ORANGE DATA MINING
STEPS for AI model to predict the penguin species .

# Step 1: Upload Dataset



**Data Acquisition**

# Step 1(a)

**This is the Welcome page that you will see when you first launch Orange.**
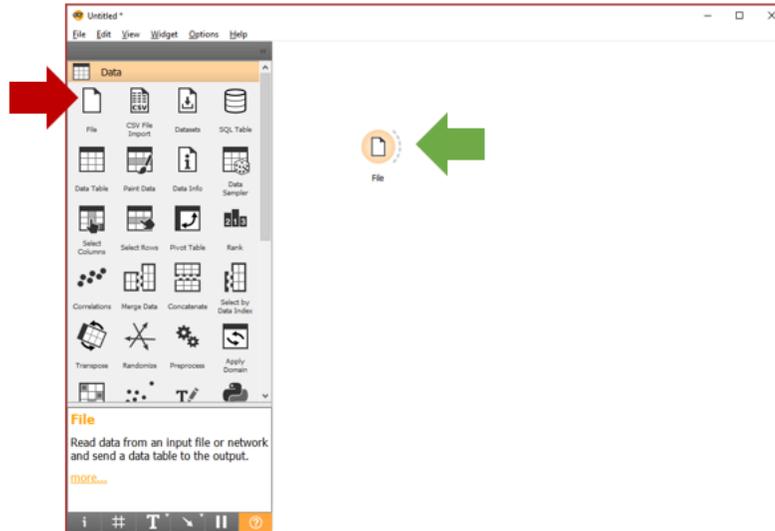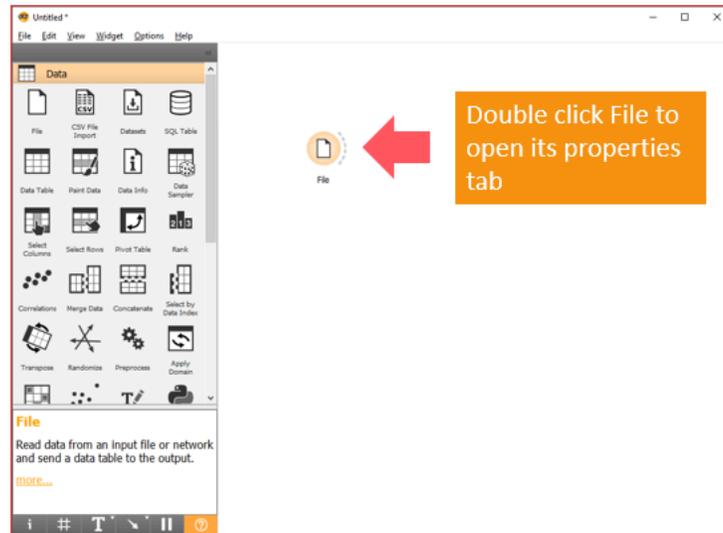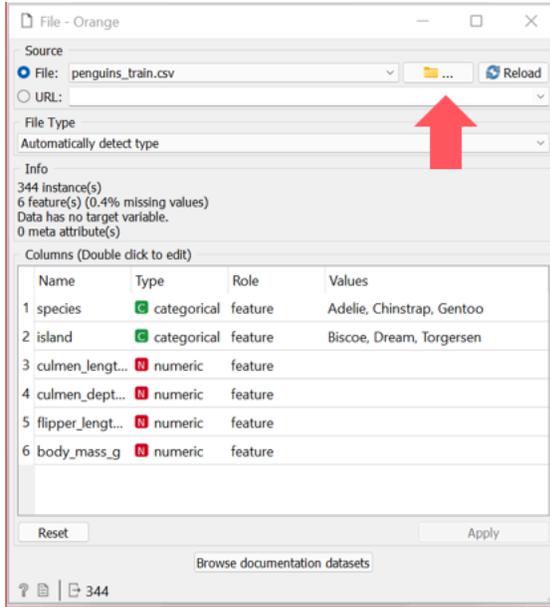
**Click on New.**

## Step 1(b)

Insert the File widget onto the canvas



## Step 1(c)

Double click File to open its properties tab

## Step 1(d)

**File - Orange**

**Source**
- File: penguins_train.csv | ... | Reload
- URL:

**File Type**
Automatically detect type

**Info**
344 instance(s)
6 feature(s) (0.4% missing values)
Data has no target variable.
0 meta attribute(s)

**Columns (Double click to edit)**

| | Name | Type | Role | Values |
|---|---|---|---|---|
| 1 | species | C categorical | feature | Adelie, Chinstrap, Gentoo |
| 2 | island | C categorical | feature | Biscoe, Dream, Torgersen |
| 3 | culmen_lengt... | N numeric | feature | |
| 4 | culmen_dept... | N numeric | feature | |
| 5 | flipper_lengt... | N numeric | feature | |
| 6 | body_mass_g | N numeric | feature | |

Reset | Apply

Browse documentation datasets

344

## Step 1(e)

**Open...**

« Desk... > Penguin E... | Search Penguin Example

Organize ▾ | New folder

- ⭐ Quick access
- OneDrive
- This PC
  - Desktop
  - Documents
  - Downloads
  - Music
  - Pictures
  - Videos
  - Local Disk (C:)
  - Local Disk (G:)

| Name | Date modified | Type |
|---|---|---|
| penguins_test.csv | 24-05-2022 12:58 AM | Microsof |
| penguins_train.csv | 24-05-2022 01:01 AM | Microsof |

