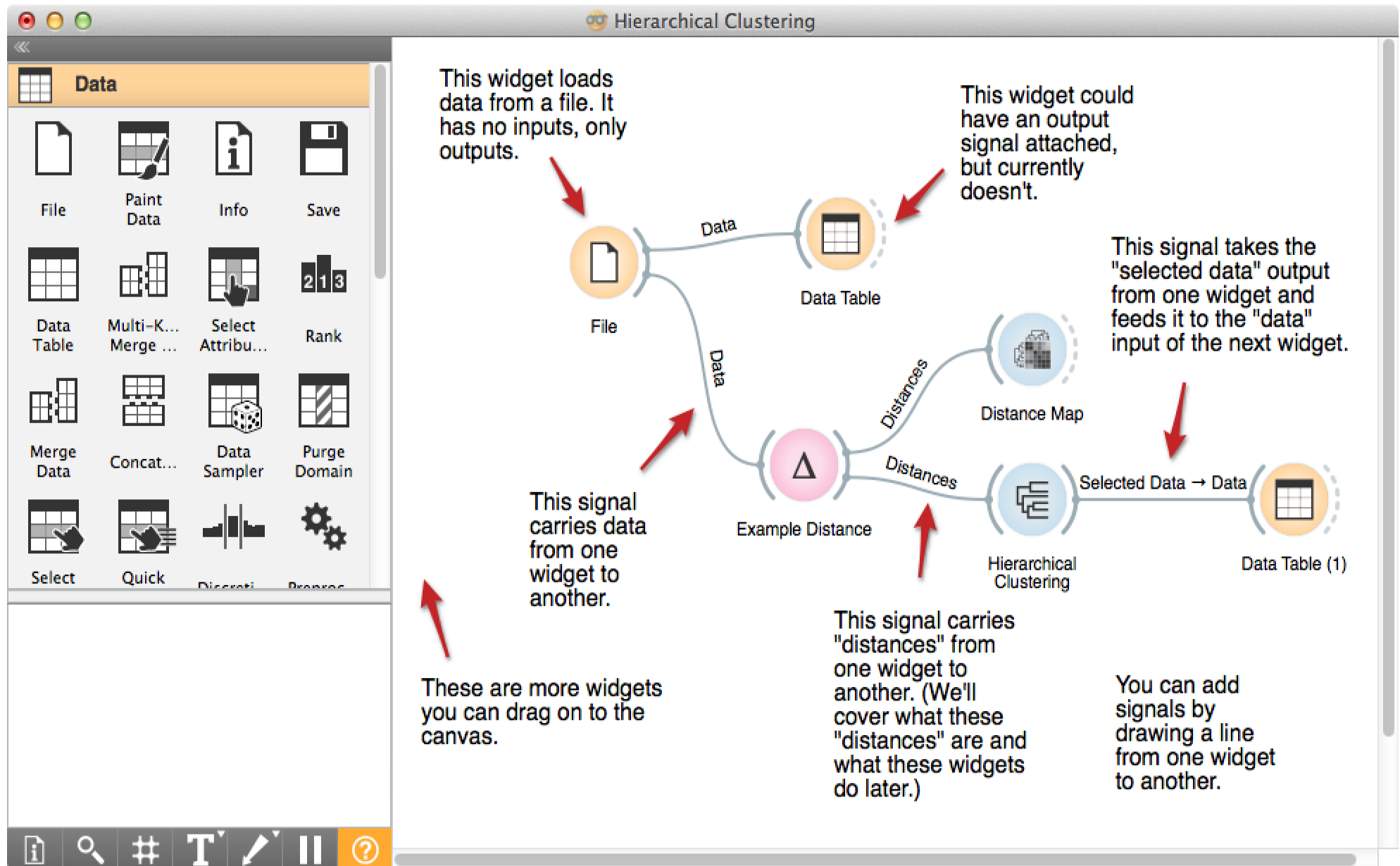# Introduction to Orange

# Introduction to Orange

- Orange is a data mining toolkit, so you don't need to be an expert in any of those subjects
- We will use Orange to:
  - load, manipulate, and save large data sets
  - visualize the relationships between variables
  - discover and quantify patterns in data
  - create rules to predict outcomes based on observed data
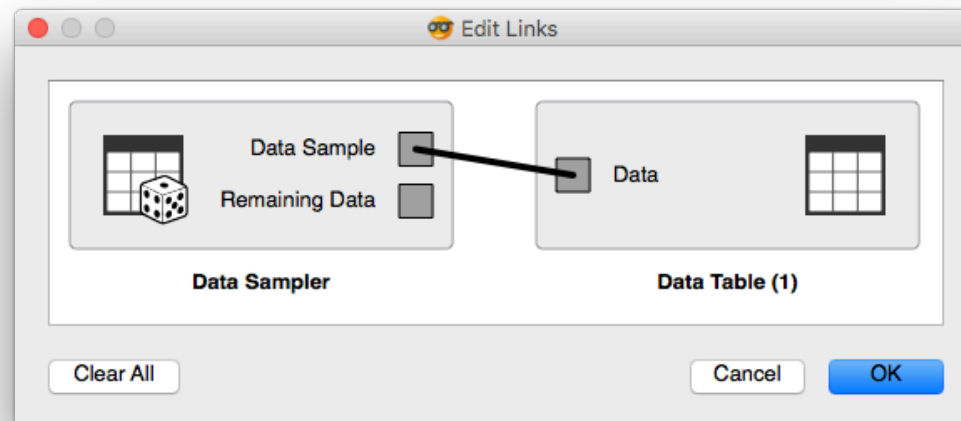
# Orange: Graphical Programming

# Using the Orange interface

- To add a widget, drag it onto the canvas from the widget panel, or just click on it in the widget panel

