

Data Analysis



Exploratory Data Analysis (EDA) Orange Data Mining



Edgar Rios Linares

Problem

▼ Problem description

Context

The original dataset contains 1000 records with 20 categorical attributes prepared by Prof. Hofmann.

In this dataset, each record represents a person taking out a loan from a bank.

Each person is classified as **good** or **bad** credit risk based on the set of attributes.

Link to the original dataset at

[UCI Machine Learning](#)

Content

The selected attributes:

- **Age** (numeric)
- **Sex** (text: male, female)
- **Job** (numeric: 0 - unskilled and non-resident, 1 - unskilled and resident, 2 - skilled, 3 - highly skilled)
- **Housing** (text: own, rent, or free)
- **Saving accounts** (text - little, moderate, quite rich, rich)
- **Checking account** (numeric, in DM - Deutsch Mark)
- **Credit amount** (numeric, in DM)
- **Duration** (numeric, in month)
- **Purpose**(text: car, furniture/equipment, radio/TV, domestic appliances, repairs, education, business, vacation/others)
- **Risk** (Value target - Good or Bad Risk)

Objective

Train a model to predict from new data whether a person applying for a loan represents a good or bad risk

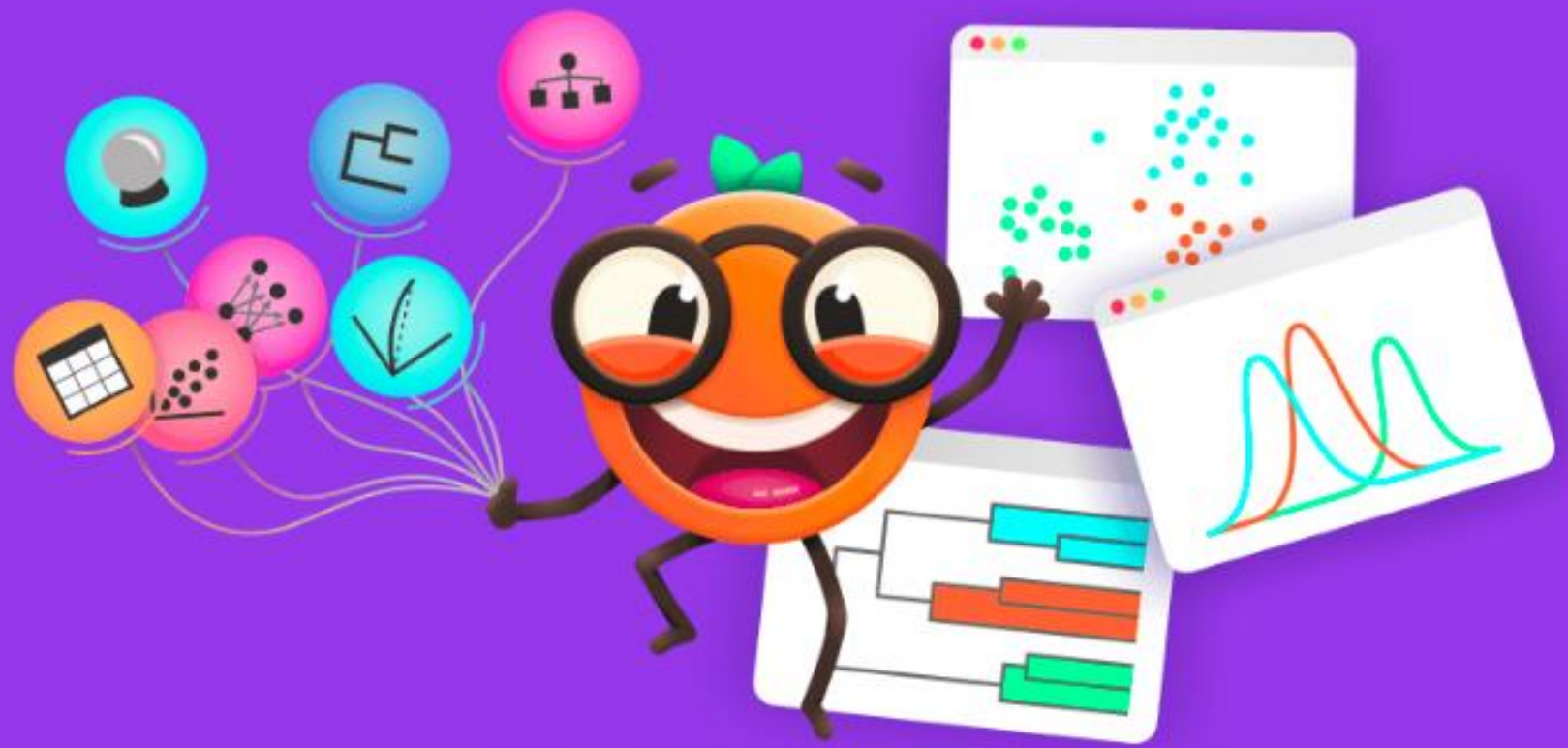
Download

Go website Orange, download & install app
<https://orangedatamining.com>

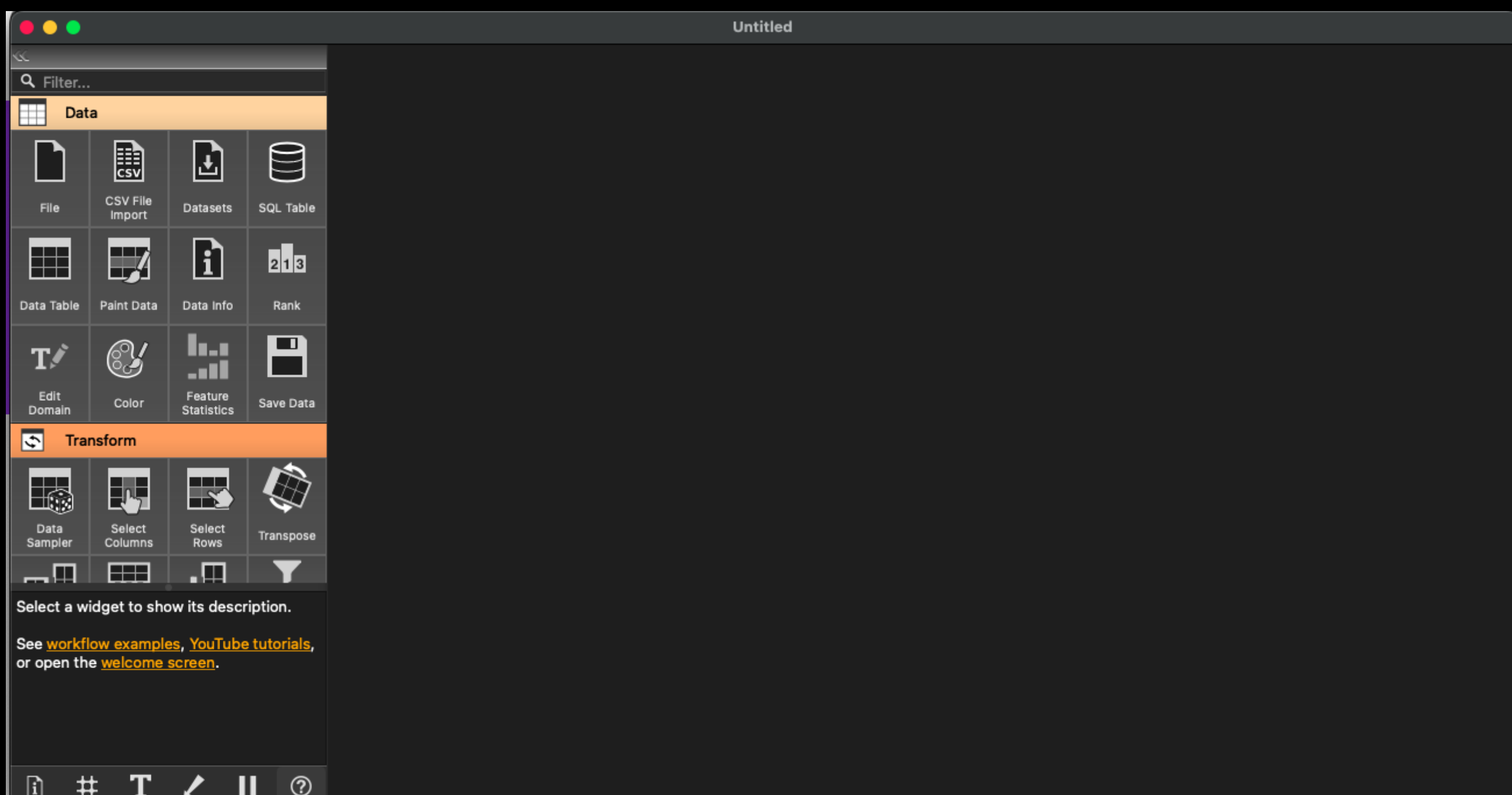
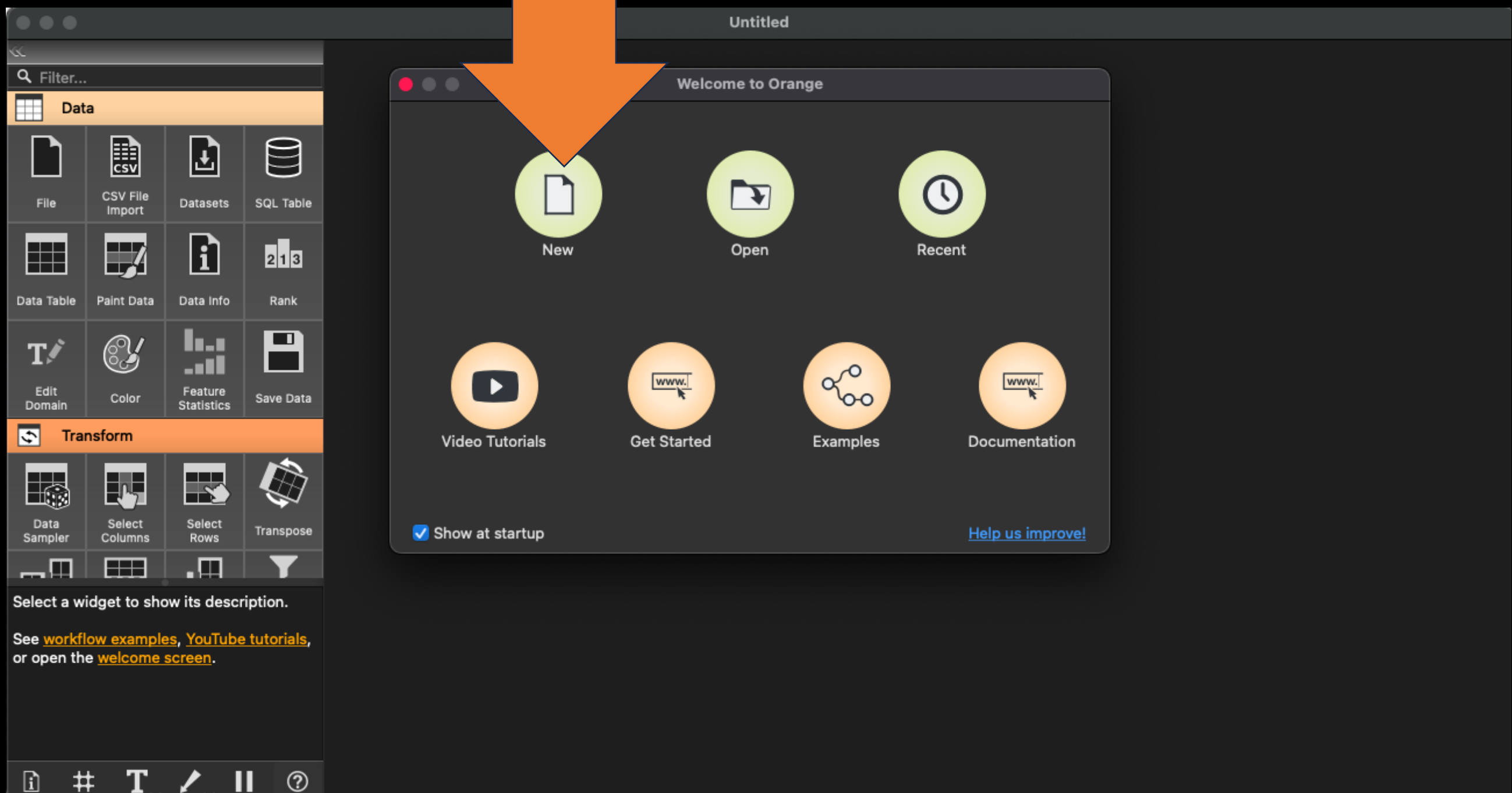
Data Mining Fruitful and Fun

Open source machine learning and data
visualization.