File name: penguins_lter.csv | All re...ble files (*.basket *.bsk

Open | Cancel

## Step 1(f)

penguins_train.csv is uploaded



## Step 1(g)

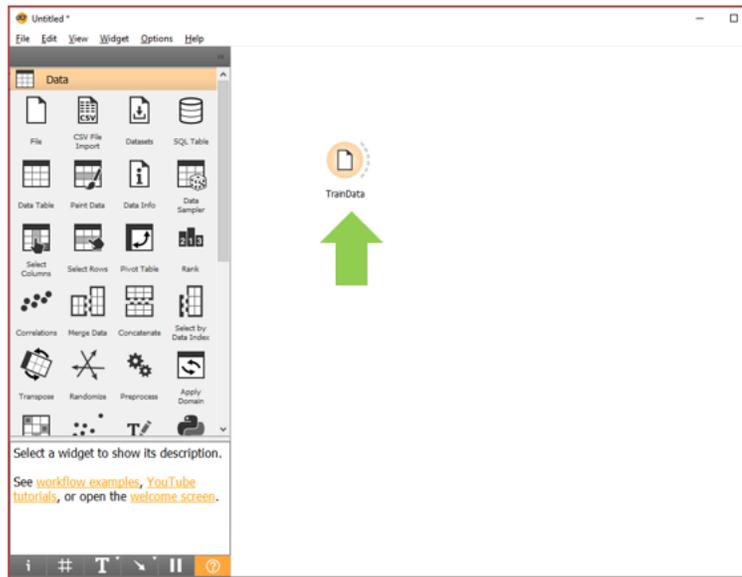Right-click on the widget File and click on Rename.
We will rename it to 'Train Data' so that we do not confuse it with the testing data later
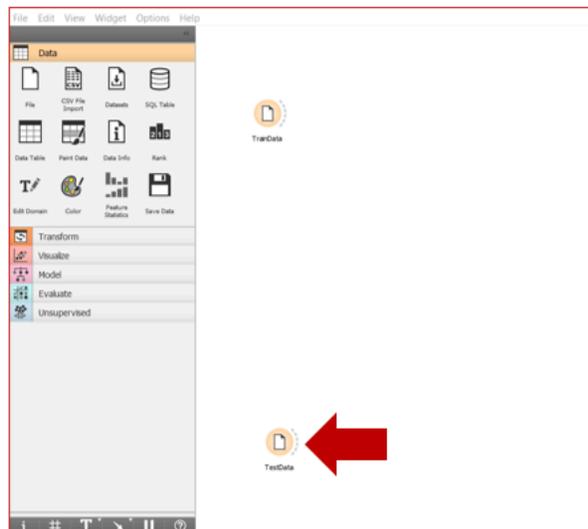
## Step 1(h)



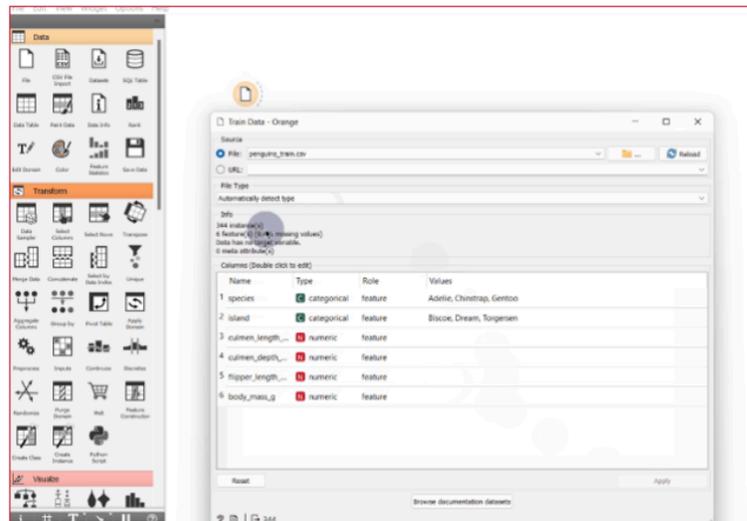The file name for the uploaded dataset has changed

## Step 1(i)



Repeat the same steps for 'Test Data' and upload testing dataset in the file widget

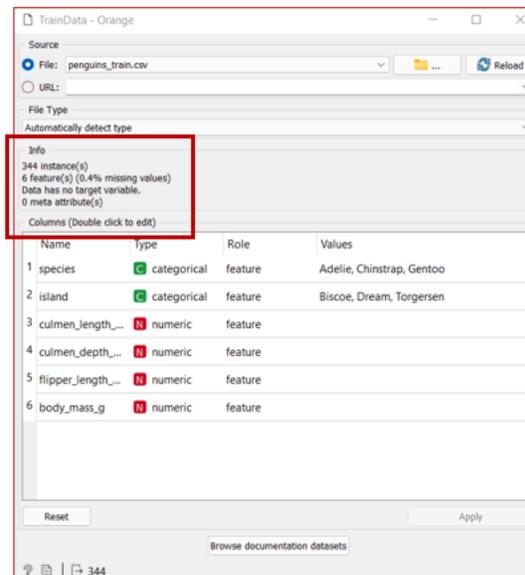After Data Acquisition, what should we do next?
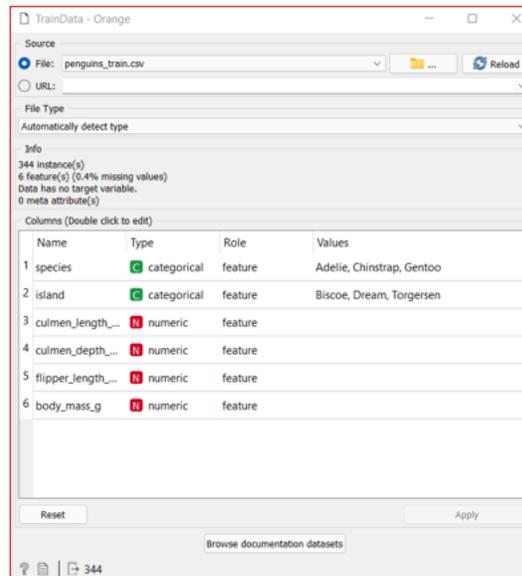
# Step 2: Clean Missing Data

## Step 2(a)

Check if there are any missing values
Notice that there are some missing values

## Step 2(b)

We will now look at another way to inspect on missing data.
Click X to close the pop up.



## Step 2(c)

Insert the widget Feature Statistics onto the canvas
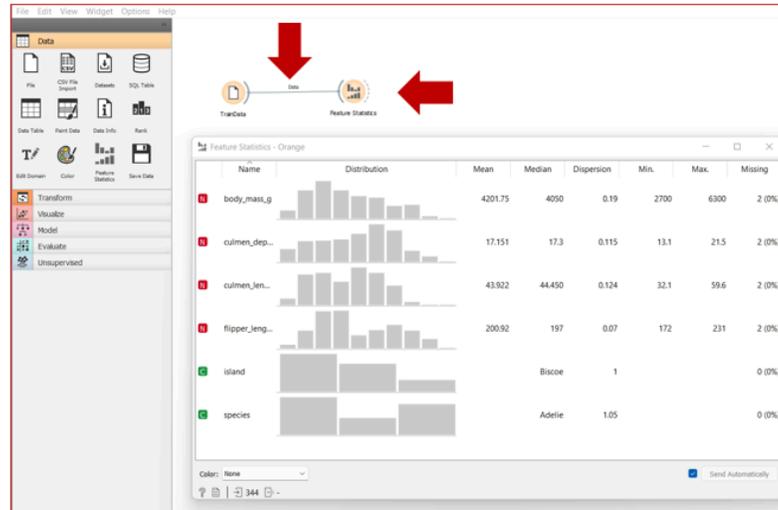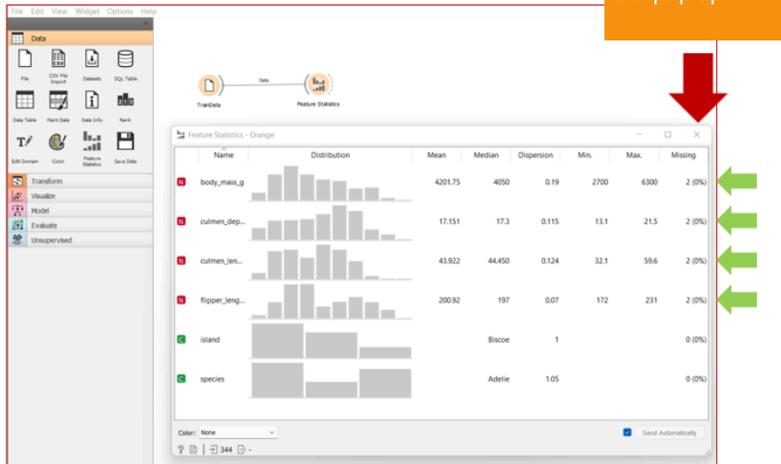
## Step 2(d)

1. Connect widget 'Train Data' to widget Feature Statistics. We can do that by dragging the output from 'Train Data' to the input of Feature Statistics.

