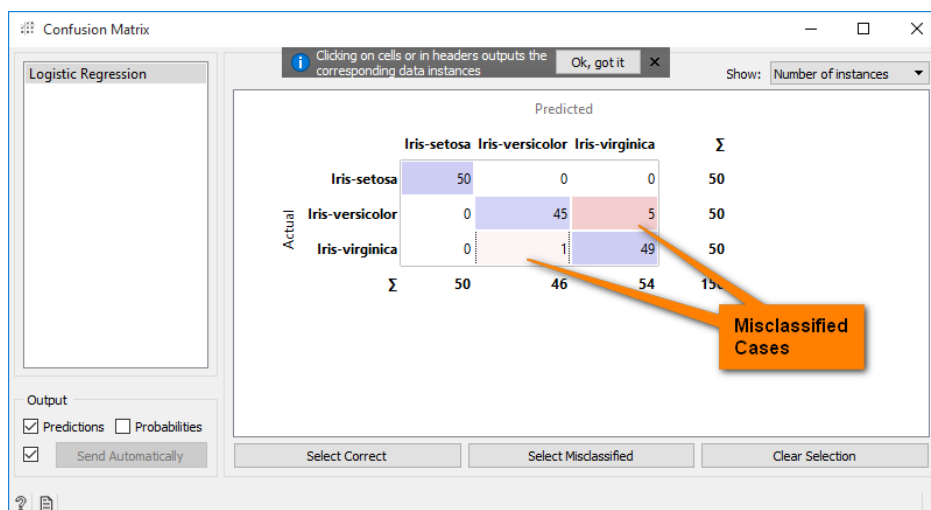| | iris | Tree | Tree (Iris-setosa) | Tree (Iris-versicolor) | Tree (Iris-virginica) | Fold | sepal length | sepal width | petal length | petal width |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 1 | 4.8 | 3.4 | 1.6 | 0.2 |
| 2 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 1 | 5.8 | 4.0 | 1.2 | 0.2 |
| 3 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 1 | 5.0 | 3.0 | 1.6 | 0.2 |
| 4 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 1 | 5.5 | 4.2 | 1.4 | 0.2 |
| 5 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 1 | 5.0 | 3.2 | 1.2 | 0.2 |
| 6 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.9 | 3.1 | 4.9 | 1.5 |
| 7 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 5.2 | 2.7 | 3.9 | 1.4 |
| 8 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.7 | 3.1 | 4.4 | 1.4 |
| 9 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.3 | 2.3 | 4.4 | 1.3 |
| 10 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.2 | 2.9 | 4.3 | 1.3 |
| 11 | Iris-virginica | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.3 | 3.3 | 6.0 | 2.5 |
| 12 | Iris-virginica | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 7.3 | 2.9 | 6.3 | 1.8 |
| 13 | Iris-virginica | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.7 | 3.3 | 5.7 | 2.1 |
| 14 | Iris-virginica | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 7.4 | 2.8 | 6.1 | 1.9 |
| 15 | Iris-virginica | Iris-versicolor | 0 | 0.5 | 0.5 | 1 | 6.7 | 3.3 | 5.7 | 2.5 |
| 16 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 2 | 4.3 | 3.0 | 1.1 | 0.1 |
| 17 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 2 | 5.7 | 4.4 | 1.5 | 0.4 |
| 18 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 2 | 5.4 | 3.4 | 1.7 | 0.2 |
| 19 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 2 | 4.7 | 3.2 | 1.6 | 0.2 |
| 20 | Iris-setosa | Iris-setosa | 1 | 0 | 0 | 2 | 4.8 | 3.0 | 1.4 | 0.3 |
| 21 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 2 | 6.0 | 2.9 | 4.5 | 1.5 |
| 22 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 2 | 6.0 | 3.4 | 4.5 | 1.6 |
| 23 | Iris-versicolor | Iris-versicolor | 0 | 0.5 | 0.5 | 2 | 5.8 | 2.6 | 4.0 | 1.2 |

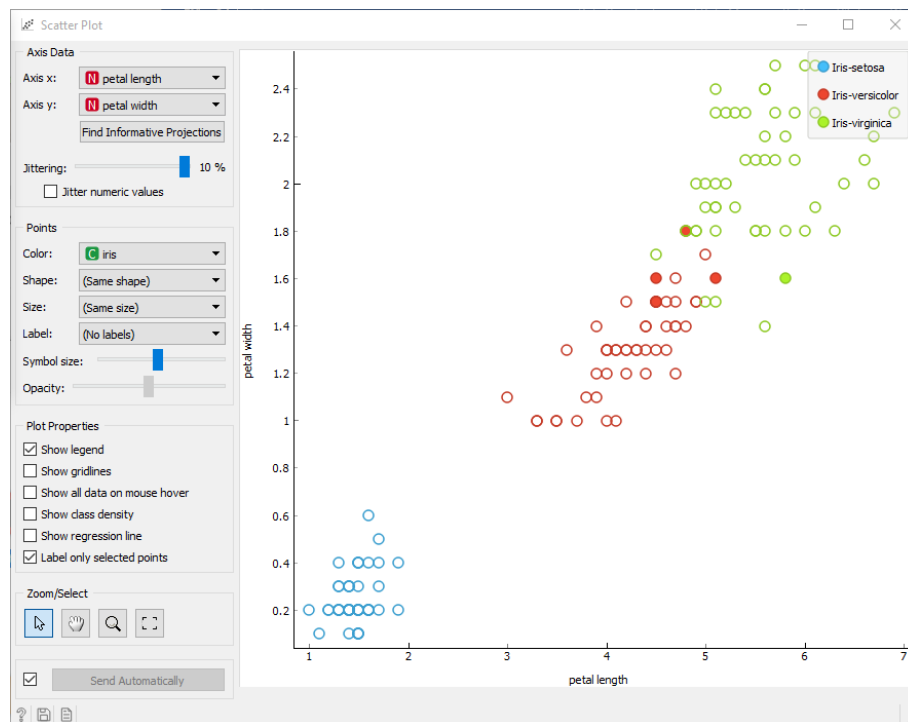Open the **Confusion Matrix** by double-clicking the widget.

