

# Test and Score

Tests learning algorithms on data.

## Inputs

- Data: input dataset
- Test Data: separate data for testing
- Learner: learning algorithm(s)

## Outputs

- Evaluation Results: results of testing classification algorithms

The widget tests learning algorithms. Different sampling schemes are available, including using separate test data. The widget does two things. First, it shows a table with different classifier performance measures, such as **classification accuracy** and **area under the curve**. Second, it outputs evaluation results, which can be used by other widgets for analyzing the performance of classifiers, such as **ROC Analysis** or **Confusion Matrix**.

The *Learner* signal has an uncommon property: it can be connected to more than one widget to test multiple learners with the same procedures.

**Test and Score**

**1 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

**2 Target Class**

(Average over classes)

**4 Model Comparison**

Area under ROC curve

☐ Negligible difference: 0.1

**3 Evaluation Results**

Model	AUC	CA	F1	Precision	Recall
Logistic Regression	0.732	0.776	0.762	0.770	0.776
Naive Bayes	0.700	0.776	0.762	0.770	0.776
SVM	0.500	0.531	0.546	0.581	0.531
Tree	0.737	0.783	0.749	0.816	0.783

**5 Model Comparison by AUC**

	Tree	SVM	Naive Bay...	Logistic R...
Tree		1.000	0.993	0.802
SVM	0.000		0.000	0.000
Naive Bayes	0.007	1.000		0.015
Logistic Regression	0.198	1.000	0.985	

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

1. The widget supports various sampling methods.

- Cross-validation** splits the data into a given number of folds (usually 5 or 10). The algorithm is tested by holding out examples from one fold at a time; the model is induced from other folds and examples from the held out fold are classified. This is repeated for all the folds.
- Cross validation by feature** performs cross-validation but folds are defined by the selected categorical feature from meta-features.
- Random sampling** randomly splits the data into the training and testing set in the given proportion (e.g. 70:30); the whole procedure is repeated for a specified number of times.

- **Leave-one-out** is similar, but it holds out one instance at a time, inducing the model from all others and then classifying the held out instances. This method is obviously very stable, reliable... and very slow.
  - **Test on train data** uses the whole dataset for training and then for testing. This method practically always gives wrong results.
  - **Test on test data:** the above methods use the data from *Data* signal only. To input another dataset with testing examples (for instance from another file or some data selected in another widget), we select *Separate Test Data* signal in the communication channel and select Test on test data.
2. For classification, *Target class* can be selected at the bottom of the widget. When *Target class* is (Average over classes), methods return scores that are weighted averages over all classes. For example, in case of the classifier with 3 classes, scores are computed for class 1 as a target class, class 2 as a target class, and class 3 as a target class. Those scores are averaged with weights based on the class size to retrieve the final score.
  3. The widget will compute a number of performance statistics. A few are shown by default. To see others, right-click on the header and select the desired statistic.

■ Classification

**Test and 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)

**Model Comparison**

Area under ROC curve

☐ Negligible difference: 0.1

**Evaluation Results**

Model	AUC	CA	F1	Precision	Recall
Logistic Regression	0.732	0.776	0.762	0.770	0.776
Naive Bayes	0.700	0.776	0.762	0.770	0.776
SVM	0.500	0.531	0.546	0.581	0.531
Tree	0.737	0.783	0.749	0.816	0.783

**Model Comparison by AUC**

	Tree	SVM	Naive Bay...	Logistic R...
Tree		1.000	0.993	0.802
SVM	0.000		0.000	0.000
Naive Bayes	0.007	1.000		0.015
Logistic Regression	0.198	1.000	0.985	

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

- **Area under ROC** is the area under the receiver-operating curve.
- **Classification accuracy** is the proportion of correctly classified examples.
- **F-1** is a weighted harmonic mean of precision and recall (see below).
- **Precision** is the proportion of true positives among instances classified as positive, e.g. the proportion of *Iris virginica* correctly identified as *Iris virginica*.
- **Recall** is the proportion of true positives among all positive instances in the data, e.g. the number of sick among all diagnosed as sick.

- **Specificity** is the proportion of true negatives among all negative instances, e.g. the number of non-sick among all diagnosed as non-sick.
- **LogLoss** or cross-entropy loss takes into account the uncertainty of your prediction based on how much it varies from the actual label.
- Train time - cumulative time in seconds used for training models.
- Test time - cumulative time in seconds used for testing models.

## ■ Regression

**Test and 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

**Evaluation Results**

Model	MSE	RMSE	MAE	R2
Constant	84.644	9.200	6.662	-0.003
kNN	38.676	6.219	4.352	0.542
SGD	24.751	4.975	3.514	0.707
SVM	36.665	6.055	3.710	0.566

**Model Comparison by MSE**

	Constant	kNN	SGD	SVM
Constant		1.000	1.000	1.000
kNN	0.000		0.964	0.626
SGD	0.000	0.036		0.001
SVM	0.000	0.374	0.999	

