

UNPHU

Universidad Nacional
Pedro Henríquez Ureña

Dirección de
Investigación

Semana Dominicana de Ciencia y Tecnología 2023 y XVIII Congreso Internacional de Investigación Científica (XVIII CIC-2023)



Taller: Implementando algoritmos de Aprendizaje Automático sin la necesidad de programar

Ing. Christian López, PhD

Profesor de Ciencias Computacionales, con afiliación en Ingeniería Mecánica, Lafayette College, USA

sites.lafayette.edu/lopezbec

lopezbec@Lafayette.edu

[@C_LopezB](https://twitter.com/C_LopezB)

github.com/lopezbec

UNPHU

Universidad Nacional
Pedro Henríquez Ureña



LAFAYETTE
COLLEGE

INTELIGENCIA ARTIFICIAL

El estudio de agentes “inteligentes”

- Visión Artificial
- Procesamiento de Lenguaje Natural
- Investigación de Operaciones
- ...

APRENDIZAJE AUTOMÁTICO (Machine Learning)

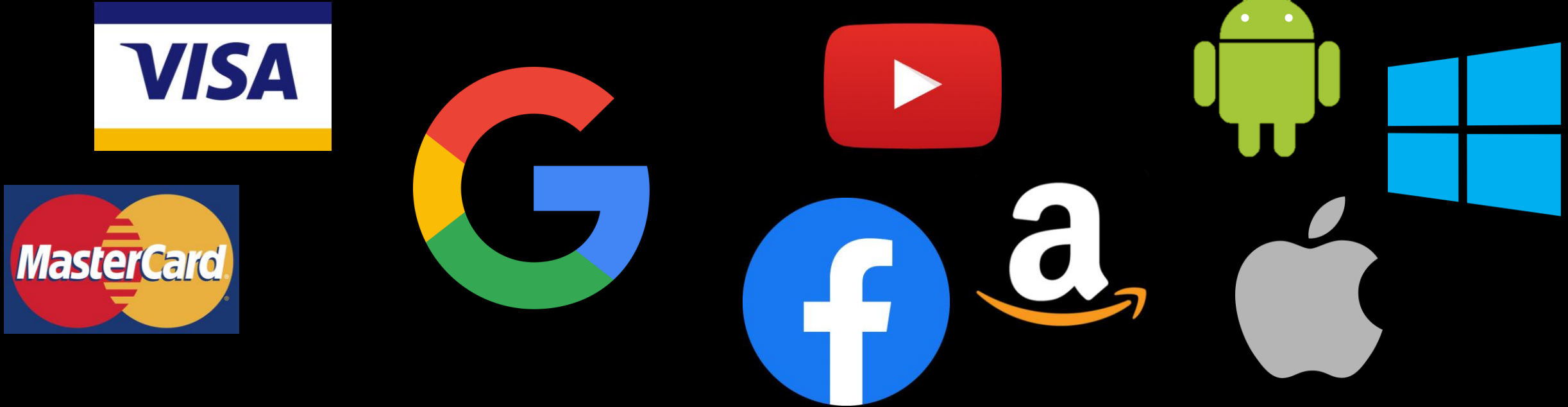
Algoritmos que desempeñan mejor al ser expuestos a mas datos

APRENDIZAJE NEURONAL (Neural Networks)

Algoritmos de aprendizaje automático compuestos de redes/sistemas con múltiples capas interconectadas

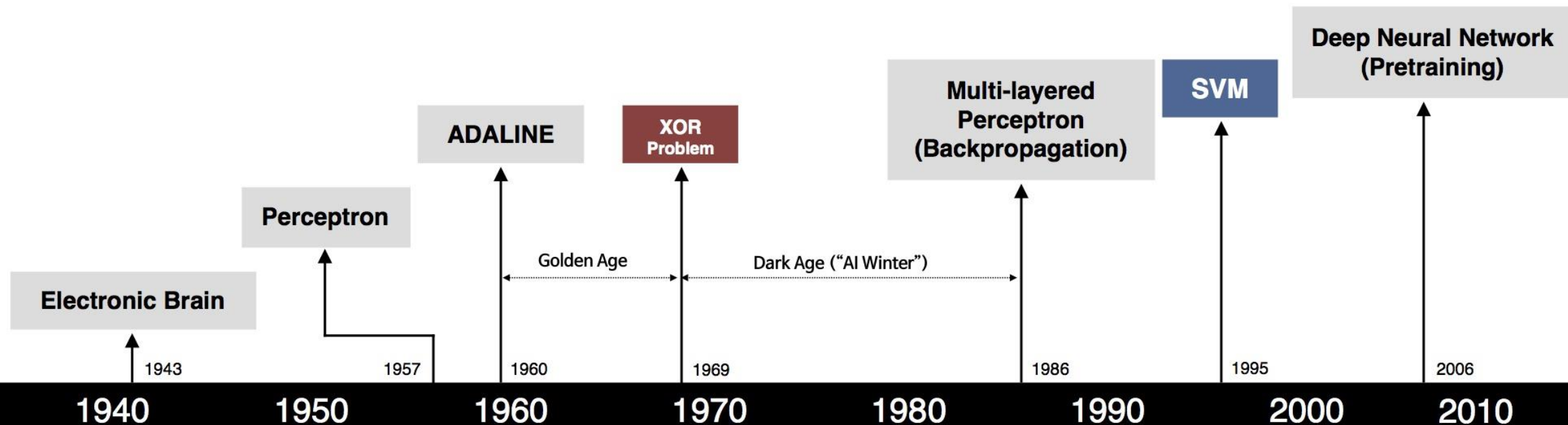
APRENDIZAJE PROFUNDO NEURONAL (Deep Learning)

¿En esta semana, has interactuado con algo que usa o implementa algoritmos de Aprendizaje Automático?

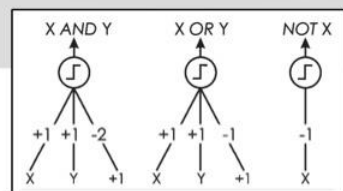


*“No necesitas un Doctorado para **implementar** algoritmos de Aprendizaje Automático”*

Aprendizaje Automático ha existido durante bastante tiempo, y su interés ha fluctuado a lo largo de las décadas



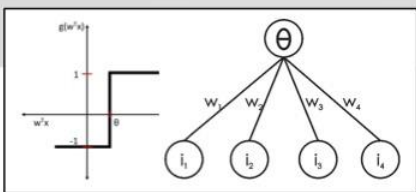
S. McCulloch – W. Pitts



- Adjustable Weights
- Weights are not Learned



F. Rosenblatt



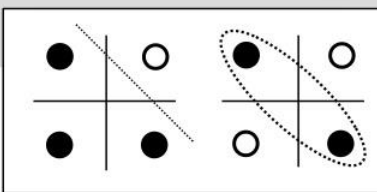
- Learnable Weights and Threshold



B. Widrow – M. Hoff



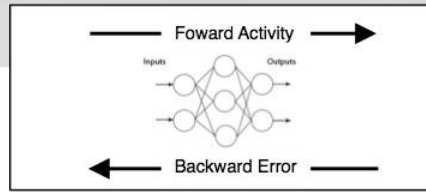
M. Minsky – S. Papert



- XOR Problem



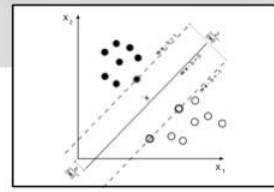
D. Rumelhart – G. Hinton – R. Williams



- Solution to nonlinearly separable problems
- Big computation, local optima and overfitting



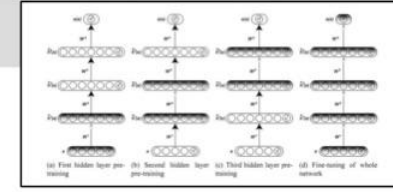
V. Vapnik – C. Cortes



- Limitations of learning prior knowledge
- Kernel function: Human Intervention



G. Hinton – S. Ruslan

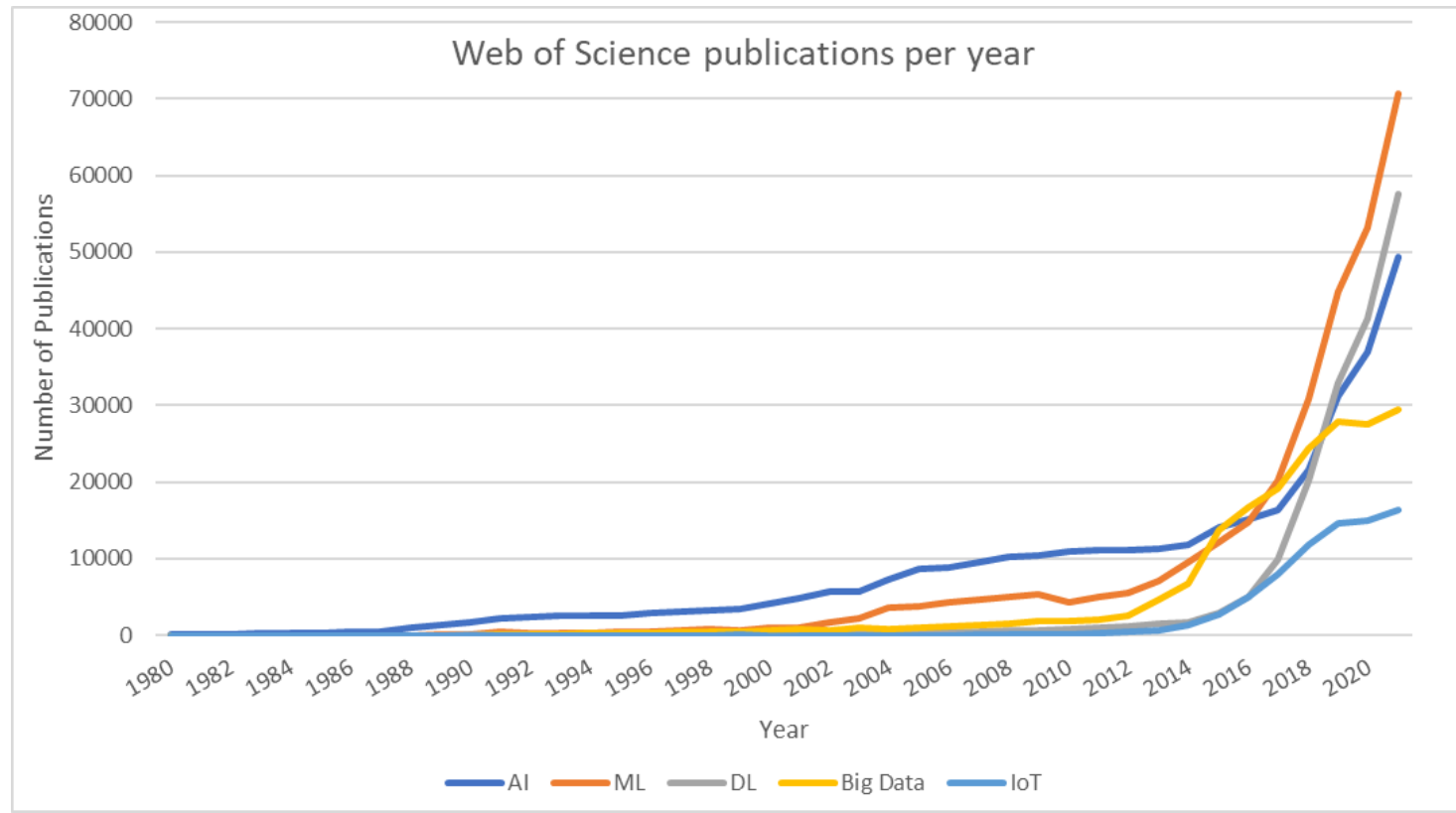


- Hierarchical feature Learning

Pero en la última década hemos visto un crecimiento en su popularidad, la cual continúa creciendo

“Programa informático que aprende de la experiencia E con respecto a una tarea T con rendimiento P , si su desempeño en la tarea T , medido por P , mejora con la experiencia E .”