- To add a signal, click on the signal attachment point on a widget and drag from it to the signal attachment point on another widget
  - Input signals come in from the left, output signals go out to the right
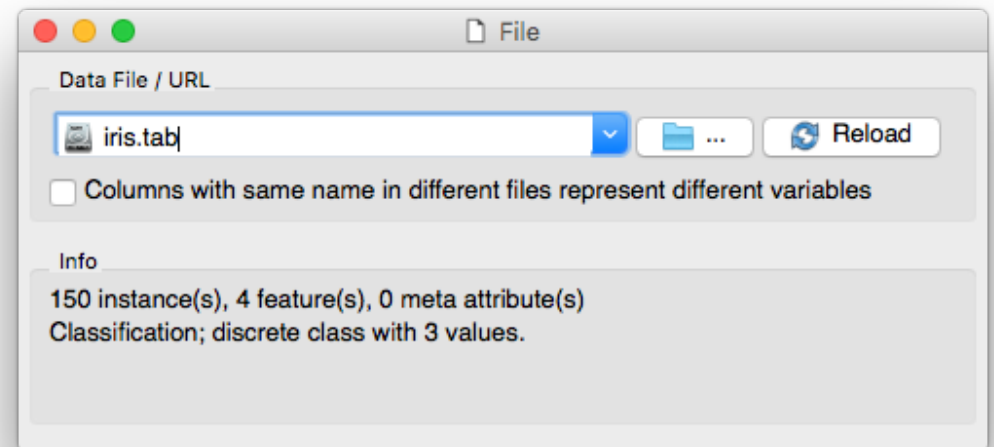
# Using the Orange interface

- Some widgets have multiple possible input and output ports
    - Orange tries to guess which one you mean
    - If it guesses wrong, double click on the signal to select which inputs and outputs you are using
    - You can also temporarily disconnect or delete signals by right-clicking on them