2. Double-click on Feature Statistics to see the results.



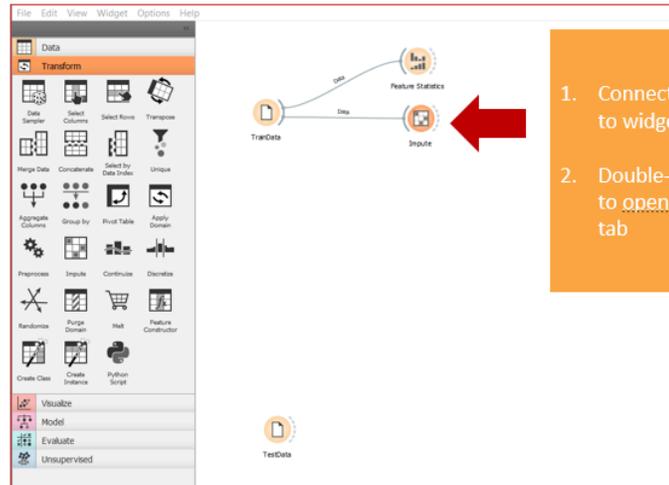What is the mean value of culmen_length feature?

# Step 2(f)



Click X to close the pop up

Notice that arrows are pointing to the features with missing values here
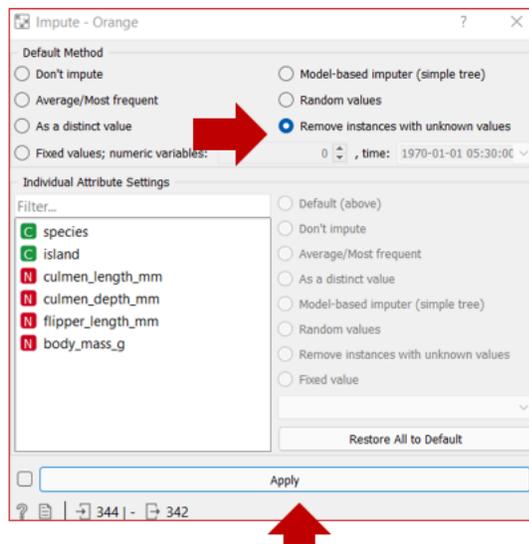
# Step 2(g)



Insert the impute widget onto the canvas

## Step 2(h)



1. Connect widget TrainData to widget Impute

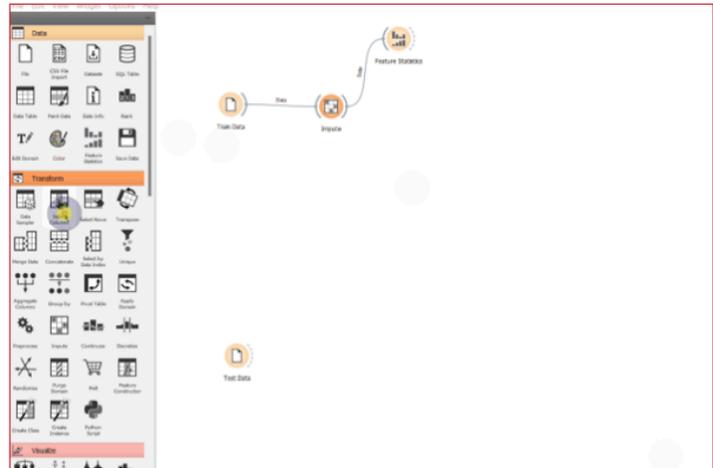2. Double-click on Impute to open up the properties tab

## Step 2(i)

# Step 2(j)



1. Connect the output of Impute to the input of the existing Feature Statistics (the previous connection between TrainData and Feature Statistics has been removed because only accepts one input)

2. Double-click on the Feature Statistics to see the output.

Now that the data is clean and without any missing values, what next?

# Step 3: Select Target Label

## Step 3(a)
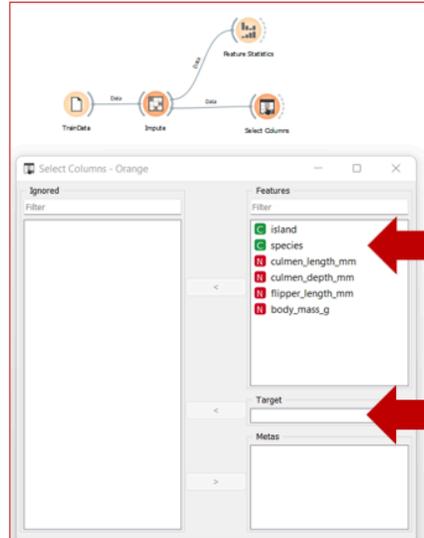
Insert the Select columns widget onto the canvas

1. Connect with Impute widget
2. Double click on the widget



From TrainData, you would have noticed that the Feature Type for most of the columns is Numeric Feature. In supervised learning models, we have both the features and the labels. The labels are the output. Therefore, we need to define an output for our Palmer Penguin model. We will assign species as our label since that is what we want to identify.