Mitchell, T. [1997].



A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion - fuelled by internet of things and the use of connected devices - are hard to comprehend, particularly when looked at in the context of one day

500m

tweets are sent every day
Twitter



4PB

of data created by Facebook, including

350m photos
100m hours of video watch time

Facebook Research

294bn

billion emails are sent

Radicati Group

320bn

emails to be sent each day by 2021

306bn

emails to be sent each day by 2020

3.9bn

people use emails

4TB

of data produced by a connected car

Intel

ACCUMULATED DIGITAL UNIVERSE OF DATA

4.4ZB

PwC

2013

44ZB

2020

DEMYSTIFYING DATA UNITS

From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

Unit	Value	Size
b bit	0 or 1	1/8 of a byte
B byte	8 bits	1 byte
KB kilobyte	1,000 bytes	1,000 bytes
MB megabyte	1,000 ² bytes	1,000,000 bytes
GB gigabyte	1,000 ³ bytes	1,000,000,000 bytes
TB terabyte	1,000 ⁴ bytes	1,000,000,000,000 bytes
PB petabyte	1,000 ⁵ bytes	1,000,000,000,000,000 bytes
EB exabyte	1,000 ⁶ bytes	1,000,000,000,000,000,000 bytes
ZB zettabyte	1,000 ⁷ bytes	1,000,000,000,000,000,000,000 bytes
YB yottabyte	1,000 ⁸ bytes	1,000,000,000,000,000,000,000,000 bytes

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.

65bn

messages sent over WhatsApp and two billion minutes of voice and video calls made

Facebook

Searches made a day

5bn

Searches made a day from Google

3.5bn

Smart Insights

463EB

of data will be created every day by 2025

ISC

95m

photos and videos are shared on Instagram

Instagram Business

28PB

to be generated from wearable devices by 2020

Statista

Quick, Draw



Can a neural network learn to recognize doodling?

Help teach it by adding your drawings to the world's largest doodling data set, shared publicly to help with machine learning research.

Let's Draw!

This is an
A.I.
Experiment

Made with
some friends from
Google

English

Privacy & Terms

<https://quickdraw.withgoogle.com/#>

Existen muchos tipos de sistemas de Aprendizaje Automático, por lo tanto, es útil clasificarlos en diferentes categorías

- Si están o no entrenado con supervisión humana (**datos utilizados**)
 - Supervisado (*Supervised*)
 - No Supervisado (*Unsupervised*)
 - Semi-supervisado (*Semi-Supervised*)
 - Aprendizaje reforzado (*Reinforcement Learning*)

La clasificación más común de sistemas de Aprendizaje Automático es basado en los "datos de entrenamiento"

Datos Etiquetados

$(x^1, y^1), \dots$
 $, (x^m, y^m)$

Supervisado

No Supervisado

Datos No Etiquetados

x^1, x^2, \dots, x^m

Supervisado es el tipo de algoritmos más común

Datos Etiquetados

$(x^1, y^1), \dots$
 $, (x^m, y^m)$

Objetivo :

$x \rightarrow Y$
 $y = f(x)$

Supervisado

```
graph TD; A[Supervisado] --> B[Continua<br/>Variable Objetivo]; A --> C[Categórica<br/>Variable Objetivo]; B --> D[Regresión]; C --> E[Clasificación]; D --> F[Tarea/Problema]; E --> F;
```

Continua
Variable Objetivo

Regresión

- *Predicción del precio de casas*
- *Predicciones mercado de valores*
- *Predicciones de uso de la red eléctrica*

Categórica
Variable Objetivo

Clasificación

➔ **Tarea/Problema**

- *Imagenes medicas*
- *Detección de objetos*
- *Predicciones mercado de valores*

En la mayoría de los casos, los datos no etiquetados son más predominantes que los datos etiquetados

Datos Etiquetados

$(x^1, y^1), \dots$
 $, (x^m, y^m)$

Supervisado

No Supervisado

Datos No Etiquetados

x^1, x^2, \dots, x^m

*No existe una
Variable Objetivo*

Agrupación

- *Segmentación de clientes*

Asociación

- *Análisis de la canasta de mercado*

Detección de anomalías

- *Detecciones de fraudes*

Reducción de dimensionalidad

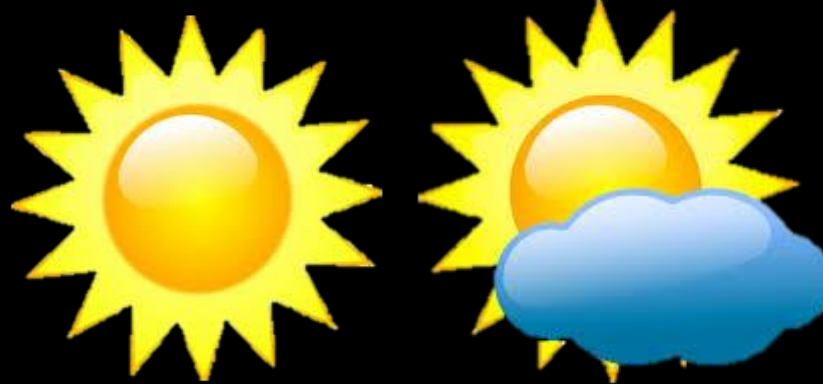
- *Visualización*
- *Ingeniería de Atributos*

¿Qué tipo de sistemas de ML puedo usar para responder las siguientes preguntas y qué tipo de tarea es?

Supervisado

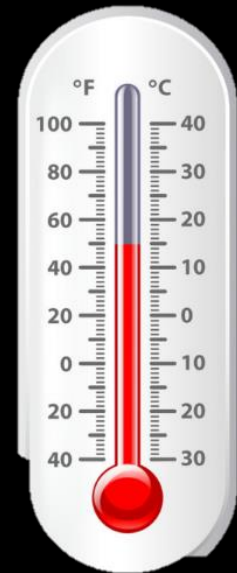
Clasificación

1. ¿Va a estar soleado o no el día de mañana?

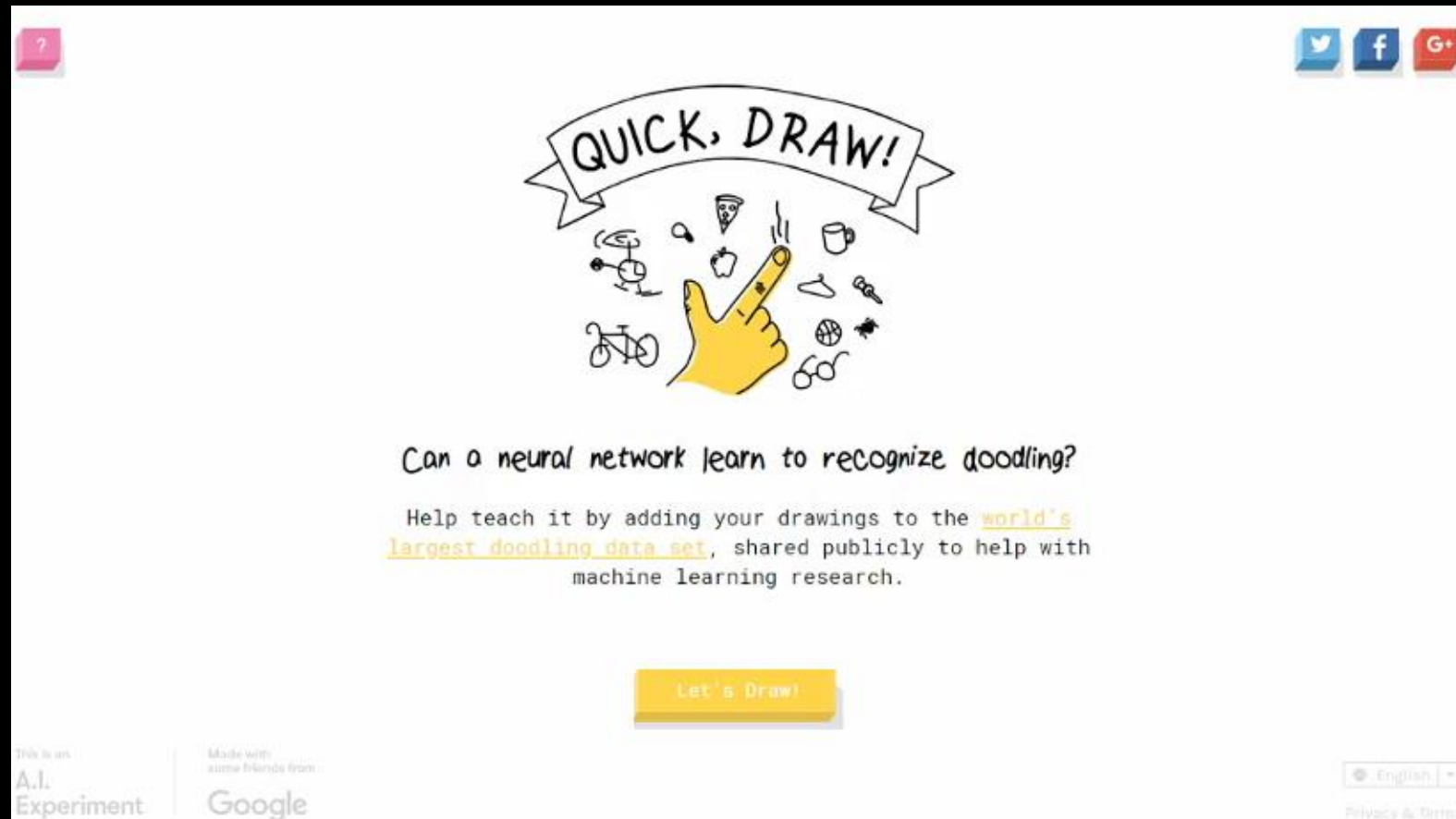


Regresión

2. ¿Cuál sería la temperatura máxima el día de mañana?