# File Widget

- Loads data from a file
- Many different file types are supported
  - Recommended: tab-delimited text
- *iris.tab* is an example dataset that comes with Orange, and contains 150 iris flowers from three species
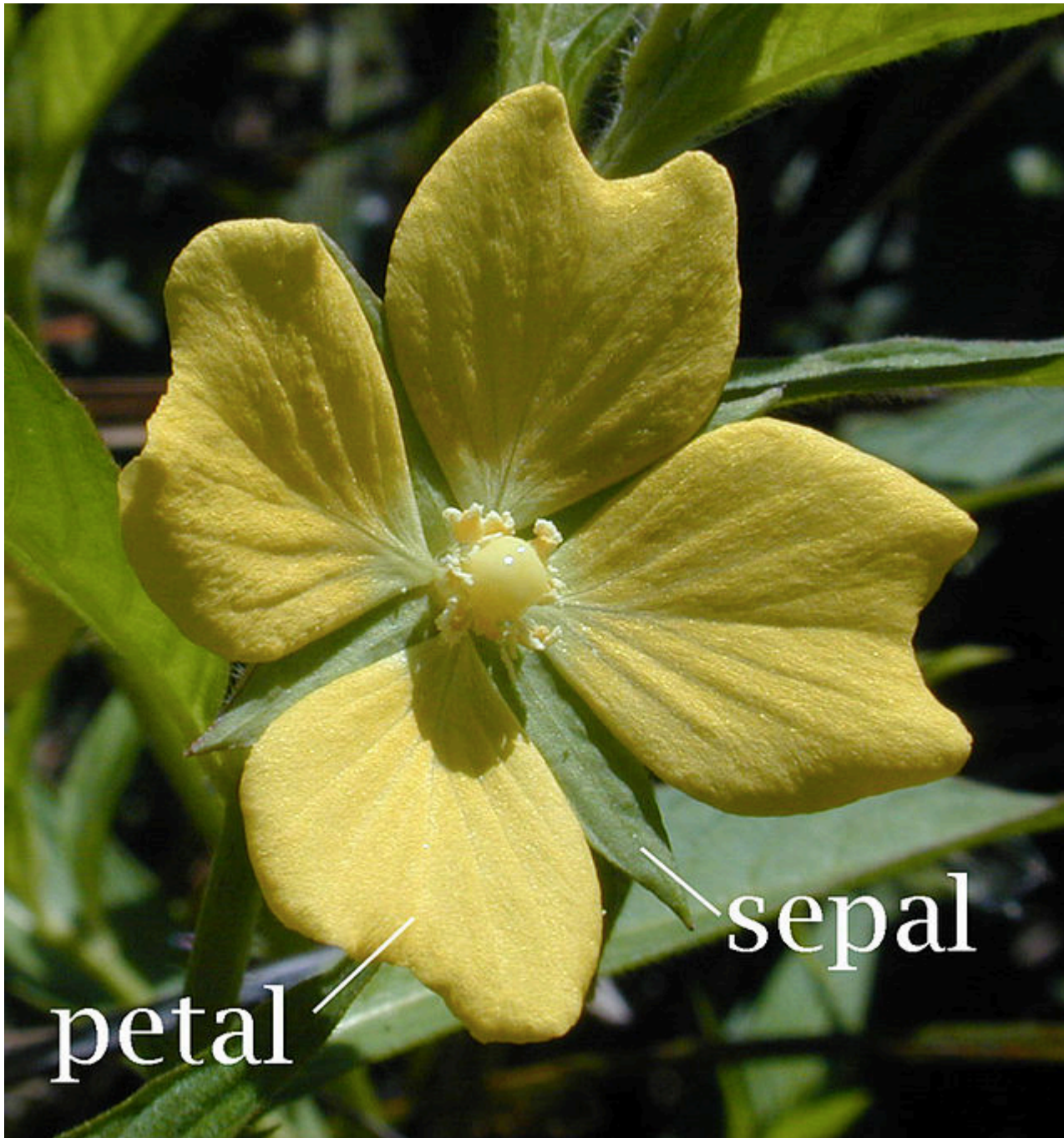
# Data Table Widget



- Lists rows in a dataset, sort by clicking on the column heading

- Each value has a bar showing how big it is

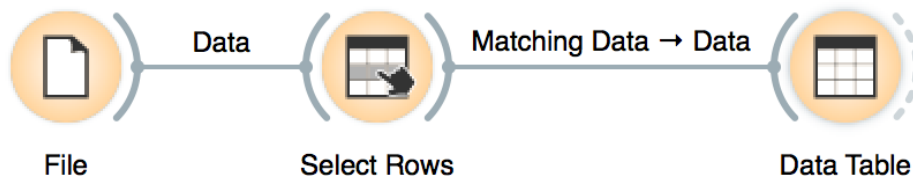- First column is assumed to be a category (in this case, species)

For each of the 150 flowers in the dataset, there is a value for:

