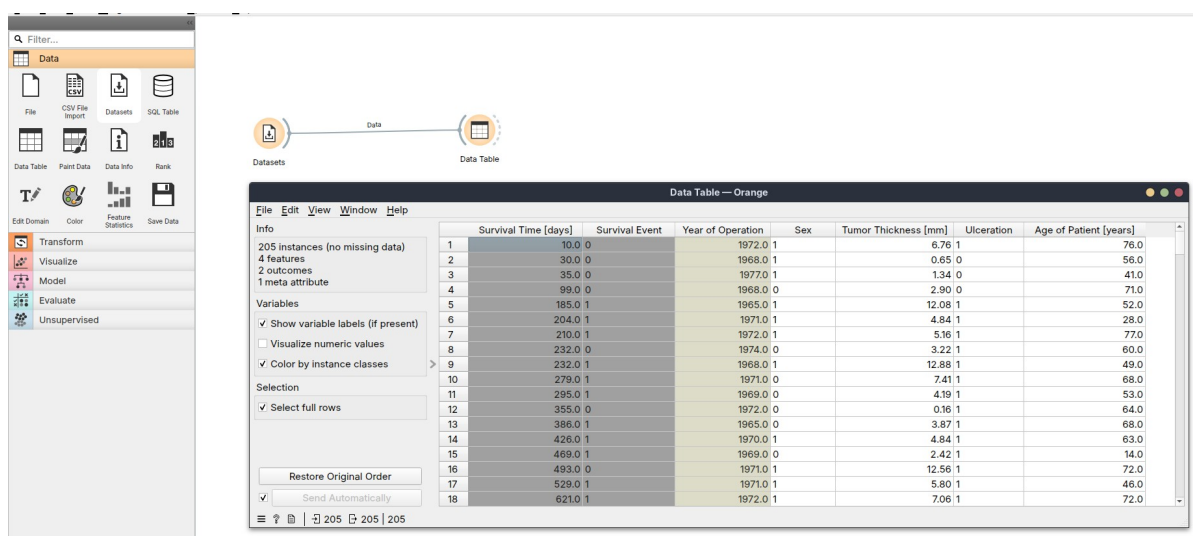# No code lab

## Data visualization

1. From the [Data] section select the **Datasets** operator. Double-click it and select the *Melanoma: Survival from Melignant Melanoma* dataset. Double-click on the name of the dataset so that a green dot appears on the left side of the dataset (indicating that the dataset has been succesfully downloaded).
2. Put the **Data Table** operator in the panel. Drag the data connection between **Datasets** and **Data Table** operators. Then, double-click on the **Data Table** operator to see the data.
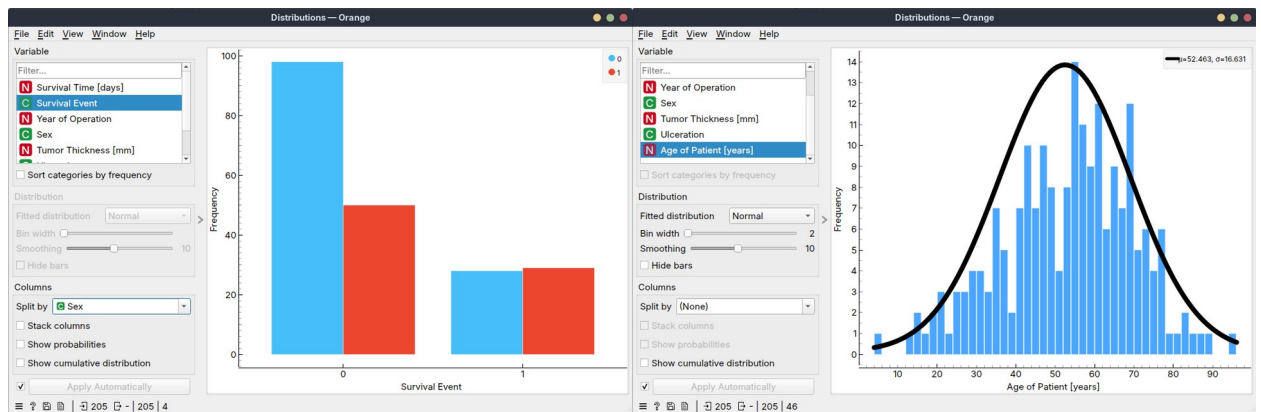


As you can see, the dataset presents the on patients being treated from melanoma, with attributes representing the year of operation, the sex of a patient (0=female, 1=male), the size of the tumor and its potential ulceration, the age of a patient, and two special attributes representing the survival time and the survival event (outcome of the operation, 0=survived, 1=died).

3. From the [Visualize] section drag the **Distributions** operator and send the data to the operator. To do this, simply drag a line from the output of the **Datasets** operator (the grey arc) to the **Distributions** operator. Open the operator and compare the display of continuous (numerical) attributes with the display of categorical (discrete) attributes.
    1. Display the *Survival Event*, are patients are more likely to survive the operation?
    2. Split the *Survival Event* by *Sex* (look for "Split by" option in the lower left part of the window). Are women more likely to survive the operation then men?
    3. Display the *Age of Patient* and see what happens when you slide the "Bin width" slider. Are there more patients in the age range between 40-50 or 50-60?
    4. Fit the normal distribution to the display (look for "Fitted distribution" dropdown list). What is the mean age of patients?

