

Installation for Mac OS X

Orange is a comprehensive data mining suite that includes components in C++, interface to scripting language Python, a number of data mining modules developed in Python, and a set of graphical user interface components called Orange widgets. The guiding ideas behind Orange are flexibility through scripting and visual programming, component-based design and ease of use.

Orange is developed as an Open Source platform. Thanks to portable code and usage of portable solutions such as Python, Qt, PyQt, Qwt, Numerical, and other can be run on Windows, Linux and Mac OS X. This document is about the installation of Orange under Mac OS X, which was beta tested under Mac OS X 10.2 and 10.3.

Installation

There are two installation files (disk images: Orange Install.dmg and Orange Doc.dmg). Put them on the desktop, double click them and drag the contents to whatever location you wish. In this way, you will copy an Orange application (bundle) and a documentation folder containing also several data sets.

Orange installation is a fairly large file (about 15 MBytes) and besides Orange application and visual programming environment Orange Canvas contains a number of precompiled libraries and programs, including Python, Qt, Qwt, Numeric, and other.

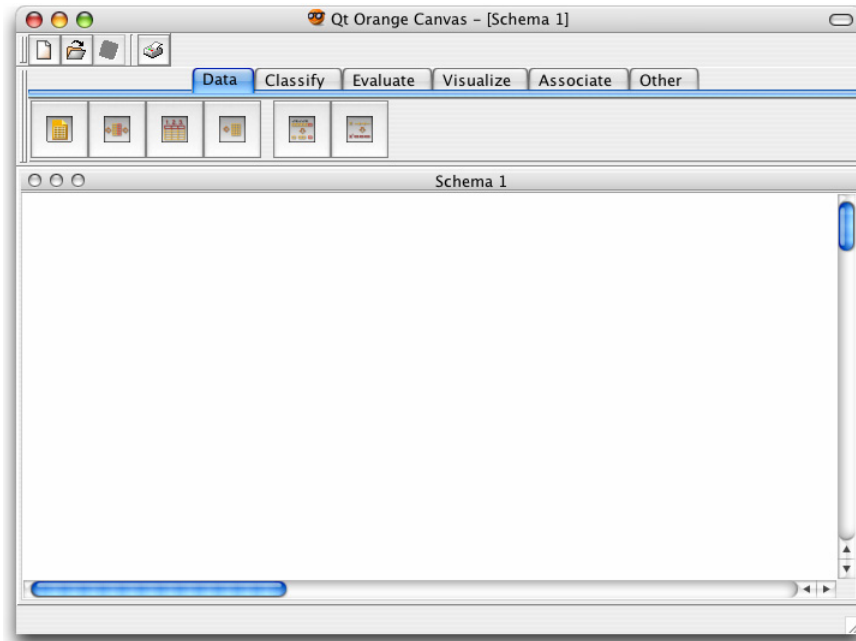
The installation package also contains three white papers that, besides this document, include a white paper on Orange and a white paper on Orange Widgets.



Orange Widgets most often include a part where user sets the parameters for particular data analysis or visualization method, and a part with visualization.

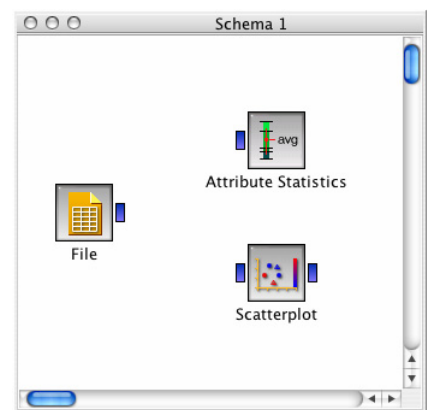
Starting With Orange: Few Examples

To run Orange, double click on Orange icon. This should open the following window.