- Petal Length
- Petal Width
- Sepal Length
- Sepal Width

# Select Rows Widget



- Filters data according to simple rules

- For example: exclude all irises with short petals

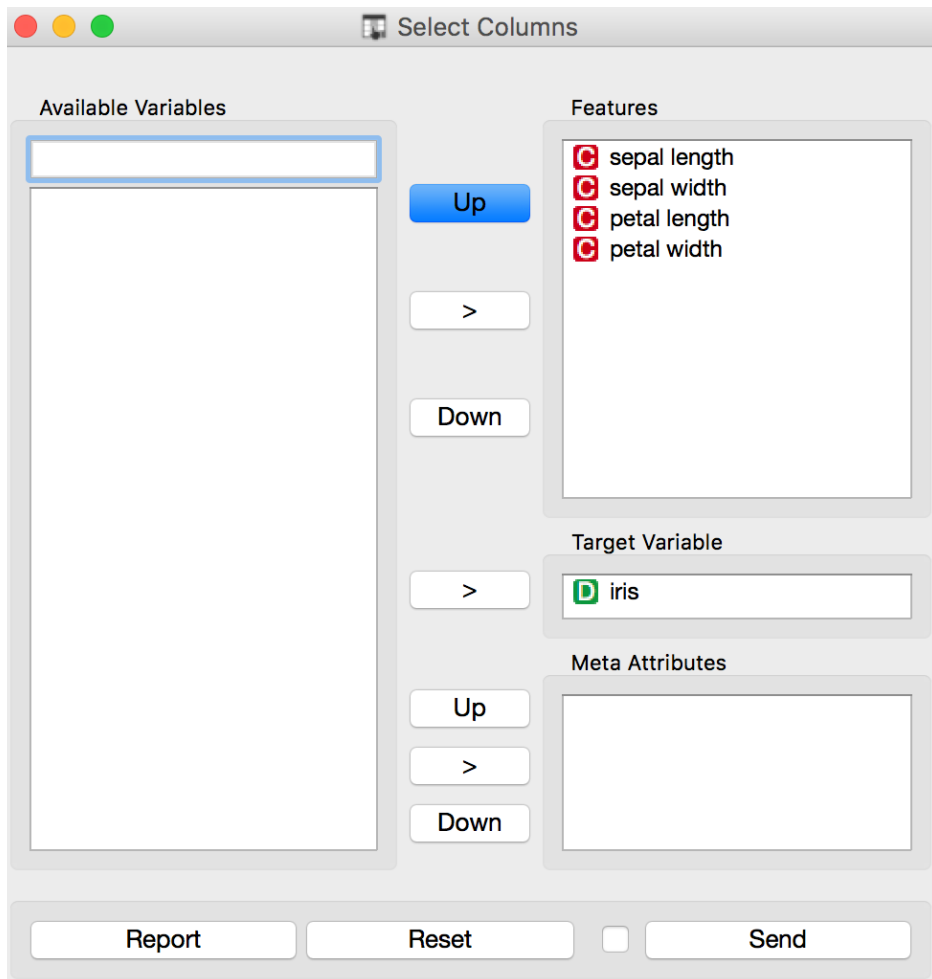- Select an attribute and a condition and press "Add" to add it to the filter

# Data Selection Results

| | iris | sepal length | sepal width | petal length | petal width |
|---|---|---|---|---|---|
| 1 | Iris-versicolor | 7.000 | 3.200 | 4.700 | 1.400 |
| 2 | Iris-versicolor | 6.400 | 3.200 | 4.500 | 1.500 |
| 3 | Iris-versicolor | 6.900 | 3.100 | 4.900 | 1.500 |
| 4 | Iris-versicolor | 5.500 | 2.300 | 4.000 | 1.300 |
| 5 | Iris-versicolor | 6.500 | 2.800 | 4.600 | 1.500 |
| 6 | Iris-versicolor | 5.700 | 2.800 | 4.500 | 1.300 |
| 7 | Iris-versicolor | 6.300 | 3.300 | 4.700 | 1.600 |
| 8 | Iris-versicolor | 4.900 | 2.400 | 3.300 | 1.000 |
| 9 | Iris-versicolor | 6.600 | 2.900 | 4.600 | 1.300 |
| 10 | Iris-versicolor | 5.200 | 2.700 | 3.900 | 1.400 |
| 11 | Iris-versicolor | 5.000 | 2.000 | 3.500 | 1.000 |
| 12 | Iris-versicolor | 5.900 | 3.000 | 4.200 | 1.500 |
| 13 | Iris-versicolor | 6.000 | 2.200 | 4.000 | 1.000 |
| 14 | Iris-versicolor | 6.100 | 2.900 | 4.700 | 1.400 |
| 15 | Iris-versicolor | 5.600 | 2.900 | 3.600 | 1.300 |
| 16 | Iris-versicolor | 6.700 | 3.100 | 4.400 | 1.400 |
| 17 | Iris-versicolor | 5.600 | 3.000 | 4.500 | 1.500 |
| 18 | Iris-versicolor | 5.800 | 2.700 | 4.100 | 1.000 |
| 19 | Iris-versicolor | 6.200 | 2.200 | 4.500 | 1.500 |
| 20 | Iris-versicolor | 5.600 | 2.500 | 3.900 | 1.100 |
| 21 | Iris-versicolor | 5.900 | 3.200 | 4.800 | 1.800 |

- The "petal length" column now only contains values longer than 3 cm

- The blue category, iris-setosa, is now completely absent.