Therefore, we will change the Feature Type for species, from Categorical Feature to Categorical Label. To do that, we will be using Select Columns.
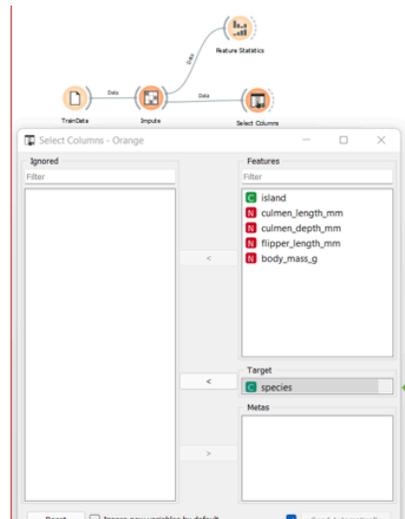
## Step 3(b)



A window displaying all the features will appear

Drag the 'species' feature to the 'Target' box

## Step 3(b)
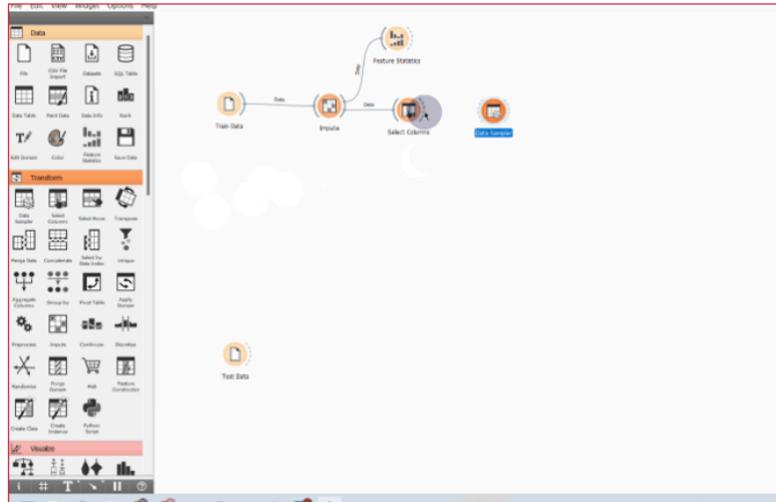


'species' is the Target label

After choosing a target label, we need to split the data
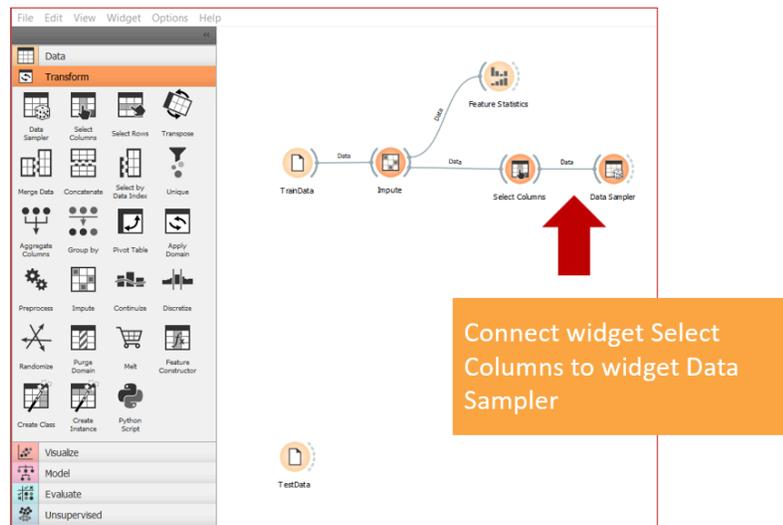
# Step 4: Data Sampler

## Step 4(a)

Insert the Data Sampler widget onto the canvas

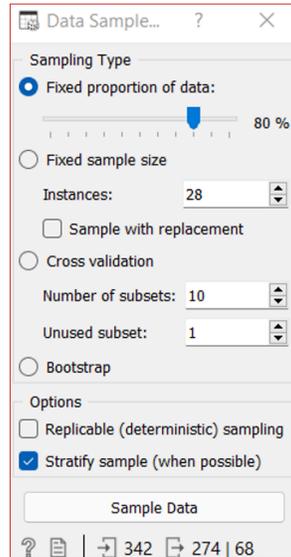

## Step 4(b)

Connect widget Select Columns to widget Data Sampler



Connect widget Select Columns to widget Data Sampler. We can do that by dragging the output from Select Columns to the input of Data Sampler.

After the connection is made, double-click on Data Sampler to open the properties tab.
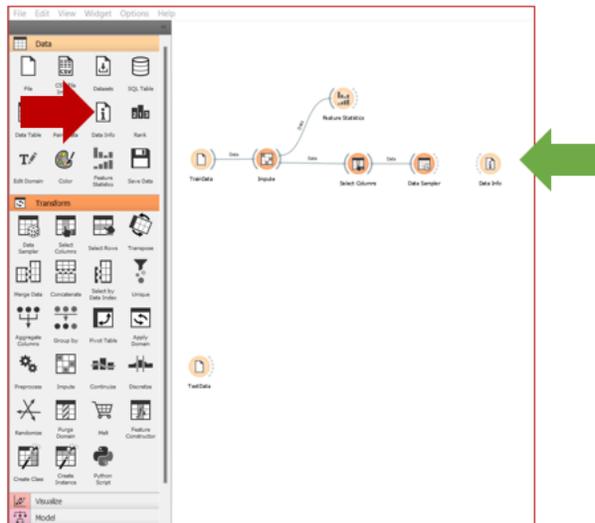
## Step 4(c)



1. Set slider to 80%
2. Click on Sample Data to effect the changes
3. Click X to close the pop-up

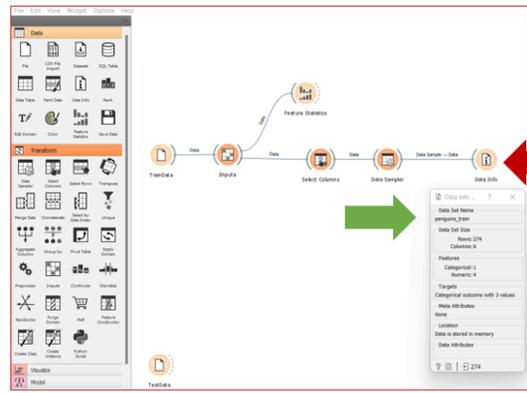How do we know if the data is actually split or not?

## Step 4(d)

Insert the Data Info widget onto the canvas



Let's inspect on how the data is being split through Data Sampler.
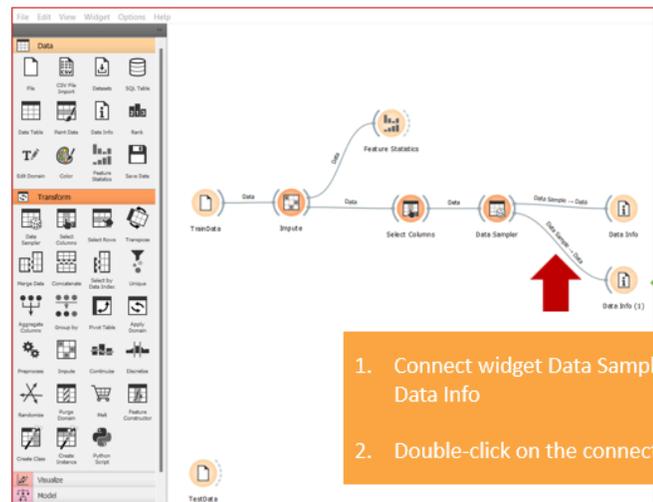
We will be using Data Info.