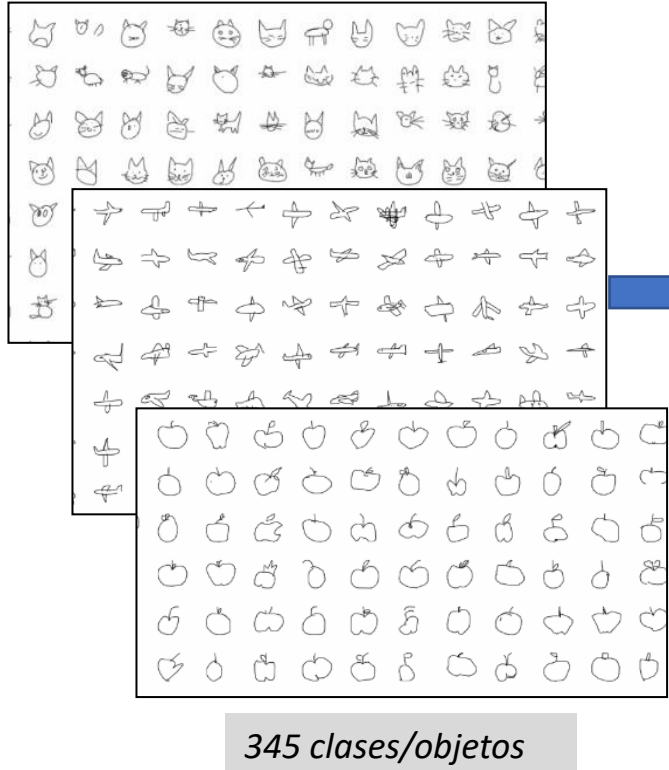
¿Qué tipo de sistemas de ML utilizaron, qué tipo de tareas, y qué tipo de datos utilizaron?



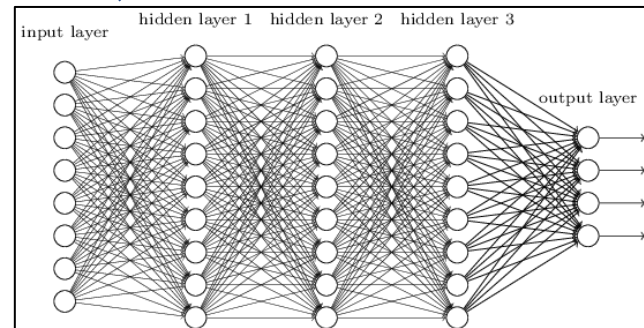
Supervisado

Clasificación

“Humanos provén el conocimiento”

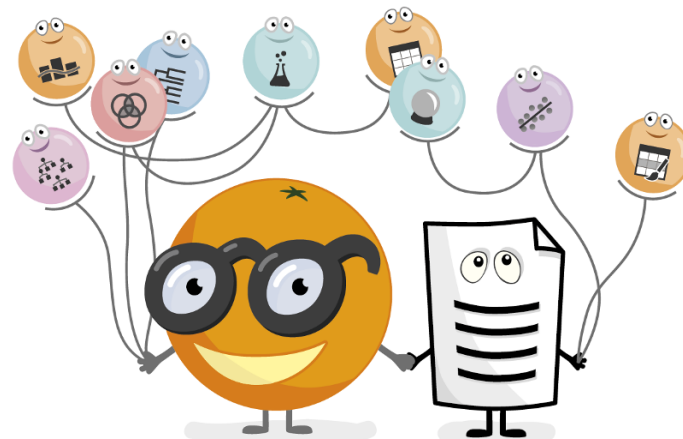
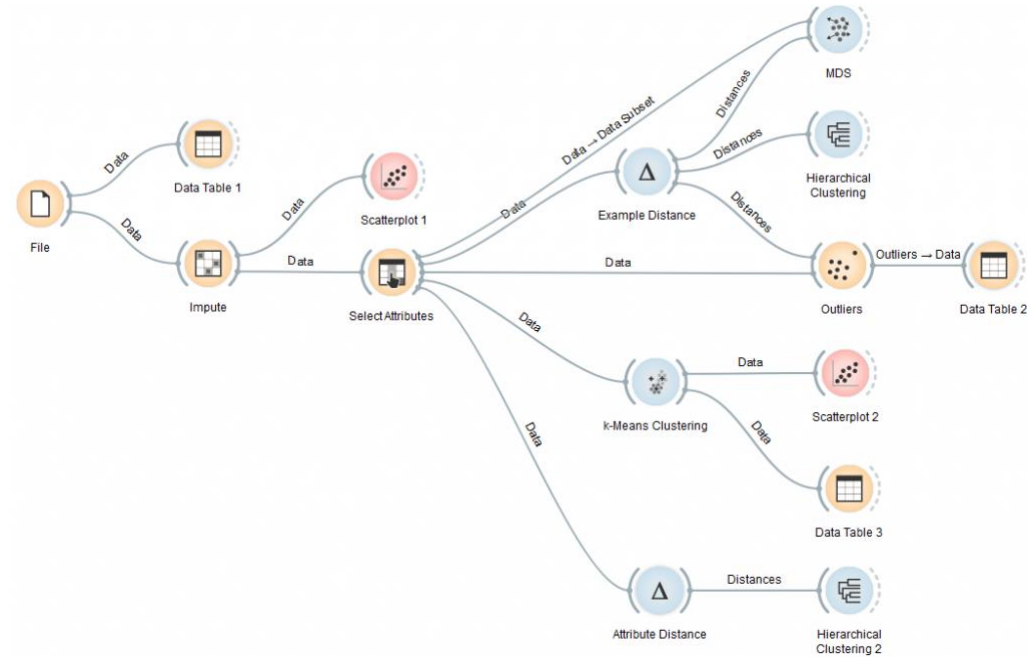


Algoritmos identifican patrones en la data



Despliegue del producto

Implementando algoritmos de Aprendizaje Automático sin la necesidad de programar*



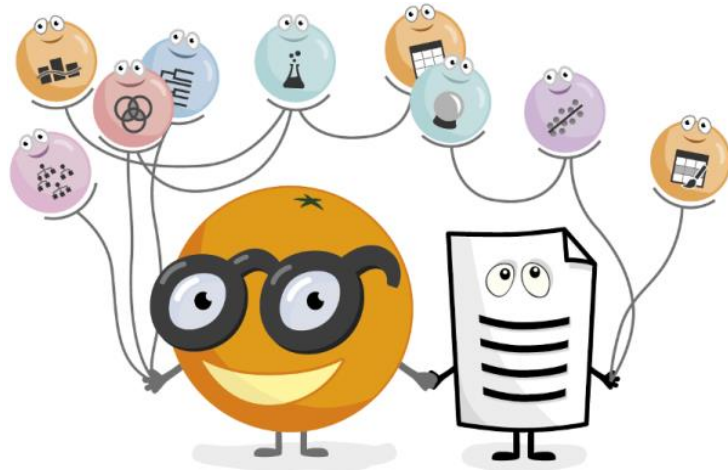
*Lenguaje de
Programación visual*

1er paso: Instalar programa Orange



“Orange Data mining”

