

Topic 2

Data Preprocessing

+ Low-code AI tool

IT8302 APPLIED MACHINE LEARNING

Learning Outcomes

- ❑ Understand principles of data representation and feature engineering
 - Understanding representation of categorical variables
 - Explain binning, discretization
 - Understand feature selection
 - Understand how to utilize expert knowledge

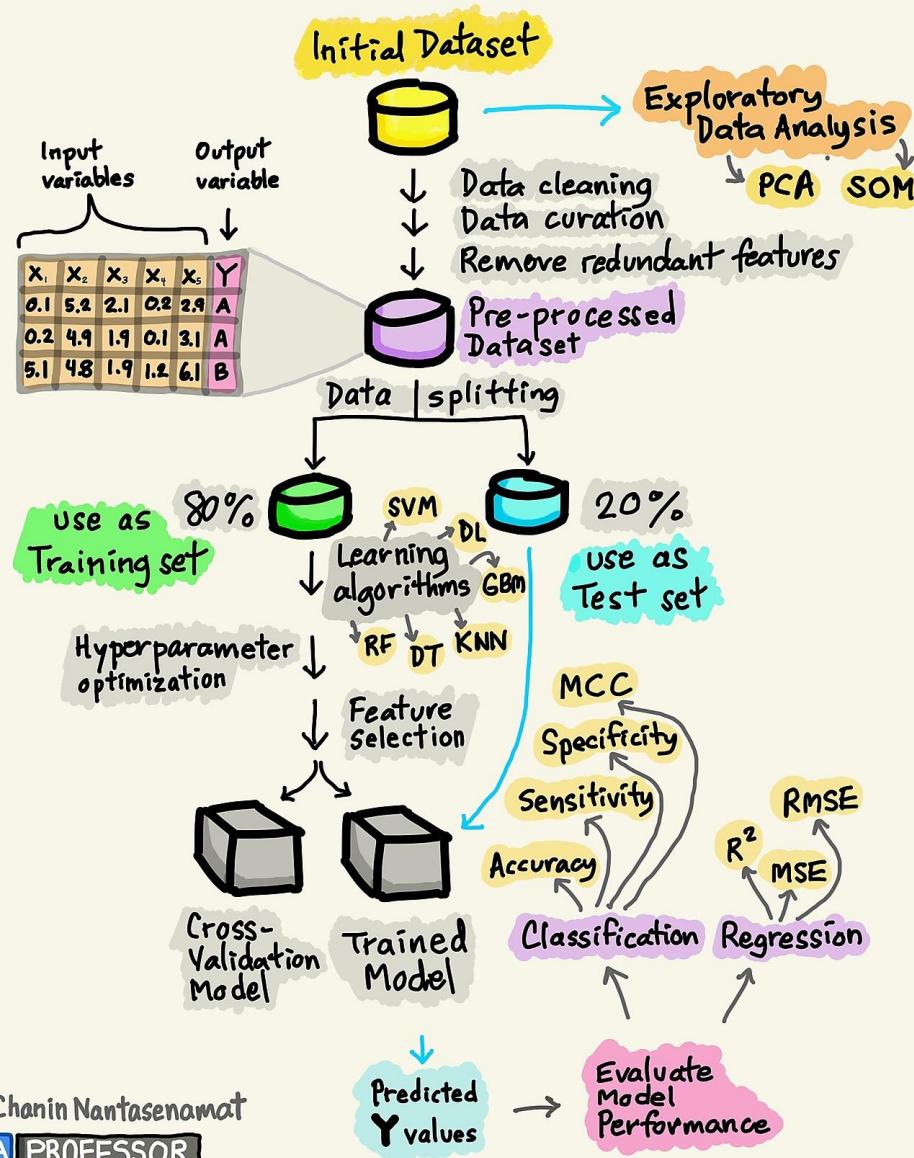
- ❑ Use of Machine Learning Tools
 - Use Orange Data Mining to carry out simple machine learning training and predictions.

Machine Learning Process

BUILDING THE

MACHINE LEARNING

MODEL



By: Chanin Nantasenamat

DATA PROFESSOR



<http://youtube.com/dataprofessor>

January 1, 2020

Machine Learning Process

- **The primary goal of the process is to identify a 'Model'.** The Model is the main thing that applications can submit requests to in order to gain insight on new data. A person working as the role of a Data Scientist performs the Machine Learning process and will ultimately decide on the right model to use.
- **The process starts with a question;** what are you trying to learn from your Machine Learning experiment? For example, in the case of recommendations, the question might be "*identify most commonly sold products for each product in the inventory*"

Machine Learning Process

- **The next step is to provide 'prepared data'.** Prepared data is one or more data sets that have been pre-processed (formatted, cleaned and sampled) in readiness to apply Machine Learning algorithms to. Preparing the data means that the data is in the best shape to draw scientific conclusions from and is not skewed in any way.
- Once you have your prepared data, you apply one or more Machine Learning algorithms to it with a view to producing a **Model**. This is an iterative process and you may loop around testing various algorithms until you have a Model that sufficiently answers your question.
- Once you have produced your chosen model, it will typically be exposed via some kind of API.

Choosing Learning Algorithms

Choosing a ML algorithm

- Every machine learning algorithm has its own style or *inductive bias*. For a specific problem, several algorithms may be appropriate and one algorithm may be a better fit than others. But it's not always possible to know beforehand which is the best fit. In cases like these, several algorithms are listed together in the cheat sheet. An appropriate strategy would be to try one algorithm, and if the results are not yet satisfactory, try the others.

Choosing a ML algorithm

There are three main categories of machine learning: **supervised learning**, **unsupervised learning**, and **reinforcement learning**.

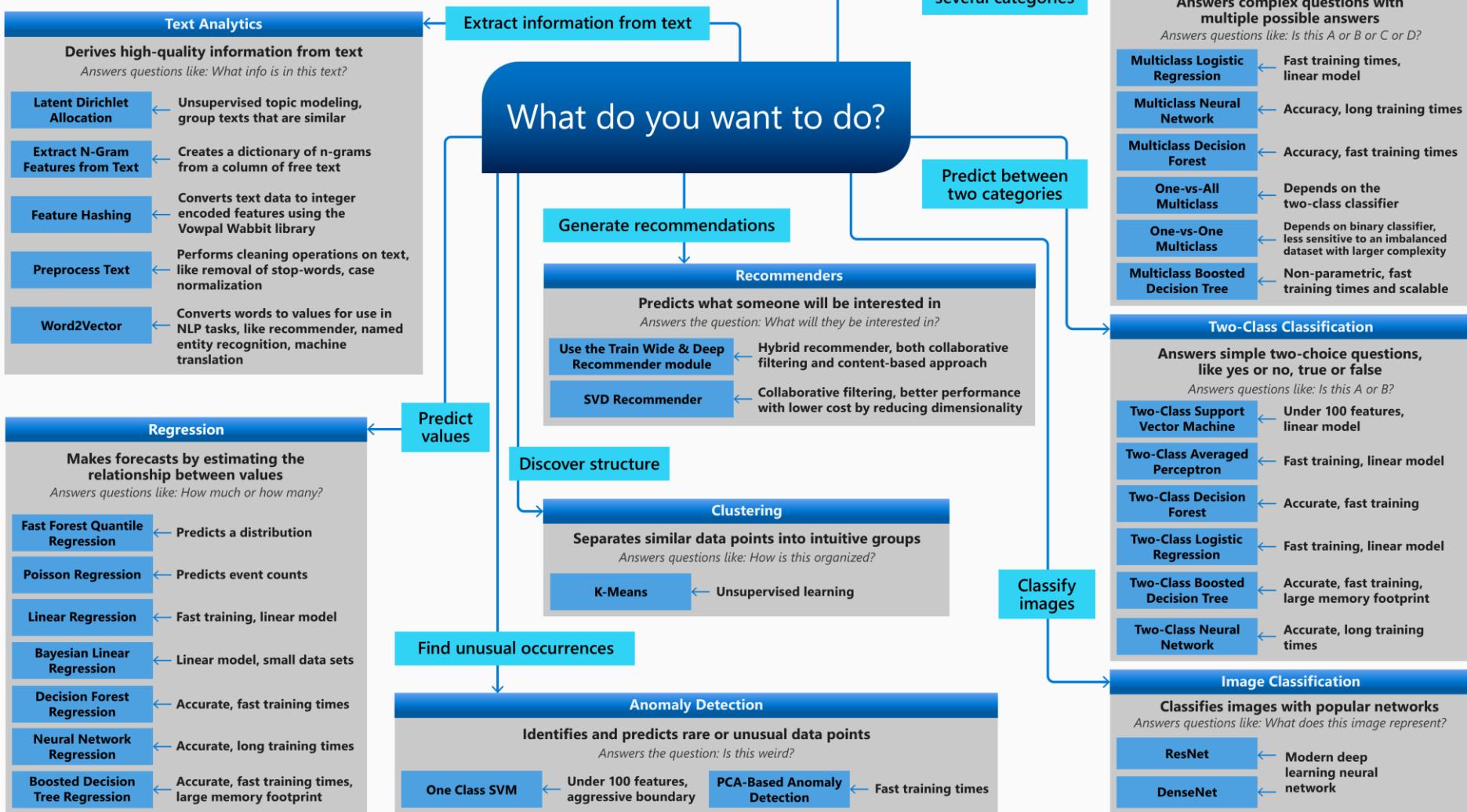
- In **supervised learning**, each data point is **labelled** or associated with a **category** or **value** of interest. An example of a categorical label is assigning an image as either a ‘cat’ or a ‘dog’. An example of a value label is the sale price associated with a used car. The goal of supervised learning is to study many labelled examples like these, and then to be able to make predictions about future data points. For example, identifying new photos with the correct animal or assigning accurate sale prices to other used cars. This is a popular and useful type of machine learning.

Choosing a ML algorithm

- In **unsupervised learning**, data points have no labels associated with them. Instead, the goal of an unsupervised learning algorithm is to organize the data in some way or to describe its structure. This can mean grouping it into clusters, as K-means does, or finding different ways of looking at complex data so that it appears simpler.
- In **reinforcement learning**, the algorithm gets to choose an action in response to each data point. It is a common approach in robotics, where the set of sensor readings at one point in time is a data point, and the algorithm must choose the robot's next action. It's also a natural fit for Internet of Things applications. The learning algorithm also receives a reward signal a short time later, indicating how good the decision was. Based on this, the algorithm modifies its strategy in order to achieve the highest reward.

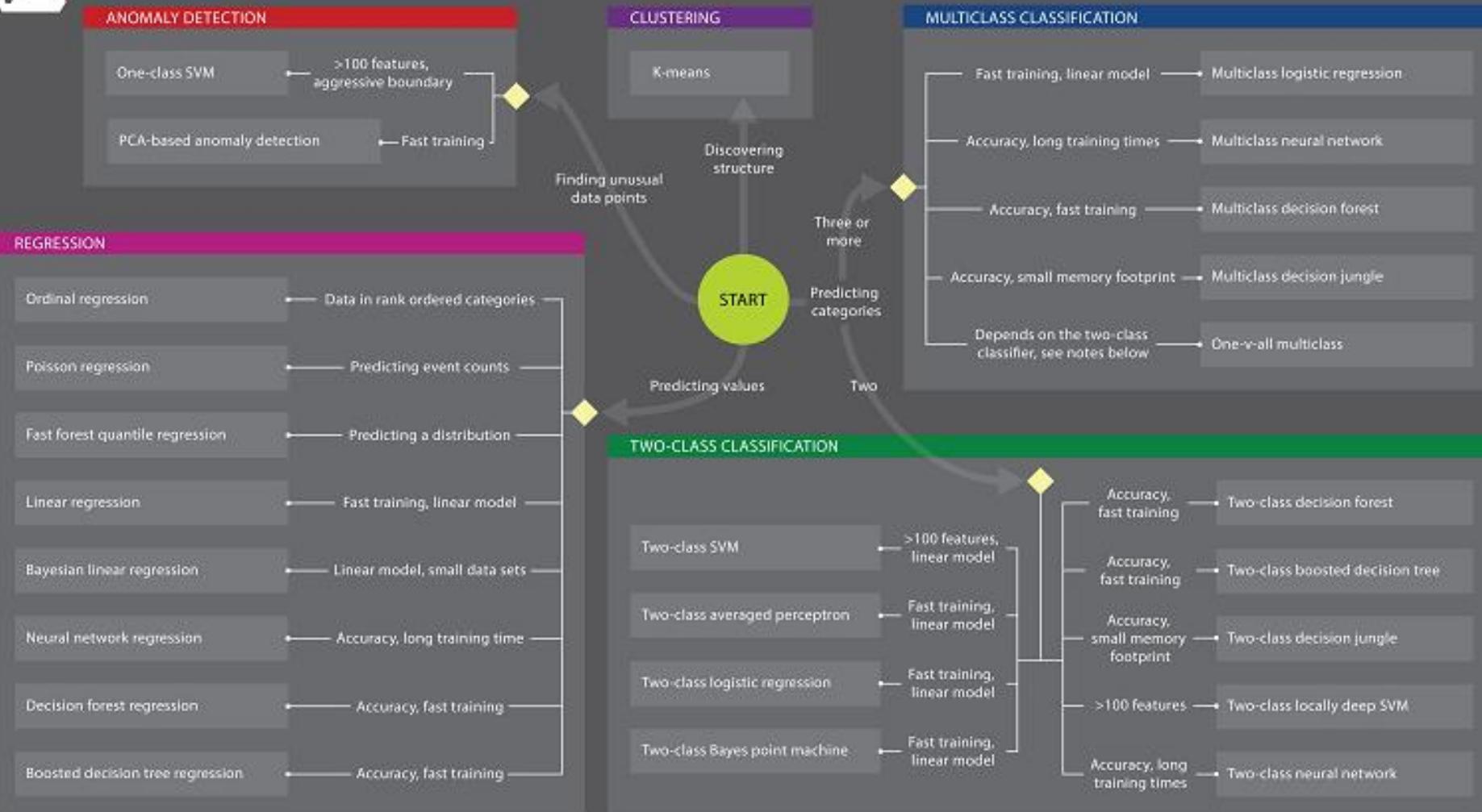
Machine Learning Algorithm Cheat Sheet

This cheat sheet helps you choose the best machine learning algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the goal you want to achieve with your data.



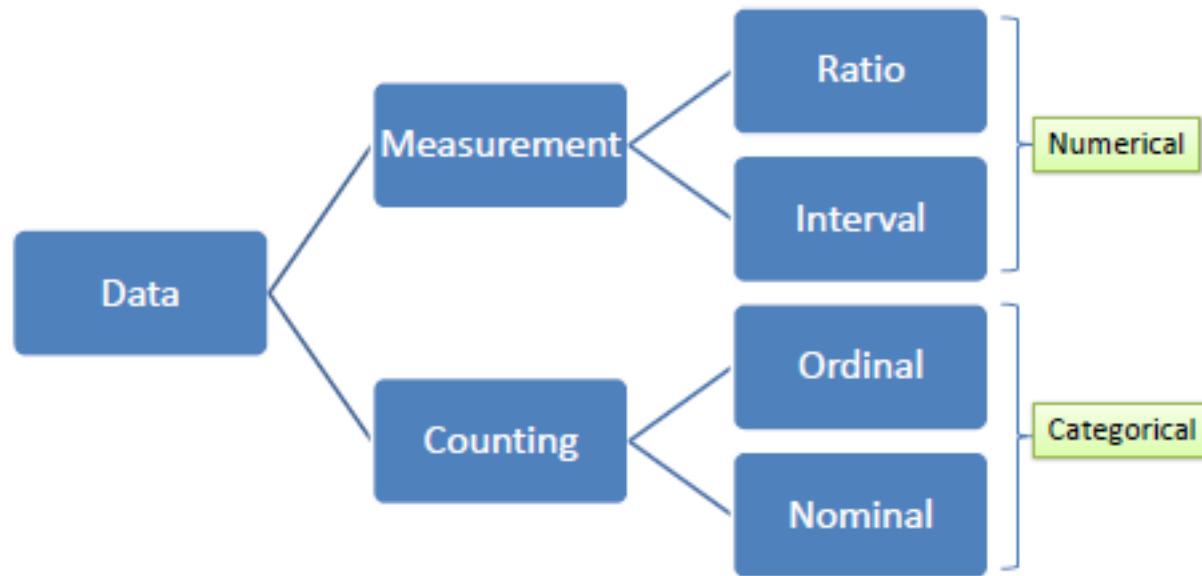
Microsoft Azure Machine Learning: Algorithm Cheat Sheet

This cheat sheet helps you choose the best Azure Machine Learning Studio algorithm for your predictive analytics solution. Your decision is driven by both the nature of your data and the question you're trying to answer.



Understand Principles of Data Representation

Data Types



source http://www.saedsayad.com/data_preparation.htm

Data Types

Data is information typically the results of measurement (numerical) or counting (categorical). Variables serve as placeholders for data. There are two types of variables, numerical and categorical.

- A **numerical** or **continuous variable** is one that can accept any value within a finite or infinite interval (e.g., height, weight, temperature, blood glucose, ...). There are two types of numerical data, **interval** and **ratio**. Data on an interval scale can be added and subtracted but cannot be meaningfully multiplied or divided because there is no true zero. For example, we cannot say that one day is twice as hot as another day. On the other hand, data on a ratio scale has true zero and can be added, subtracted, multiplied or divided (e.g., weight).
- A **categorical** or **discrete variable** is one that can accept two or more values (categories). There are two types of categorical data, **nominal** and **ordinal**. Nominal data does not have an intrinsic ordering in the categories. For example, "gender" with two categories, male and female. In contrast, ordinal data does have an intrinsic ordering in the categories. For example, "level of energy" with three orderly categories (low, medium and high).