4. From the [Visualize] section drag the **Scatter Plot** operator. Drag the mouse from the **Datasets** to the **Scatter Plot** operator. Try to answer the following questions:
    1. How does age of a patient and survival time relate to the probability of surviving the melanoma? *Hint*: plot *Age of Patient* on X-axis, plot *Survival Time* on Y-axis, and use *Survival Event* to set the color and shape of points.
    2. Is there a difference between men and women when it comes to surviving the illness with respect to the age of the patient? *Hint*: plot *Age of Patient* on X-axis, plot *Sex* on Y-axis, and use *Survival Event* to set the color and shape of points. Since *Sex* is a discrete attribute, in order to make the figure more readable you will have to use the jittering option. You may also want to see what happens when you check the "Show color regions" checkbox.
    3. How do ulceration and tumor thickness influence the survival of the disease? *Hint*: plot *Tumor Thickness* on X-axis, plot *Ulceration* on Y-axis, and use *Survival Event* to set the color and shape of points.
5. Click on the button "Find Informative Projections" in the upper left corner of the window. The tool will propose the combinations of features that will produce the most informative displays of the data.
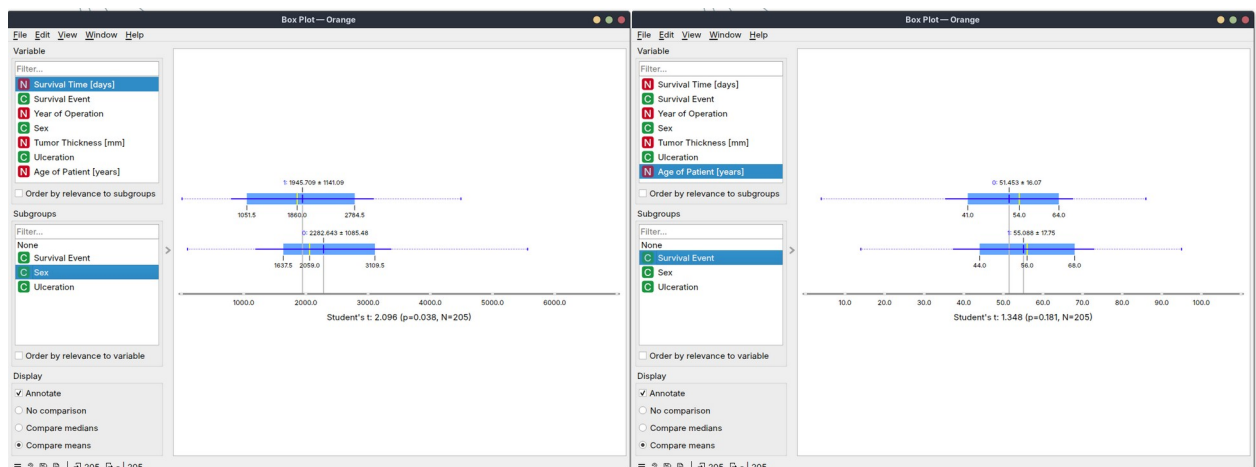


6. Send the data to the **RadVis** operator. Use color to distinguish between patients who survived and who did not. Use size of a data point to represent the tumor thickness. Display color regions to better understand, how these 6 features together influence the survival probability. Think of each patient as a point connected to features with rubber bands. The larger the value of a given attribute, the stronger the rubber band pulls in the direction of the attribute. Experiment with adding or removing features (simply click on a feature in the upper left part of the window). Do you find this type of visualization useful?

7. Drag the **Box Plot** operator to the panel and send the data to the operator. Box Plot is useful for visualizing simple data statistics. The blue line in the middle represents the mean, the values at the ends of the blue shaded area are the values of the 1st and 3rd quartile (25% and 50% of data), the yellow blue line represents the median, and the horizontal blue line is the standard deviation around the mean.

   Try to answer the following questions:

   1. Is there a significant difference in survival time between men and women? *Hint*: use *Survival Time [days]* as the main variable in the upper left part of the window. Mark *Sex* as the Subgroups variable in the lower left part of the window. Select the "Compare means" radiobox in the lower left part of the window and make sure that "Annotate" checkbox is checked. You can see that the average survival for men is 1945 days, and 2282 days for women. In addition, the difference is statistically significant, because the p-value of the Student's t-test is below 0.05.

   2. Is there a significant difference in patients age between patients who survived and those who did not? Hint: use *Age of Patient [years]* as the main variable, and *Survival Event* as the subgroup variable. Is there enough data to say that the difference of two years (51.6 vs 53.9) is statistically significant?



## Assignment

Use the **Datasets** operator to download the *German Credit Data*. Using visualization tools try to answer the following questions:

- Does age influence the credit score rating of customers?
- Is there a relationship between sex/marital status and credit score rating?
- What can be said about credit score rating from the credit history and status of existing checking account?