<https://orangedatamining.com/download/>



Hoy vamos a trabajar con la base de datos

Iris de Ronald Fisher



Iris Versicolor



Iris Setosa

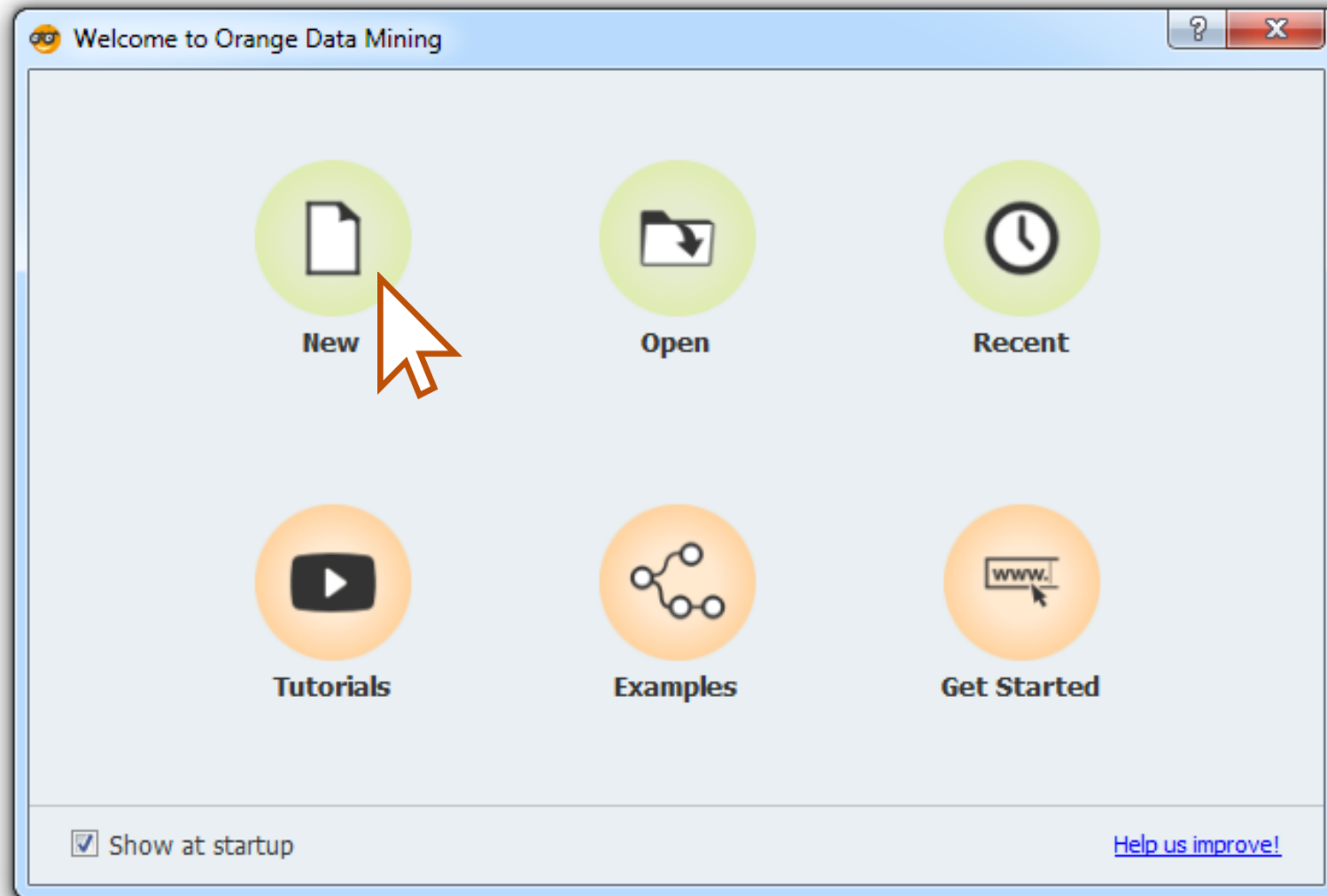


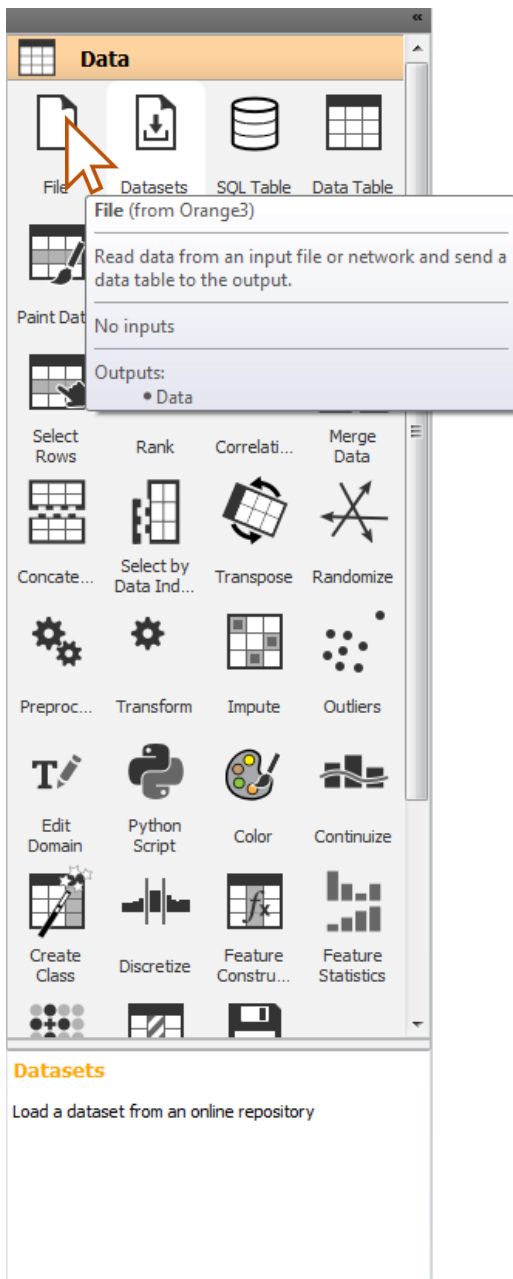
Iris Virginica

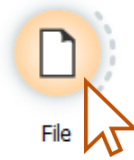
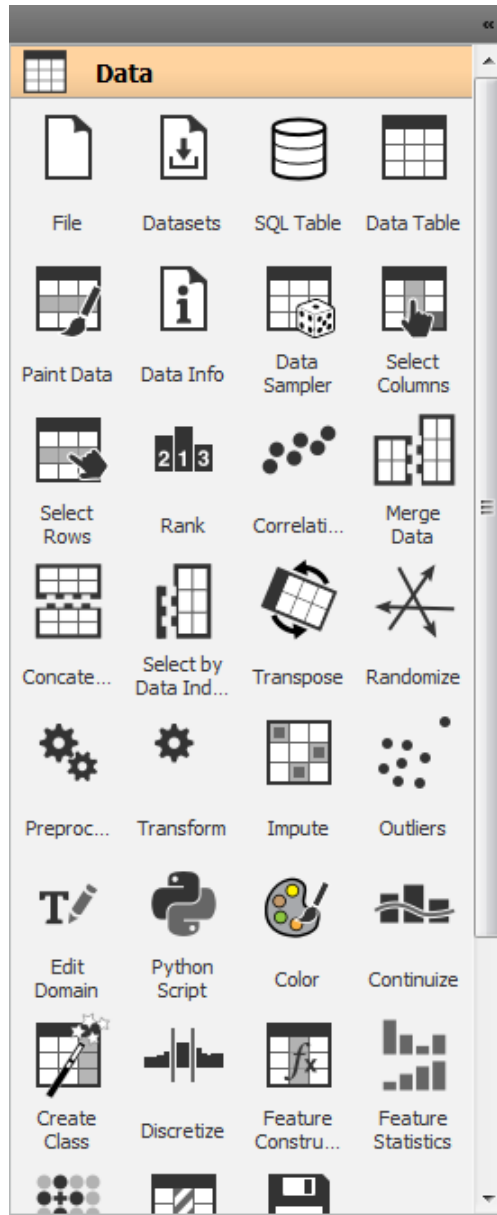
	Sepal Length	Sepal Width	Petal Length	Petal Width	Species
1	5.1	3.5	1.4	0.2	setosa
2	4.9	3	1.4	0.2	setosa
...
51	7	3.2	4.7	0.2	versicolor
...
150	5.9	3	5.1	1.8	virginica

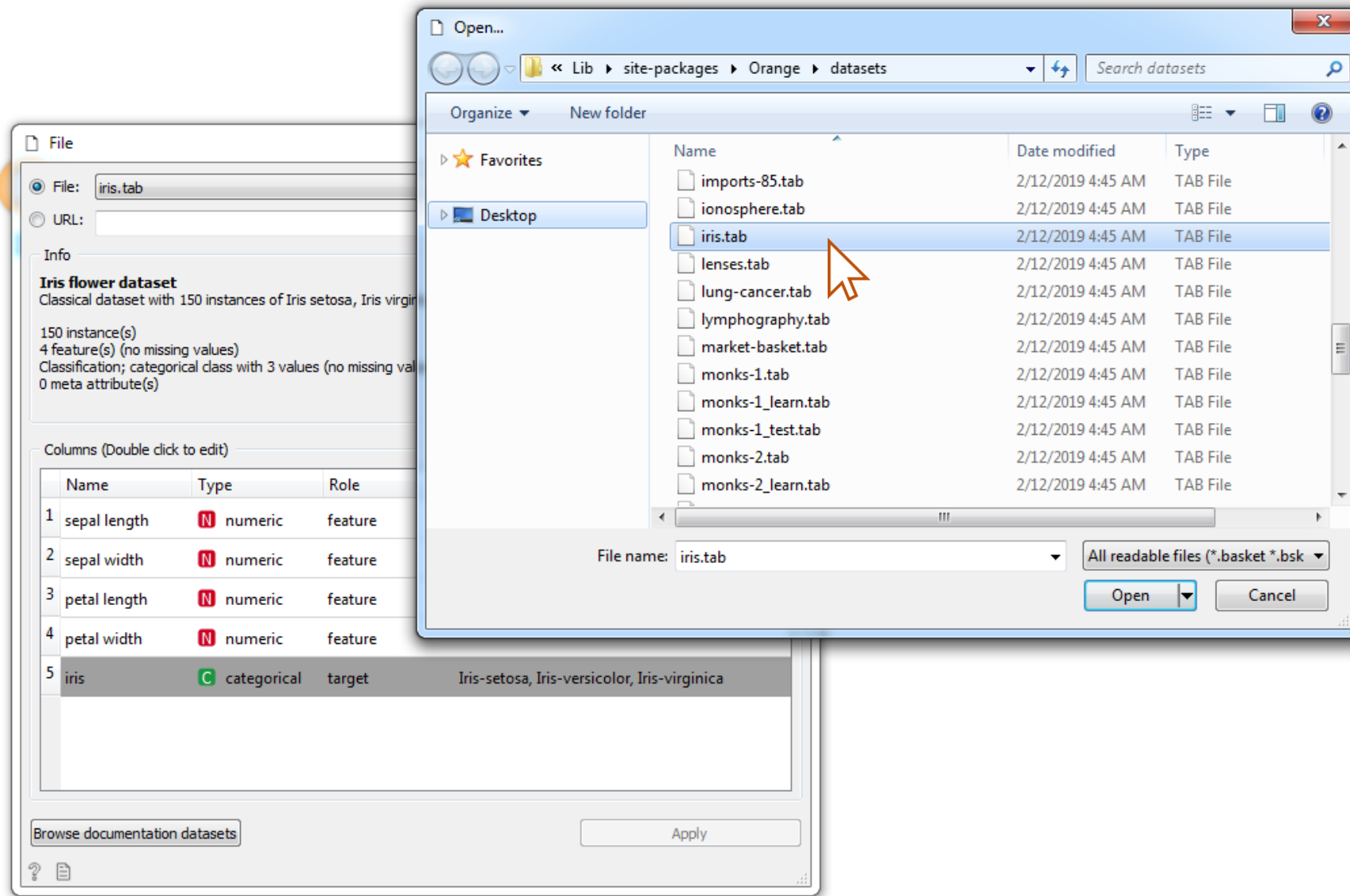
150 Flores medidas
(*tuplas=filas*)
 \times
4 atributos
(*variables=columnas*)