ML algorithms require the inputs variables to be numerical

- No problem for data that is already a real number or integer
- Categorical variables need to be converted or encoded into a number

Encoding categorical variables

- Use of one-hot encoding for categorical variables

Original values (Direction)	I0	I1	I2	I3
NORTH	1	0	0	0
EAST	0	1	0	0
SOUTH	0	0	1	0
WEST	0	0	0	1

One-hot encoding of variables = number of values

In the case 4 possible values = 4 (I0,I1,I2,I3 variables)

The original column is replaced by 4 new columns

Encoding categorical variables

- Use of dummy variables encoding for categorical variables

Original values (Direction)	I0	I1	I2
NORTH	1	0	0
EAST	0	1	0
SOUTH	0	0	1
WEST	0	0	0

Number of dummy/indicator variables = number of values – 1
In the case 4 possible values – 1 = 3 (I0,I1,I2 dummy variables)
The original column is replaced by 3 new columns

Why not just use a value encoding (e.g. LabelEncoder)?

The problem with value encoding or LabelEncoder is that it leads to undesirable side effects if the ML algorithm uses the numerical values. For example, Average of East and West = $(1+3)/2 = 2$ (South).

Direction	Value
NORTH	0
EAST	1
SOUTH	2
WEST	3

This is nonsense as the average of East and West has no meaning.

Discretization/Binning

Discretization (otherwise known as quantization or binning) provides a way to partition continuous features into discrete values.

Certain datasets with continuous features may benefit from discretization, because discretization can transform the dataset of continuous attributes to one with only nominal attributes.

Time	TimeOfDay
00:00	0
07:30	1
08:00	1
12:00	2
13:00	2
14:00	2
17:00	2

We convert the 24H into 4 time period bins
00:00-06:00 is the bin 0 for early morning

Dataset

Dataset is a collection of data, usually presented in a tabular form. Each **column** represents a particular **variable (attribute)**, and each **row** corresponds to a given **member (observation)** of the data.

In predictive modeling, **predictors** or **attributes** are the **input variables** and **target** or **class attribute** is the **output variable** whose value is determined by the values of the predictors and function of the predictive model.

Diagram illustrating a dataset structure:

The dataset consists of 14 rows and 6 columns. The columns are labeled: ID, Outlook, Temp, Humidity, Windy, and Play Golf.

Annotations explain the structure:

- Columns:** Points to the column headers.
- Rows:** Points to the first row of data.
- Values:** Points to the value 'Yes' in the 'Play Golf' column for row 5.

ID	Outlook	Temp	Humidity	Windy	Play Golf
1	Rainy	85	82	False	No
2	Rainy	80	88	True	No
3	Overcast	83	86	False	Yes
4	Sunny	70	80	False	Yes
5	Sunny	68	?	False	Yes
6	Sunny	65	58	True	No
7	Overcast	64	62	True	Yes
8	Rainy	72	95	?	No
9	Rainy	?	70	False	Yes
10	Sunny	75	72	False	Yes
11	Rainy	75	74	True	Yes
12	?	72	78	True	Yes
13	Overcast	81	66	False	Yes
14	Sunny	71	79	True	No

Dataset

There are some alternative terms for columns, rows and values.

- Columns, Fields, Attributes, Variables
- Rows, Records, Objects, Cases, Instances, Examples, Vectors
- Values, Data

Feature Engineering and Feature Selection

Feature Engineering

A **feature** is an **attribute** or **property** shared by all of the independent units on which analysis or prediction is to be done. Any attribute could be a feature, as long as it is useful to the model.

Feature engineering is the process of using domain knowledge of the data to create features that make machine learning algorithms work.

- **feature engineering:** This process attempts to create additional relevant features from the existing raw features in the data, and to increase the predictive power of the learning algorithm.
- **feature selection:** This process selects the key subset of original data features in an attempt to reduce the dimensionality of the training problem.

Normally **feature engineering** is applied first to generate additional features, and then the **feature selection** step is performed to eliminate irrelevant, redundant, or highly correlated features.

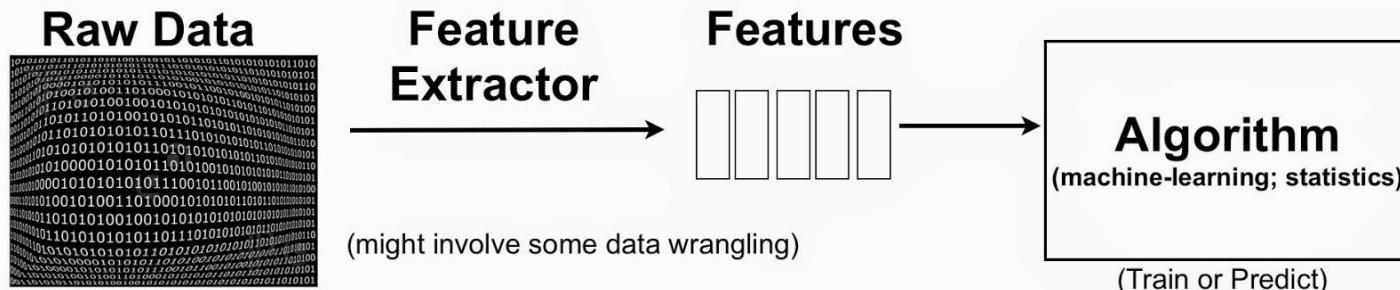
Feature engineering is fundamental to the application of machine learning, and is both difficult and expensive.

source <https://docs.microsoft.com/en-us/azure/machine-learning/team-data-science-process/create-features>

Feature Engineering

The **training data** used in machine learning can often be enhanced by **extraction of features** from the raw data collected. An example of an engineered feature in the context of learning how to classify the images of handwritten characters is creation of a bit density map constructed from the raw bit distribution data. This map can help locate the edges of the characters more efficiently than simply using the raw distribution directly.

What kind of features should be created to enhance the dataset when training a model? Engineered features that enhance the training provide information that better differentiates the patterns in the data. The new features are expected to provide additional information that is not clearly captured or easily apparent in the original or existing feature set. But **this process is something of an art**. Sound and productive decisions often require some domain expertise.



Feature Engineering

- Engineered and selected features increase the efficiency of the training process, which attempts to extract the key information contained in the data.
- They also improve the power of these models to classify the input data accurately and to predict outcomes of interest more robustly.
- Feature engineering and selection can also combine to make the learning more computationally tractable. It does so by enhancing and then reducing the number of features needed to calibrate or train a model. Mathematically speaking, the features selected to train the model are a minimal set of independent variables that explain the patterns in the data and then predict outcomes successfully.

Feature Selection vs Feature Extraction

In **feature selection** we try to find the best subset of the input feature set

In **feature extraction** we create new features based on the transformation or combination of the original feature set

No Code / Low Code Machine Learning with Orange Data Mining

PRACTICAL EXERCISE



INTRODUCTION TO ORANGE DATA MINING

<https://www.orangedatamining.com>

The screenshot shows the official website for Orange Data Mining. The header features the "orange" logo in orange and black, followed by navigation links for Screenshots, Workflows, Download, Blog, Docs, Workshops, a search icon, and a Donate button. Below the header is a main section with the heading "Data Mining" and the tagline "Fruitful and Fun". It describes Orange as an open-source machine learning and data visualization tool, emphasizing visual workflow building and a diverse toolbox. A large, friendly orange character wearing glasses and holding a document is the central illustration. A prominent orange "Download Orange" button is located below the heading. The footer contains three news items: "Explainable AI Project Meeting" (Nov 26, 2021), "Characterizing Clusters with a Box Plot" (Oct 21, 2021), and "Semantic Analysis of Documents" (Sep 17, 2021). Each news item includes a thumbnail image, a brief description, and a "Read more >" link.

ORANGE DATA MINING FEATURES

- Designed to be an easy to use interface for applied machine learning
- Simple to use drag-and-drop interface
- Pre-built algorithms (no need to code anything)
- Sample data sets
- Based on Python (for more advanced users)
- Supports Windows, macOS and Linux
- For more information about what you can do: <https://orangedatamining.com/>

DOCUMENTATION

Documentation is available online

- <https://orangedatamining.com/docs/>

Orange Data Mining: Python Library Reference

- <https://orange3.readthedocs.io/projects/orange-data-mining-library/en/latest/>

WHAT CAN YOU DO WITH ORANGEDM?

- Supervised Learning (Classification)
 - Answers simple, two-choice questions, like yes-or-no, true-or-false.

Is this tweet positive?

Will this customer renew their service?

Which of two coupons draws more customers

WHAT CAN YOU DO WITH ORANGEDM?

- Supervised Learning (Regression)

- Forecast the future by estimating the relationship between variables.

Estimate product demand

Predict sales figures

Determine equipment servicing priorities

WHAT CAN YOU DO WITH ORANGEDM?

- Unsupervised Learning (Clustering)
 - Separate similar data points into intuitive groups.