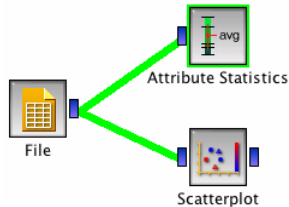
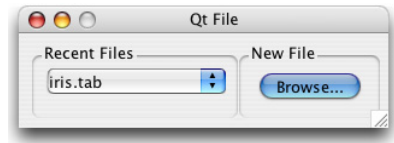


Some Simple Data Exploration

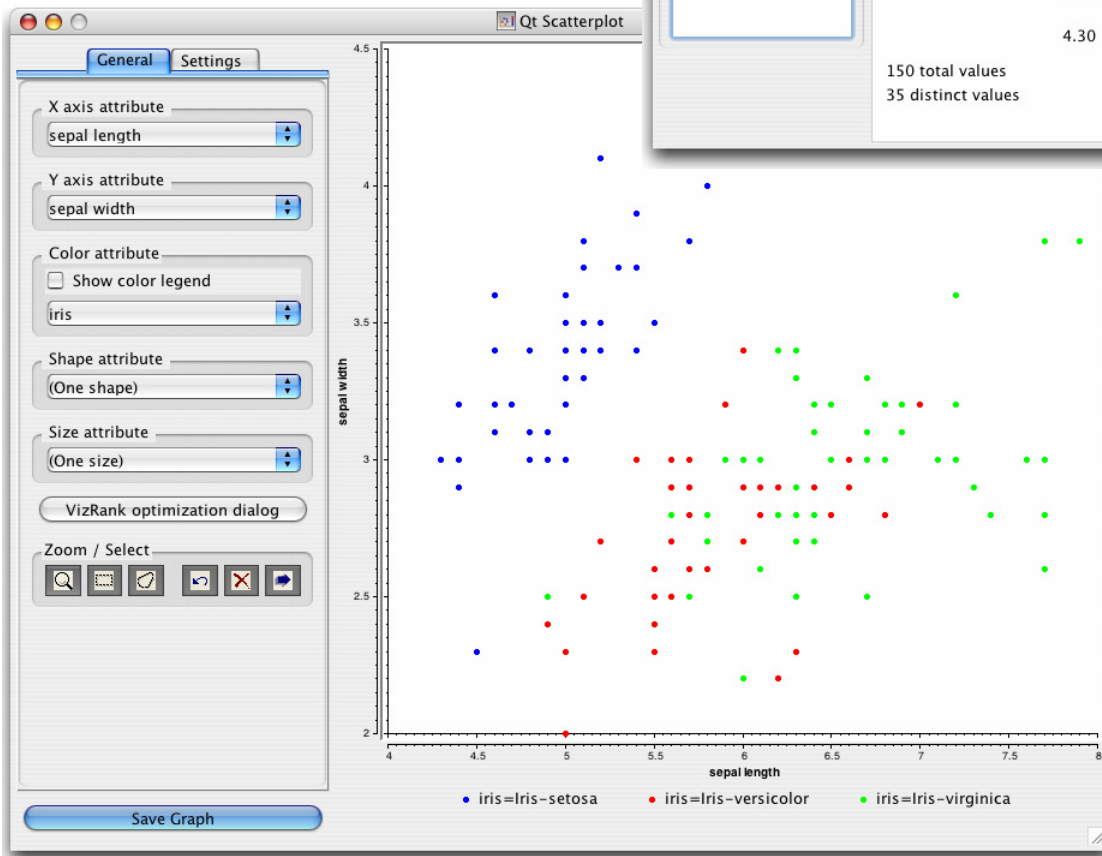
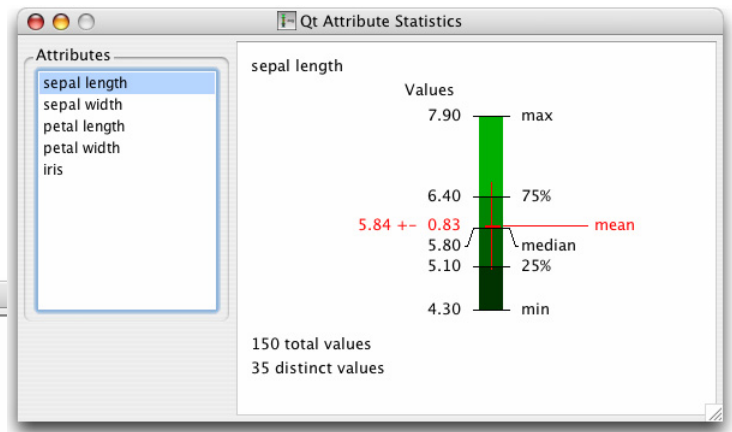
Let us now create an example schema. Put a File, Attribute Statistics, and Scatterplot widget on a canvas (finding these widgets in Data and Visualize tabs and clicking on them once will do the job). With some appropriate placing (click-and-drag) you should get something similar to what we show here.