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

**Model Comparison**

Mean square error  
☐ Negligible difference: 0.1

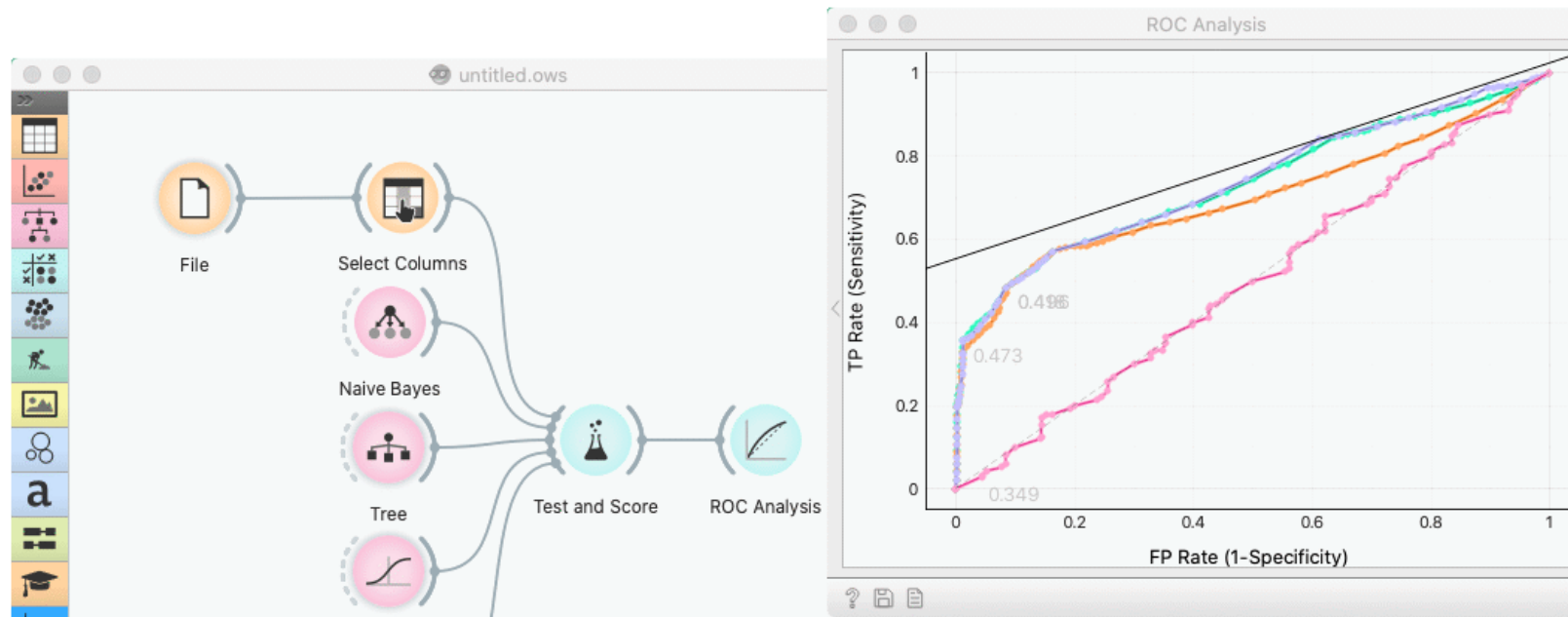
- **MSE** measures the average of the squares of the errors or deviations (the difference between the estimator and what is estimated).
- **RMSE** is the square root of the arithmetic mean of the squares of a set of numbers (a measure of imperfection of the fit of the estimator to the data)
- **MAE** is used to measure how close forecasts or predictions are to eventual outcomes.

- **R<sup>2</sup>** is interpreted as the proportion of the variance in the dependent variable that is predictable from the independent variable.
  - **CVRMSE** is RMSE normalized by the mean value of actual values.
  - Train time - cumulative time in seconds used for training models.
  - Test time - cumulative time in seconds used for testing models.
4. Choose the score for pairwise comparison of models and the region of practical equivalence (ROPE), in which differences are considered negligible.
  5. Pairwise comparison of models using the selected score (available only for cross-validation). The number in the table gives the probability that the model corresponding to the row is better than the model corresponding to the column. If negligible difference is enabled, the smaller number below shows the probability that the difference between the pair is negligible. The test is based on the **Bayesian interpretation of the t-test** ([shorter introduction](#)).
  6. Get help and produce a report.

## Example

In a typical use of the widget, we give it a dataset and a few learning algorithms and we observe their performance in the table inside the **Test & Score** widget and in the **ROC**. The data is often preprocessed before testing; in this case we did some manual feature selection (**Select Columns** widget) on *Titanic* dataset, where we want to know only the sex and status of the survived and omit the age.

In the bottom table, we have a pairwise comparison of models. We selected that comparison is based on the *area under ROC curve* statistic. The number in the table gives the probability that the model corresponding to the row is better than the model corresponding to the column. We can, for example, see that probability for the tree to be better than SVM is almost one, and the probability that tree is better than Naive Bayes is 0.001. Smaller numbers in the table are probabilities that the difference between the pair is negligible based on the negligible threshold 0.1.



Logistic Regression

Select Columns

Available Variables

Filter

age

Up

>

Down

Features

Filter

status

sex

Up

>

Down

Target Variable

survived

Up

>

Down

Meta Attributes

Up

>

Down

Reset

Send Automatically

Test and 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)

Model Comparison

Area under ROC curve

☒ Negligible difference: 0.1

Evaluation Results

Model	AUC	CA	F1	Precision	Recall
Tree	0.737	0.783	0.749	0.816	0.783
SVM	0.500	0.531	0.546	0.581	0.531
Naive Bayes	0.700	0.776	0.762	0.770	0.776
Logistic Regression	0.732	0.776	0.762	0.770	0.776

Model Comparison by AUC

	Tree	SVM	Naive Bay...	Logistic R...
Tree		0.999 0.001	0.001 0.999	0.000 1.000
SVM	0.000 0.001		0.000 0.006	0.000 0.002
Naive Bayes	0.000 0.999	0.994 0.006		0.000 0.999
Logistic Regression	0.000 1.000	0.998 0.002	0.001 0.999	

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

Another example of using this widget is presented in the documentation for the [Confusion Matrix](#) widget.