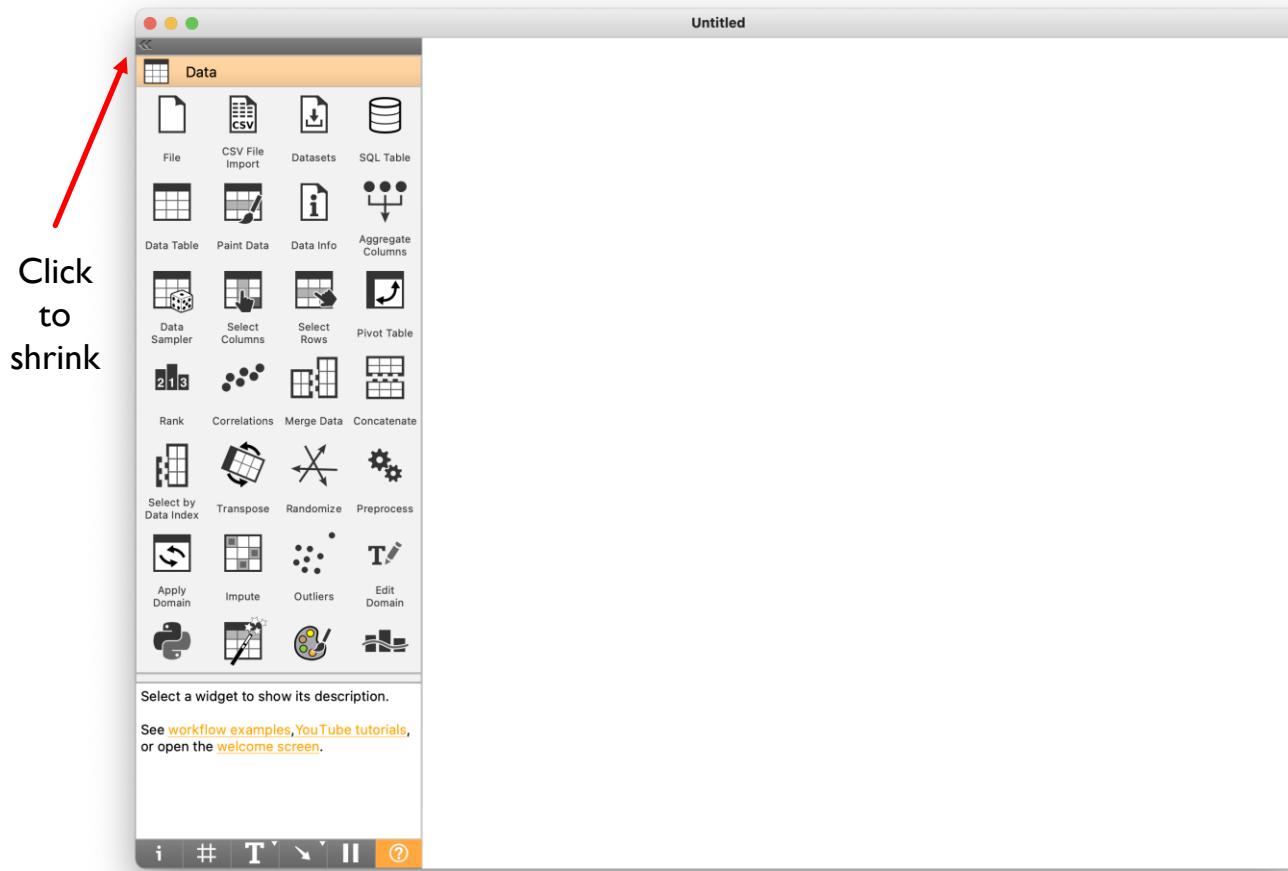
Perform customer segmentation

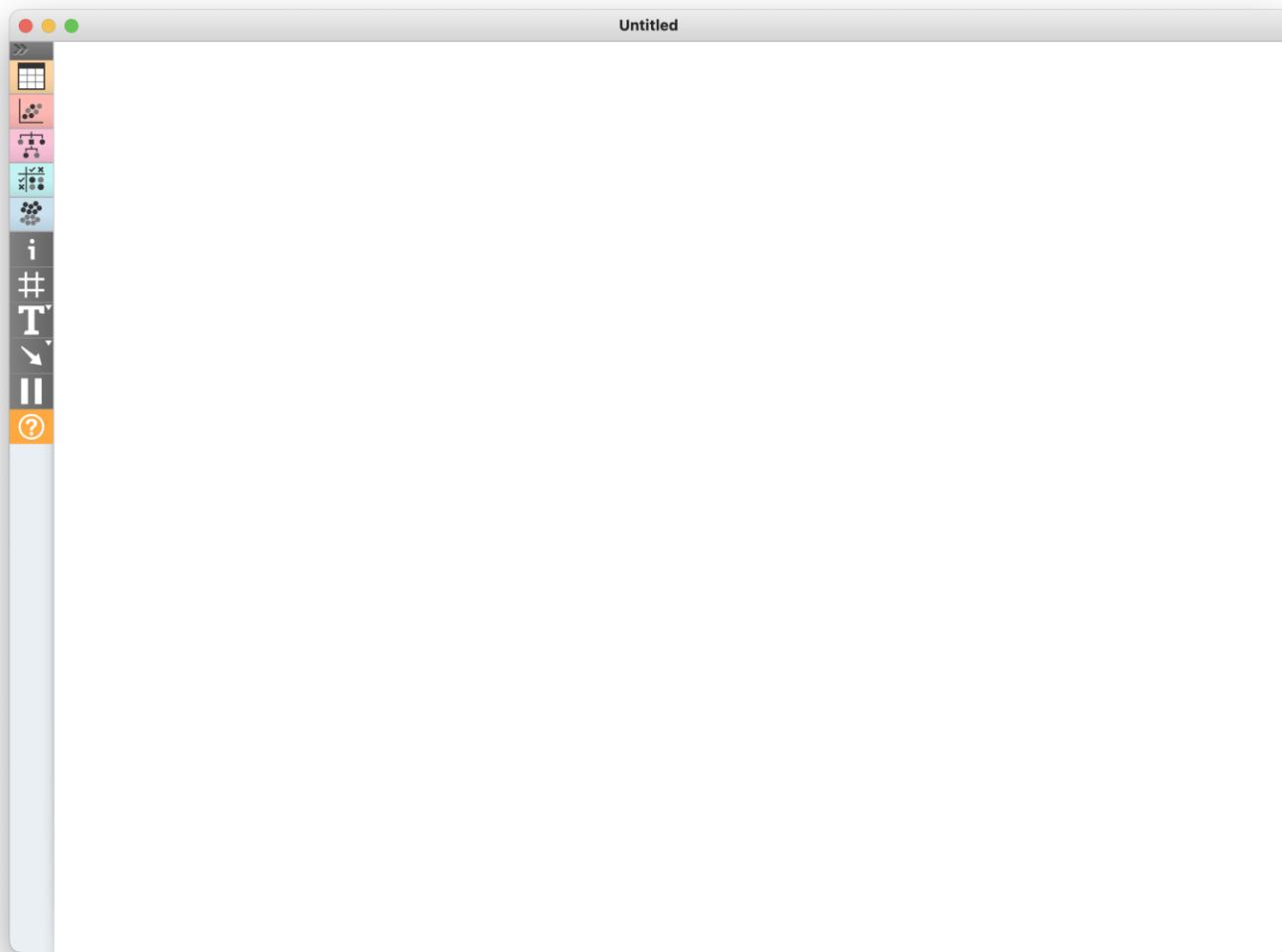
Predict customer tastes

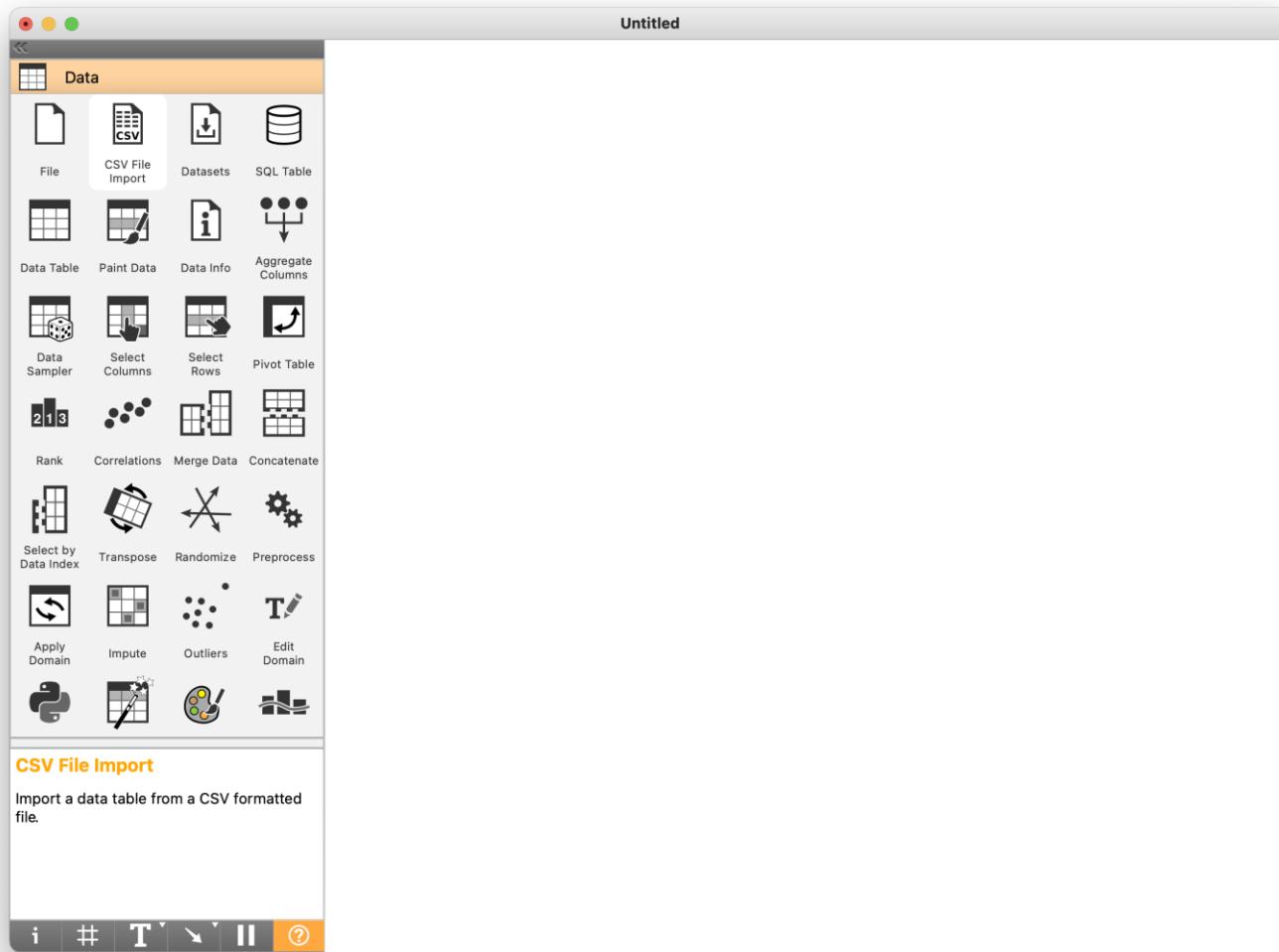
Determine which products fail the same way

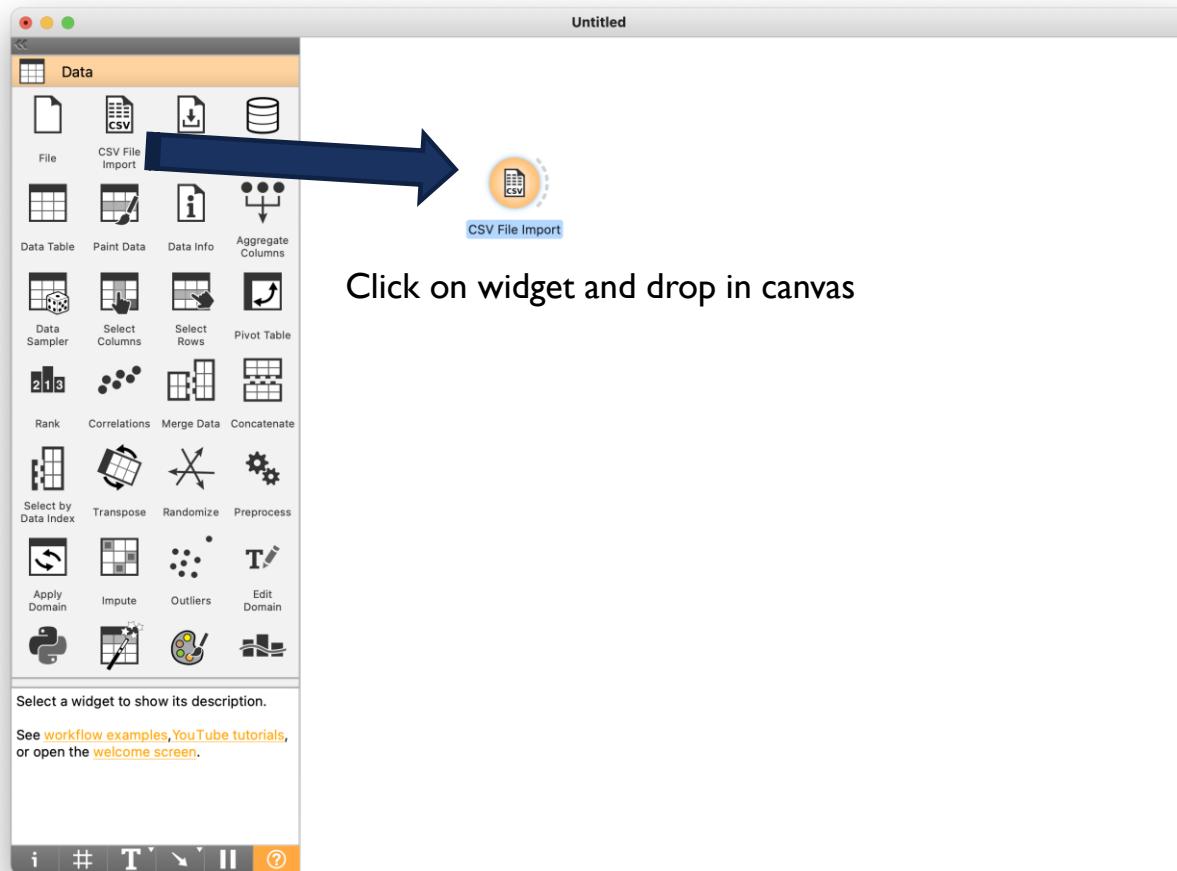
ACTIVITY

- Navigating around OrangeDM





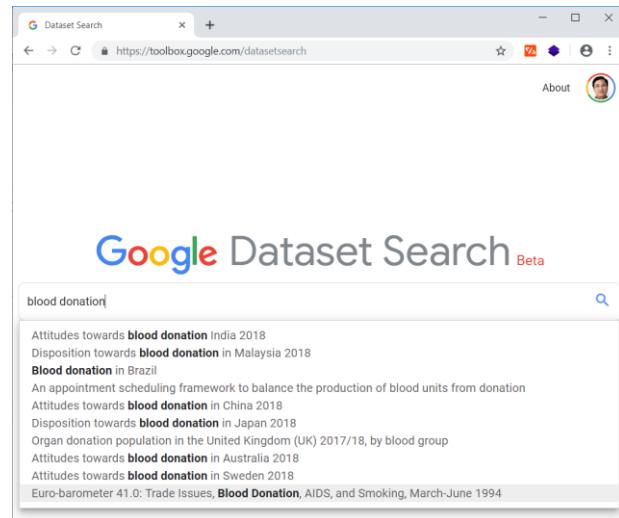




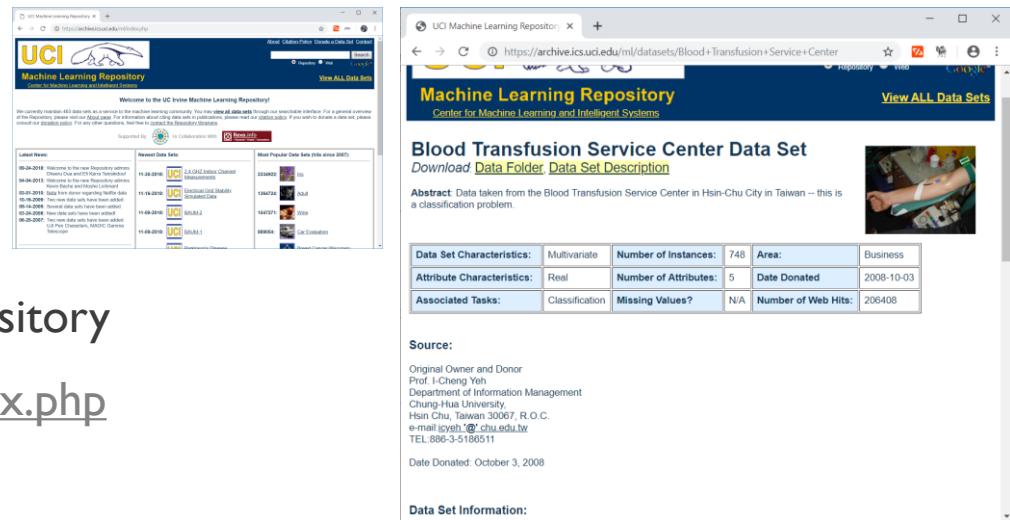
IMPORTING DATA

SEARCH FOR THE “BLOOD DONATION” DATASET

- Use Google Dataset Search Toolbox
- <https://toolbox.google.com/datasetsearch>

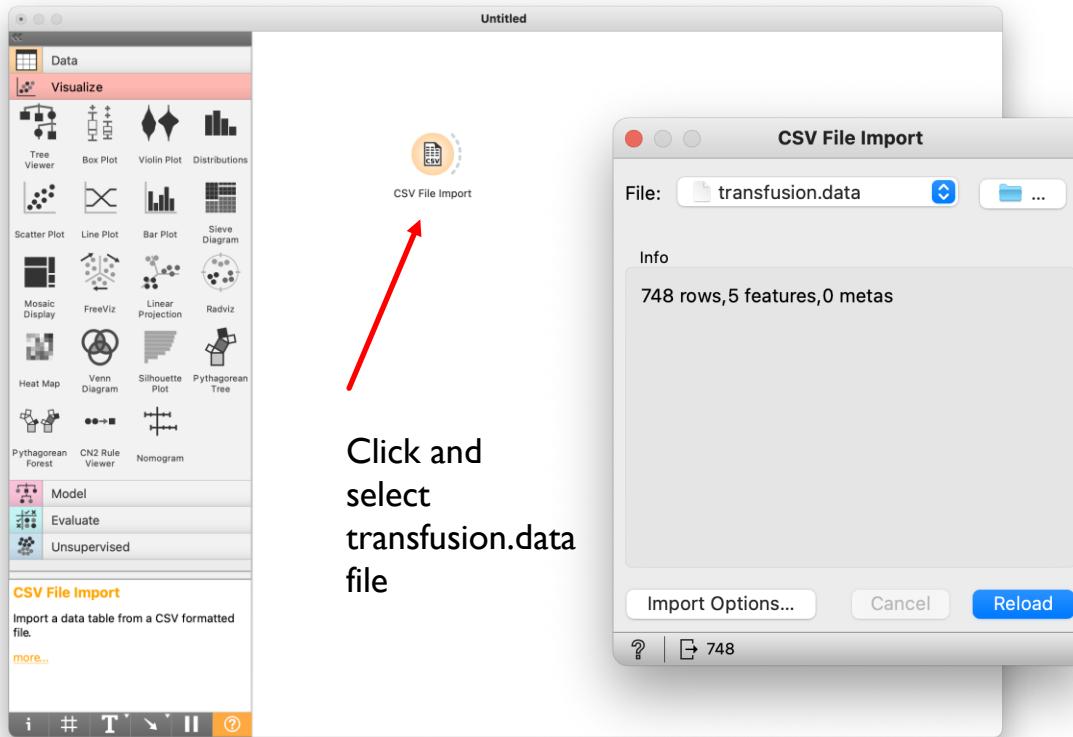


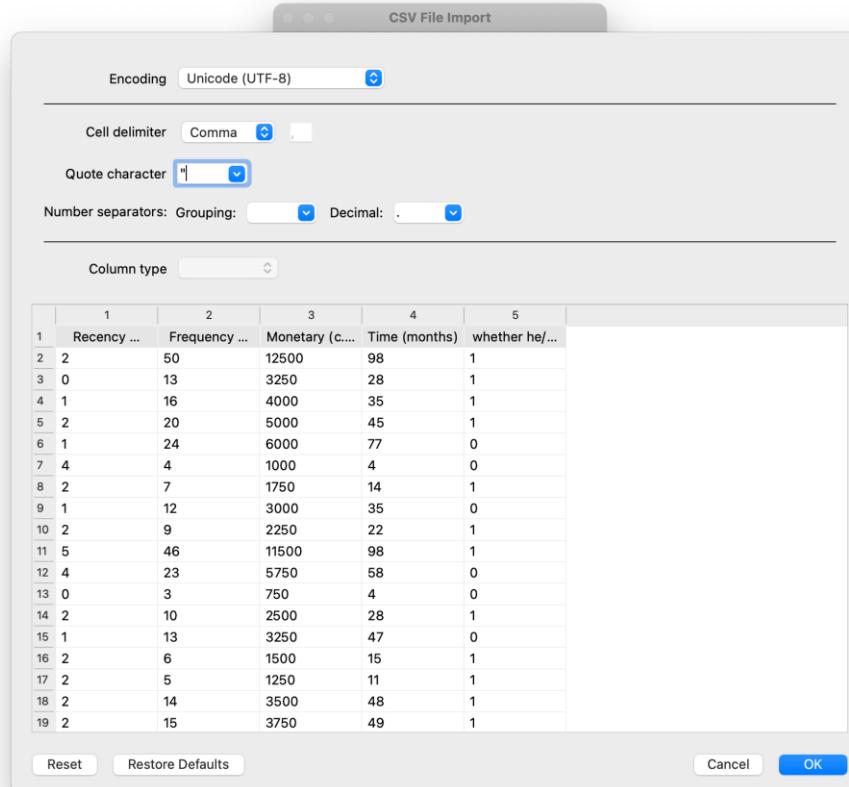
SEARCH FOR THE “BLOOD DONATION” DATASET



- Use UCI Machine Learning Repository
 - <https://archive.ics.uci.edu/ml/index.php>