- Apparently all iris-setosa flowers have petals shorter than 3 cm.
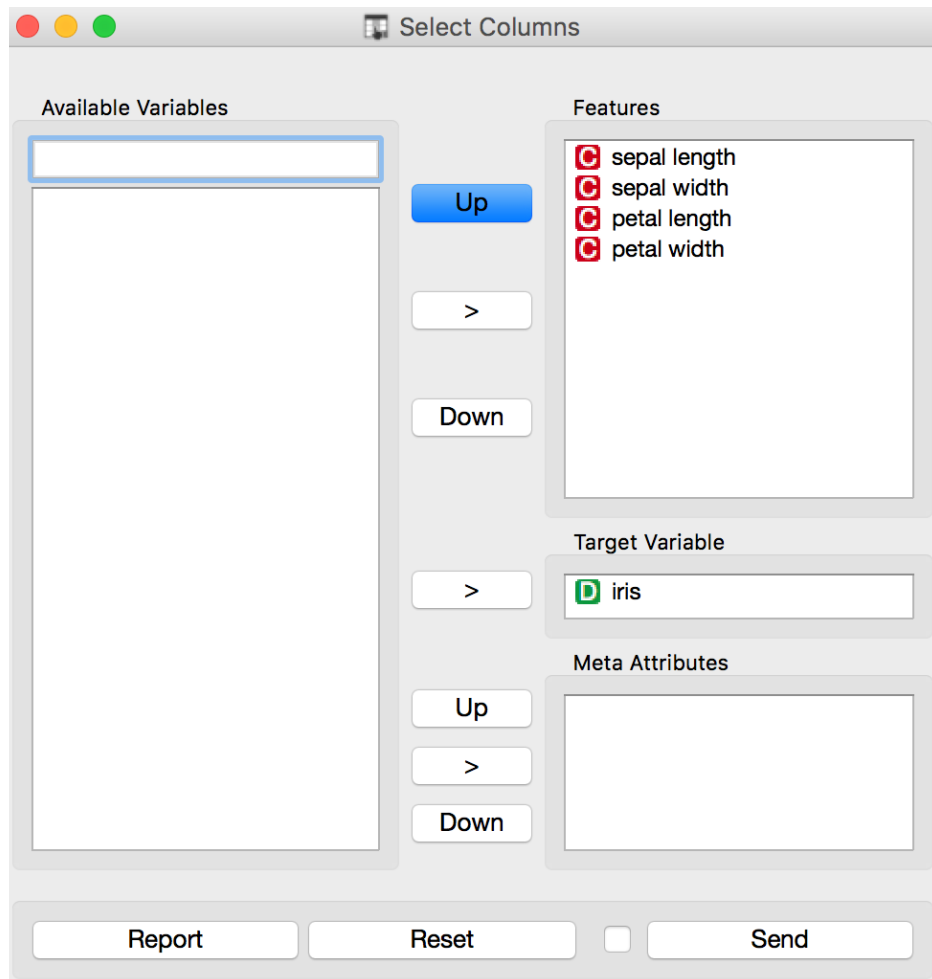
# Select Columns Widget (1)