Step 4(e)



1. Connect widget Data Sampler to widget Data Info

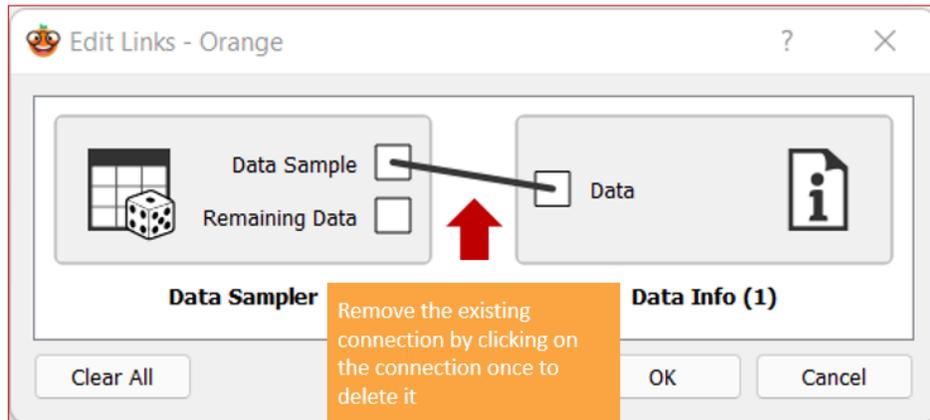2. Double-click on Data Info to open the properties tab

Step 4(f)



1. Connect widget Data Sampler to widget Data Info
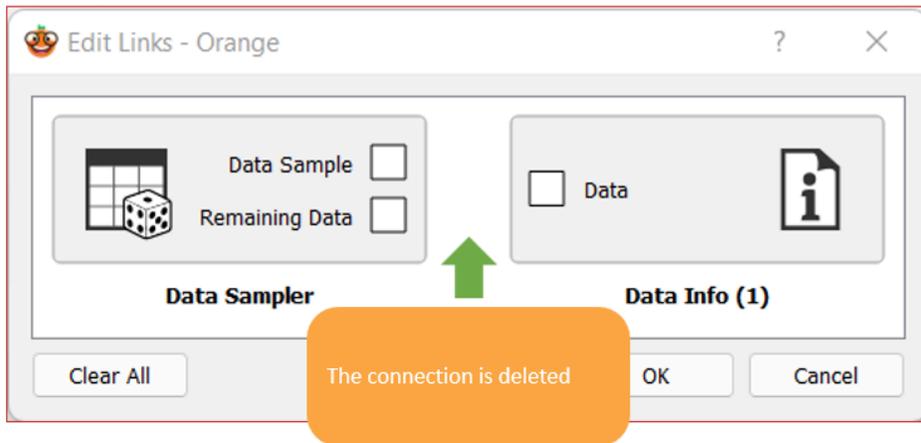
2. Double-click on the connection

Connect widget Data Sampler to the second widget Data Info. We can do that by dragging the output from Data Sampler to the input of the second Data Info.

Take note of the connection name. We will change this. Double-click on the connection.
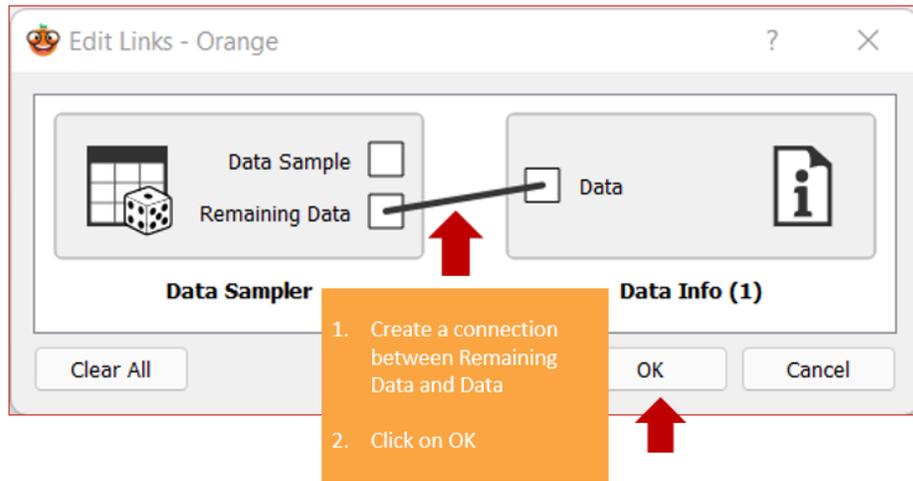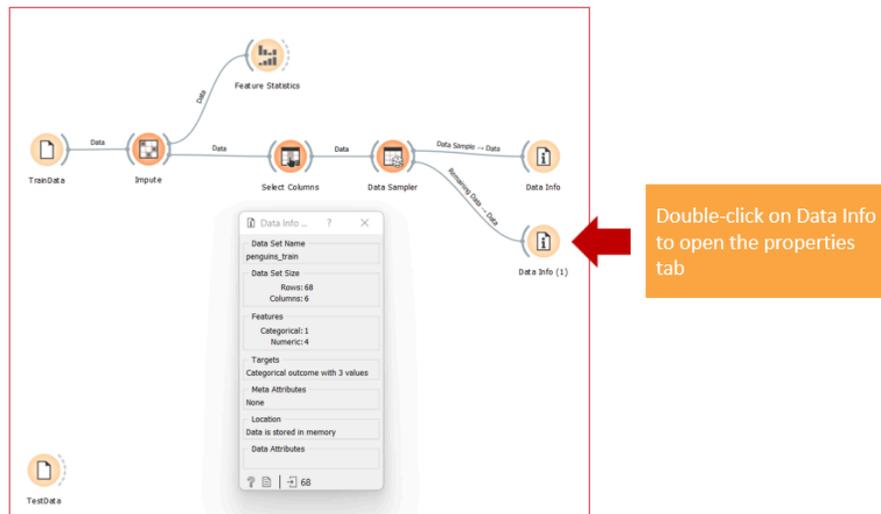
## Step 4(g)



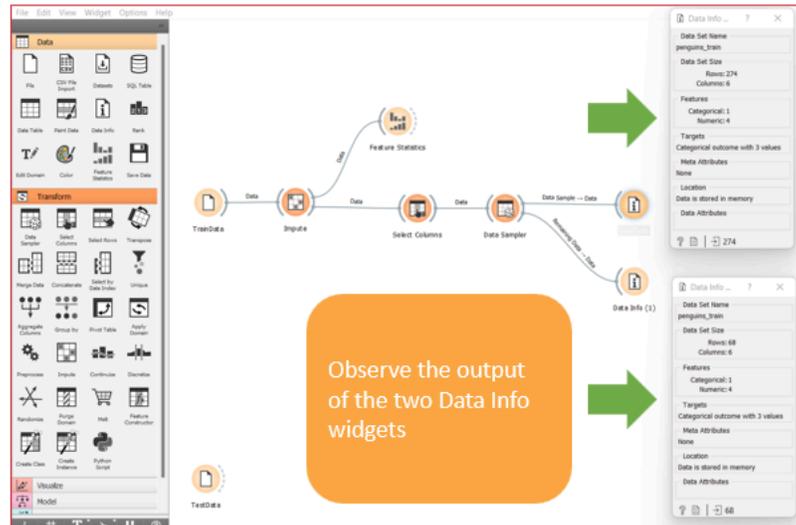## Step 4(g)