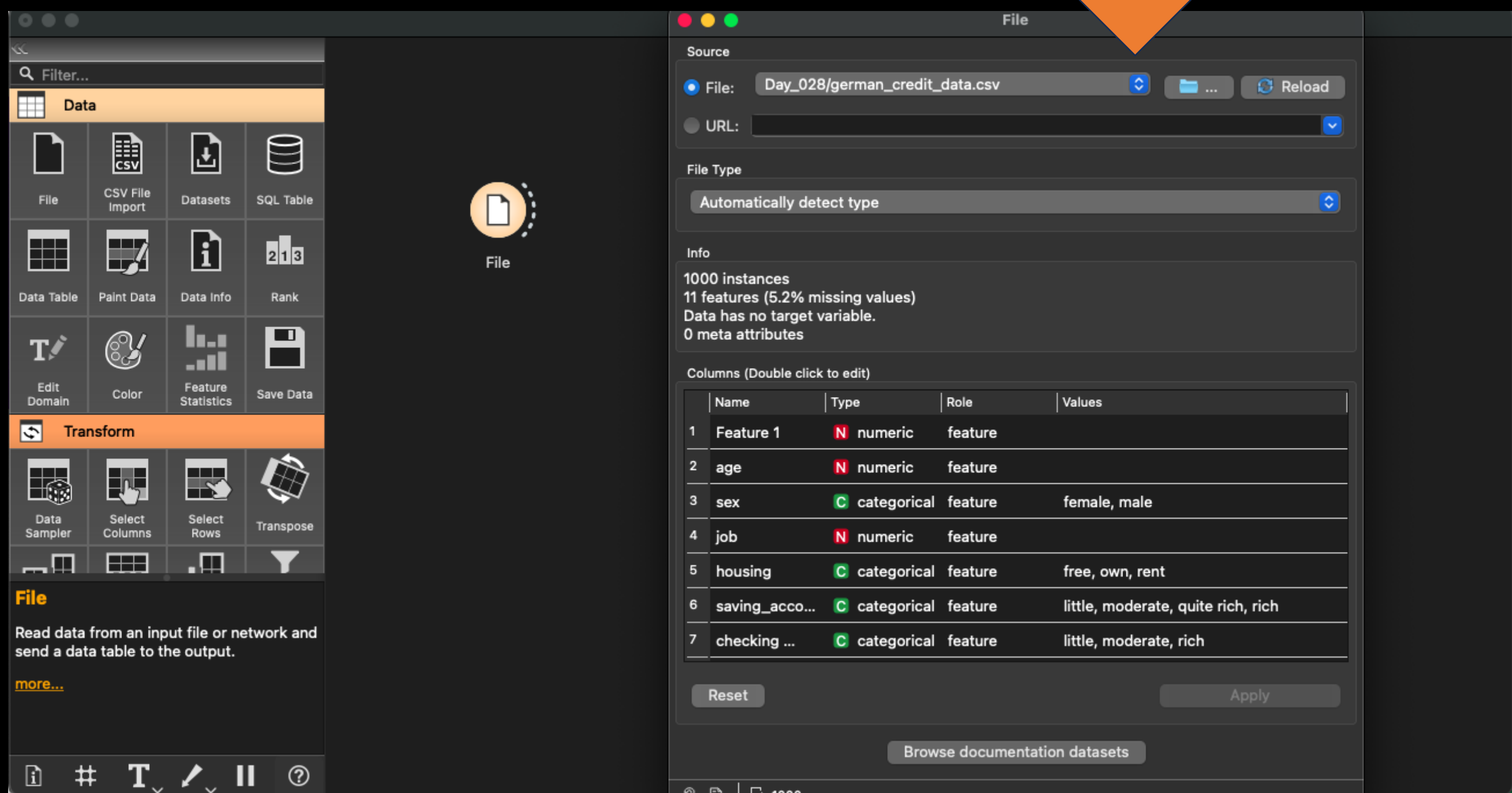
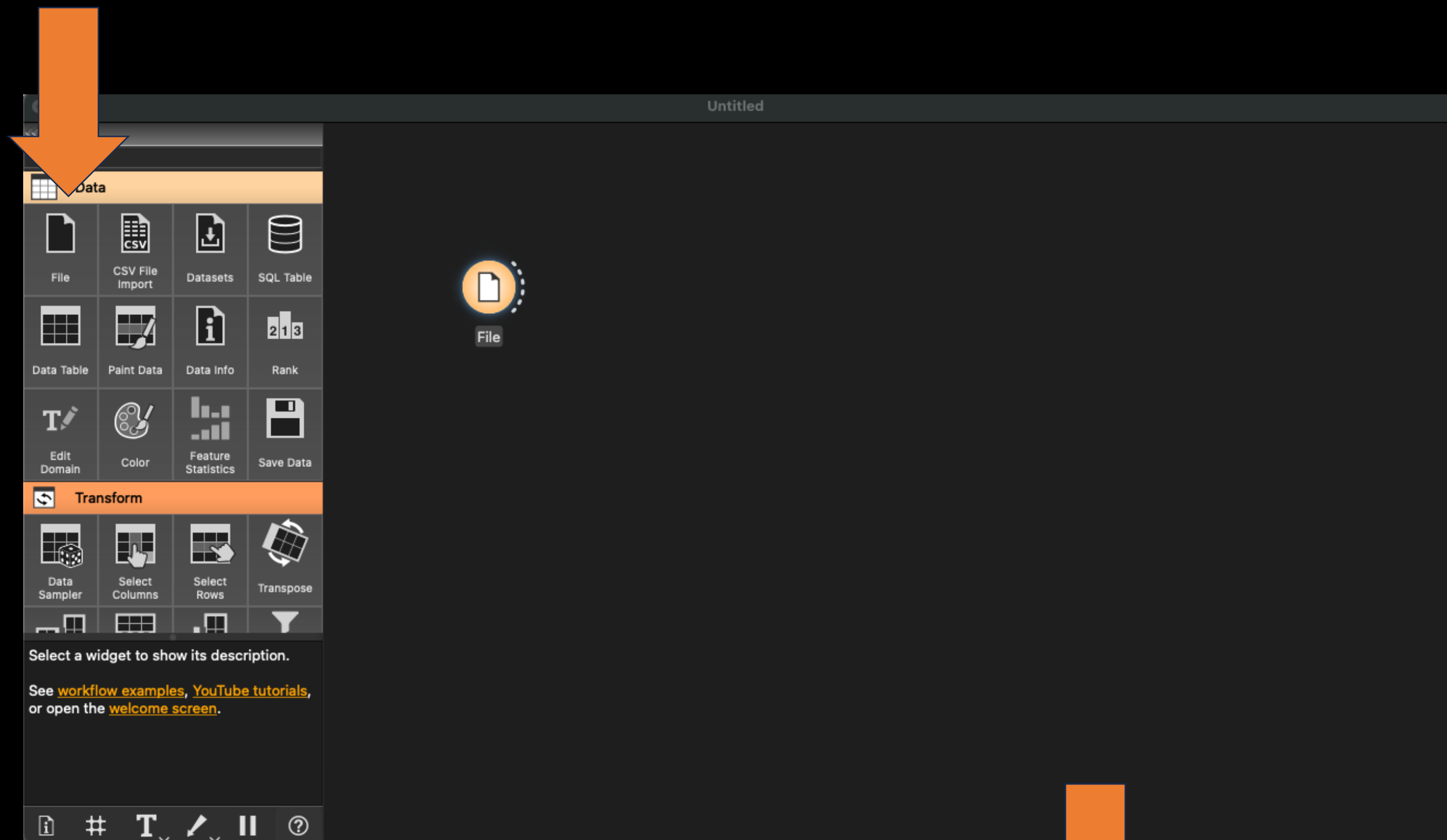
Download Orange 3.38.1