Refine the workflow as follows.



To further explore the data, select the **misclassified cases** on the confusion matrix shown by holding "**Ctrl**" and click on the misclassified cases which are **highlighted in red** in the confusion matrix. We can view the data subset (the misclassified cases) in the data table and scatter plot.
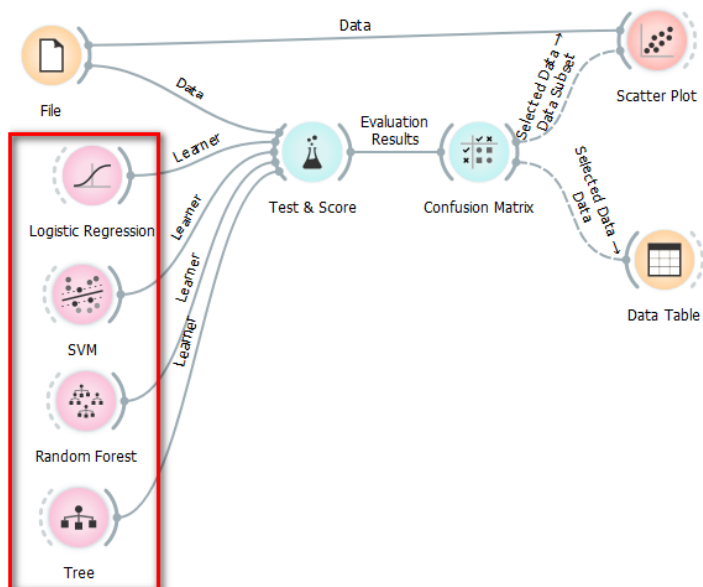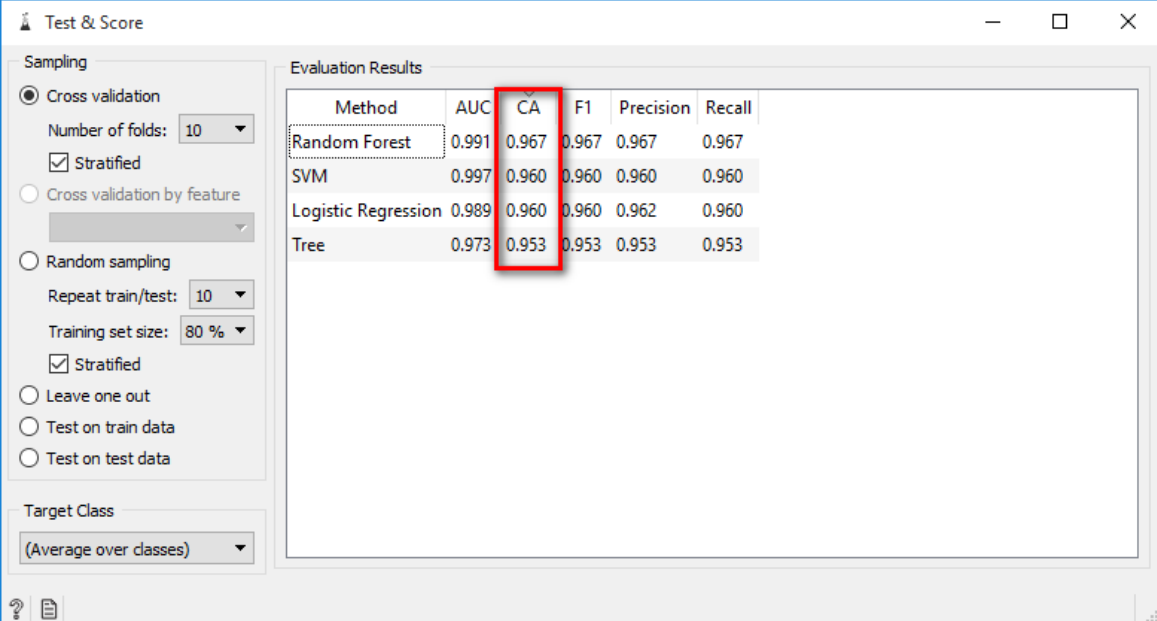
In the scatter plot above, the misclassified flowers are shown as filled circles.

Comparing different models

You may compare the performance of logistic regression with other classification algorithms (SVM, Decision Tree, Random Forest) by refining the workflow as follows.

Open the **Test & Score** to and compare the **Classification Accuracy** (CA) and other metrics of the various algorithms.



Select the option **Test on train data**. Compare the classification accuracy of the model.

Saving the model

Save the model as "**iris cross validation.ows**".

# 7. References

1. Orange Documentation
   https://orange.biolab.si/docs/
2. Getting Started with Orange
   https://www.youtube.com/playlist?list=PLmNPvQr9Tf-ZSDLwOzxpvY-HrE0yv-8Fy
3. Machine Learning Jargon
   https://orangedatamining.com/blog/2022/2022-02-01-machine-learning-jargon
4. The Kaggle titanic dataset
   https://www.kaggle.com/c/titanic