

# Data Sampler

Selects a subset of data instances from an input dataset.

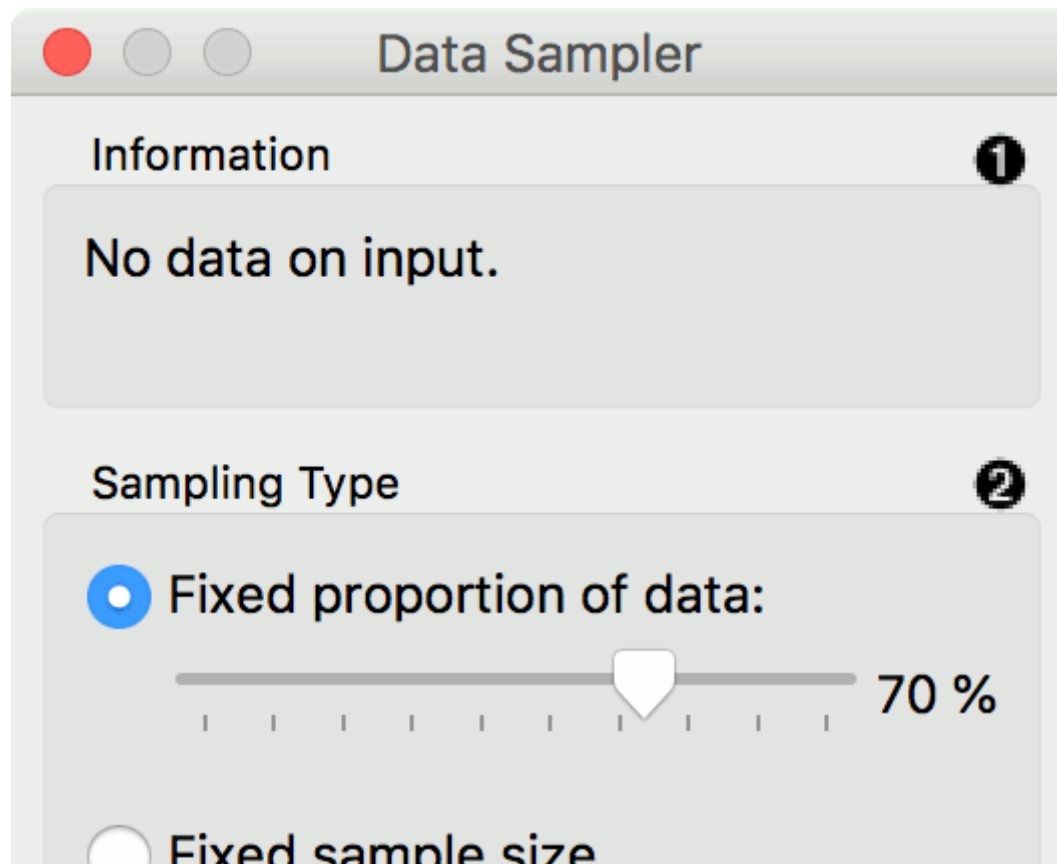
## Inputs

- Data: input dataset



## Outputs

- Data Sample: sampled data instances
- Remaining Data: out-of-sample data

The **Data Sampler** widget implements several data sampling methods. It outputs a sampled and a complementary dataset (with instances from the input set that are not included in the sampled dataset). The output is processed after the input dataset is provided and *Sample Data* is pressed.







**Fixed sample size**

Instances:   


☐ Sample with replacement

☐ Cross validation

Number of folds:   


Selected fold:   



☐ Bootstrap

**Options** 

☐ Replicable (deterministic) sampling

☐ Stratify sample (when possible)

**Sample Data** 

1. Information on the input and output dataset.
2. The desired sampling method:

- **Fixed proportion of data** returns a selected percentage of the entire data (e.g. 70% of all the data)
  - **Fixed sample size** returns a selected number of data instances with a chance to set *Sample with replacement*, which always samples from the entire dataset (does not subtract instances already in the subset). With replacement, you can generate more instances than available in the input dataset.
  - **Cross Validation** partitions data instances into the specified number of complementary subsets. Following a typical validation schema, all subsets except the one selected by the user are output as Data Sample, and the selected subset goes to Remaining Data. (Note: In older versions, the outputs were swapped. If the widget is loaded from an older workflow, it switches to compatibility mode.)
  - **Bootstrap** infers the sample from the population statistic.
3. *Replicable sampling* maintains sampling patterns that can be carried across users, while *stratify sample* mimics the composition of the input dataset.
  4. Press *Sample Data* to output the data sample.

If all data instances are selected (by setting the proportion to 100 % or setting the fixed sample size to the entire data size), output instances are still shuffled.

## Examples

First, let's see how the **Data Sampler** works. We will use the *iris* data from the **File** widget. We see there are 150 instances in the data. We sampled the data with the **Data Sampler** widget and we chose to go with a fixed sample size of 5 instances for simplicity. We can observe the sampled data in the **Data Table** widget (Data Table (in-sample)). The second **Data Table** (Data Table (out-of-sample)) shows the remaining 145 instances that weren't in the sample. To output the out-of-sample data, double-click the connection between the widgets and rewire the output to *Remaining Data* → *Data*.