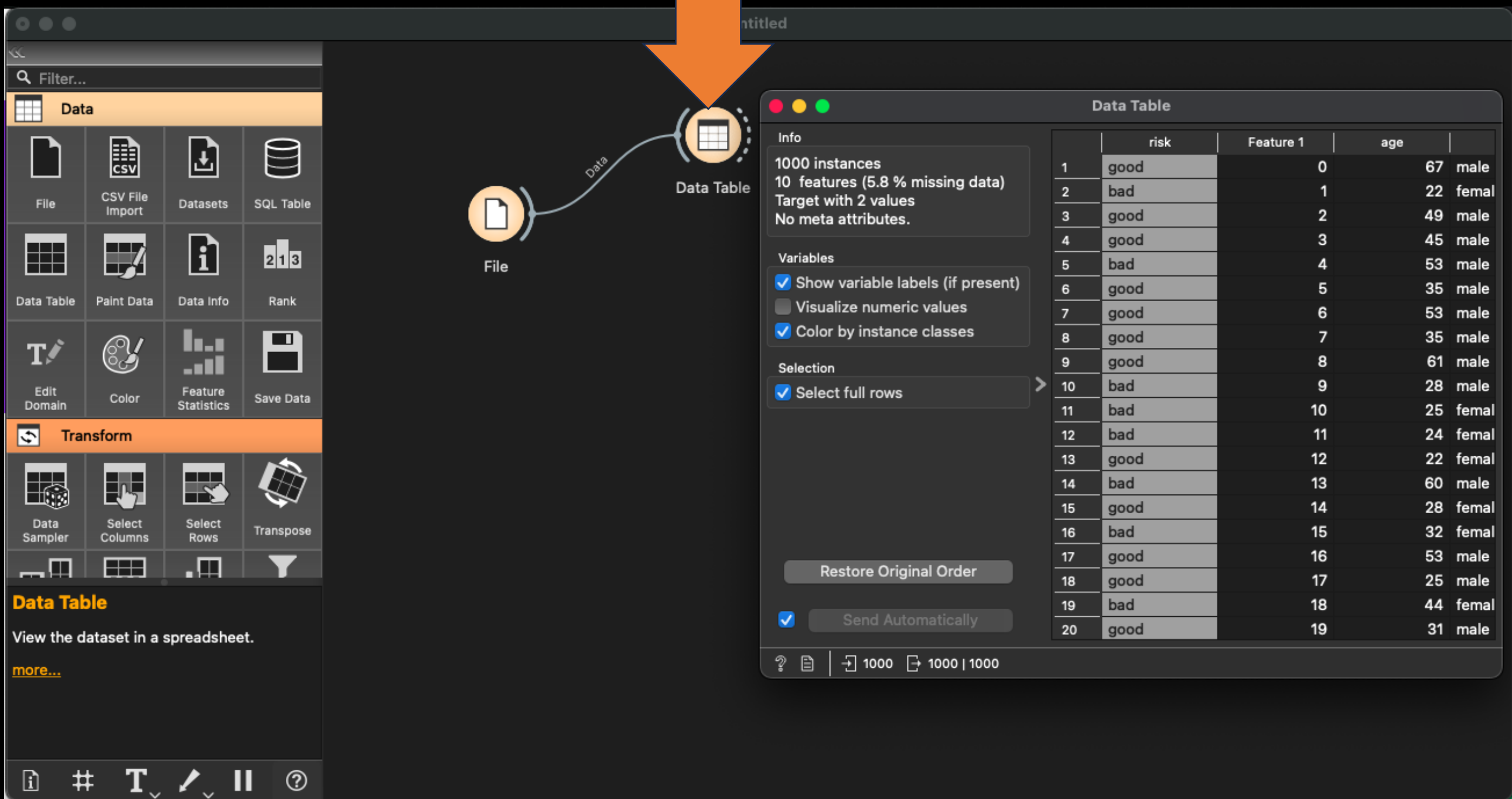
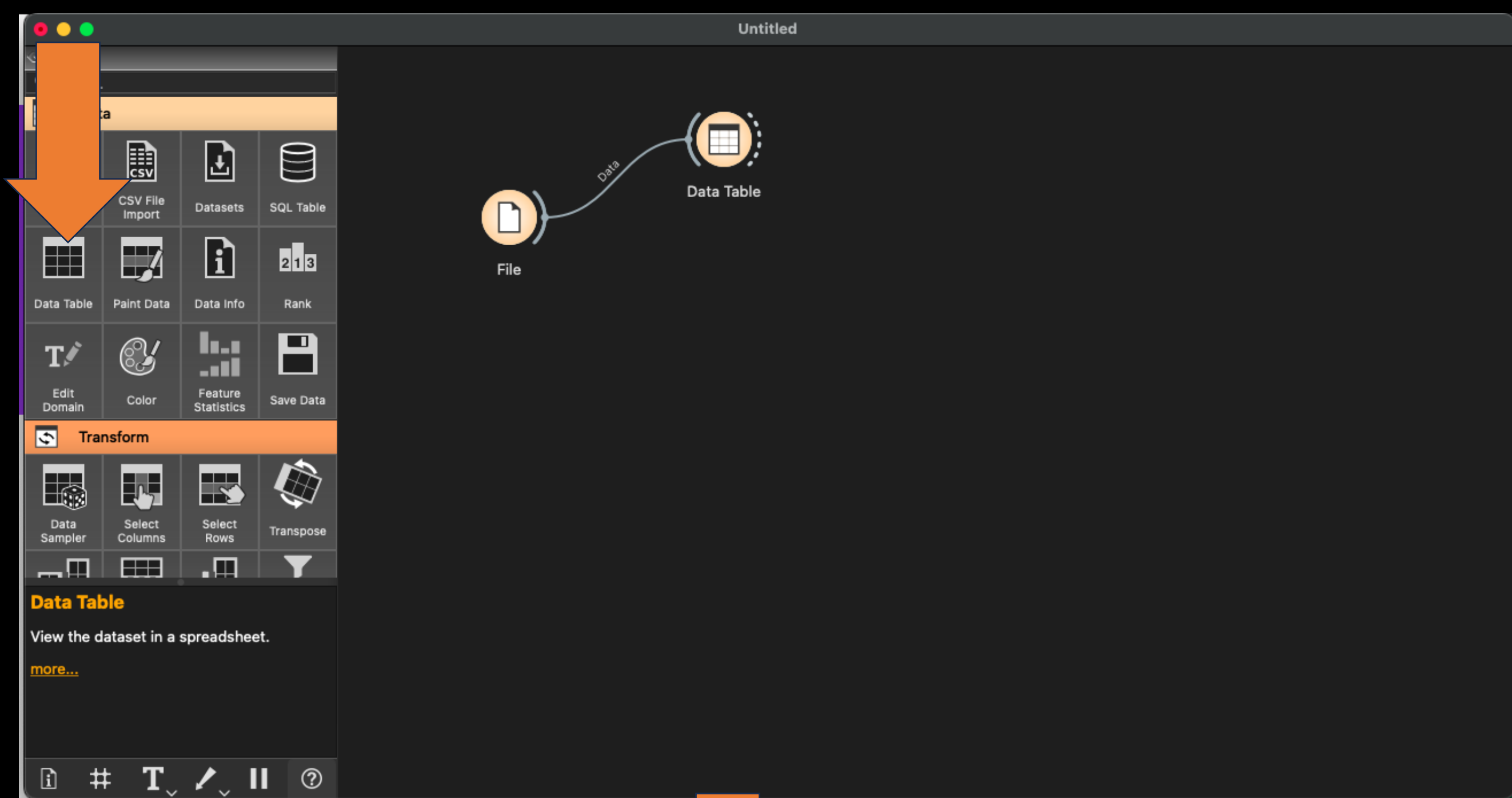
New Project