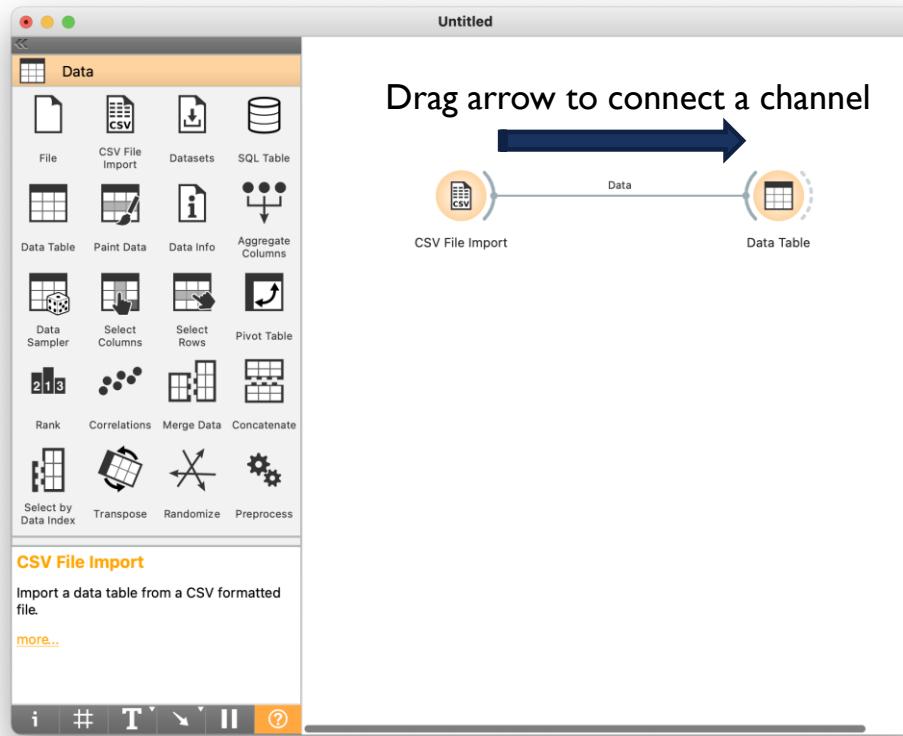
USE THE CSV FILE IMPORT





DATA PREPARATION

ADD A MODULE TO CHECK THE DATA

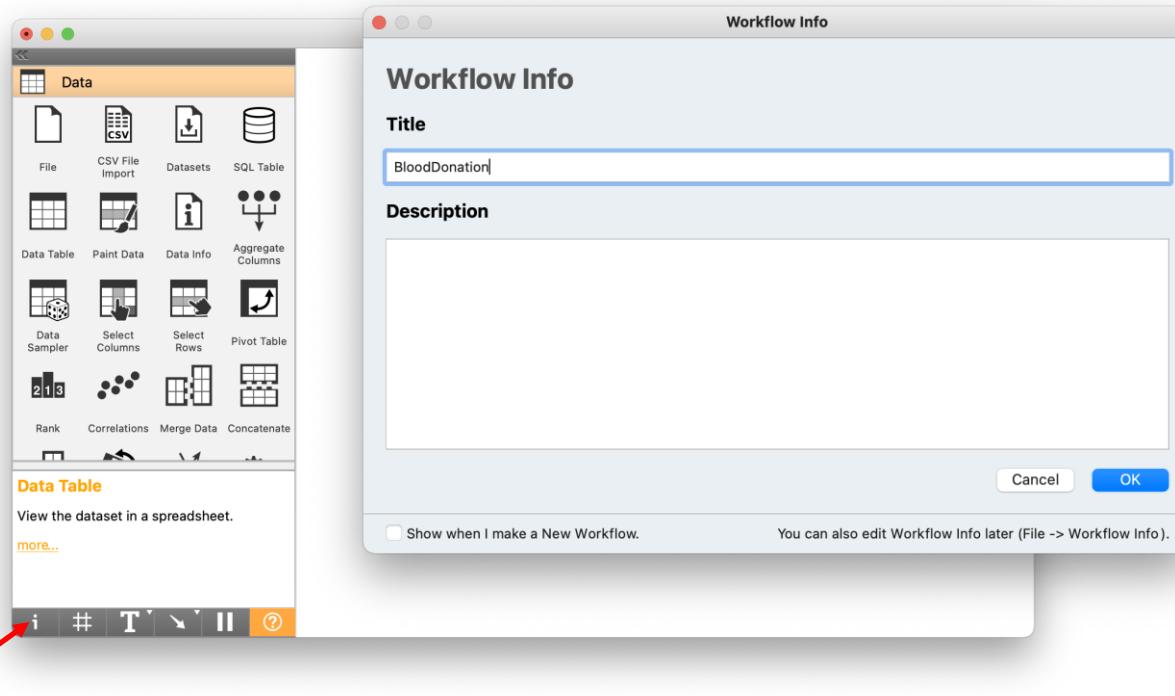


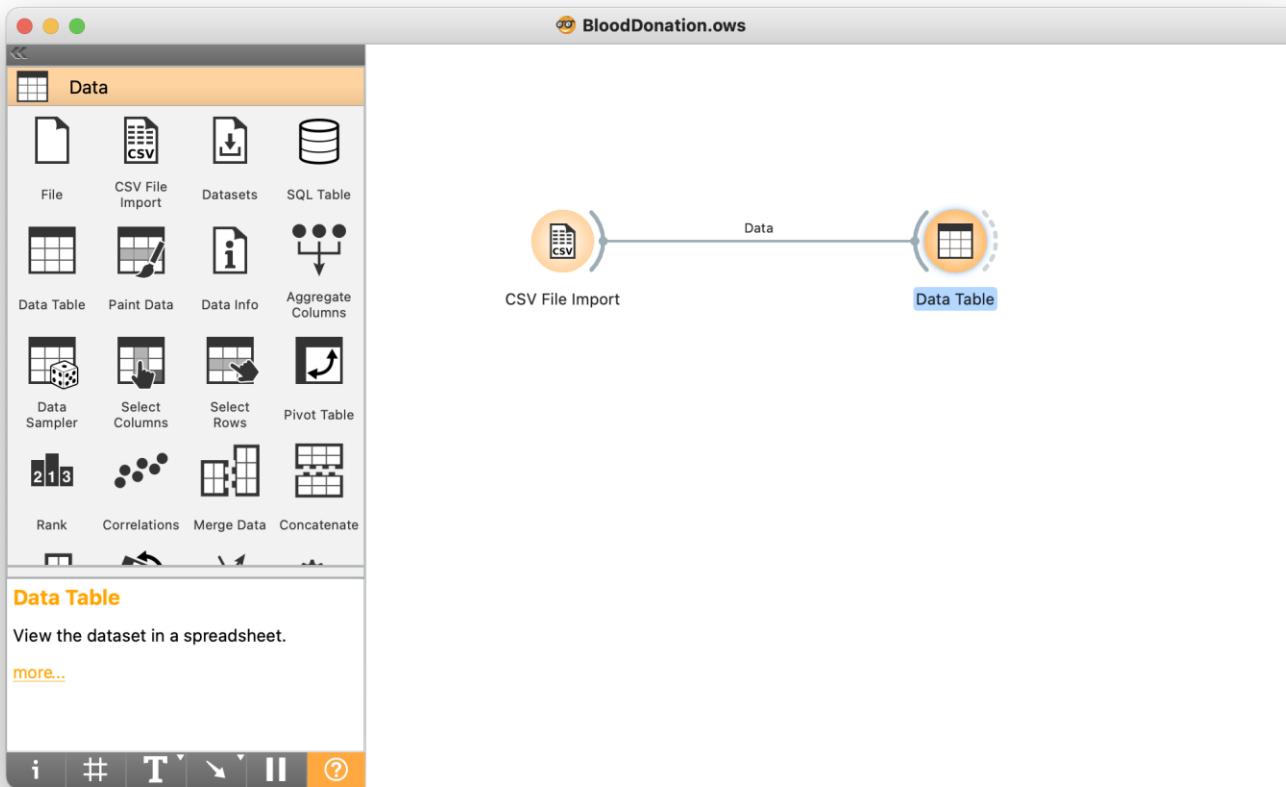
ADD A MODULE TO CHECK THE DATA

The screenshot shows the KNIME Data Explorer interface. On the left, the 'Data' palette is open, displaying various data processing modules like 'File', 'CSV File Import', 'Datasets', 'SQL Table', 'Data Table', 'Paint Data', etc. In the center, a workflow titled 'Untitled' is shown with a 'CSV File Import' node connected to a 'Data Table' node. A red arrow points from the text 'Double click or Right-click > Open' to the 'Data Table' node. Below the workflow, a 'Data Table' window is open, showing the following data:

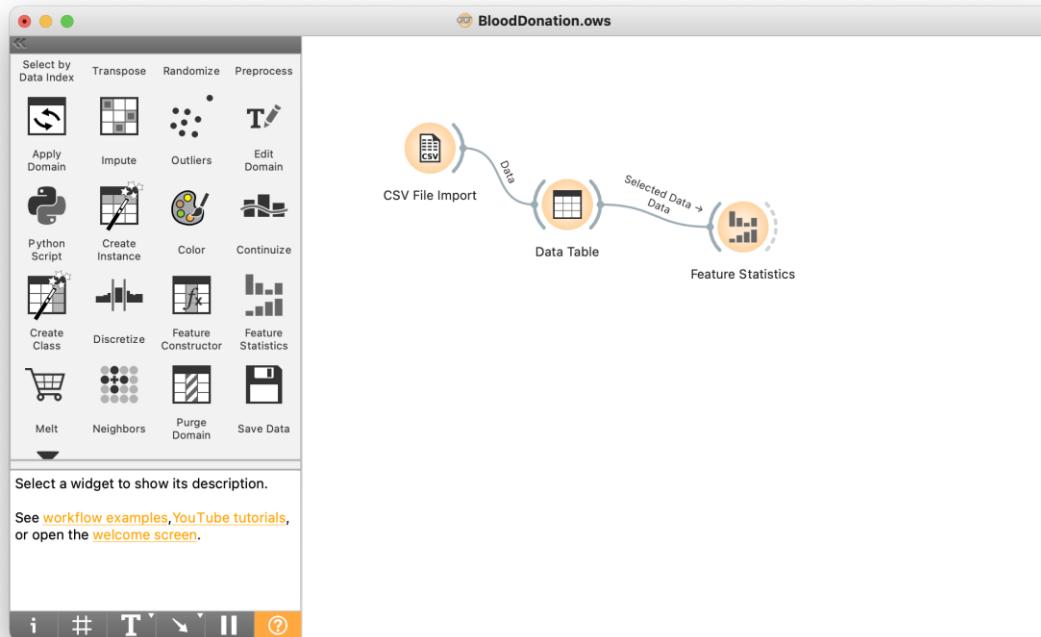
	Recency (months)	Frequency (times)	monetary (c.c.blo...)	Time (months)	#e donated blood i
1	2	50	12500	98	1
2	0	13	3250	28	1
3	1	16	4000	35	1
4	2	20	5000	45	1
5	1	24	6000	77	0
6	4	4	1000	4	0
7	2	7	1750	14	1
8	1	12	3000	35	0
9	2	9	2250	22	1
10	5	46	11500	98	1
11	4	23	5750	58	0
12	0	3	750	4	0
13	2	10	2500	28	1
14	1	13	3250	47	0
15	2	6	1500	15	1
16	2	5	1250	11	1
17	2	14	3500	48	1
18	2	15	3750	49	1
19	2	6	1500	15	1
20	2	3	750	4	1

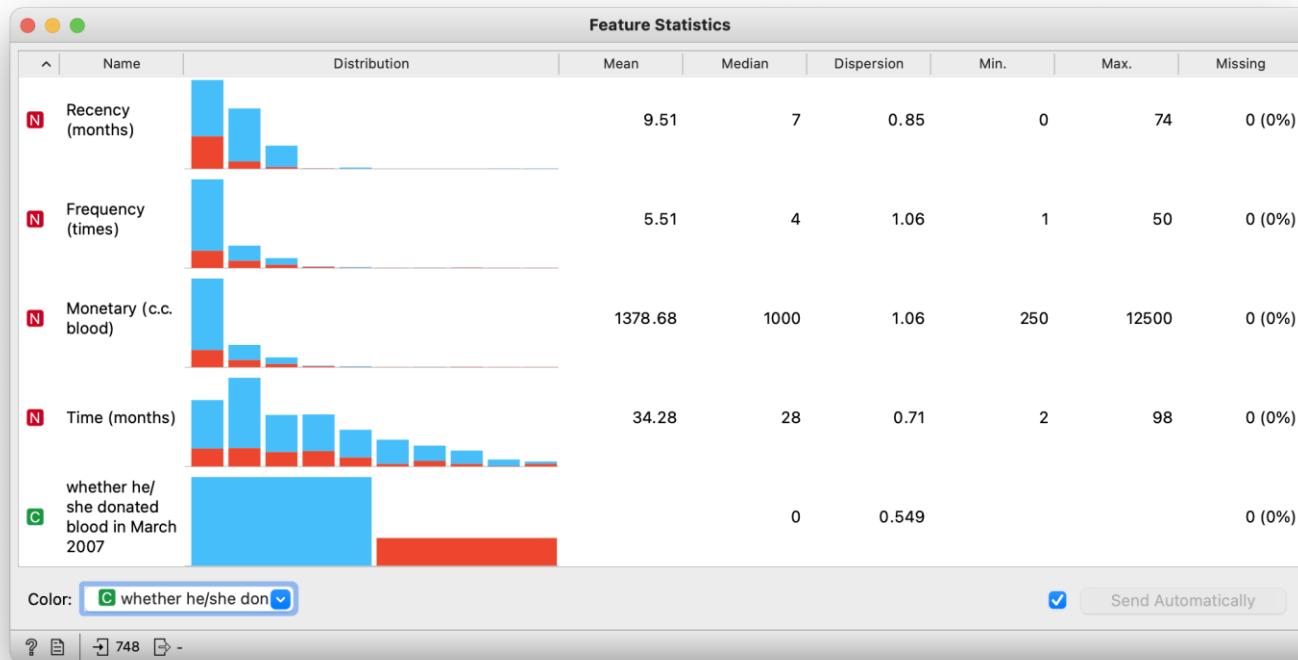
CHANGE WORKFLOW INFO TO BLOODTRANSFUSION. SAVE FILE.