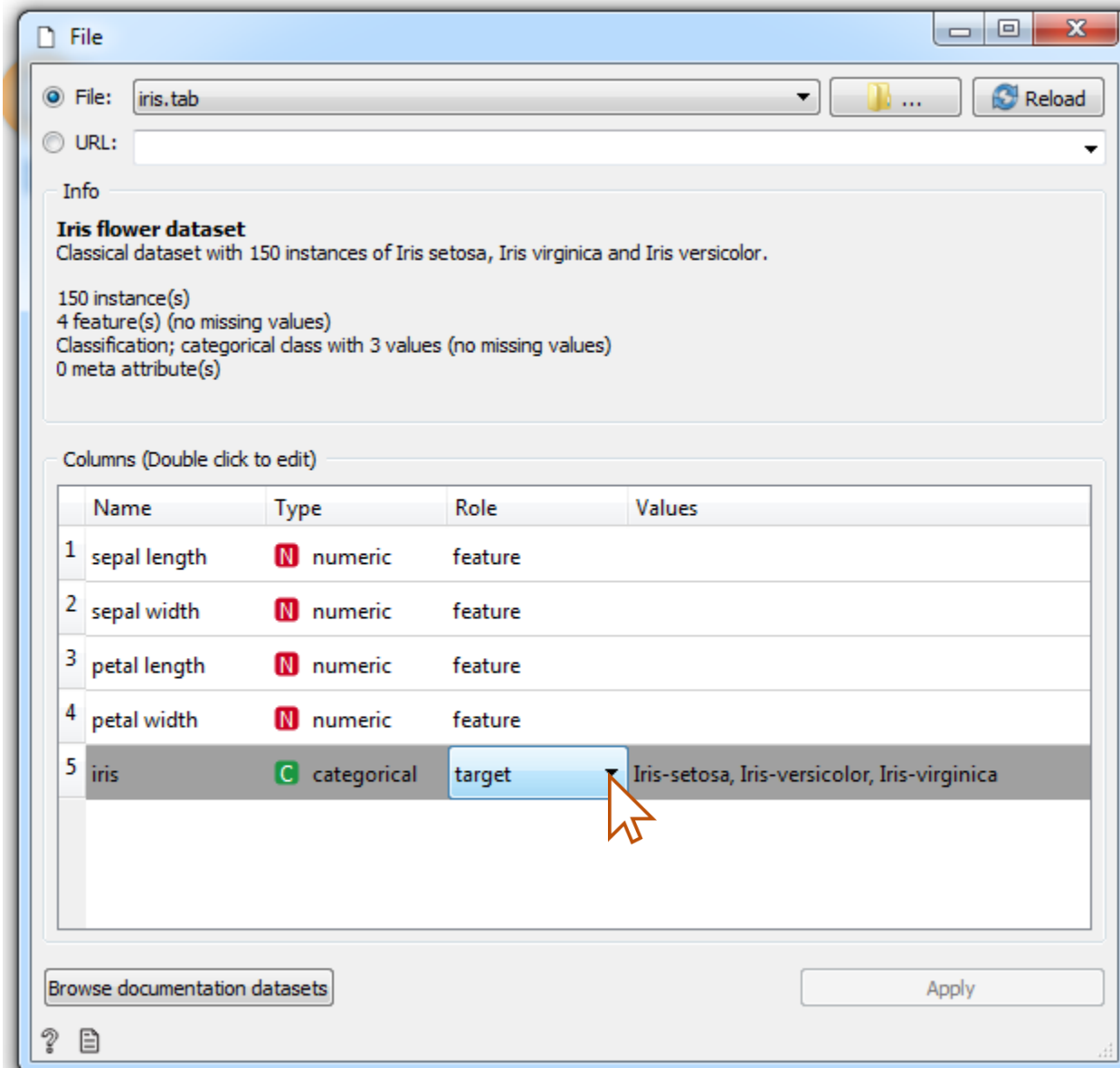
2do paso: Cargar los datos



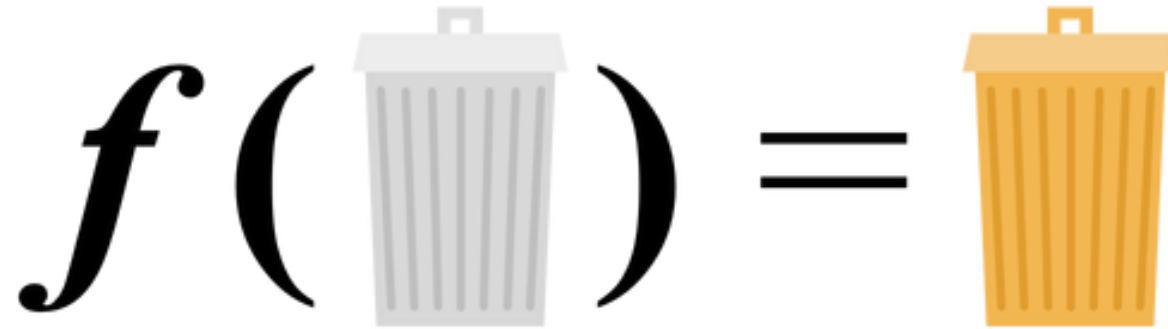








Antes de iniciar cualquier análisis hay que entender bien la data

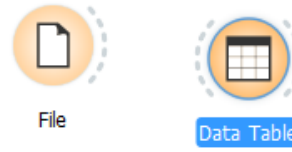
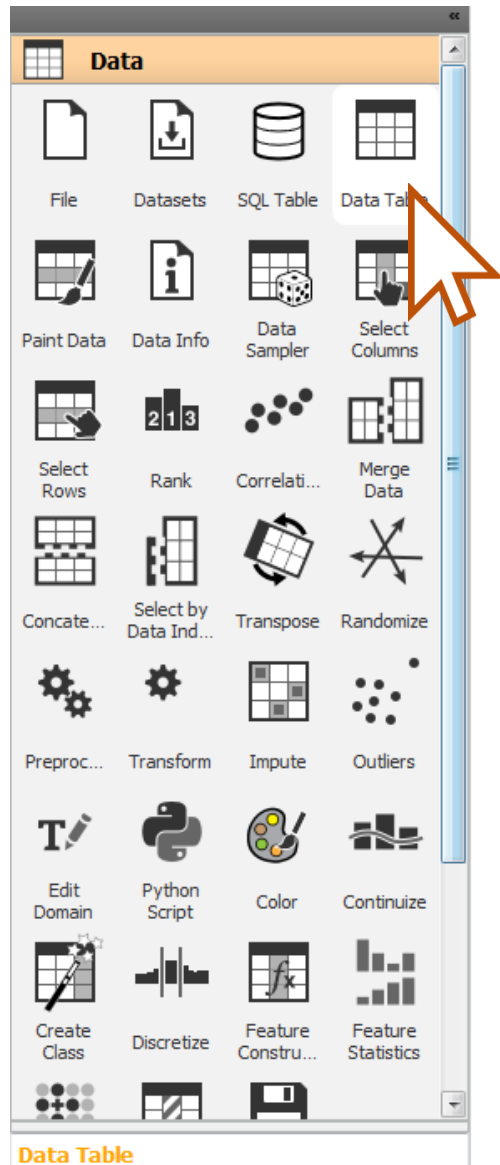


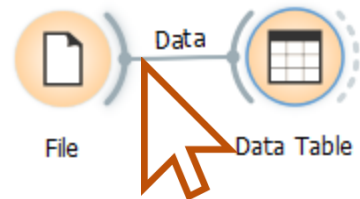
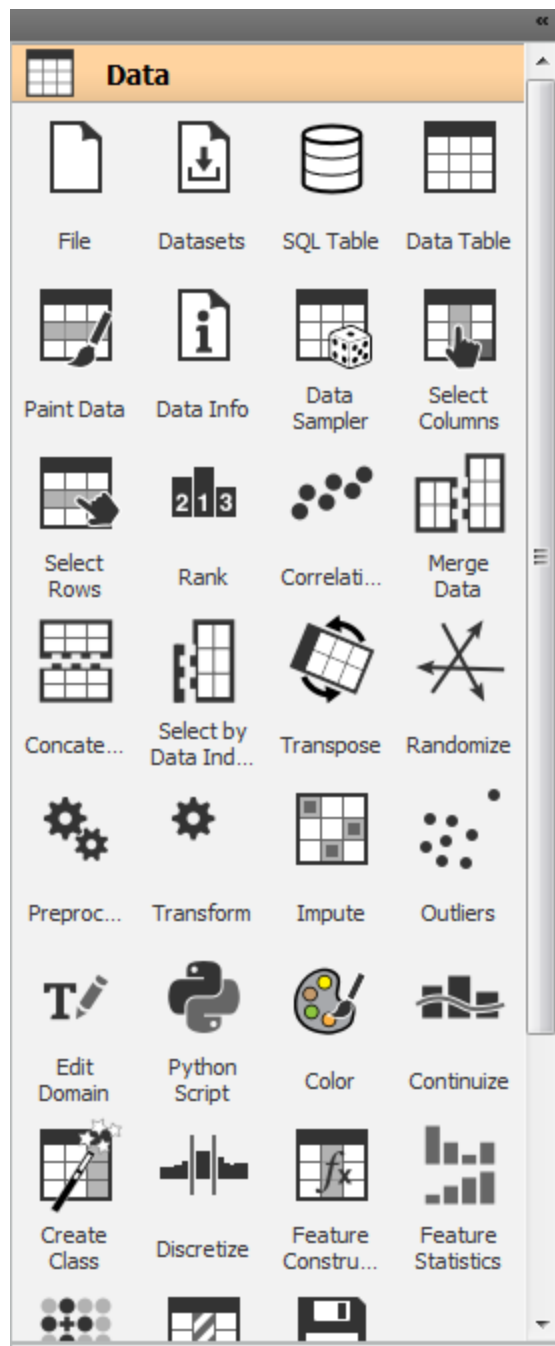
“Garbage in, Garbage out”

- **Cuántas variables y tuplas hay?**
(filas y columnas)
- **Qué tipo de variables tengo?**
(nominal, ordinal, intervalo, proporción, categóricas)
- **Qué información me falta ?**
(NAs, celdas vacías)
- **Cuál es la distribución de mis variables?**
(promedio, rango, moda,...)

...

3er paso: Explorar los datos





Data

File Datasets SQL Table Data Table

Paint Data Data Info Data Sampler Select Columns

Select Rows Rank Correlati... Merge Data

Concatenate... Select by Data Ind... Transpose Randomize

Preproc... Transform Impute Outliers

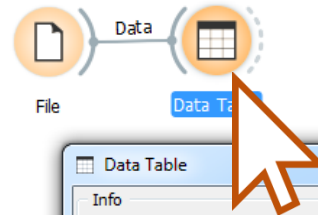
Edit Domain Python Script Color Continuize

Create Class Discretize Feature Constr... Feature Statistics

Data Table

View the dataset in a spreadsheet.

[more...](#)



Data Table

Info

150 instances (no missing values)
4 features (no missing values)
Discrete class with 3 values (no missing values)
No meta attributes

Variables

☒ Show variable labels (if present)
☒ Visualize numeric values
☒ Color by instance classes

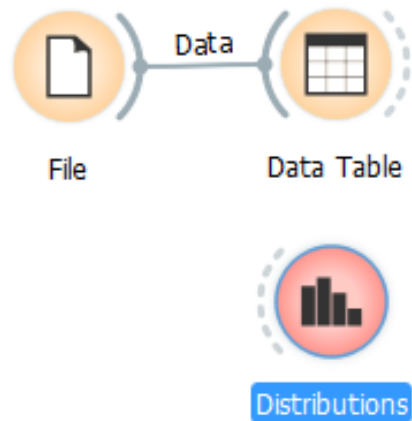
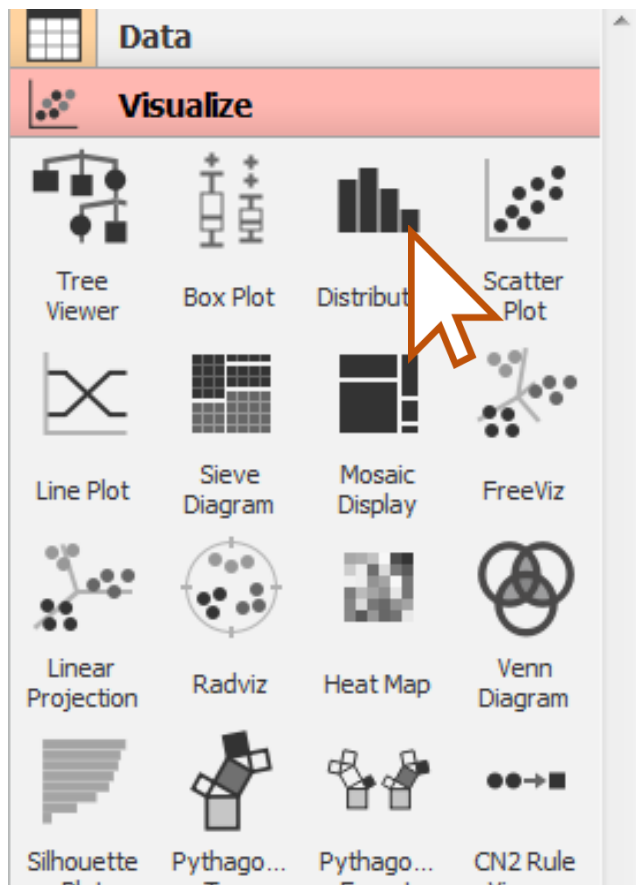
Selection