- Choose which columns go in the dataset
  - "Attributes" are data values to be included in output
  - "Class" is the category of the row
  - "Meta Attributes" are descriptive attributes that are excluded from the analysis (such as a row ID)
  - "Available Attributes" are attributes available to be loaded, but ignored
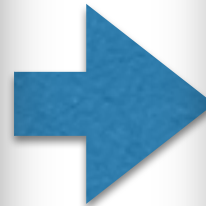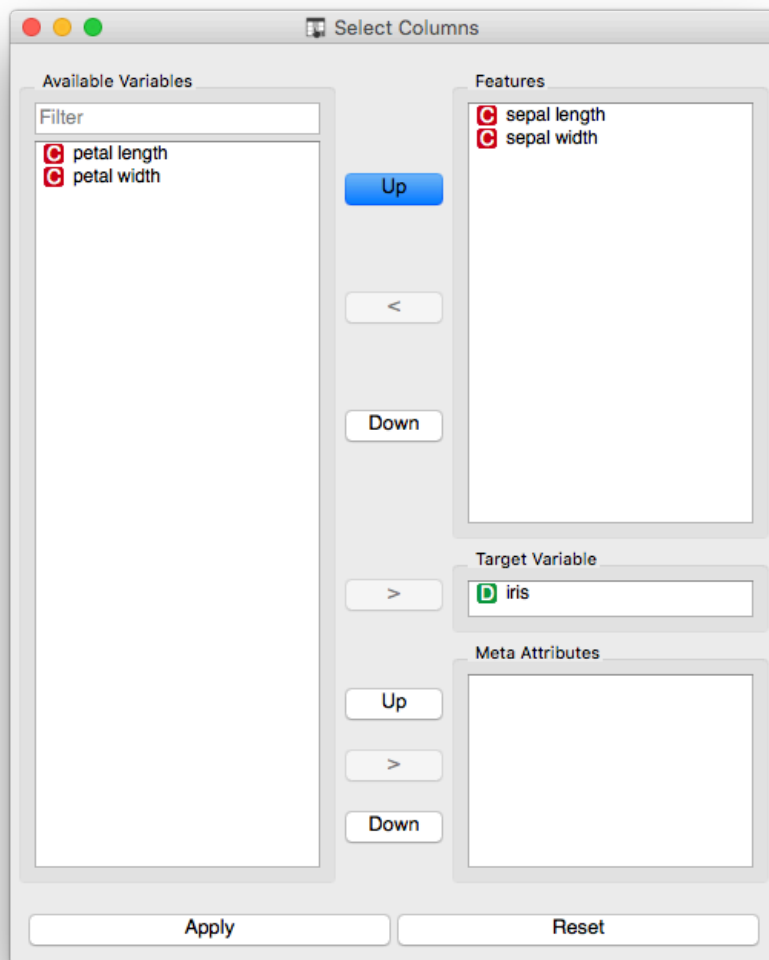
# Select Columns Widget (1)



- Drag or move variables between categories with the ">" and "<" buttons

- Each variable is marked "C" for continuous (numerical values) or "D" for discrete (categorical values)

- You may need to click "Apply" before any changes you make take effect

# Select Columns in action

- Suppose we were only interested in sepals, not petals.

# Feature Constructor Widget



- Defines new attributes (i.e. columns) based on the values of existing attributes
  - Type a formula and click "Add" to add a new feature
  - Select fields using "(all attributes)" and "(all functions)"
- Widget outputs the same data set with new attributes added
- This particular calculation is assuming petals are triangular