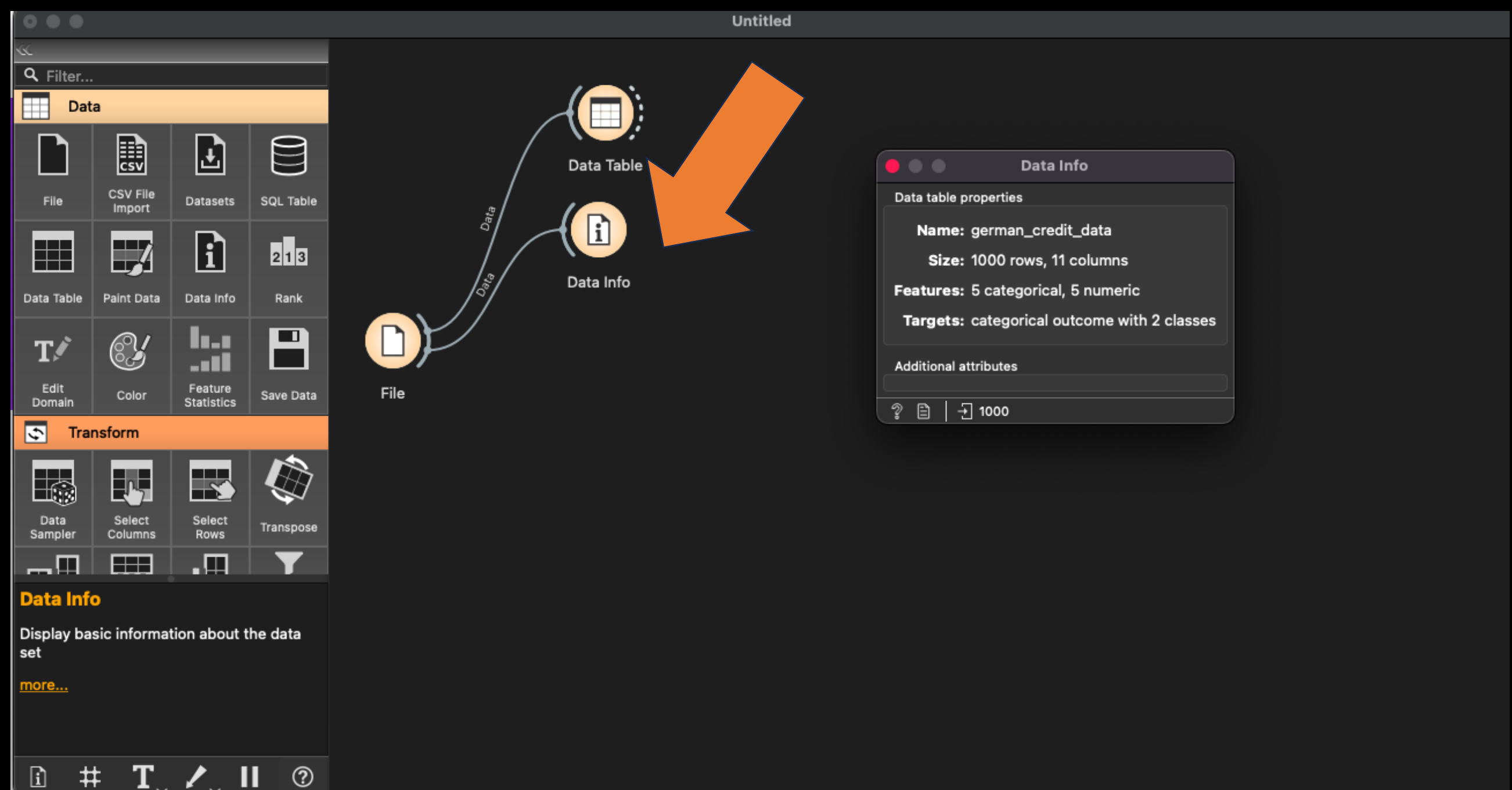
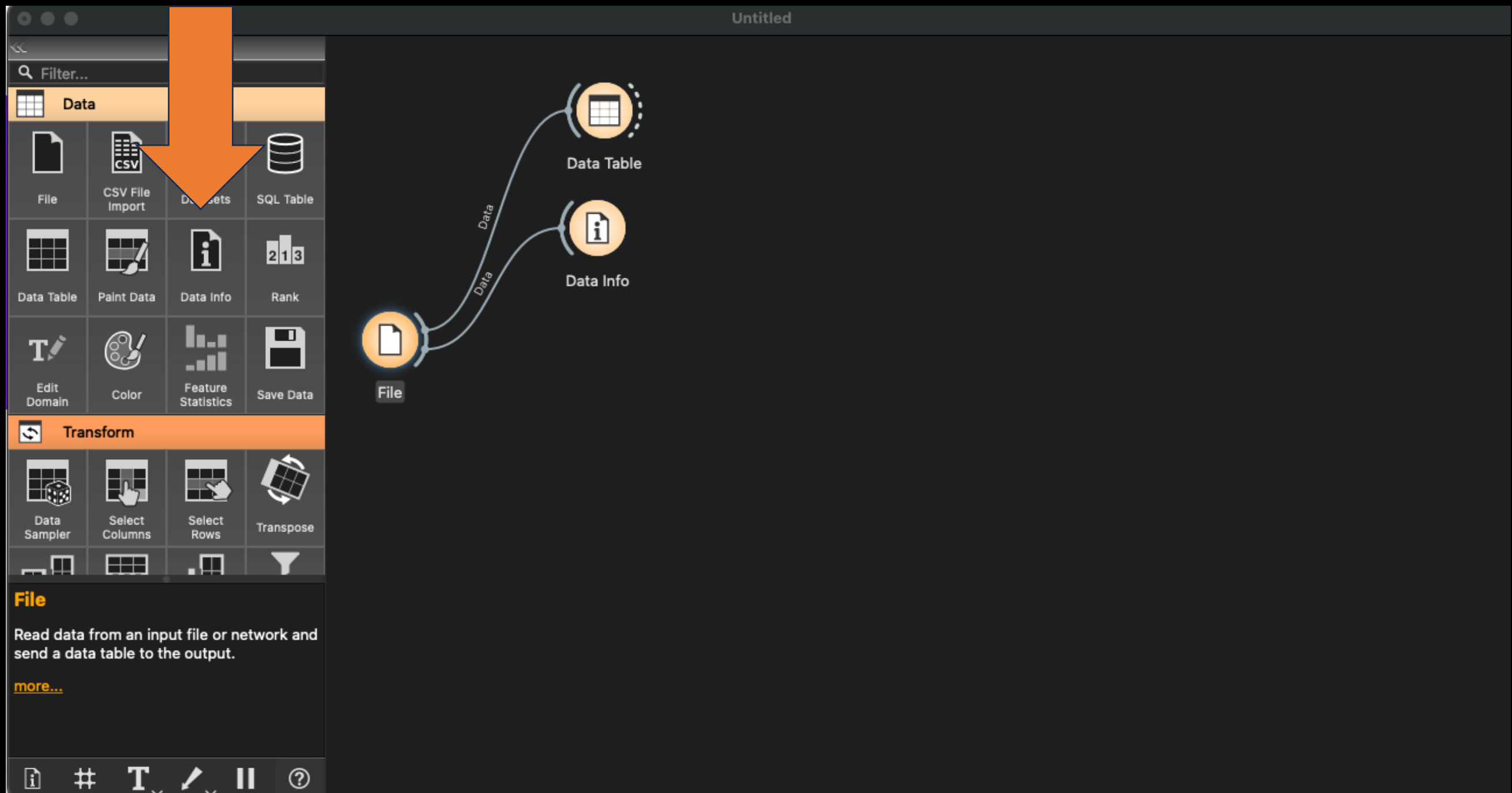
Load File