## Step 4(h)



**Edit Links - Orange**

Data Sample ☐

Remaining Data ☐ ━━━━ ☐ Data

**Data Sampler**          **Data Info (1)**

Clear All          OK          Cancel

1. Create a connection between Remaining Data and Data

2. Click on OK

## Step 4(i)



Feature Statistics

TrainData — Data — Impute — Data — Select Columns — Data — Data Sampler — Data Sample → Data — Data Info

Remaining Data → Data

Data Info (1)

**Data Info ...**

Data Set Name
penguins_train

Data Set Size
Rows: 68
Columns: 6

Features
Categorical: 1
Numeric: 4

Targets
Categorical outcome with 3 values

Meta Attributes
None

Location
Data is stored in memory

Data Attributes

68

TestData

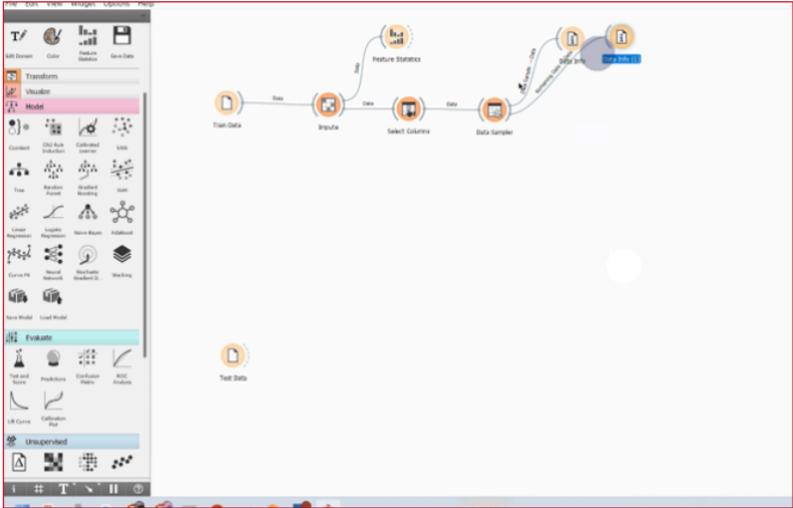Double-click on Data Info to open the properties tab

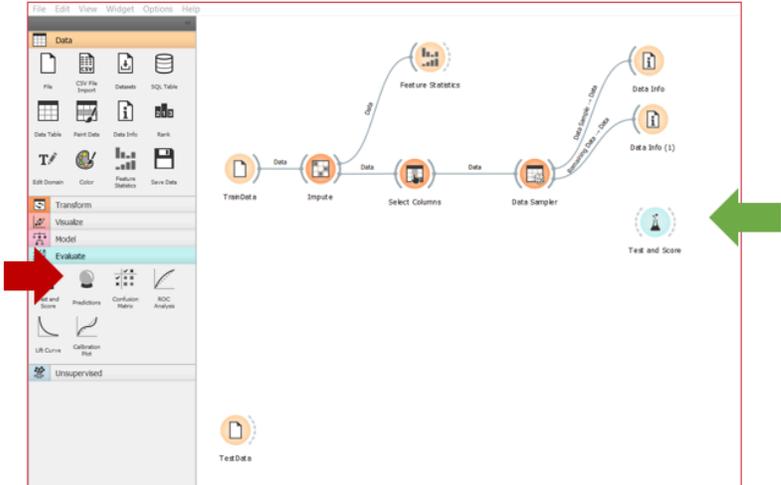# Step 4(j)



What do we do after having split the data?

# Step 5: Train Model



Modeling

# Step 5(a)

Insert the Test and Score widget onto the canvas
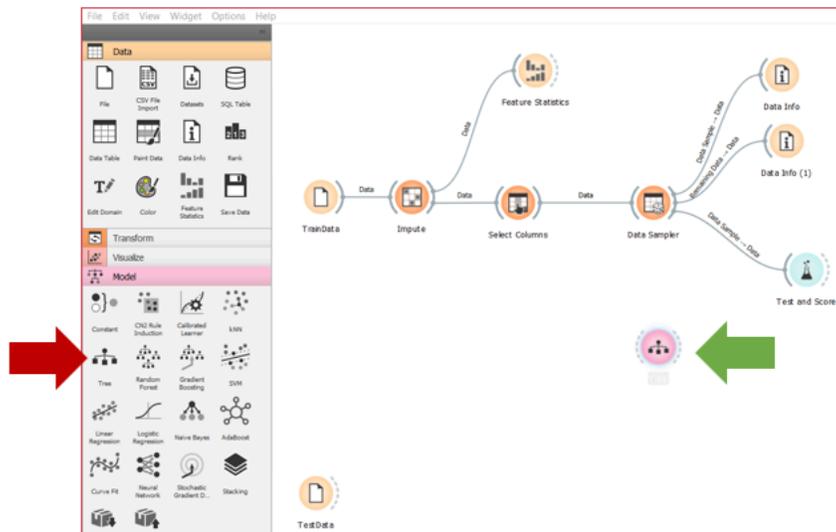
## Step 5(b)



1. Connect widget Data Sampler to widget Test and Score
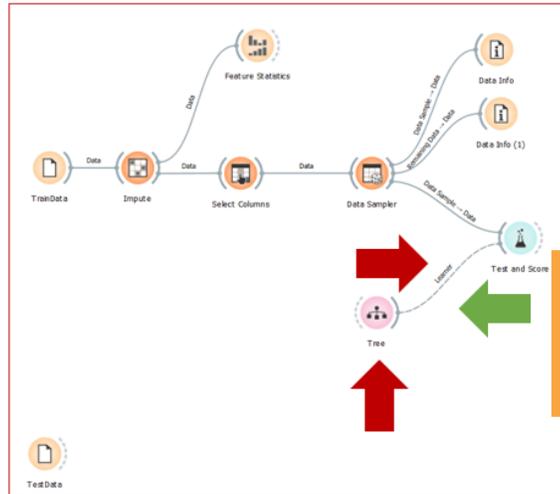
2. Make sure the connection says Data Sample -> Data