Now double click the File widget, and choose the data file called iris.tab. You will find it in Orange Doc/datasets folder. This is the data set about classification of Iris flower to three different species.



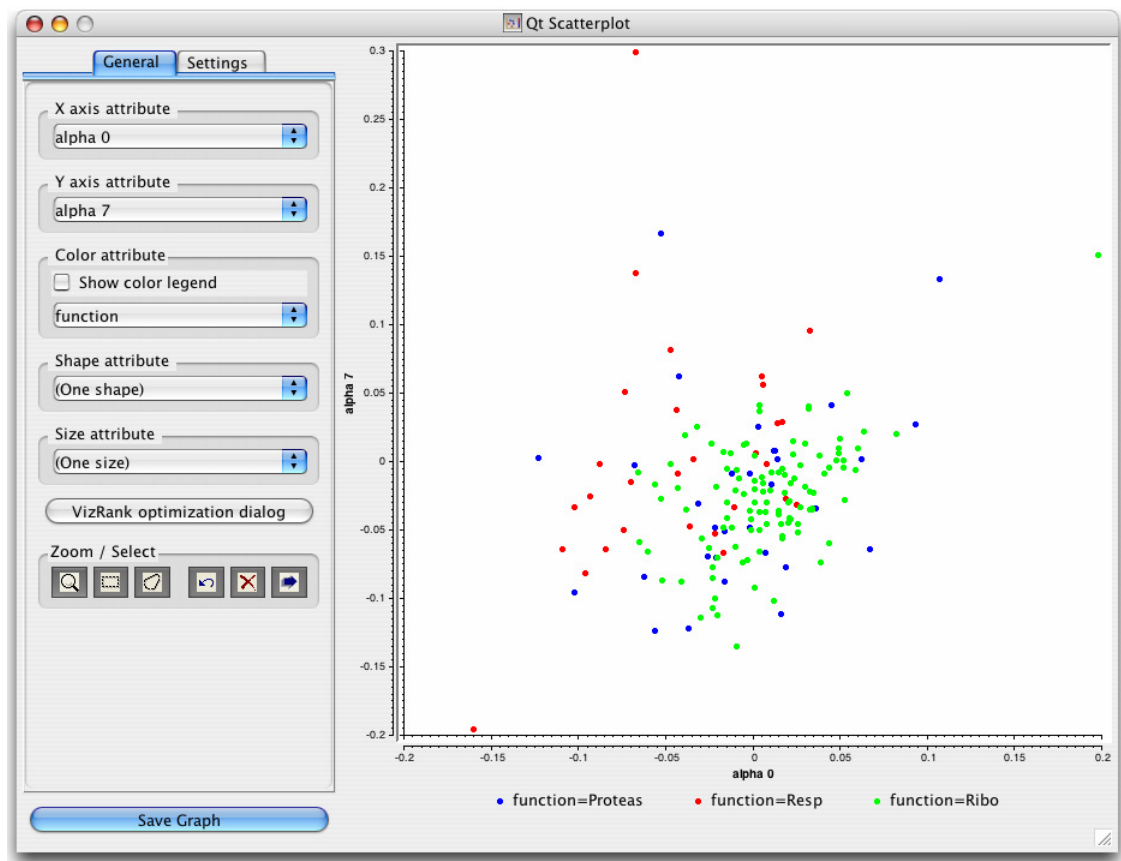
Now, to see some basic statistics and visualizations, connect the three widgets as shown below (click on right blue rectangle of File widget and draw the line to each of the other two widgets). Double click the widgets to obtain the visualizations.