Data Table



Data Info



Feature Statistics

Filter...

Data

File

CSV File Import

SQL Table

Data Table

Paint Data

Rank

Edit Domain

Color

Feature Statistics

Save Data

Transform

Data Sampler

Select Columns

Select Rows

Transpose

Select a widget to show its description.
See [workflow examples](#), [YouTube tutorials](#), or open the [welcome screen](#).

Untitled

File

Data

Data Table

Data

Data Info

Data

Feature Statistics

Feature Statistics

Inputs:

- Data

Outputs:

- Reduced Data
- Statistics

Filter...

Data

File

CSV File Import

Datasets

SQL Table

Data Table

Paint Data

Data Info

Rank

Edit Domain

Color

Feature Statistics

Save Data

Transform

Data Sampler

Select Columns

Select Rows

Transpose

Feature Statistics

Show basic statistics for data features.
[more...](#)

Untitled

File

Data

Data Table

Data

Data Info

Data

Feature Statistics

Feature Statistics

	Name	Distribution	Mean	Mode	Median
N	Feature 1		499.50	0	499.50
N	age		35.55	27	35.55
N	job		1.90	2	1.90
N	credit amount		3271.26	1258	2319.50
N	duration		20.90	24	20.90

Color: None ☒ Send Automatically

1000

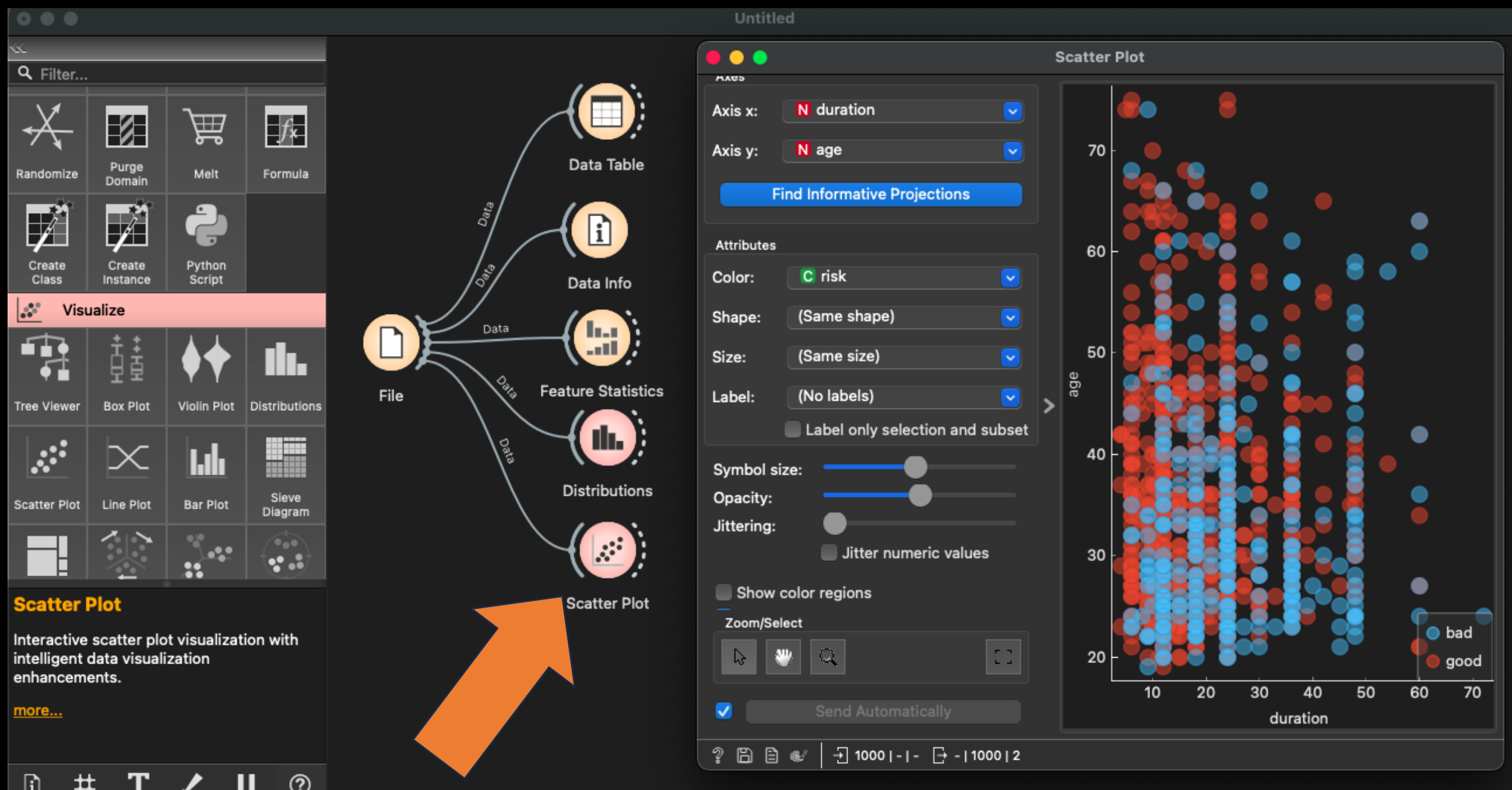
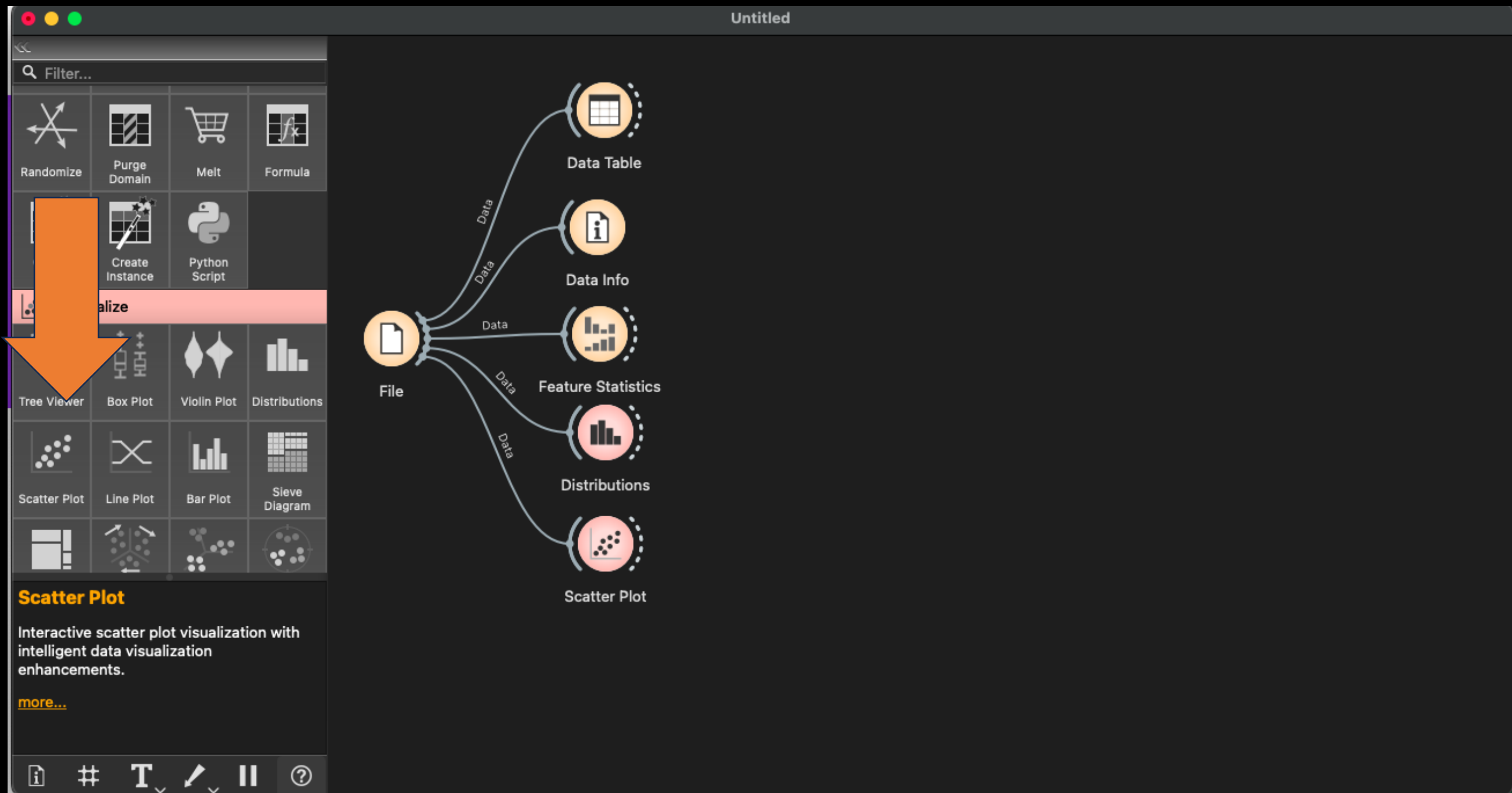
- | 11

Distributions

The screenshot shows the Orange3 software interface. On the left, the 'Visualize' widget panel is open, displaying various visualization options. A large orange arrow points to the 'Distributions' widget in this panel. Below the panel, a description for 'Distributions' is visible: 'Display value distributions of a data feature in a graph.' and a link to 'more...'. In the center, a workflow diagram shows a 'File' widget connected to four 'Data' widgets, which are then connected to 'Data Table', 'Data Info', 'Feature Statistics', and 'Distributions' widgets. The 'Distributions' widget is highlighted with a red border.