## Step 5(c)

Insert the widget Tree into the canvas and put it to the left of widget Test and Score
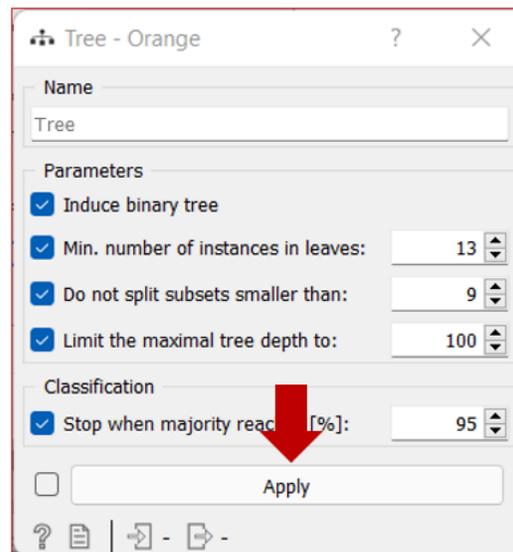
# Step 5(d)

1. Connect widget Tree to widget Test and Score

2. Double click on the model icon to open its properties

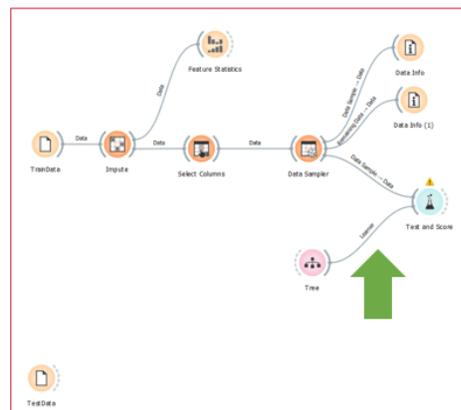If the connection is dotted it means the model has not been applied yet

## Step 5(e)

1. Click on Apply if automatically apply has not been checked

2. Click X to close the pop-up window

**Tree - Orange**

### Name
Tree

### Parameters
- [x] Induce binary tree
- [x] Min. number of instances in leaves: 13
- [x] Do not split subsets smaller than: 9
- [x] Limit the maximal tree depth to: 100

### Classification
- [x] Stop when majority reach [%]: 95

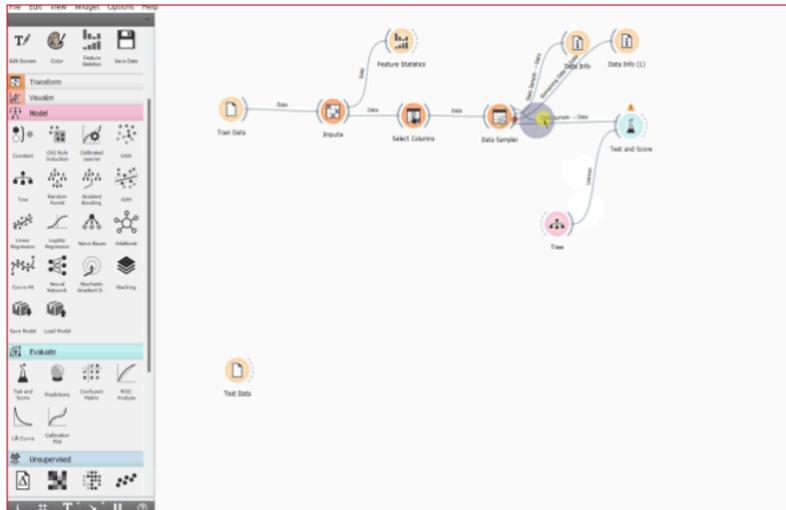Apply

## Step 5: Train Model

The model is ready

If there is a warning sign on the test and score widget, we will need to change its properties. Let's talk about it next
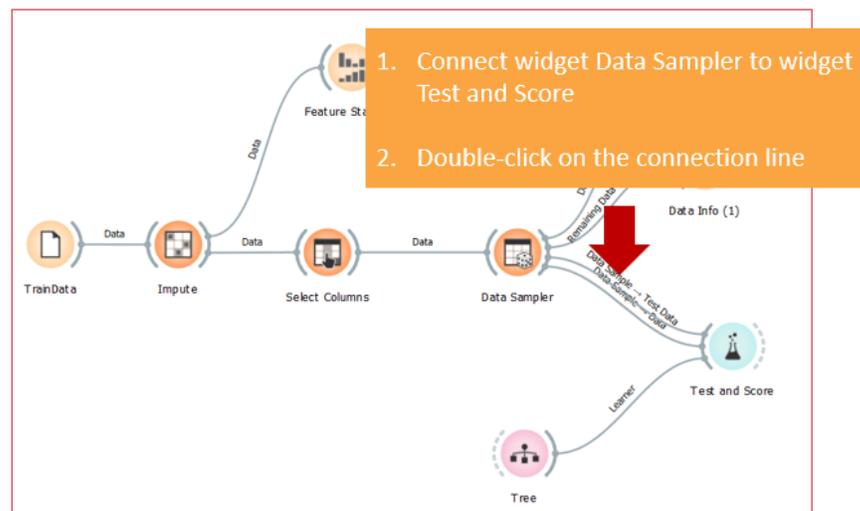
After creating a model, we need to test the model and check its accuracy
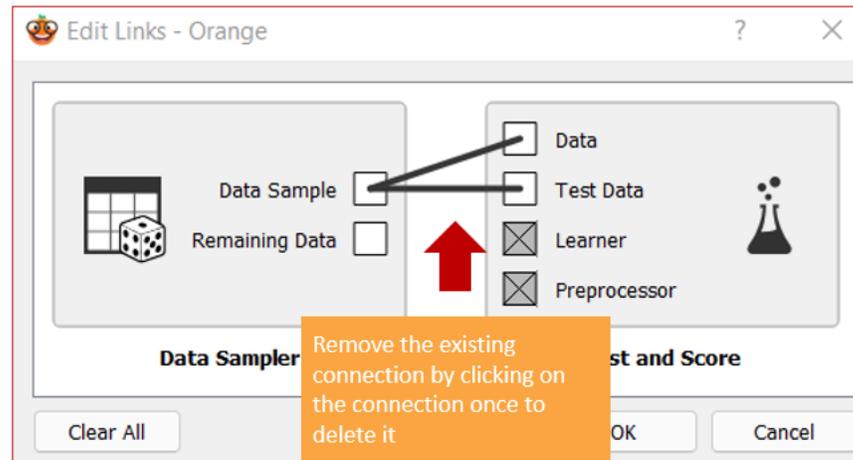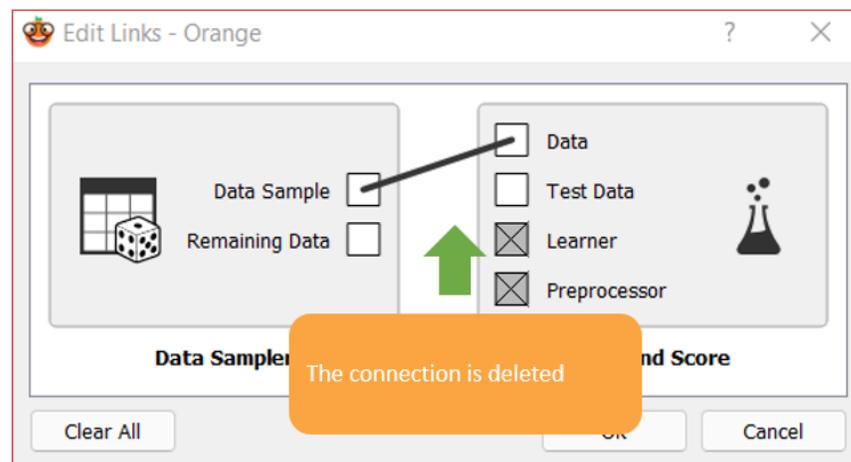
# Step 6: Evaluate Model