Play with the widgets and settings. To see how tokens propagate through the schema, select a different file in the File widget and observe the update of the other two widgets.

Finding the Best Scatterplot

For another exercise with widgets and the same schema, choose yeast-class-RPR.tab data set. This is a data set on expressions of yeast genes, which contains 79 measurements (attributes) and where each gene is labeled with one of three classes. The number of different scatterplots is large, and the one you will first see is rather messy. How to find a good one? In the Scatterplot widget choose VizRank optimization dialog, click on “Start evaluating projections” button, and then after few projections are proposed (list box on the right), stop the process and click on the first suggested projection. Better than the one on this page, right?



In Orange, Scatterplot goes beyond simple 2D visualization of the data as it includes a wizard called VizRank that can suggest which attribute pair would define an interesting projection for visualization.

Qt VizRank Optimization Dialog

Optimization Settings

- Number of neighbors (k): 12
- Number of interesting projections: 100
- Minimum examples in data set: 0
- Percent of data used in evaluation: 100

Set heuristic for attribute ordering

Measure of classification success

- Average probability assigned to the correct class

Testing method

- 10 fold cross validation

Find interesting projections...

Start evaluating projections

Evaluate current projection / classifier

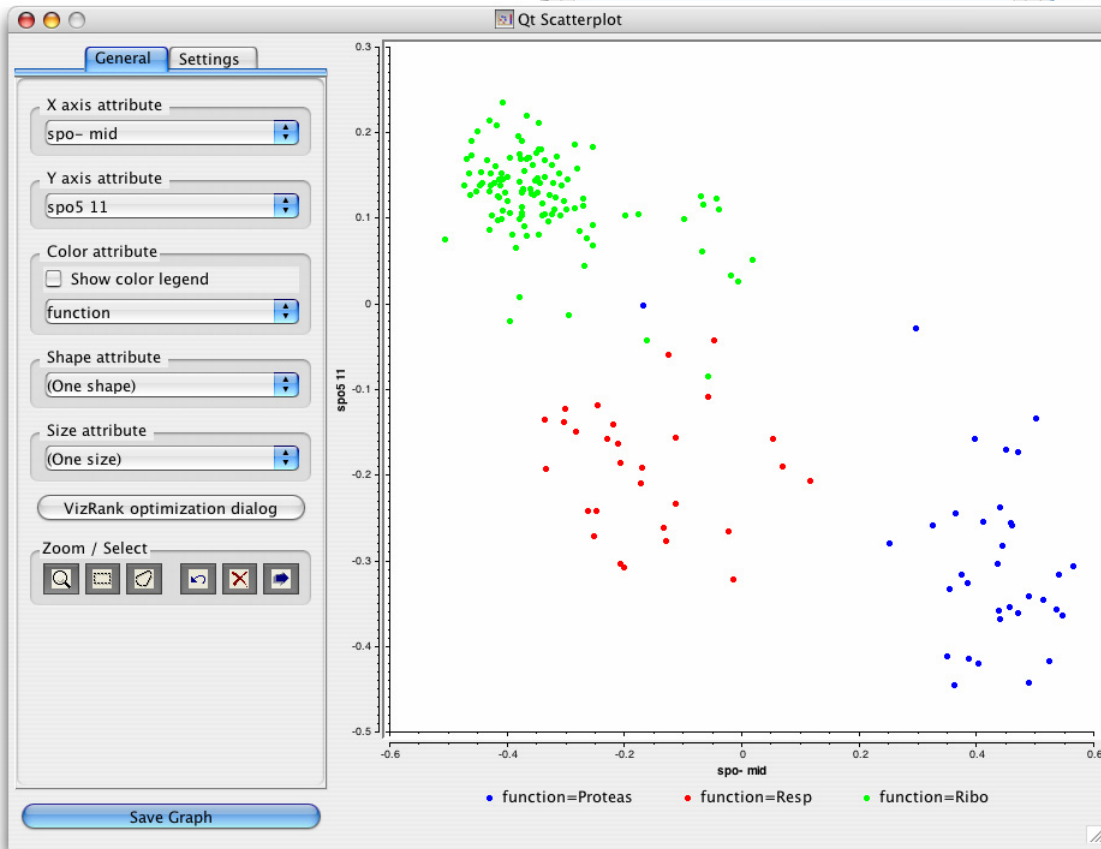
- Evaluate projection
- Save projection
- kNN correct
- kNN wrong
- Original