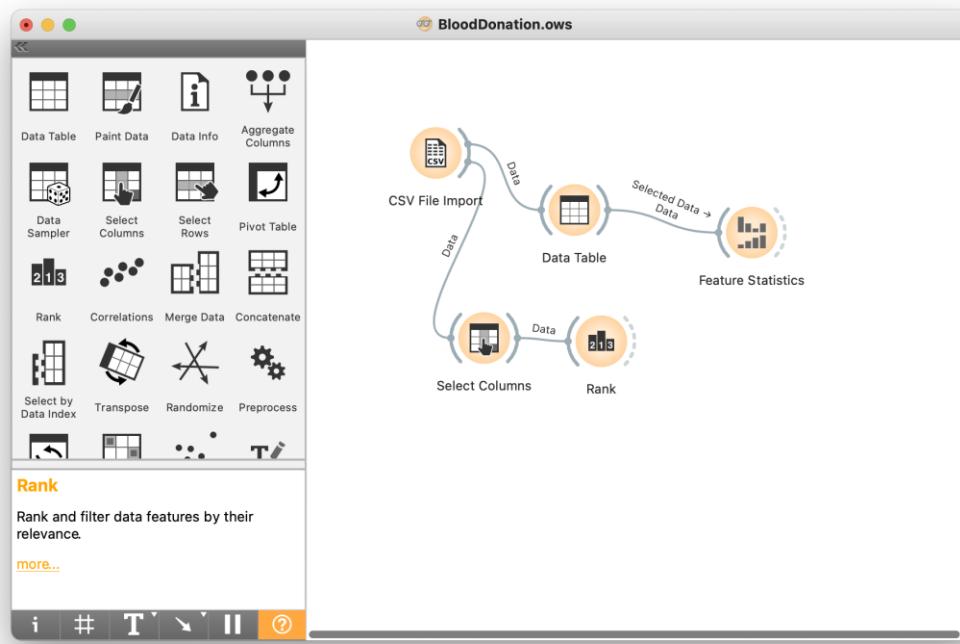
CHECK DATA QUALITY

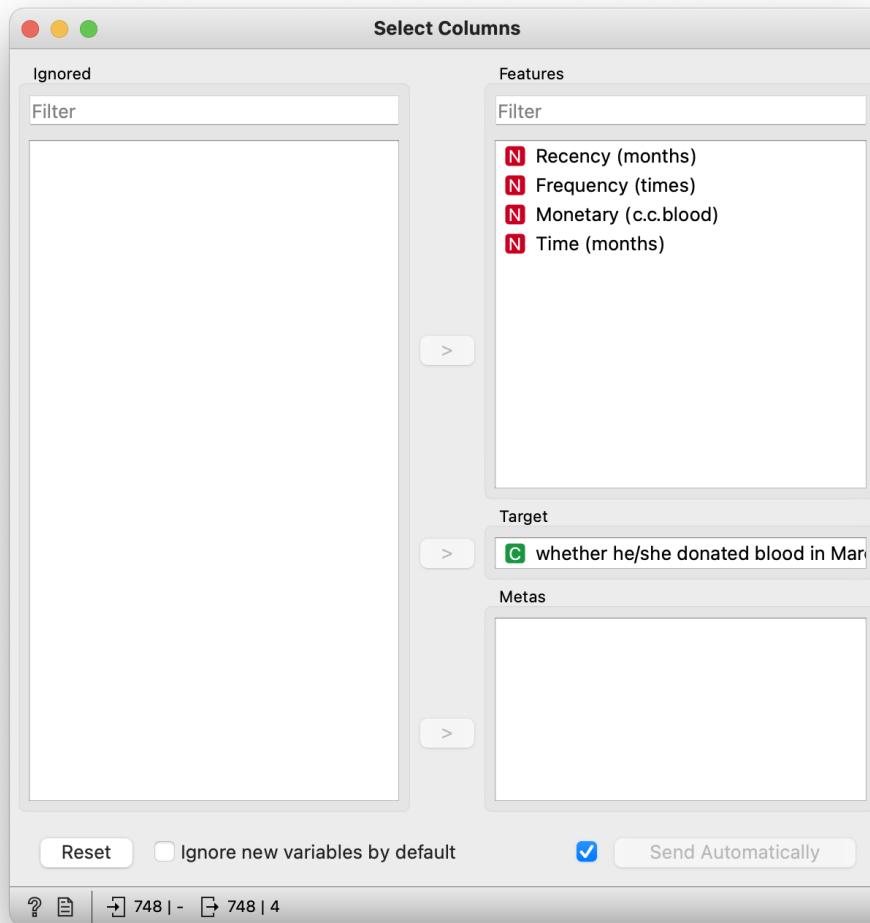


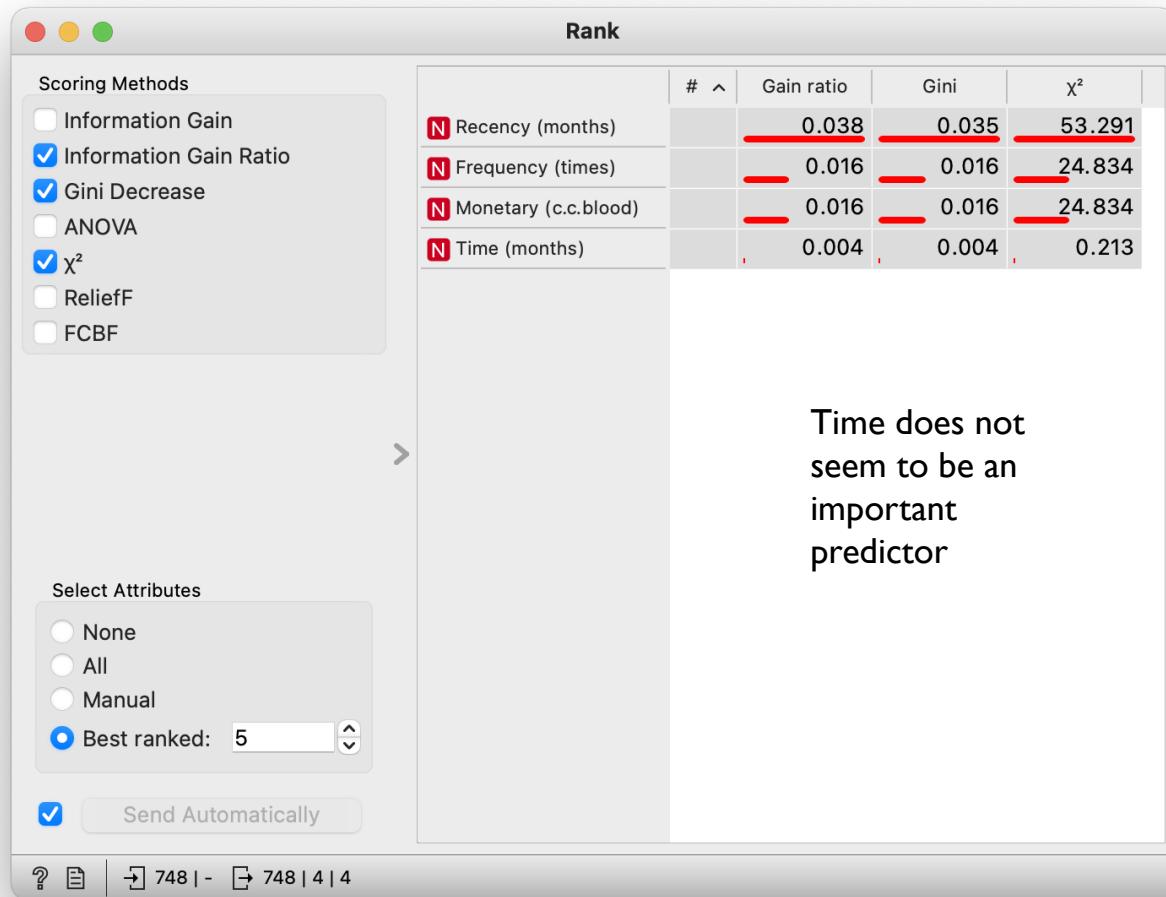


FEATURE SELECTION

FEATURE SELECTION







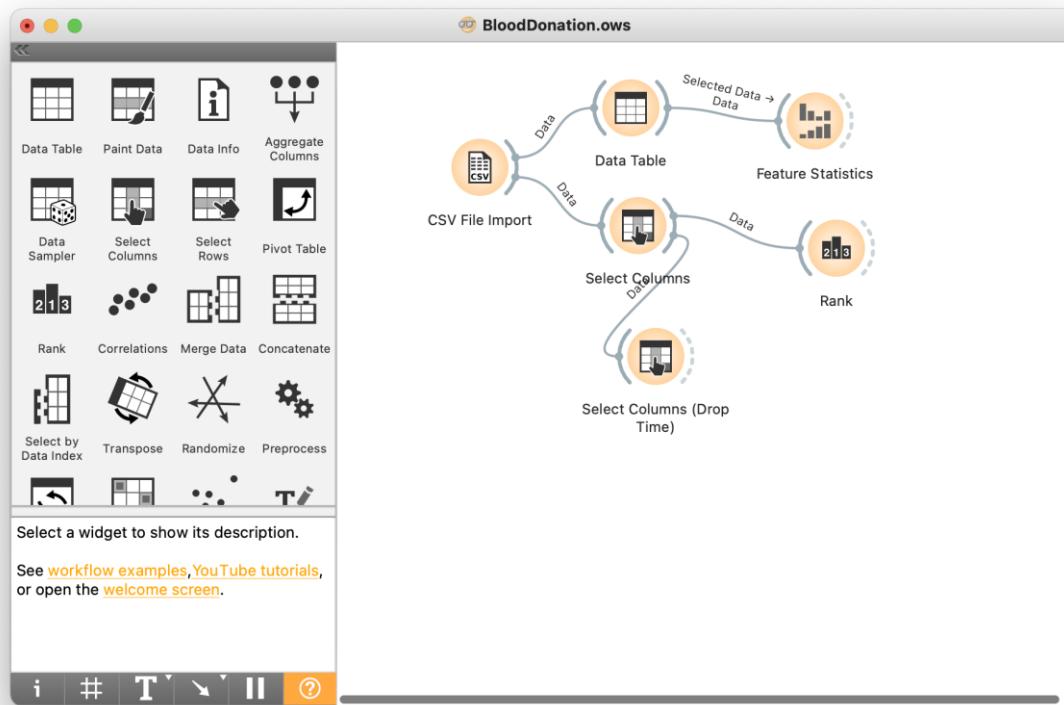
3-57

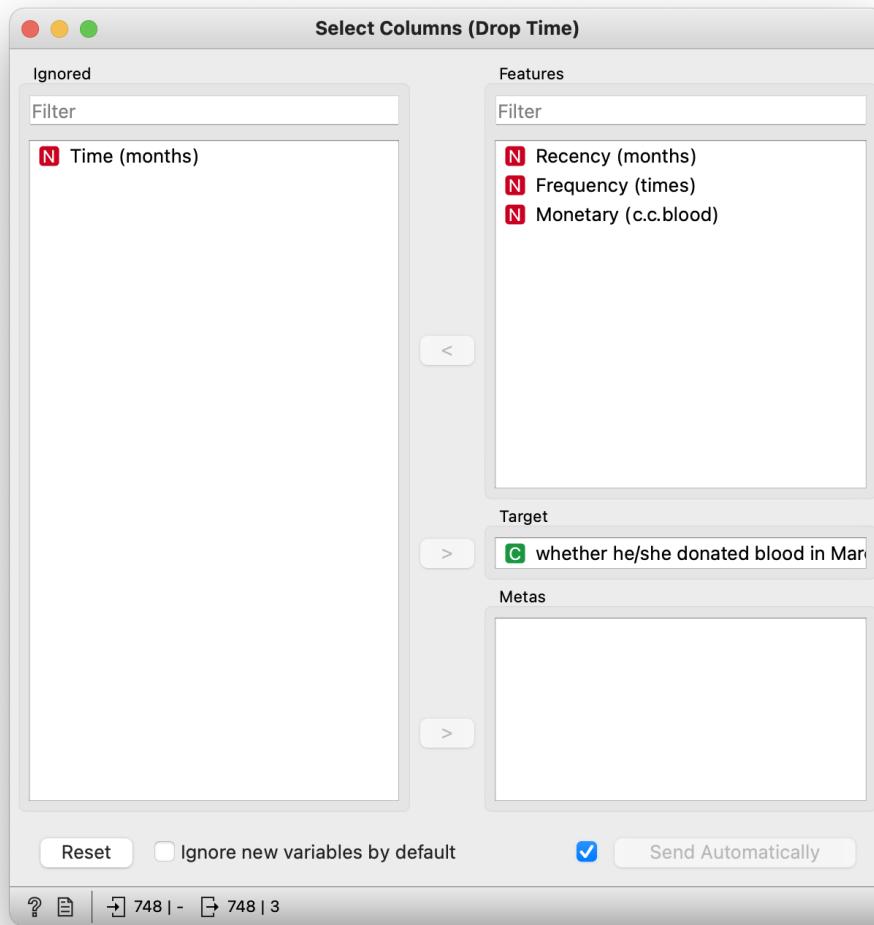
DROP TIME

IT8302 APPLIED MACHINE LEARNING



USE SECOND SELECT COLUMNS

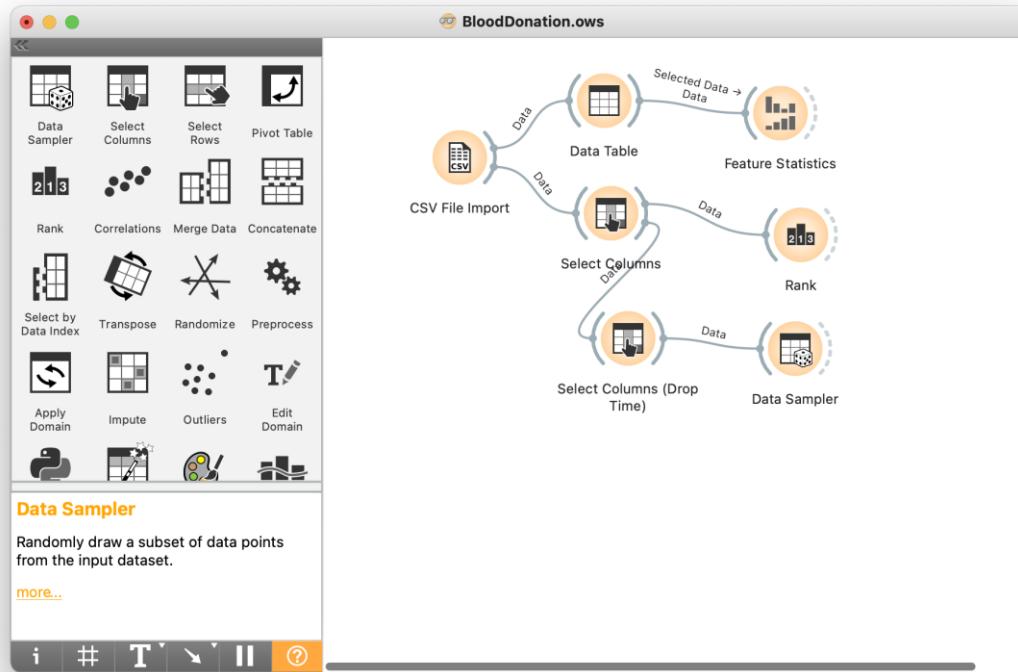


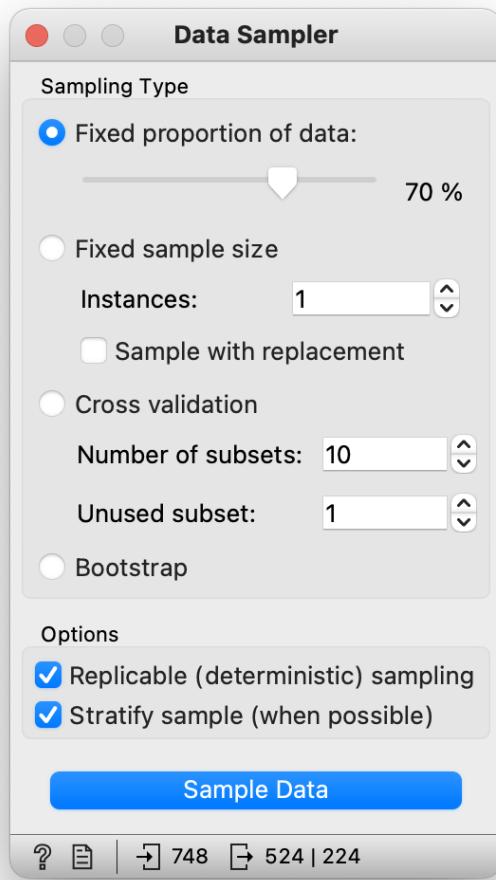




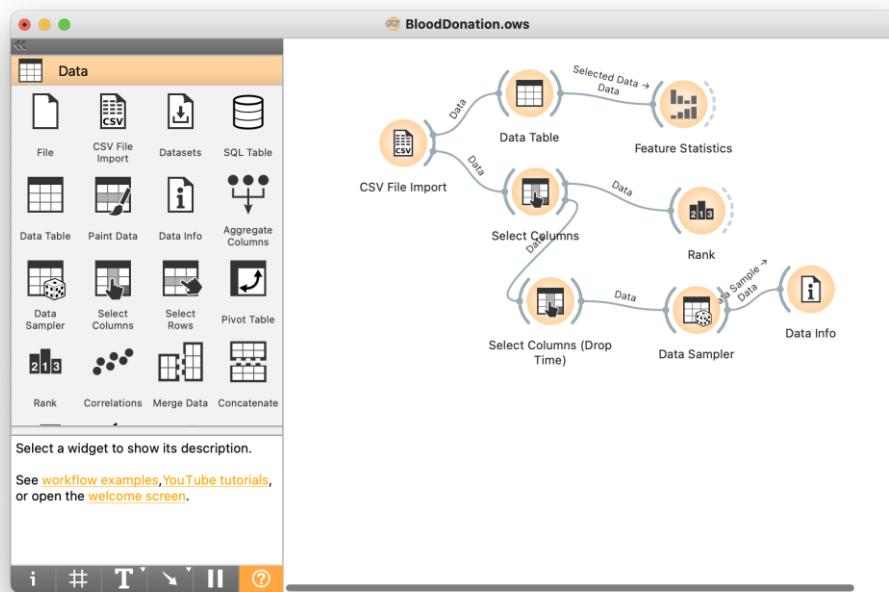
TRAIN TEST SPLIT

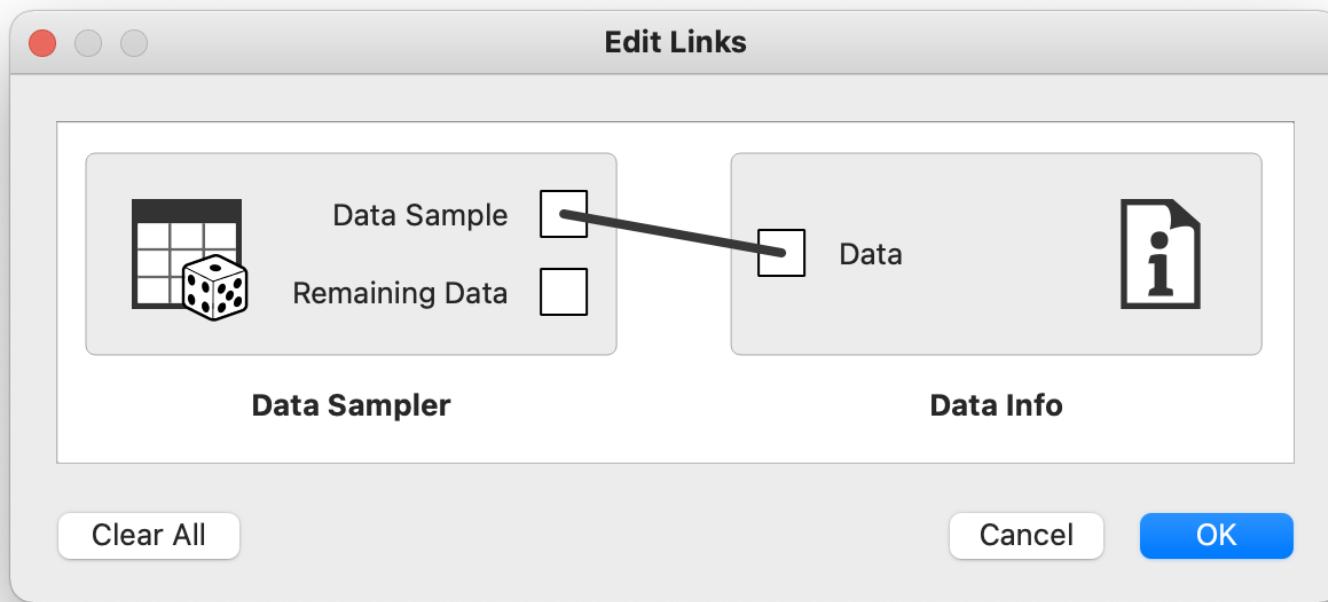