This screenshot shows the 'Distributions' widget configuration and its output. The 'Variable' list on the left includes 'checking account', 'credit amount', 'duration', and 'purpose'. The 'credit amount' variable is selected. The 'Distribution' section shows 'Fitted distribution' set to 'None', 'Bin width' set to 500, and 'Smoothing' set to 10. The 'Columns' section shows 'Split by' set to 'risk'. The 'Apply Automatically' checkbox is checked. The output is a histogram showing the frequency of 'credit amount' for two categories: 'bad' (blue bars) and 'good' (red bars). The x-axis is labeled 'credit amount' and ranges from 0 to 10000. The y-axis is labeled 'Frequency' and ranges from 0 to 140. The 'good' category shows a much higher frequency than the 'bad' category.

Scatter Plot



Correlations

The screenshot shows the Orange3 software interface. On the left, the 'Unsupervised' widget panel is visible, with an orange arrow pointing to the 'Correlations' widget. In the center, a workflow is shown in the 'Untitled' canvas. A 'File' widget is connected to several data analysis widgets: 'Data Table', 'Data Info', 'Feature Statistics', 'Distributions', 'Scatter Plot', and 'Correlations'. The 'Correlations' widget is highlighted with an orange arrow.

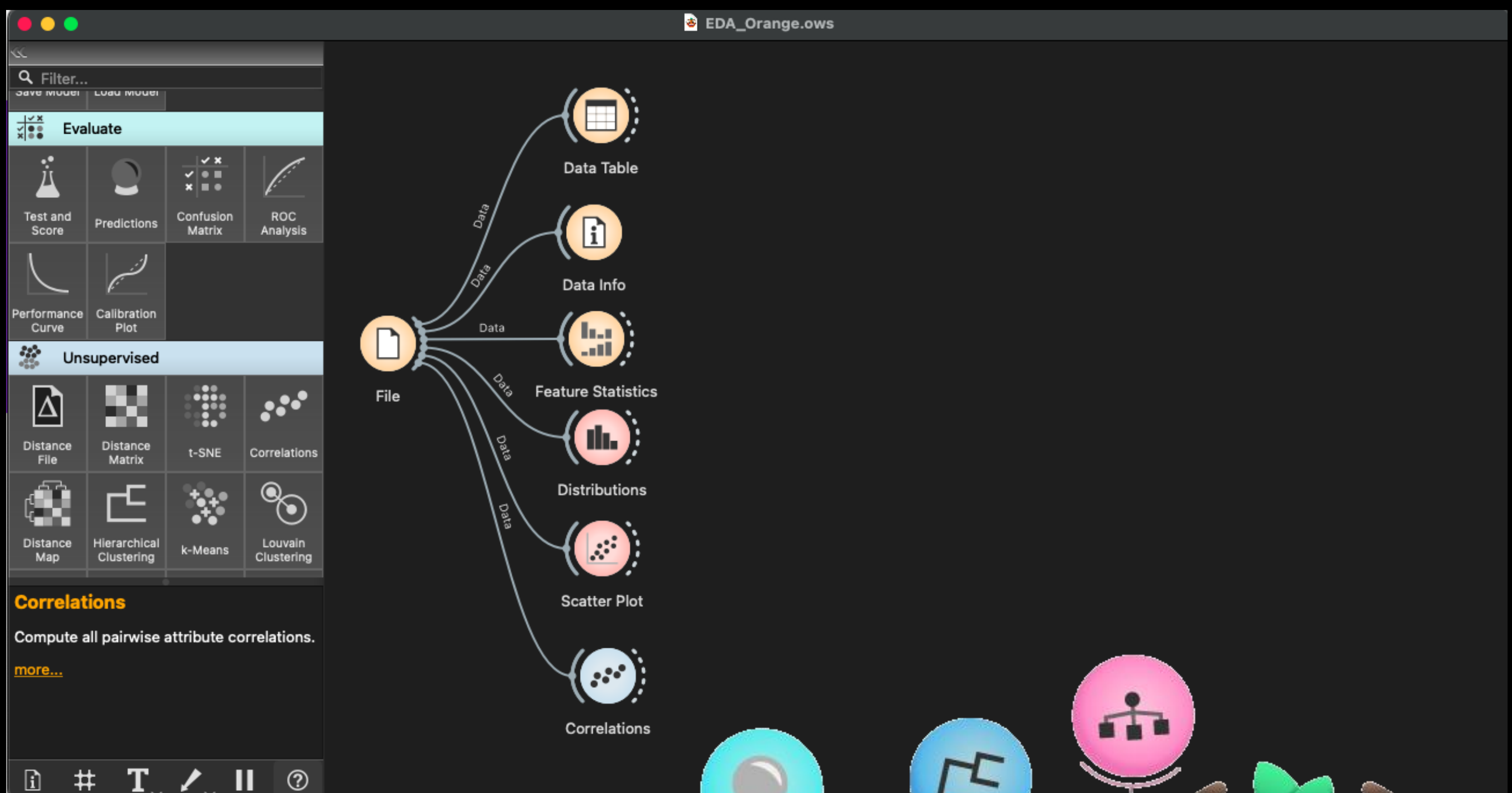
The screenshot shows the 'Correlations' widget configuration window. The 'Pearson correlation' method is selected. The 'Filter' dropdown is set to '(All combinations)'. The output table shows the following data:

		credit amount	duration
1	+0.625	credit amount	duration
2	+0.285	credit amount	job
3	+0.211	duration	job
4	-0.036	age	duration
5	+0.033	age	credit amount
6	+0.031	Feature 1	duration
7	-0.027	Feature 1	job
8	+0.016	age	job
9	+0.013	Feature 1	credit amount
10	-0.010	Feature 1	age

The 'Correlations' widget is highlighted with an orange arrow in the workflow canvas.

The work is done, great job!"

You have a basic EDA quickly



Educator in AI

**Artificial
Intelligence**

Data Engineering



Machine Learning

Data Science

📌 **Linkedin** —> <https://www.linkedin.com/in/erlinares/>

👋 **Follow us on X**: <https://x.com/erlinares>^[SEP]

💻 **GitHub**: https://github.com/erlinares/365_AI_Journey/

💬 **Discord**: <https://discord.gg/5fFM2zh8>



Edgar Rios Linares