# Feature Construction Results

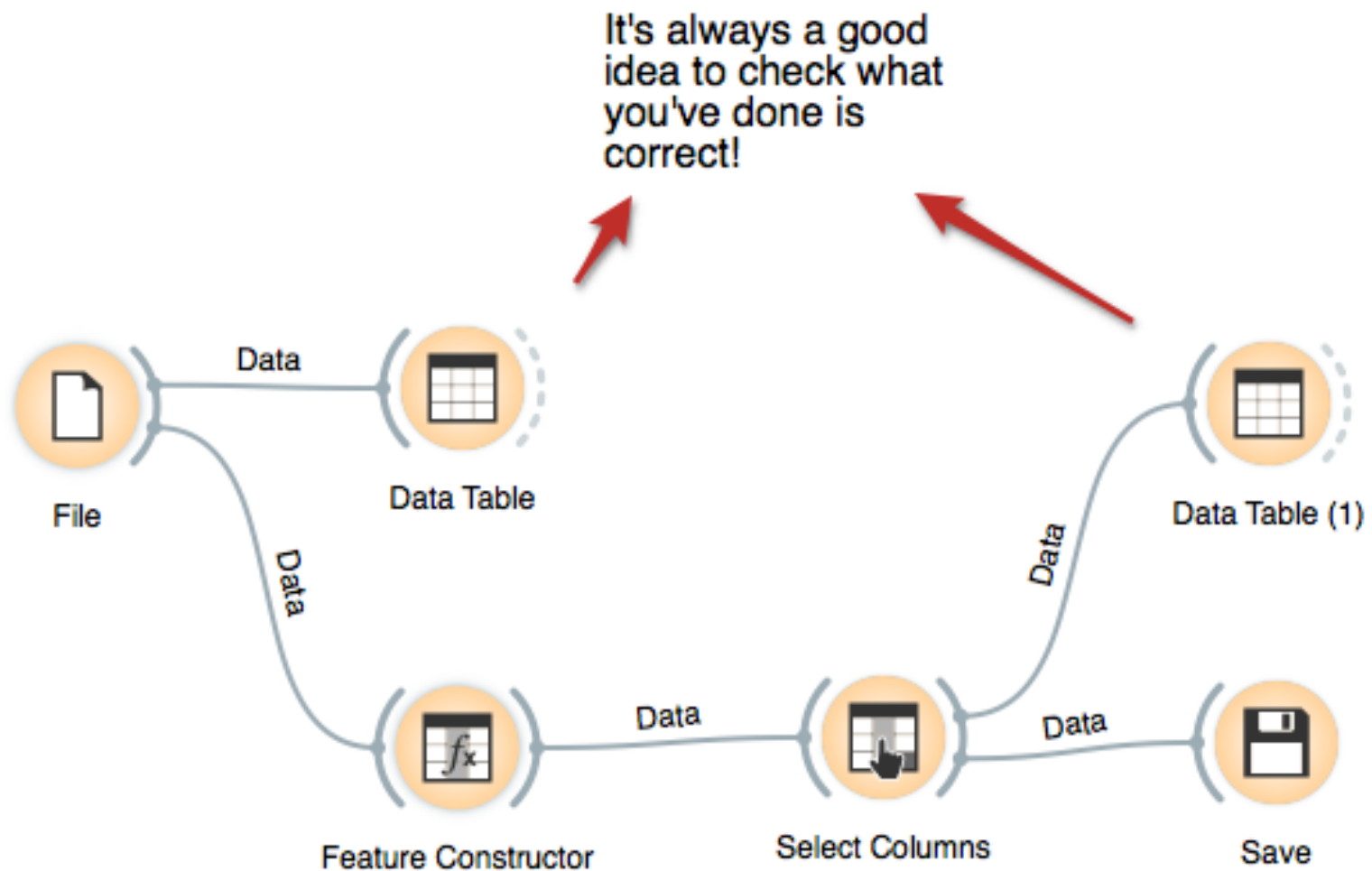- New attribute is added after existing attributes but before class