The screenshot displays the Orange Data Mining software interface. A workflow is shown starting with a **File** widget connected to a **Data Sampler** widget. The **Data Sampler** widget is configured with the following settings:

- Information:** 150 instances in input dataset. Outputting 5 instances.
- Sampling Type:**
  - ☐ Fixed proportion of data: 70 %
  - ☒ Fixed sample size: Instances: 5
  - ☐ Sample with replacement
  - ☐ Cross validation: Number of folds: 10, Selected fold: 1
  - ☐ Bootstrap
- Options:**
  - ☐ Replicable (deterministic) sampling
  - ☐ Stratify sample (when possible)
- Buttons:** Sample Data

The **Data Sampler** widget has two outputs:

- Data Sample → Data:** Connected to a **Data Table (in-sample)** widget.
- Remaining Data → Data:** Connected to a **Data Table (out-of-sample)** widget.

The **Data Table (in-sample)** widget displays the following data:

	iris	sepal length	sepal width	petal length	petal width
1	Iris-setosa	4.8	3.4	1.6	0.2
2	Iris-versicolor	4.9	2.4	3.3	1.0
3	Iris-versicolor	6.1	3.0	4.6	1.4
4	Iris-virginica	7.7	3.0	6.1	2.3
5	Iris-virginica	6.7	3.0	5.2	2.3

The **Data Table (out-of-sample)** widget displays the following data:

	iris	sepal length	sepal width	petal length	petal width
1	Iris-versicolor	5.4	3.0	4.5	1.5
2	Iris-versicolor	6.5	2.8	4.6	1.5
3	Iris-setosa	5.4	3.4	1.5	0.4
4	Iris-versicolor	5.0	2.0	3.5	1.0
5	Iris-setosa	5.0	3.5	1.3	0.3
6	Iris-virginica	6.5	3.0	5.5	1.8
7	Iris-setosa	4.3	3.0	1.1	0.1
8	Iris-versicolor	5.7	2.6	3.5	1.0
9	Iris-setosa	4.9	3.1	1.5	0.1
10	Iris-virginica	6.8	3.0	5.5	2.1
11	Iris-versicolor	5.6	2.9	3.6	1.3
12	Iris-versicolor	5.8	2.7	4.1	1.0
13	Iris-setosa	5.1	3.4	1.5	0.2
14	Iris-setosa	5.4	3.9	1.3	0.4
15	Iris-versicolor	5.6	2.5	3.9	1.1
16	Iris-versicolor	5.6	3.0	4.1	1.3

Now, we will use the **Data Sampler** to split the data into training and testing part. We are using the *iris* data, which we loaded with the **File** widget. In **Data Sampler**, we split the data with *Fixed proportion of data*, keeping 70% of data instances in the sample.

Then we connected two outputs to the **Test & Score** widget, *Data Sample → Data* and *Remaining Data → Test Data*. Finally, we added **Logistic Regression** as the learner. This runs logistic regression on the Data input and evaluates the results on the Test Data.



**Data Sampler**

Information

150 instances in input dataset.  
Outputting 105 instances.

Sampling Type

☒ Fixed proportion of data:

70 %

☐ Fixed sample size

Instances: 5

☐ Sample with replacement

☐ Cross validation

Number of folds: 10

Selected fold: 1

☐ Bootstrap

Options

☐ Replicable (deterministic) sampling

☐ Stratify sample (when possible)

**Sample Data**

**Test & Score**

Sampling

☐ Cross validation

Number of folds: 10

☒ Stratified

☐ Cross validation by feature

☐ Random sampling

Repeat train/test: 10

Training set size: 66 %

☒ Stratified

☐ Leave one out

☐ Test on train data

☒ Test on test data

Target Class

(Average over classes)

Evaluation Results

Method	AUC	CA	F1	Precision	Recall
Logistic Regression	0.992	0.956	0.956	0.962	0.956

## Over/Undersampling

**Data Sampler** can also be used to oversample a minority class or undersample majority class in the data. Let us show an example for oversampling. First, separate the minority class using a **Select Rows** widget. We are using the *iris* data from the **File** widget. The data set has 150 data instances, 50 of each class. Let us oversample, say, *iris-setosa*.

In **Select Rows**, set the condition to *iris is iris-setosa*. This will output 50 instances of the *iris-setosa* class. Now, connect *Matching Data* into the **Data Sampler**, select *Fixed sample size*, set it to, say, 100 and select *Sample with replacement*. Upon pressing *Sample Data*, the widget will output 100 instances of *iris-setosa* class, some of which will be duplicated (because we used *Sample with replacement*).

Finally, use **Concatenate** to join the oversampled instances and the *Unmatched Data* output of the **Select Rows** widget. This outputs a data set with 200 instances. We can observe the final results in the **Distributions**.