Manage projections

Number of concurrently visualized attributes:

List of interesting projections

- (97.09, 182) - spo- mid, spo5 11
- (94.60, 182) - spo- early, spo5 11
- (93.35, 180) - spo 2, spo- early
- (92.43, 180) - spo 2, spo- mid
- (91.55, 180) - spo 2, spo5 11
- (91.44, 183) - spo 7, spo5 11
- (91.39, 182) - spo 11, spo- mid
- (91.30, 182) - spo 11, spo- early
- (91.00, 180) - spo 5, spo- early
- (90.79, 180) - spo 5, spo5 11
- (90.21, 183) - spo 11, spo 7
- (89.93, 180) - spo 11, spo 5
- (89.39, 180) - spo 5, spo- mid
- (89.35, 180) - spo 11, spo 2
- (88.82, 183) - spo 9, spo- early
- (88.63, 181) - spo- early, spo5 7
- (88.58, 181) - spo- mid, spo5 7
- (88.51, 182) - spo 7, spo- early
- (88.37, 184) - heat 80, spo5 11
- (87.86, 183) - spo 11, spo5 11
- (87.78, 183) - spo5 11, spo5 2
- (87.77, 183) - spo 9, spo5 11
- (87.57, 183) - spo- early, spo5 2
- (87.40, 183) - spo 9, spo- mid
- (87.23, 182) - spo 7, spo- mid
- (86.64, 183) - spo 7, spo 9
- (86.32, 181) - spo 5, spo 9
- (86.32, 181) - spo 2, spo 9
- (86.10, 183) - spo- mid, spo5 2
- (81) - spo5 11, spo5 7
- (82) - spo- early, spo- mi
- (79) - spo 2, spo 5
- (80) - spo 2, spo5 7
- (78) - spo 0, spo5 11
- (78) - spo 0, spo- early
- (81) - spo 2, spo5 2
- (81) - spo 2, spo 7
- (81) - spo 5, spo5 2
- (82) - spo 7, spo5 7
- (78) - spo 0, spo- mid



The driving idea in Orange is to complement computationally intensive widgets, like those that build predictive models from data, with widget that can visualize the results of such data analysis.

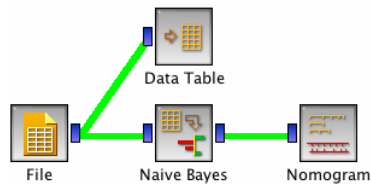
Survival on Titanic

Orange includes a number of widgets that are about induction of classification rules. In essence, these widgets take a data set where each example (data instance) is labeled with one of the classes, and build a model that is able to predict the class given some information about an example.