USE THE DATA SAMPLER

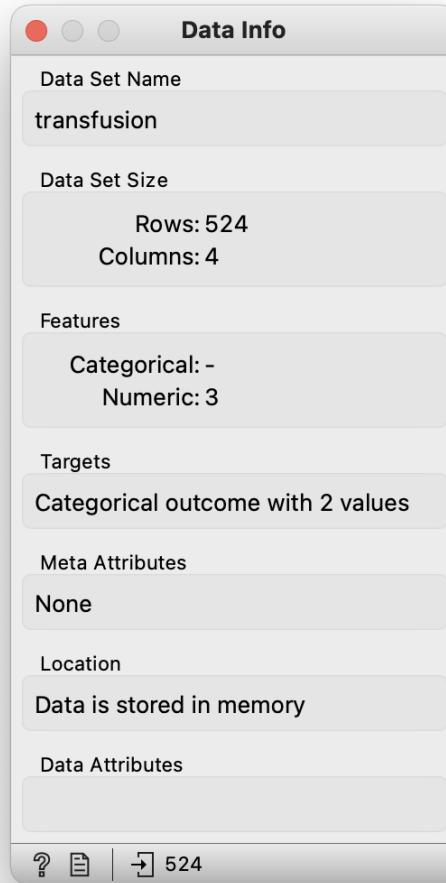




USE DATA INFO TO CHECK THE DATA SAMPLING





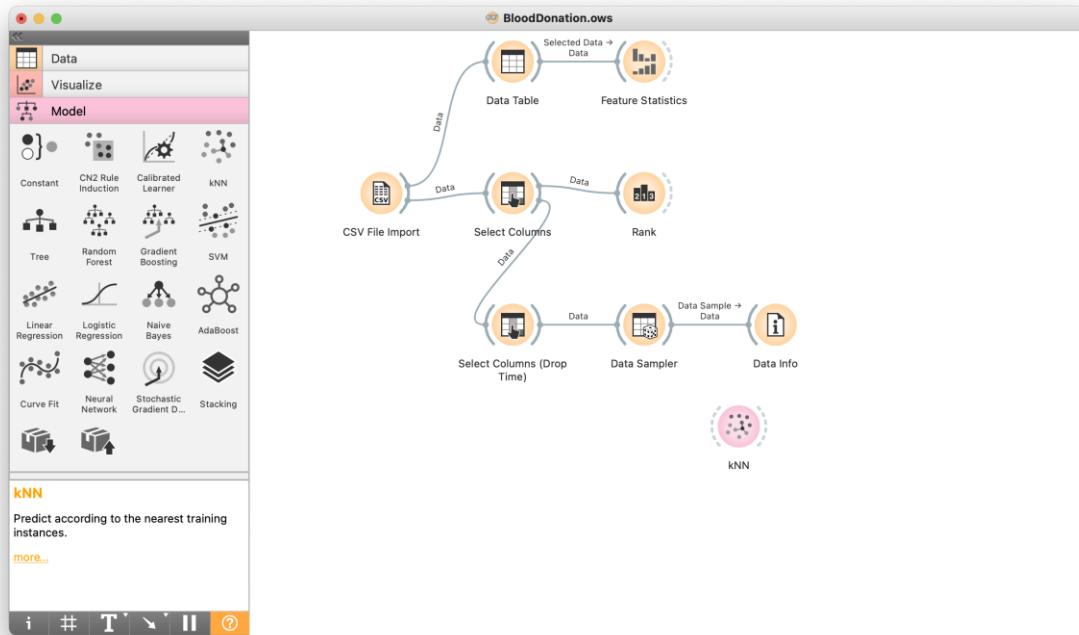




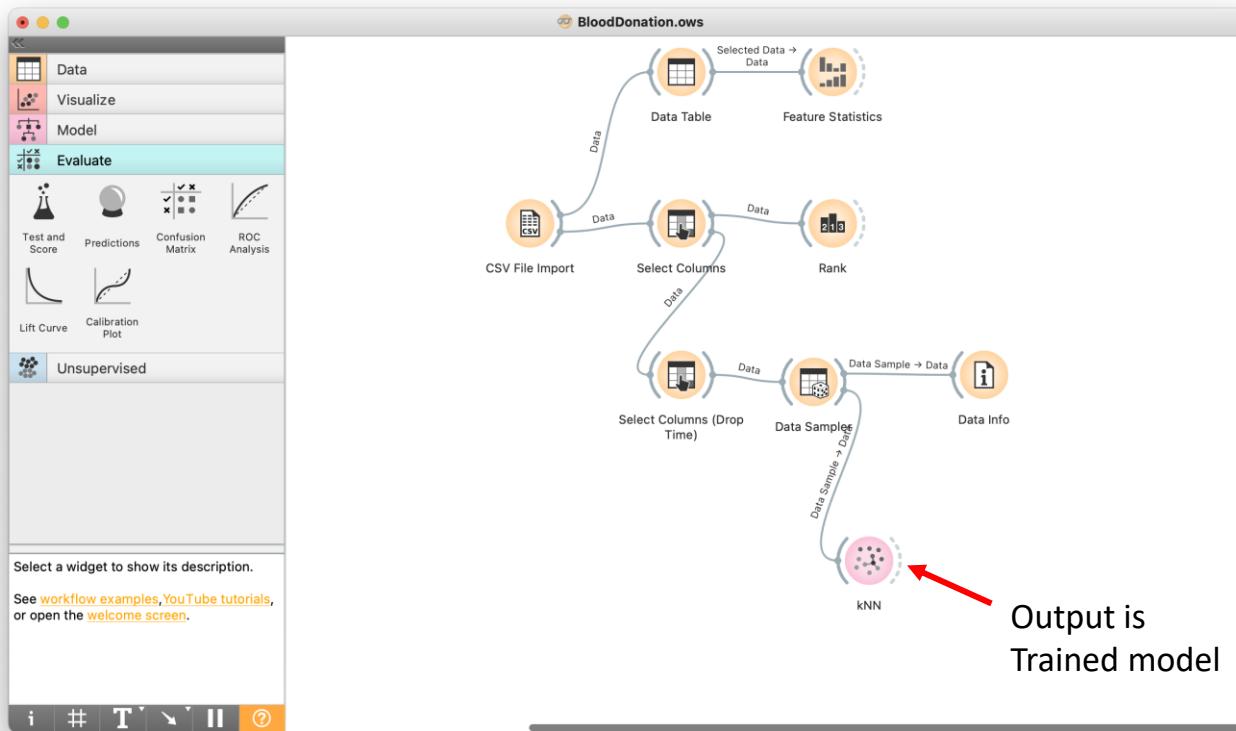
END OF DATA PREPARATION



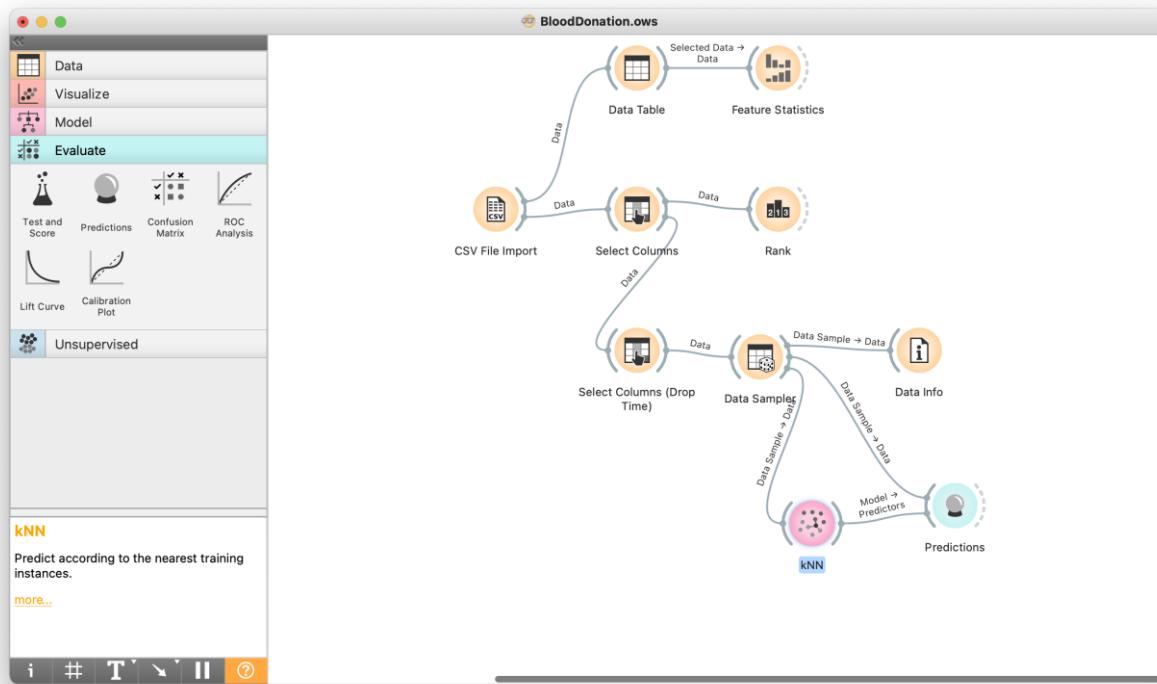
Use the KNN



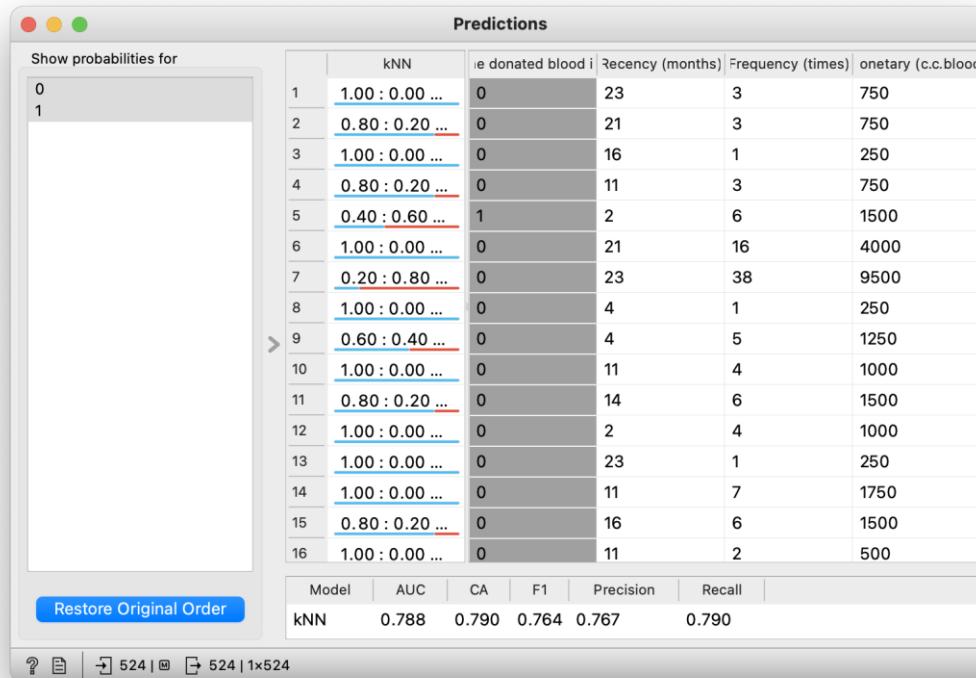
Connect the training data to KNN



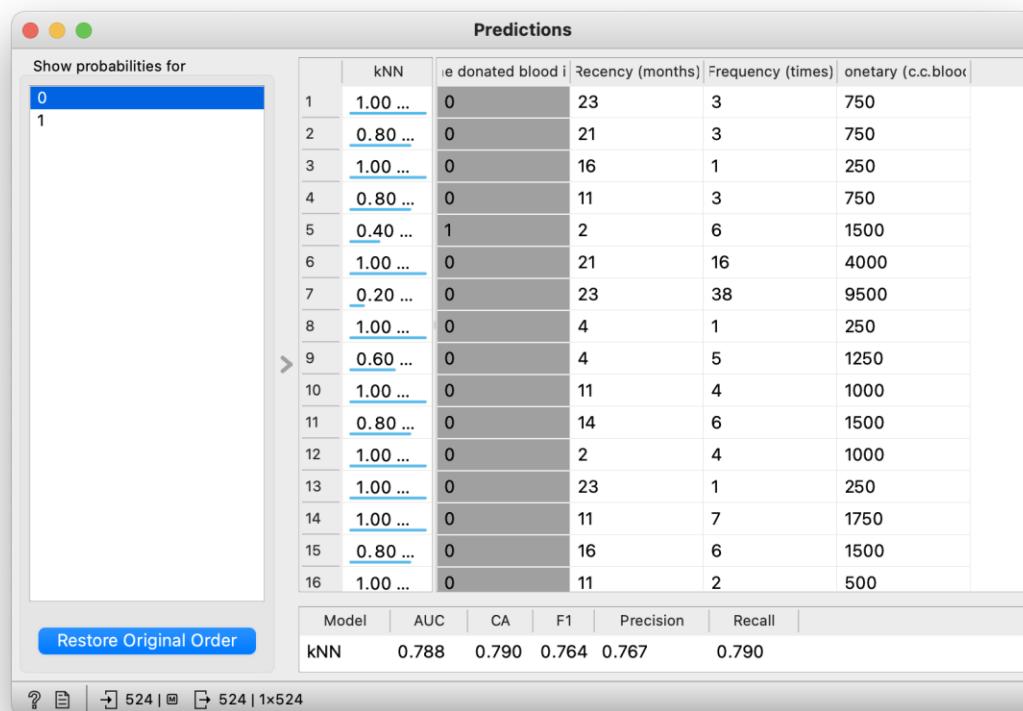
Use the Predictions Widget



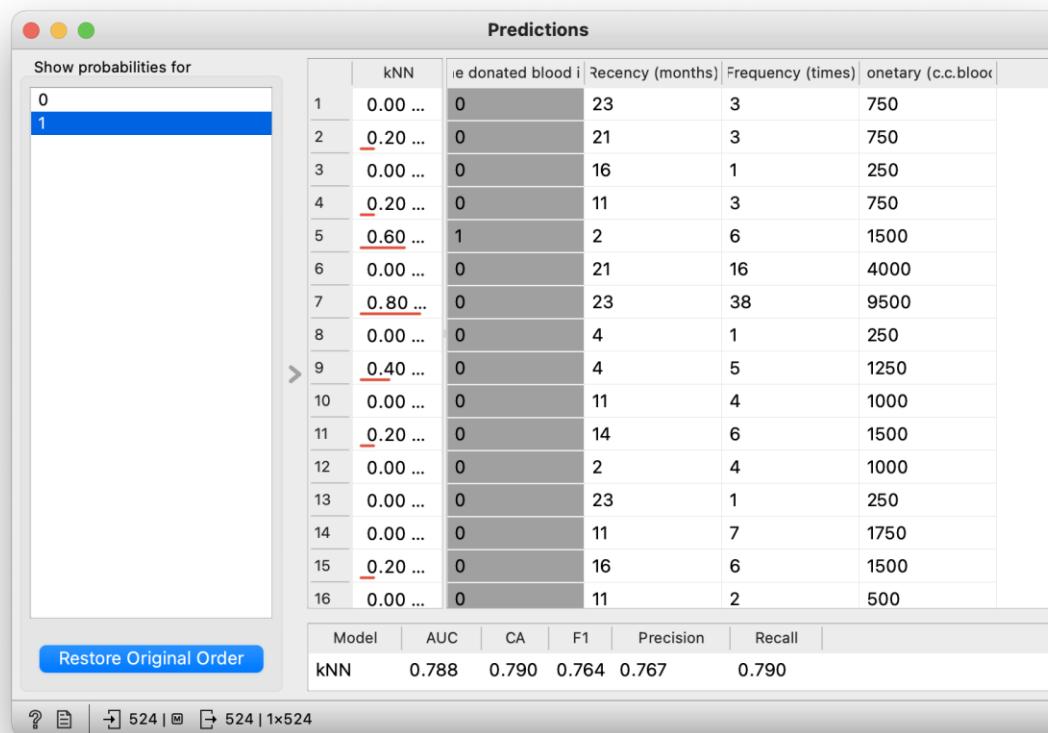
The Prediction uses the trained model from the KNN