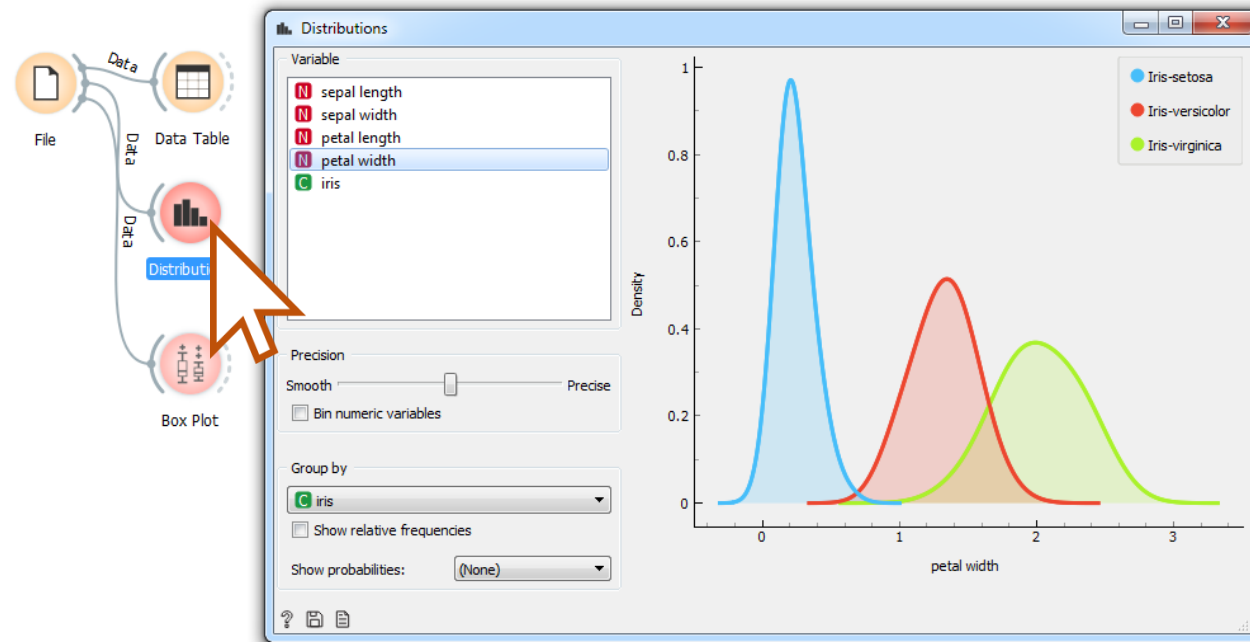
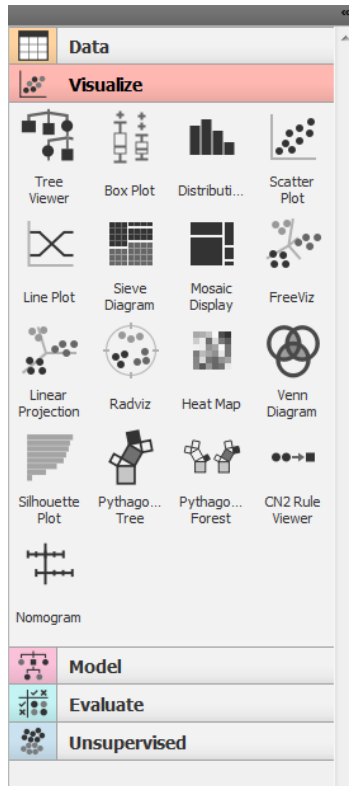
☒ Select full rows

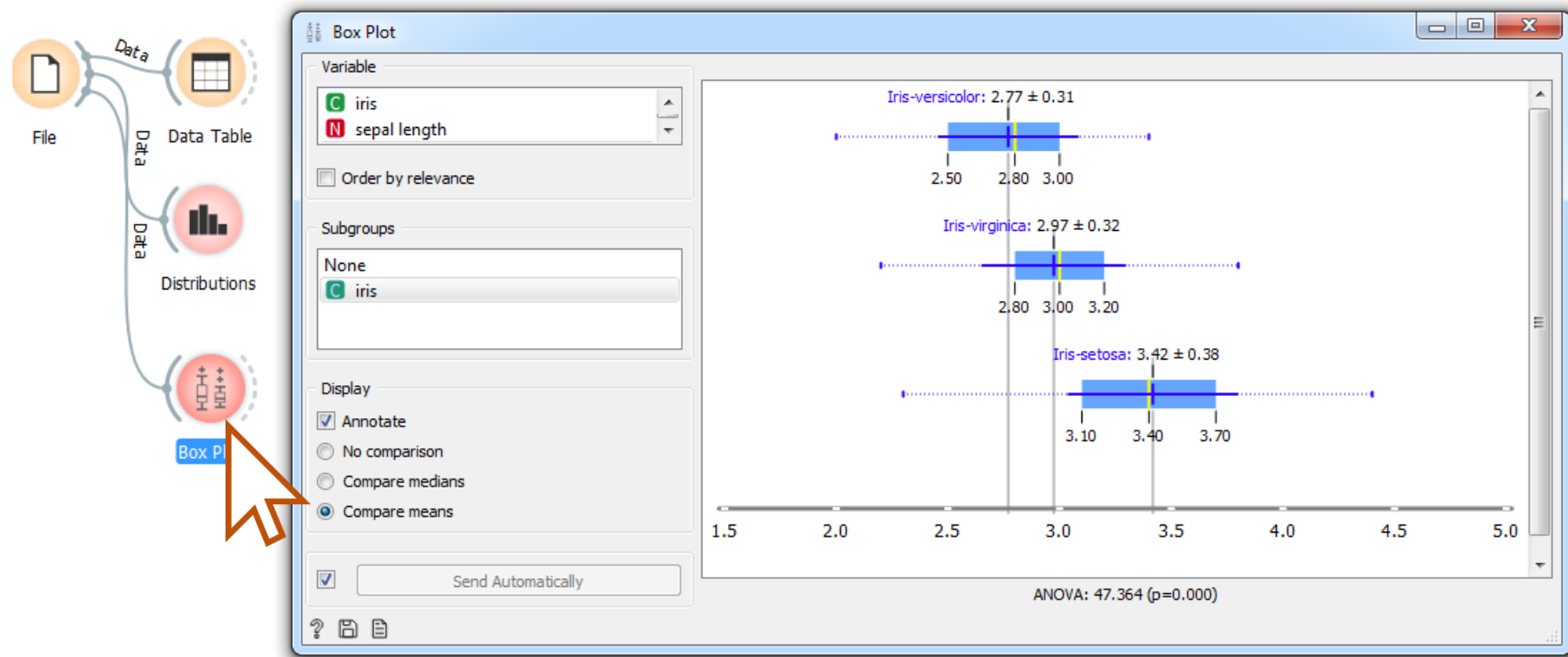
Restore Original Order

☒ Send Automatically

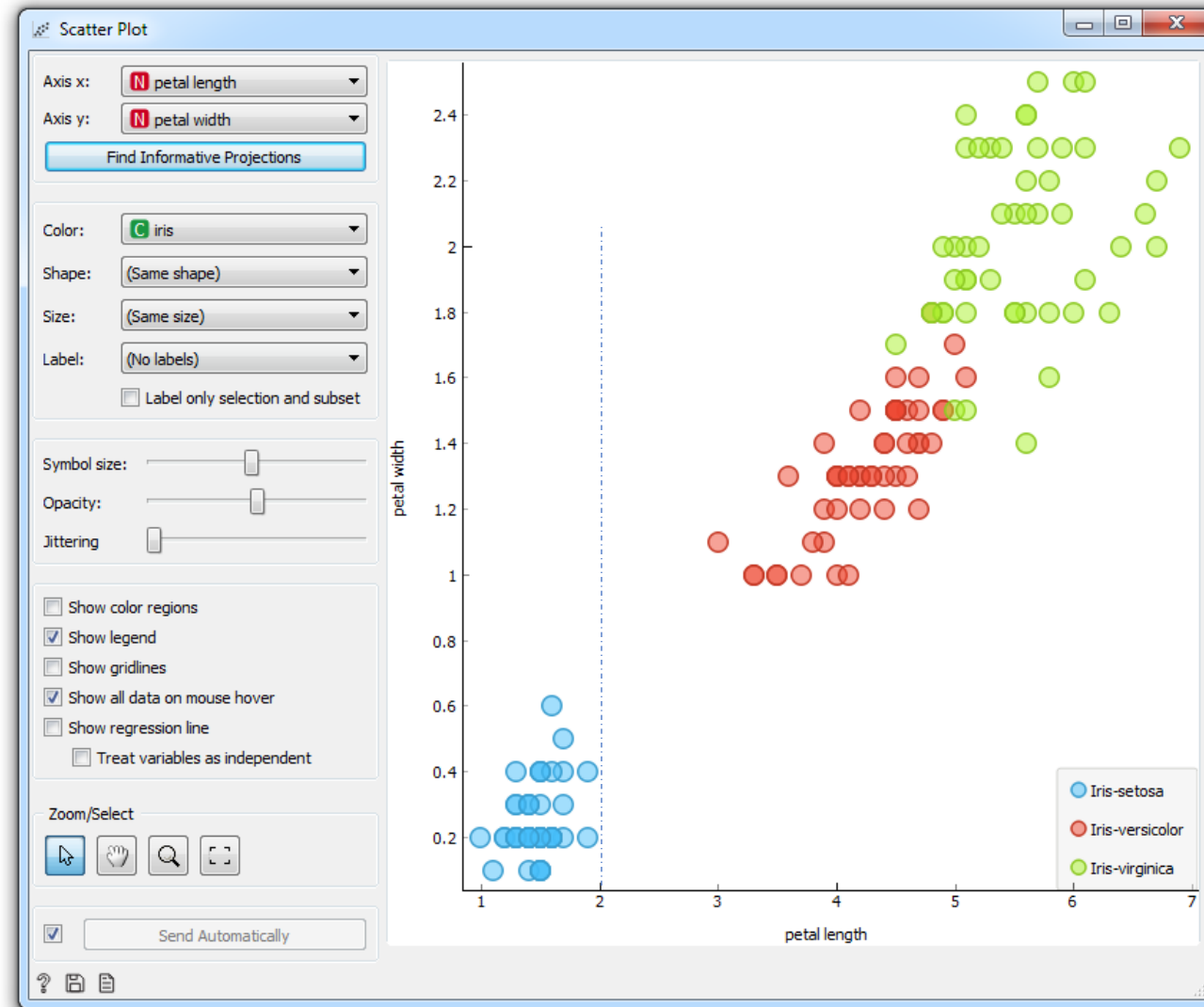
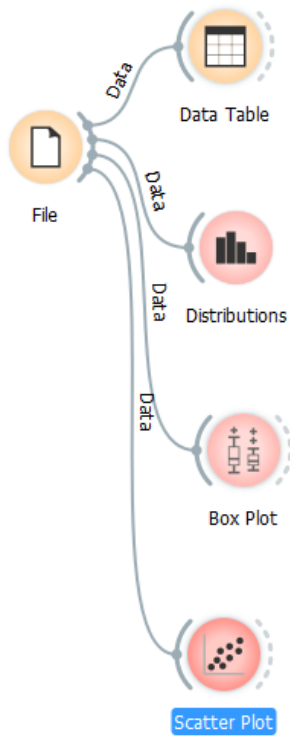
	iris	sepal length	sepal width	petal length	petal width
1	Iris-setosa	5.1	3.5	1.4	0.2
2	Iris-setosa	4.9	3.0	1.4	0.2
3	Iris-setosa	4.7	3.2	1.3	0.2
4	Iris-setosa	4.6	3.1	1.5	0.2
5	Iris-setosa	5.0	3.6	1.4	0.2
6	Iris-setosa	5.4	3.9	1.7	0.4
7	Iris-setosa	4.6	3.4	1.4	0.3
8	Iris-setosa	5.0	3.4	1.5	0.2
9	Iris-setosa	4.4	2.9	1.4	0.2
10	Iris-setosa	4.9	3.1	1.5	0.1
11	Iris-setosa	5.4	3.7	1.5	0.2
12	Iris-setosa	4.8	3.4	1.6	0.2
13	Iris-setosa	4.8	3.0	1.4	0.1
14	Iris-setosa	4.3	3.0	1.1	0.1
15	Iris-setosa	5.8	4.0	1.2	0.2
16	Iris-setosa	5.7	4.4	1.5	0.4
17	Iris-setosa	5.4	3.9	1.3	0.4
18	Iris-setosa	5.1	3.5	1.4	0.3
19	Iris-setosa	5.7	3.8	1.7	0.3
20	Iris-setosa	5.1	3.8	1.5	0.3
21	Iris-setosa	5.4	3.4	1.7	0.2
22	Iris-setosa	5.1	3.7	1.5	0.4
23	Iris-setosa	4.6	3.6	1.0	0.2