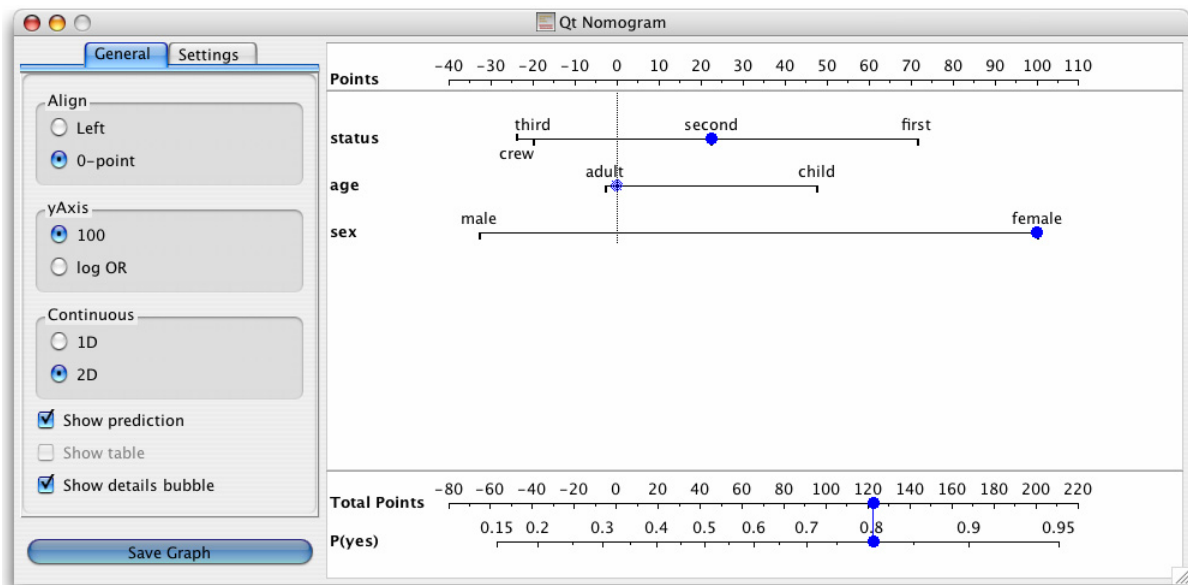
Here is an example. Consider a data set on H.M.S. Titanic (titanic.tab included in the distribution). There were 2201 passengers, each described with three attributes: status (first, second, third class or crew member), age (child or adult) and gender.

	status	age	sex	survived
2177	crew	adult	male	no
2178	crew	adult	male	no
2179	crew	adult	female	yes
2180	crew	adult	female	yes
2181	crew	adult	female	yes

The data also records whether a particular passenger survived Titanic's accident. Orange offers various methods to build models that can predict the probability of the class values. Perhaps the simplest of them is naive Bayesian classifier, which also has a nice visualization in the form of the nomogram. To explore the Titanic data, we have used the schema on the left. Notice that the data is first sent to naive Bayesian classifier, which builds the predictive model that is then sent to a Nomogram widget.



To use the nomogram for prediction of survival on Titanic, move the blue dots on the nomogram around. The snapshot shows, for instance, that women traveling in the second class had about 80% chance of survival.



What Next?

These were just a few examples to get you started. Playing with other widgets requires the users to be knowledgeable in data mining. If you are, you should go ahead: widgets are fairly easy to understand, and combining them is just like thinking about what to do in any data analysis task.

Notice that Orange is work in progress, and that this document together with an installation is in beta and prepared for Trolltech Mac Qt competition. We thank the organizers of competition for providing the Qt, which is an excellent environment, and for setting a firm deadline which pushed us towards speeding our Orange port to Mac OS X. The version submitted to the Trolltech should become available at Orange's web site (<http://magix.fri.uni-lj.si>, to become <http://www.ailab.si/orange> in the next few days when the package is released to general public).