# Step 6(a)



1. Connect widget Data Sampler to widget Test and Score

2. Double-click on the connection line

## Step 6(b)



Remove the existing connection by clicking on the connection once to delete it

## Step 6(b)



The connection is deleted

# Step 6(c)



Edit Links - Orange

Data Sample → Data

Remaining Data → Test Data

Learner

Preprocessor

**Data Sample** **Test and Score**

Clear All    OK    Cancel

1. Create a connection between Remaining Data and Test Data

2. Click on OK

Step 6(d)



We can observe the different scores for the model

Click on X to close the pop-up

Evaluation results for target (None, show average over classes)

| Model | AUC | CA | F1 | Precision | Recall |
|-------|-----|-----|-----|-----------|--------|
| Tree | 0.970 | 0.926 | 0.928 | 0.931 | 0.926 |



Through Evaluation we know
if a model is good or bad

Let's try a couple of other classification algorithms

## Step 6(h)



Insert the widget Random Forest into the canvas

## Step 6(i)



Connect the Random Forest widget to the Test and Score widget

## Step 6(j)

This pop up will appear

We can observe the individual scores for each algorithm here

Double click on the Test and Score widget

## Step 6(k)

We can choose the methods for evaluation in this pop up

**Test and Score - Orange**

- ○ Cross validation
  - Number of folds: 5
  - ☑ Stratified
- ○ Cross validation by feature
- ○ Random sampling
  - Repeat train/test: 10
  - Training set size: 66 %
  - ☑ Stratified
- ○ Leave one out
- ○ Test on train data
- ○ Test on test data

Evaluation results for target (None, show average over classes)

| Model | AUC | CA | F1 | Precision | Recall |
|---|---|---|---|---|---|
| Tree | 0.958 | 0.931 | 0.930 | 0.930 | 0.931 |
| Random Forest | 0.997 | 0.982 | 0.982 | 0.982 | 0.982 |

Compare models by: Area under ROC curve    ☐ Negligible diff.:    0.1

| | Tree | Random... |
|---|---|---|
| Tree | | 0.071 |
| Random Forest | 0.929 | |

Table shows probabilities that the score for the model in the row is higher than that of the model in the column. Small numbers show the probability that the difference is negligible.

274 | 68 | - | 274 | 2×274

## Step 6(I)



Now that we have found which model gives us the best results, we can use that one!

# Step 7: Predictions



Predictions

# Step 7(a)

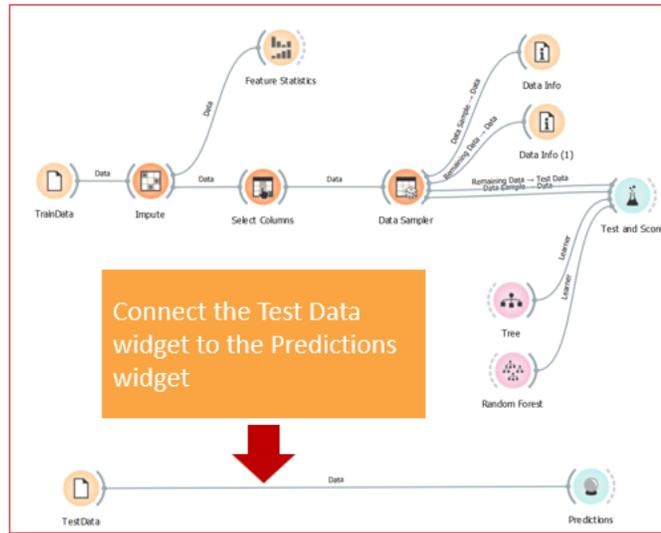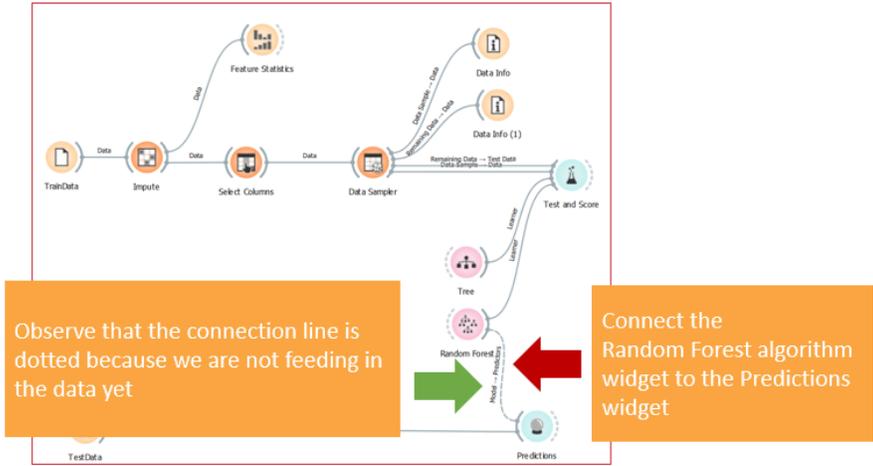Insert the widget Predictions into the canvas

# Step 7(b)



Connect the Test Data widget to the Predictions widget

## Step 7(b)



Observe that the connection line is dotted because we are not feeding in the data yet

Connect the Random Forest algorithm widget to the Predictions widget

## Step 7(c)



Connect the Data Sampler widget to the Random Forest algorithm widget

Observe that the connection line has now become solid

## Step 7(d)



Predictions pop up appears

Double click on the Predictions widget

## Step 7: Predictions



Random Forest Predictions

Actual Species

Observe that the predictions made for Chinstrap by Random Forest are false

Random Forest is classifying Chinstrap as Adelie

Since the Random Forest algorithm is not working well with one of the species, let's use another algorithm

## Step 7(e)



Connect the Data Sampler widget to the Tree widget

Now we are using multiple models at the same time

## Step 7: Predictions



Random Forest Predictions

Tree Predictions

Actual Species

Observe that the predictions made for Chinstrap by Tree are correct

This suggests that some models give better results than others