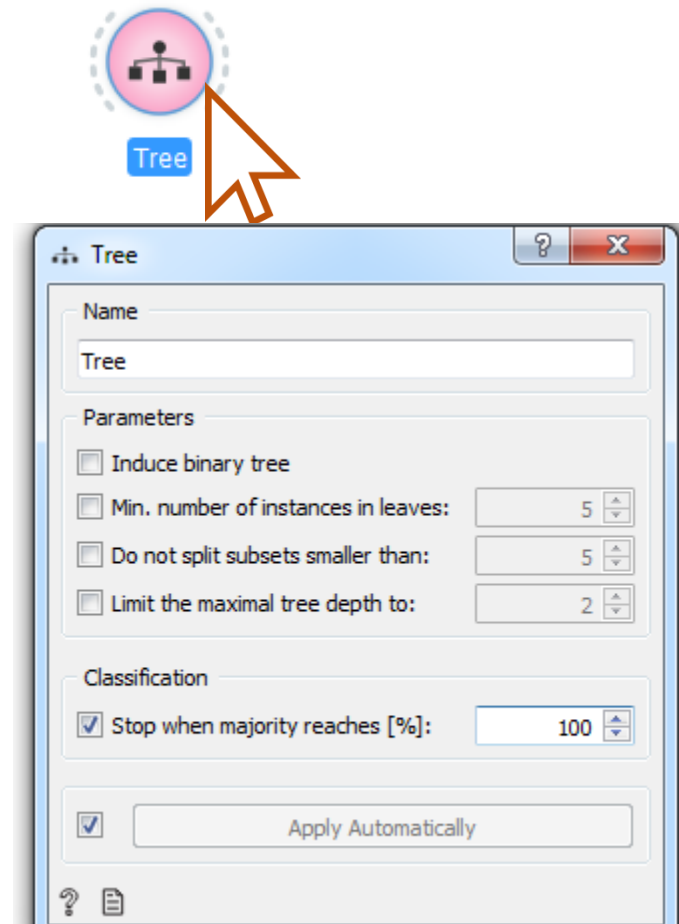
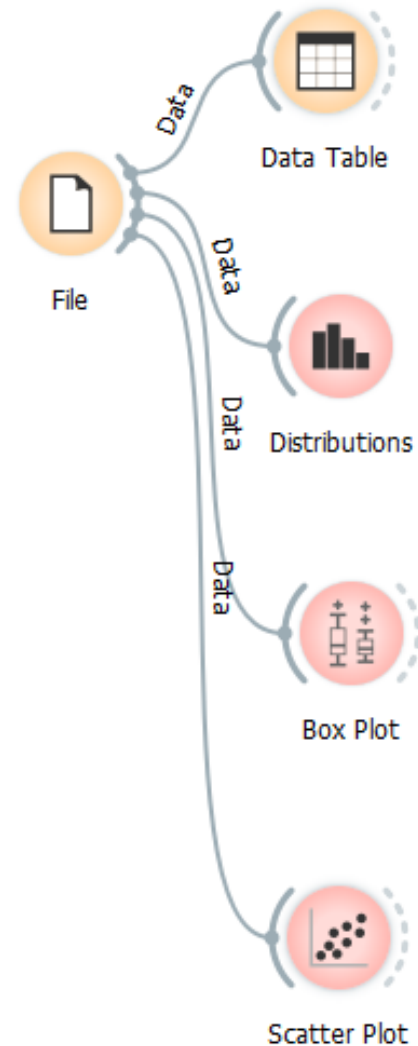
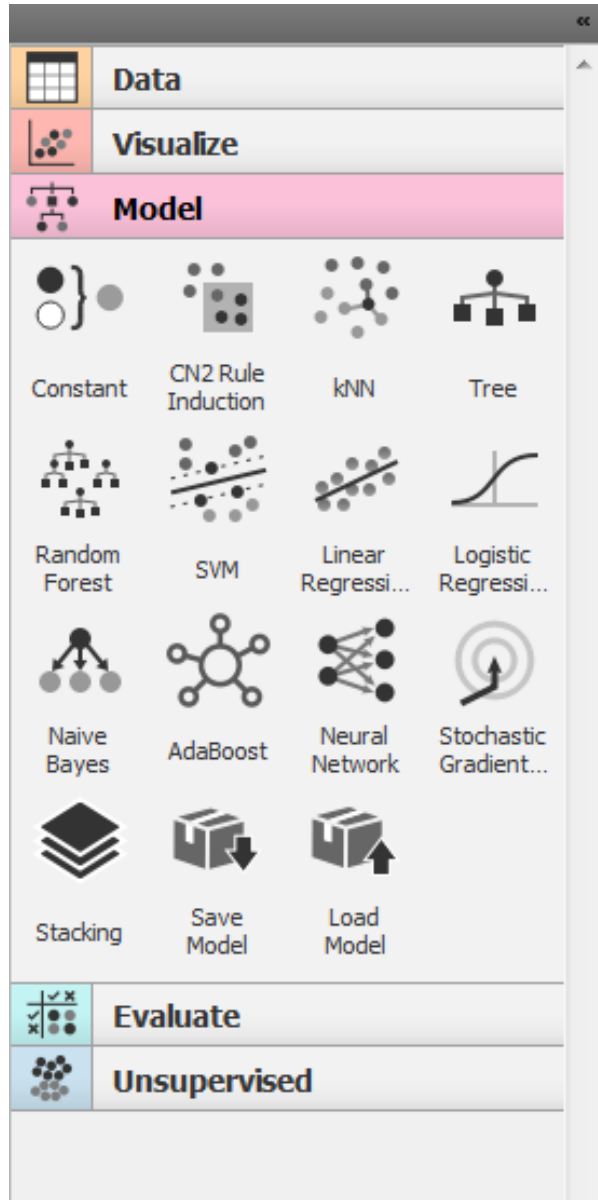


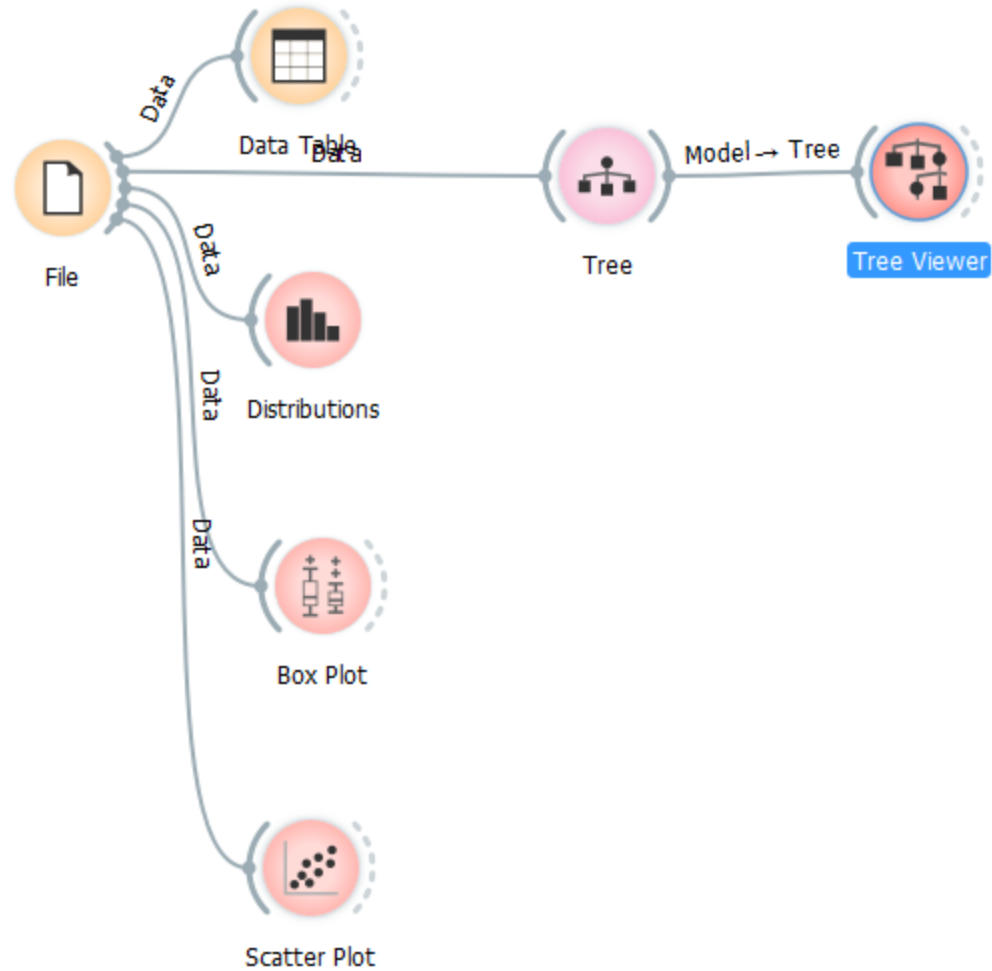
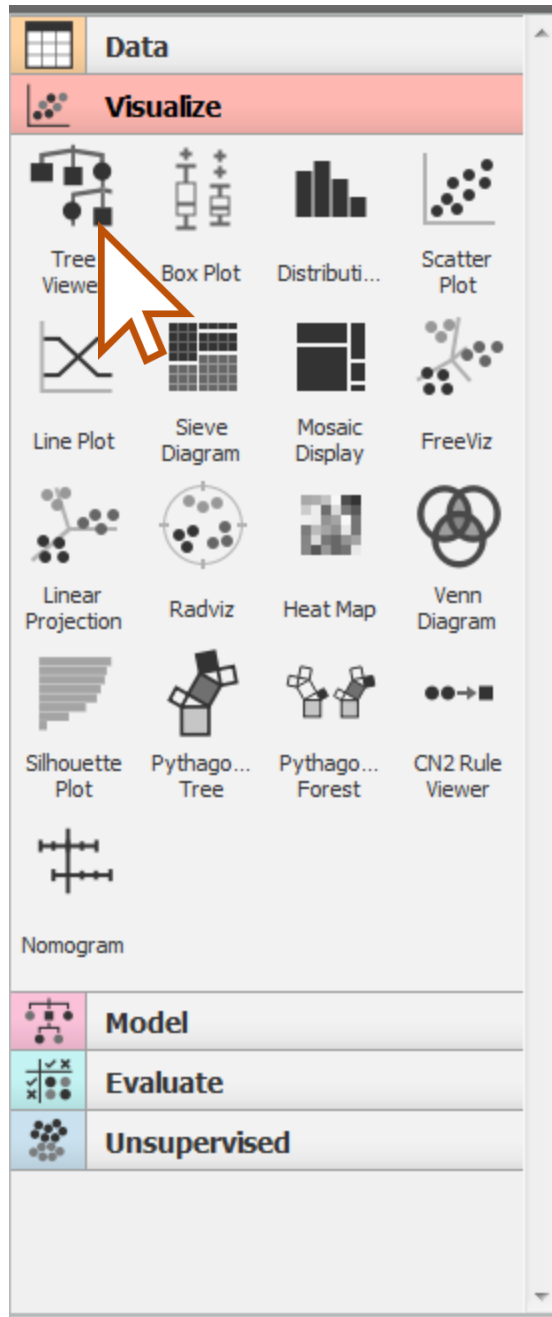