# Save Widget



- Save a modified file
- Saves whatever is going to its input
  - If you made changes elsewhere in the scheme, they will not be saved
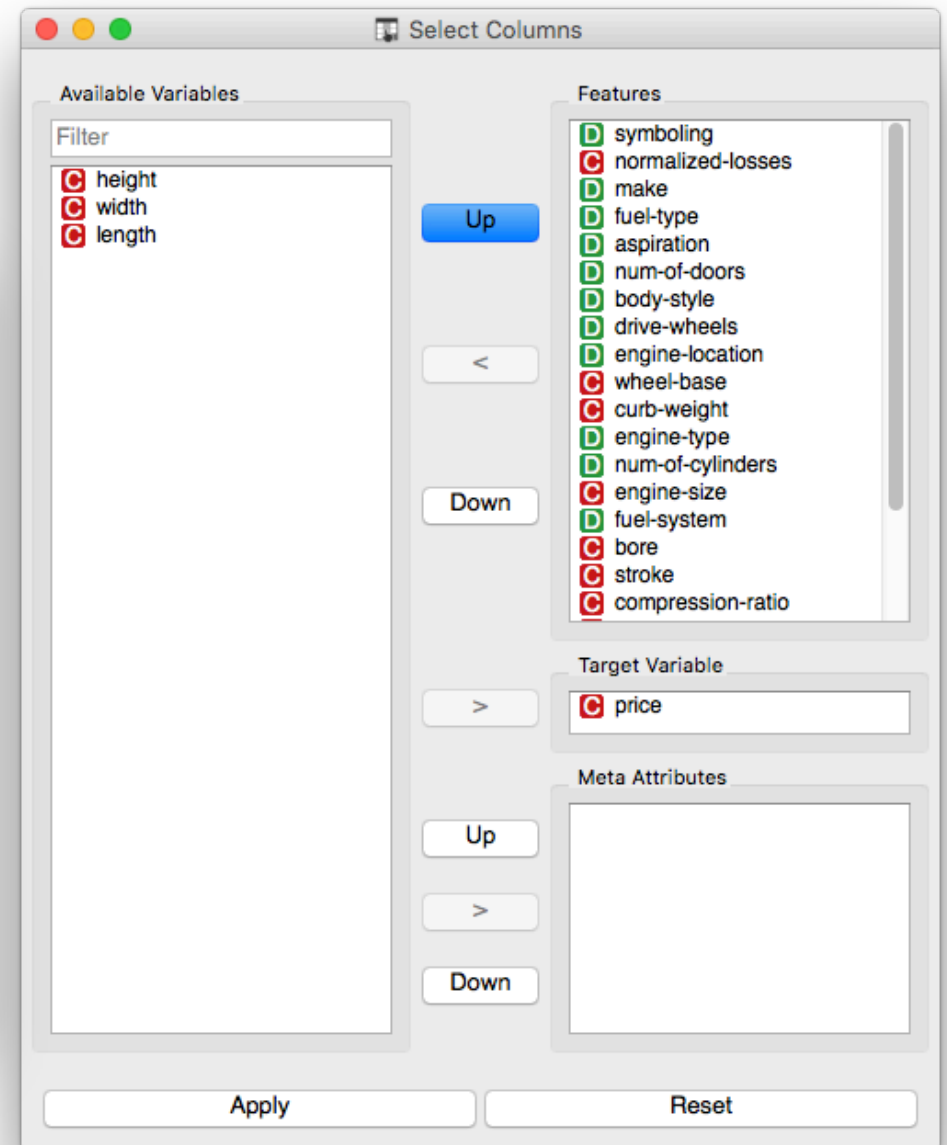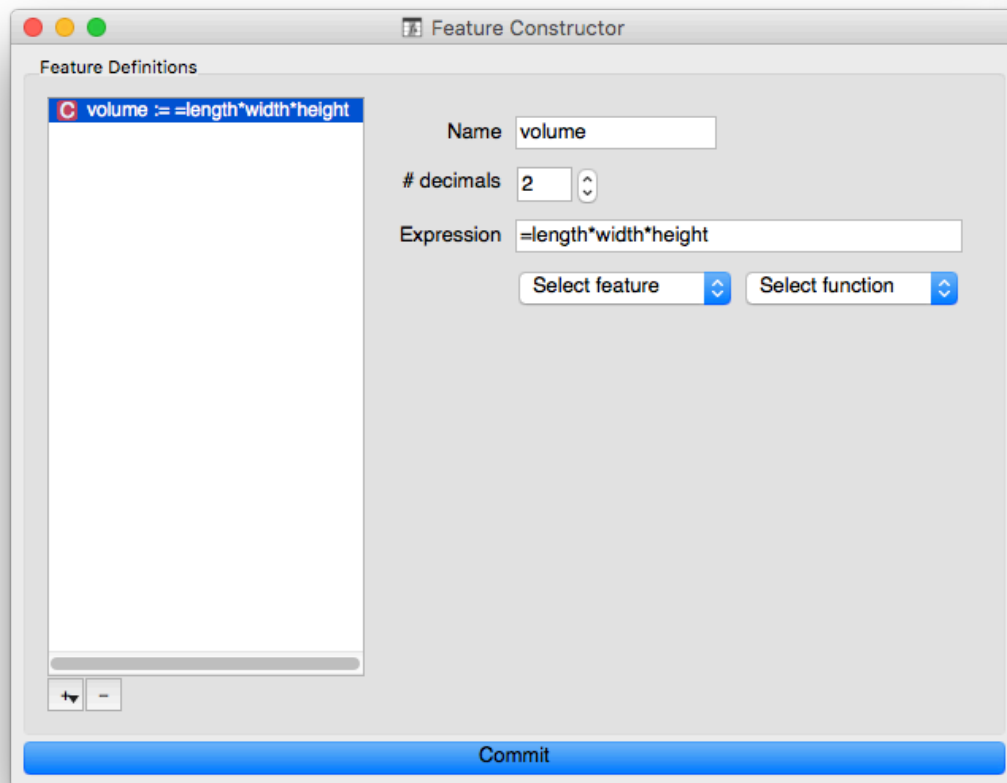- Be careful not to accidentally overwrite your input file

# Exercise: First Scheme

- Load and inspect the *imports-85.tab* data file (on course website), which contains information about various imported cars

- Add a "volume" attribute (i.e. length x width x height)

- Remove the original length, width, and height attributes

- Save the dataset using a different filename

# Solution

# Solution, continued



Remember to click "Apply" after you make changes!