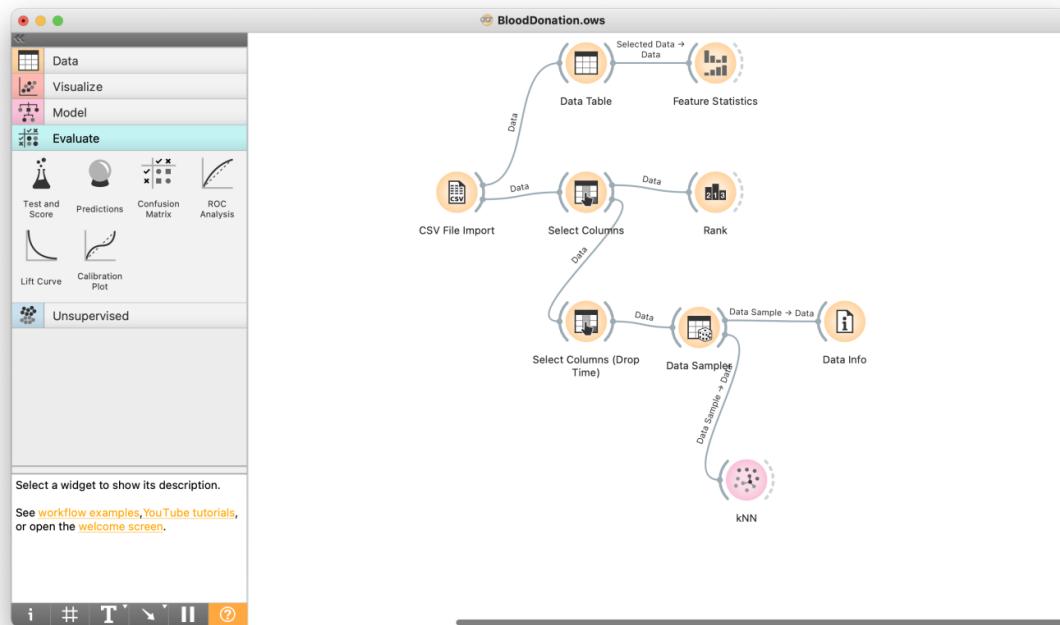
Predict 0



Predict 1



Remove the predictions/scorer



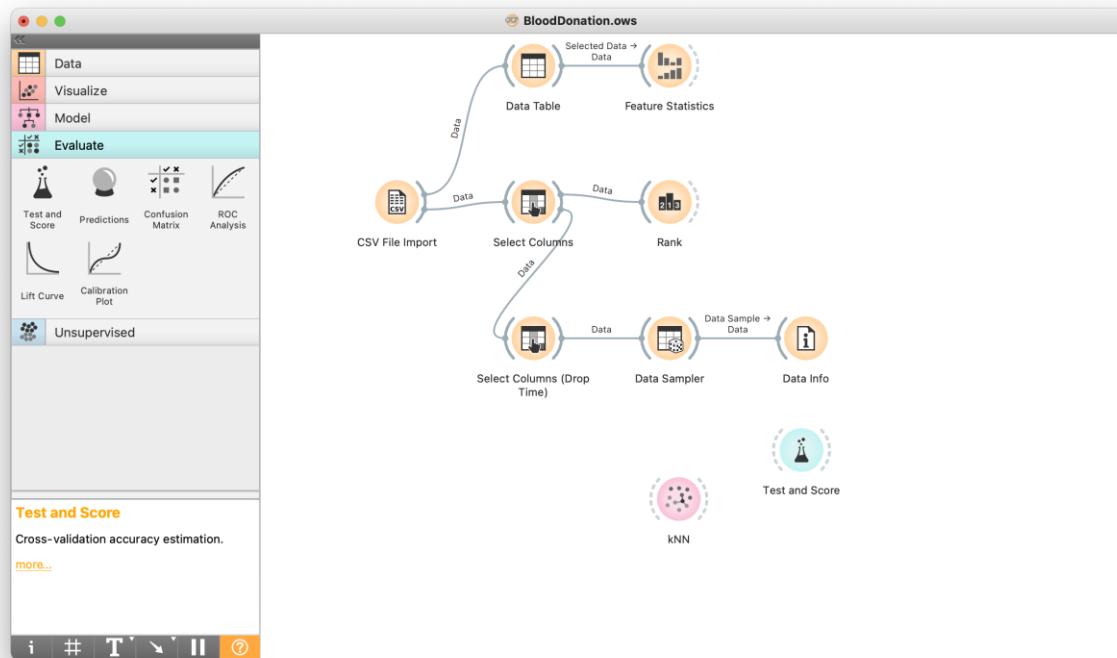
Evaluation

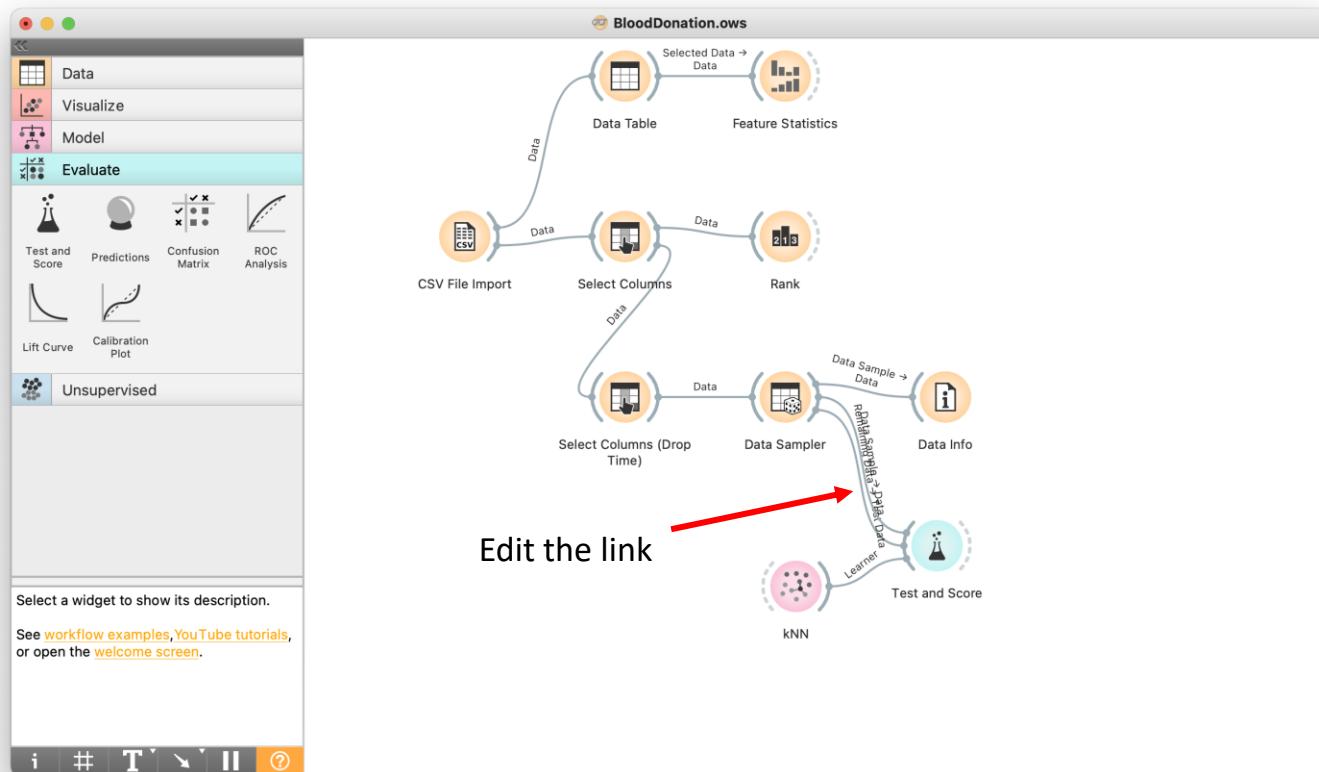


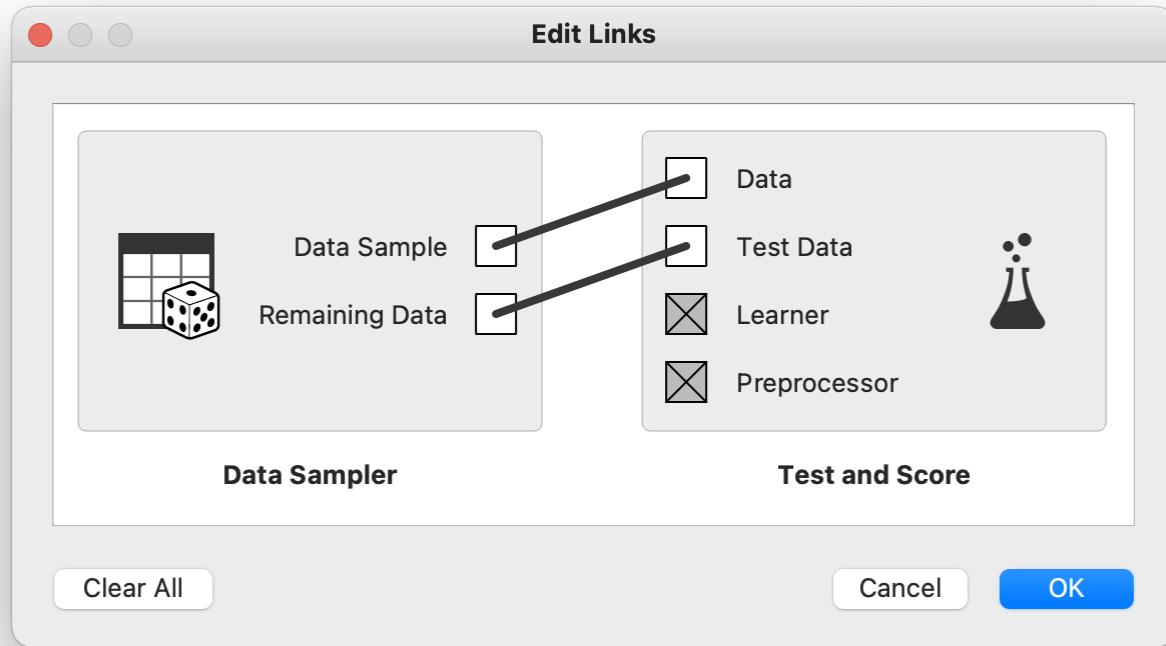
Connect
them to the
train and test
set



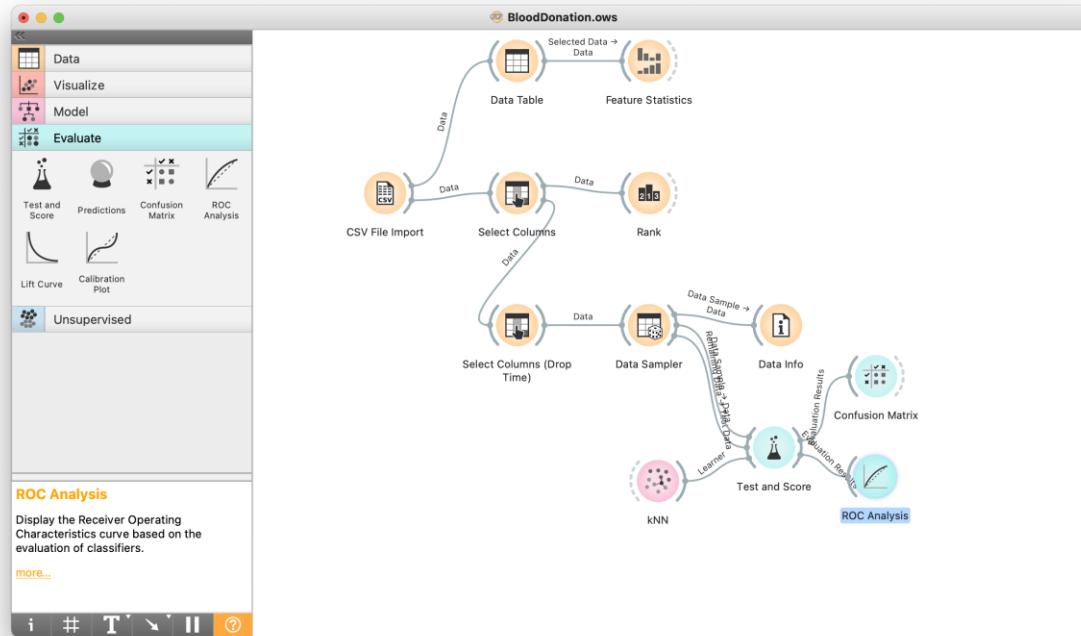
Use the Test and Score



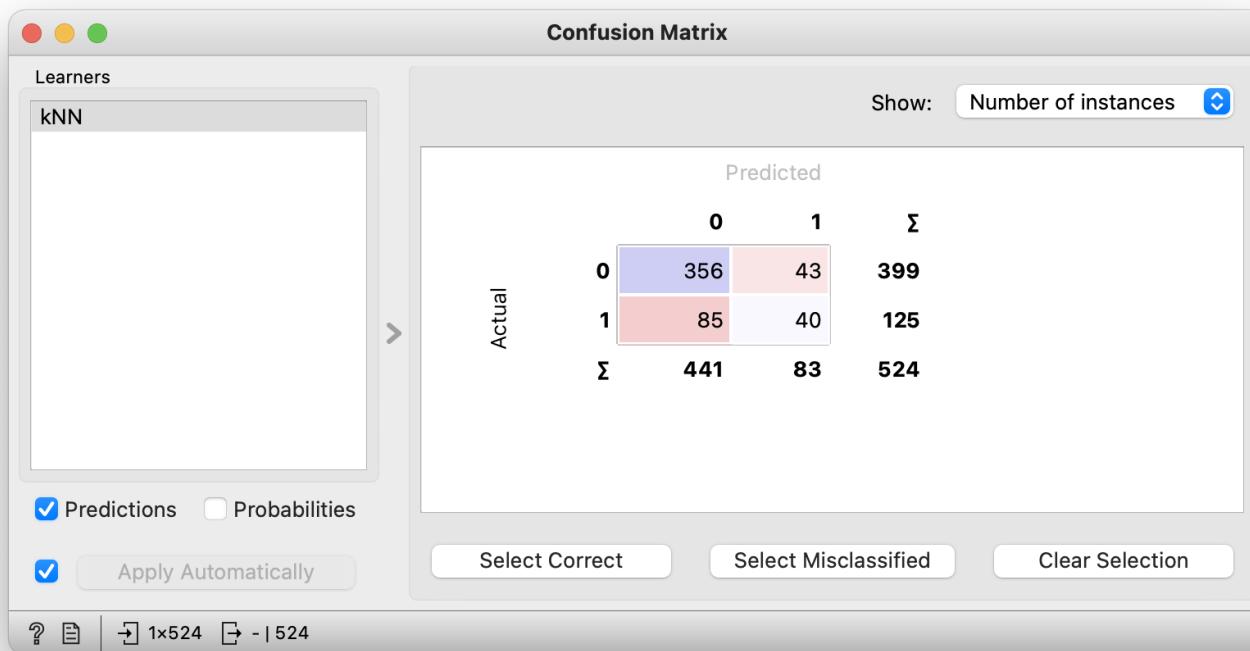




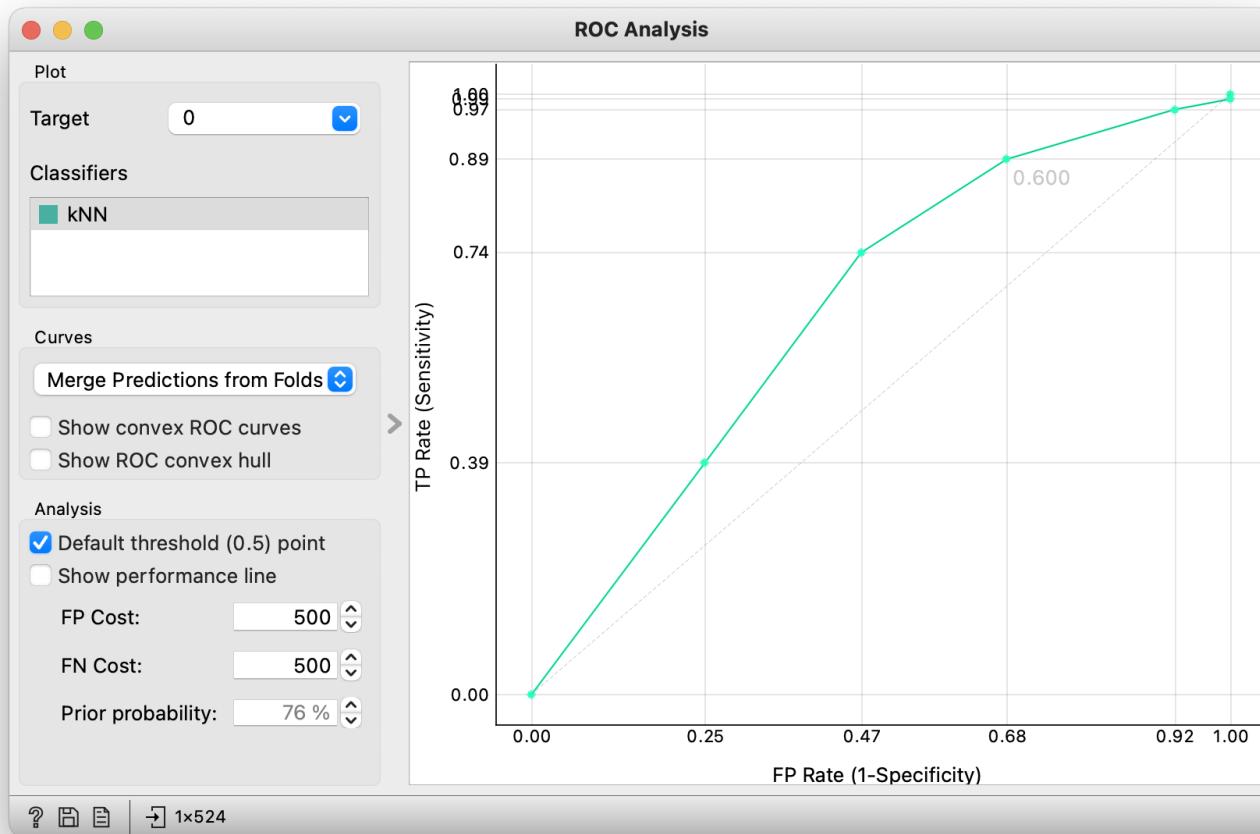
Use Confusion Matrix and ROC Analysis



Confusion Matrix

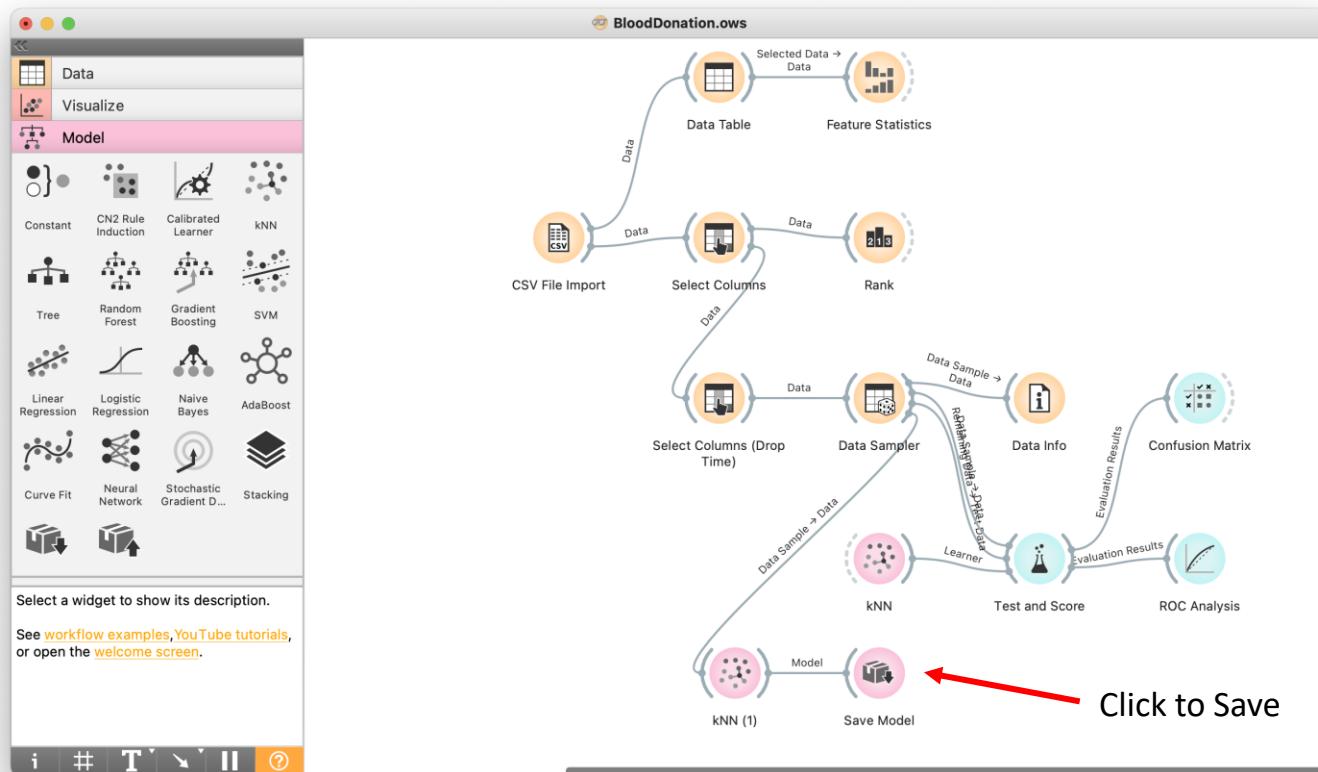


ROC

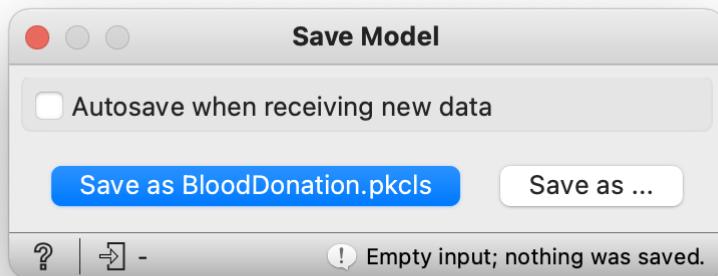


Deployment





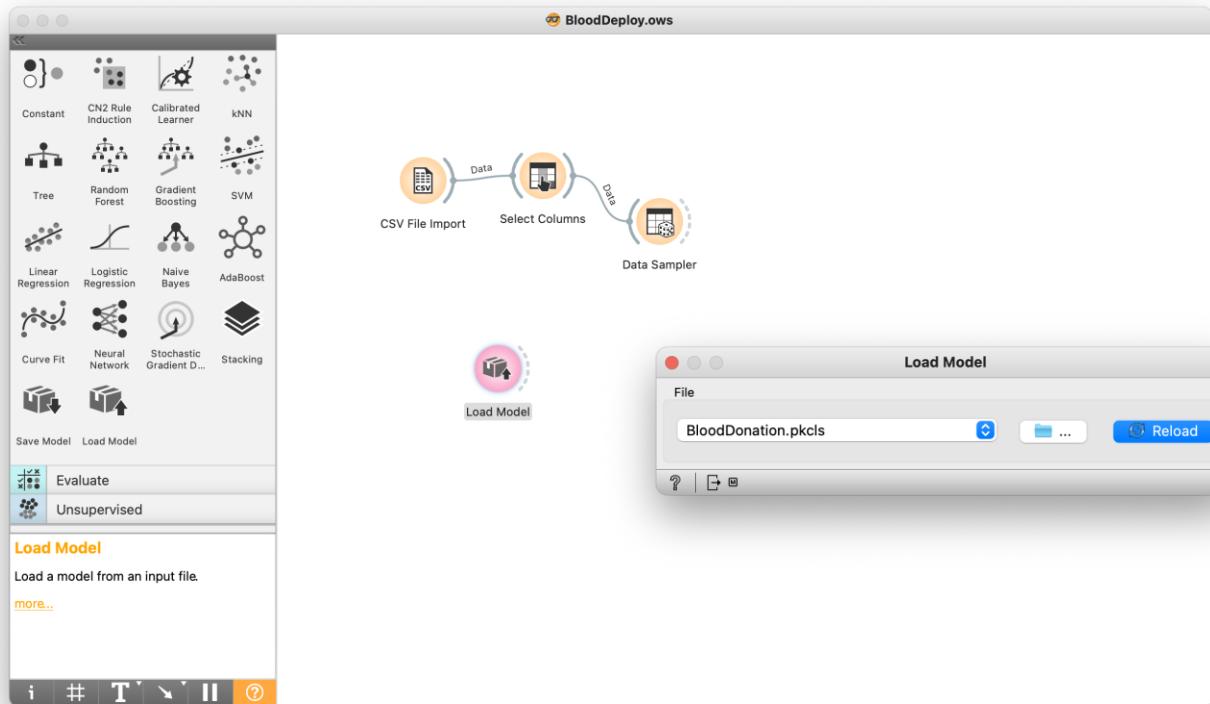
Save model



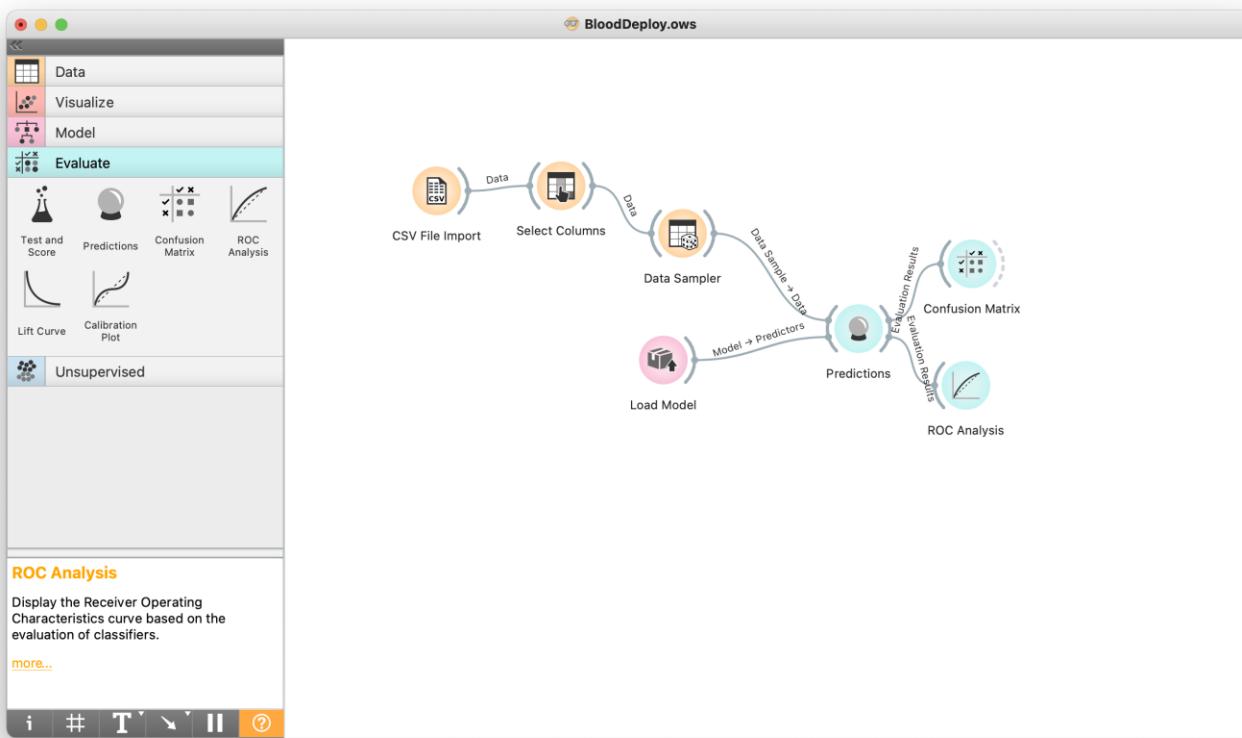


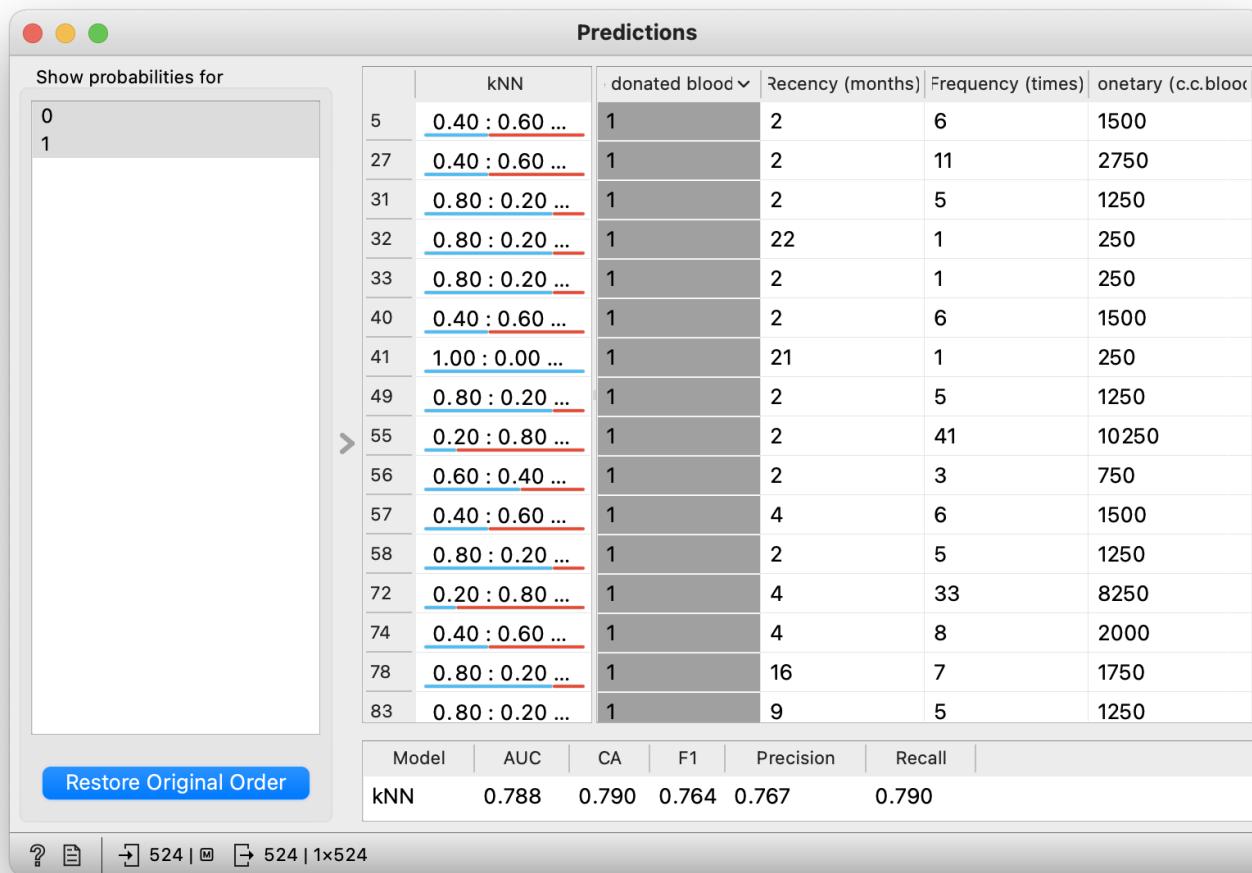
Create another
BloodDeploy
workflow

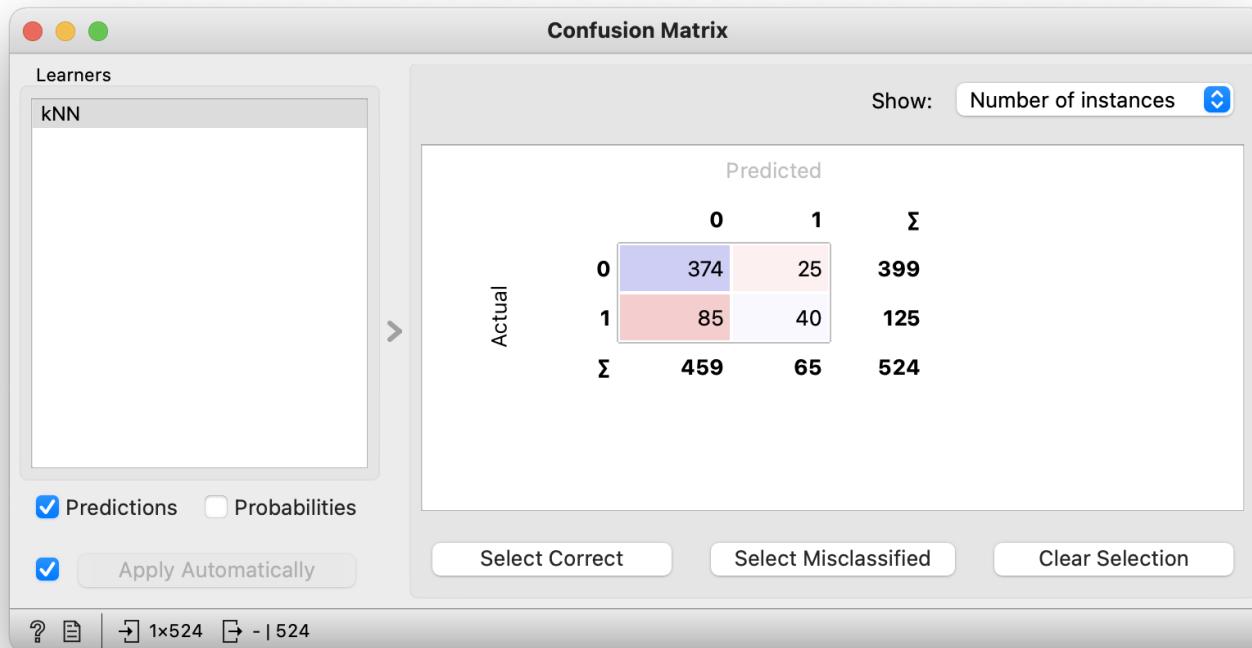
Use Load Model

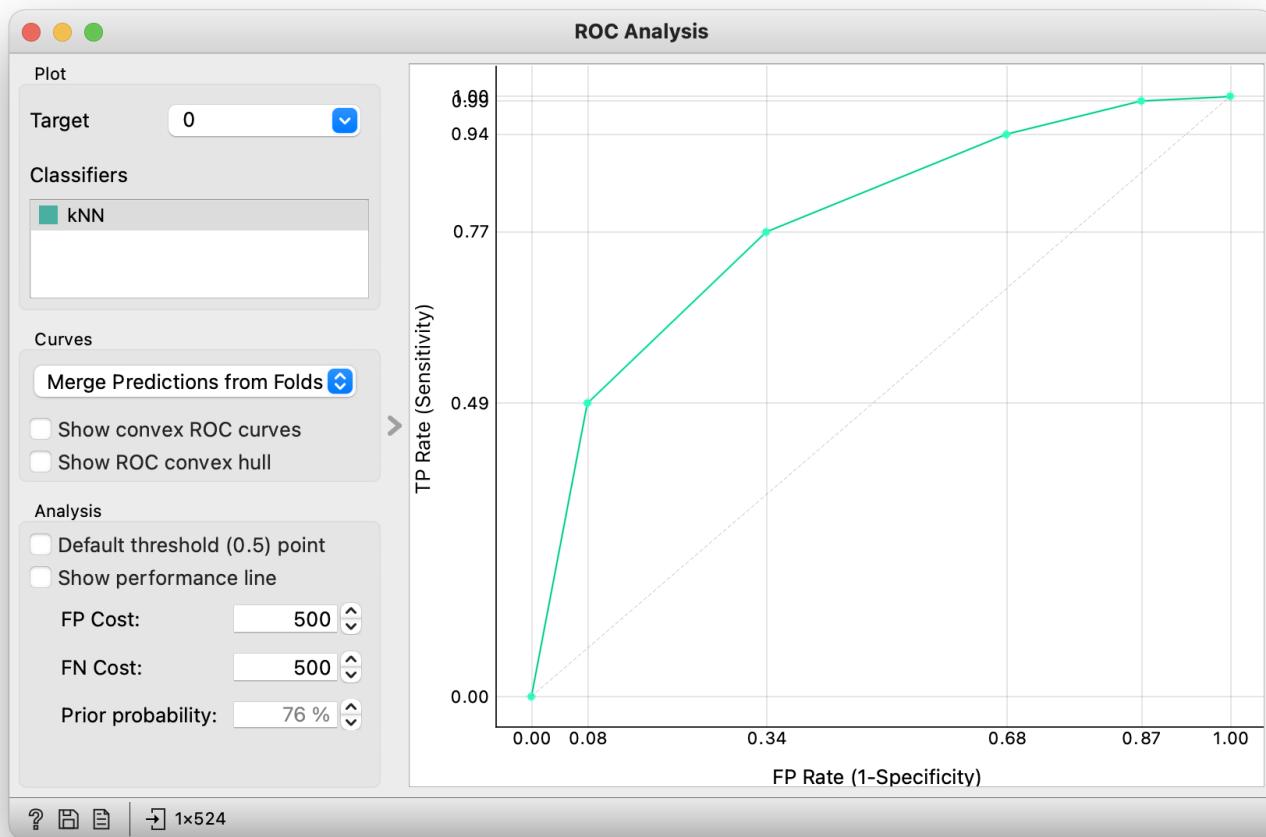


Use Predictions









Summary

We have learnt that:

- Orange Data Mining allows us to familiarize with the ML workflow and also quickly prototype some possible approaches for our modelling
- Data may need to be transformed/encoded before it can be fed to a machine algorithm
- Categorical variables can be encoded as dummy variables
- Sometimes it is useful to discretize numeric values
- Feature engineering is the part that needs real intelligence