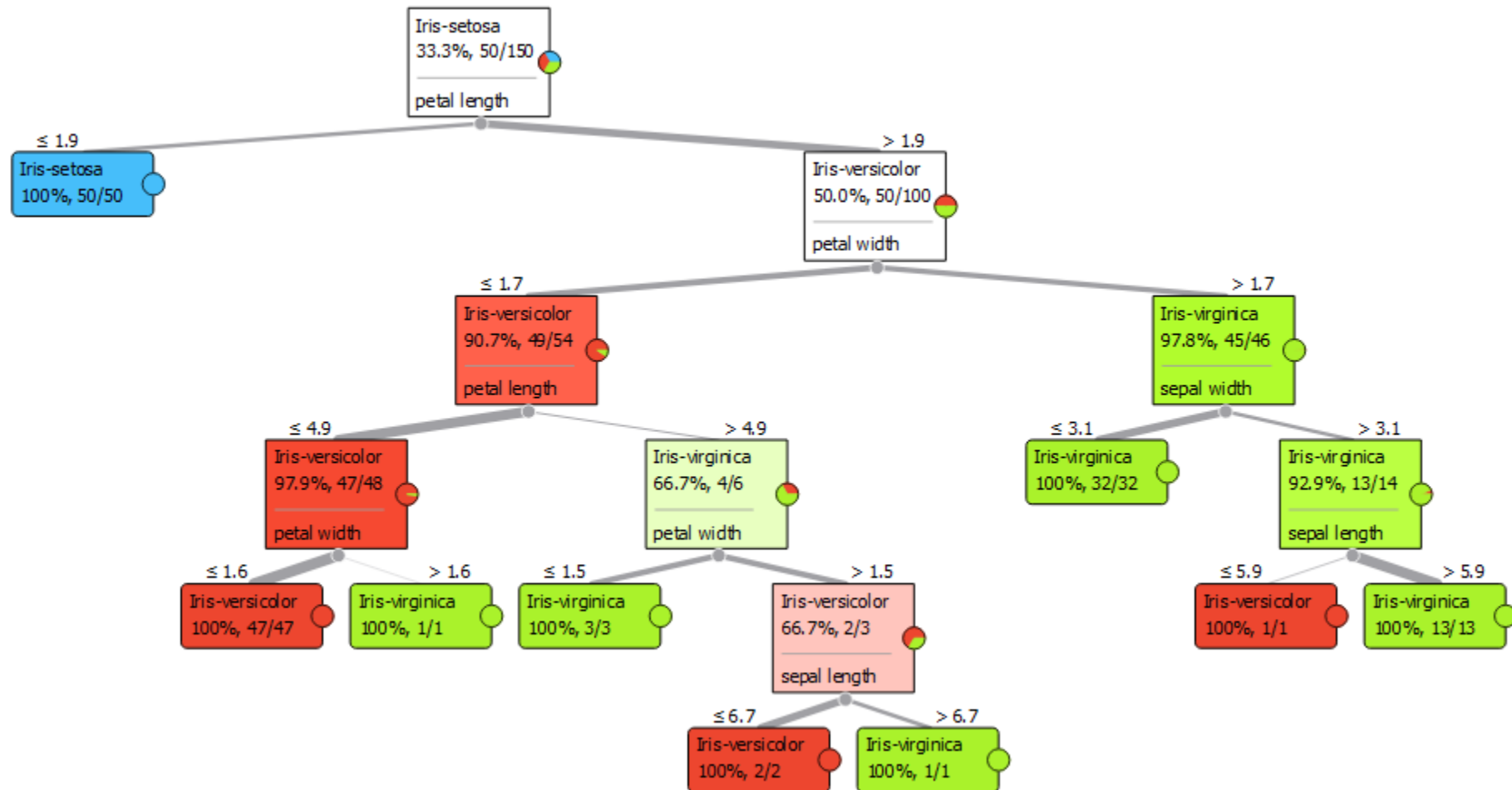
Qué tipo de reglas podemos crear para diferenciar entre las clases de Iris ?



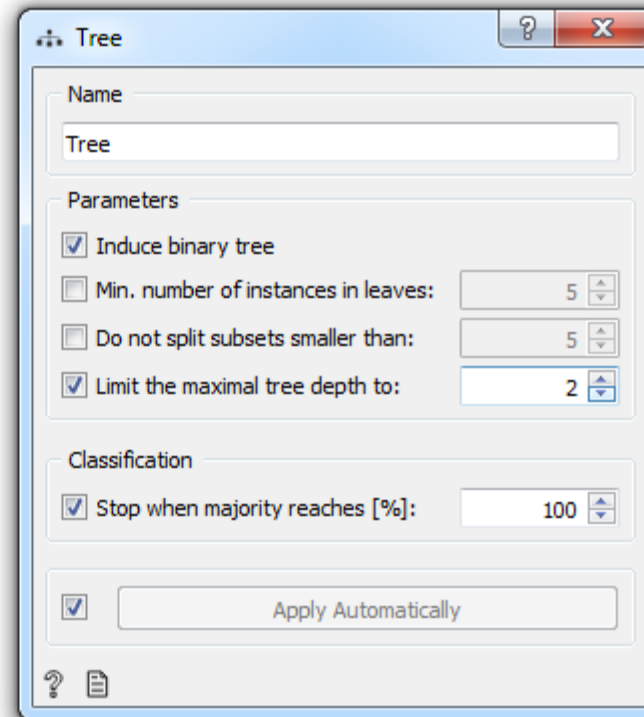
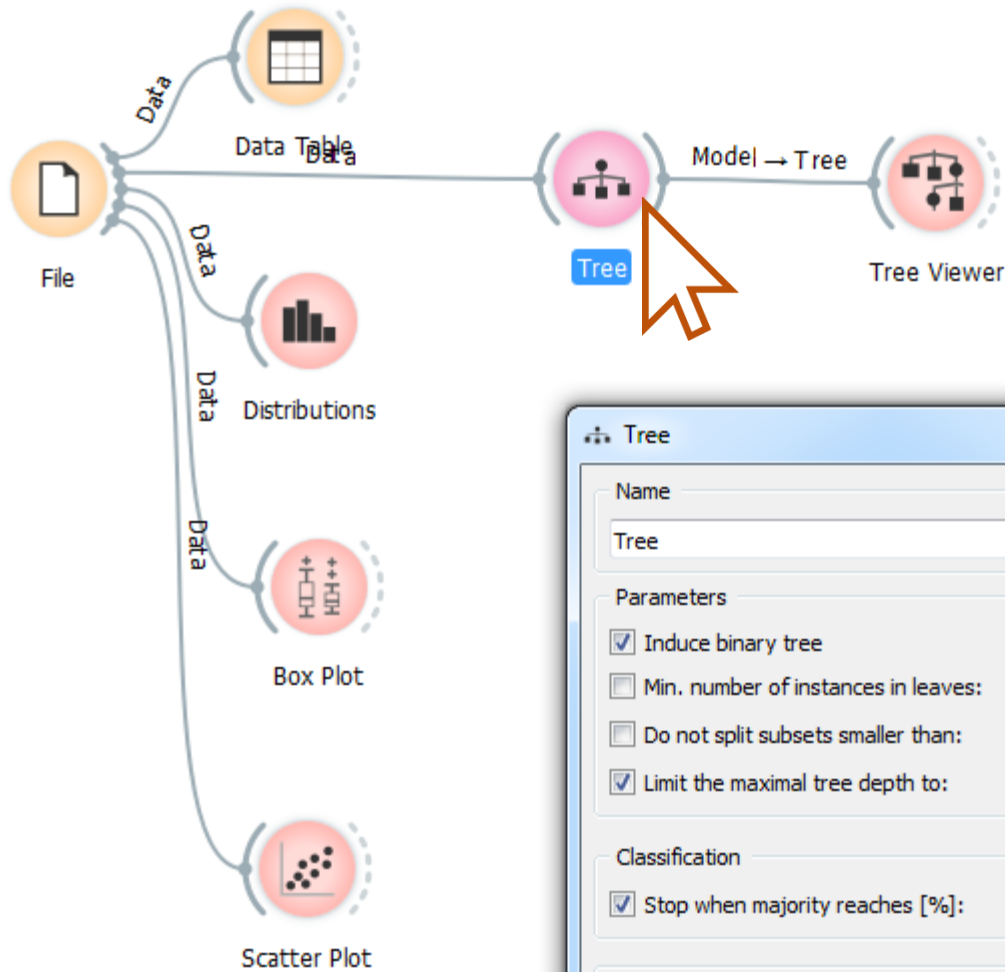
Puede un algoritmo generar estas reglas automáticamente?

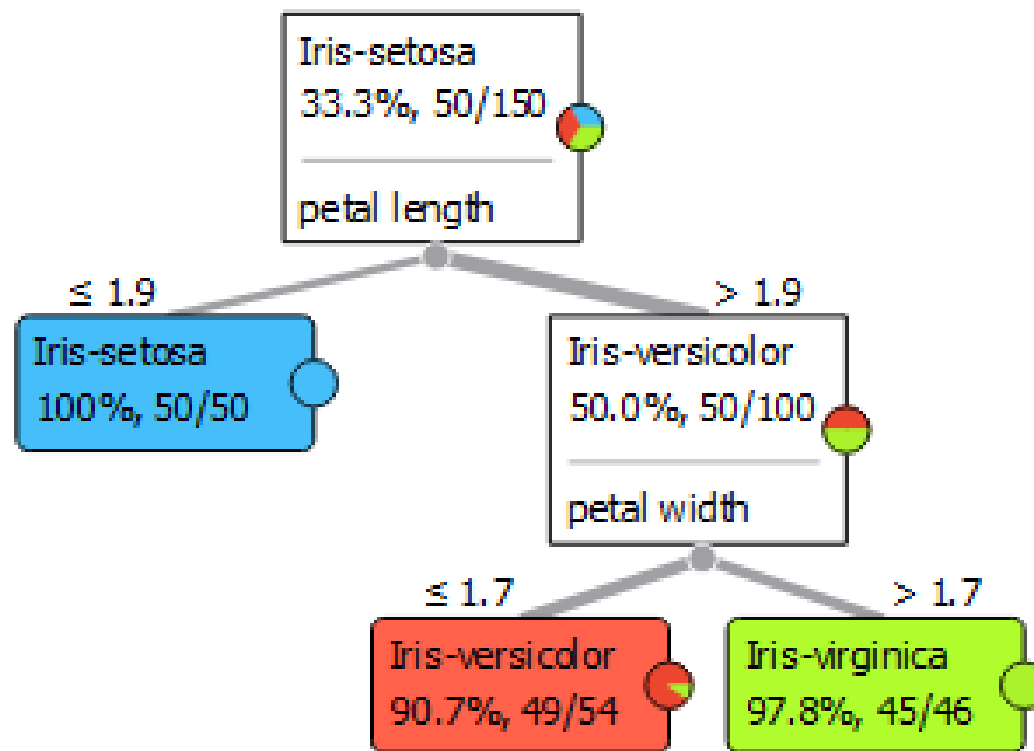






Tenemos que entender los parámetros de los diferentes algoritmos

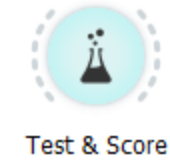
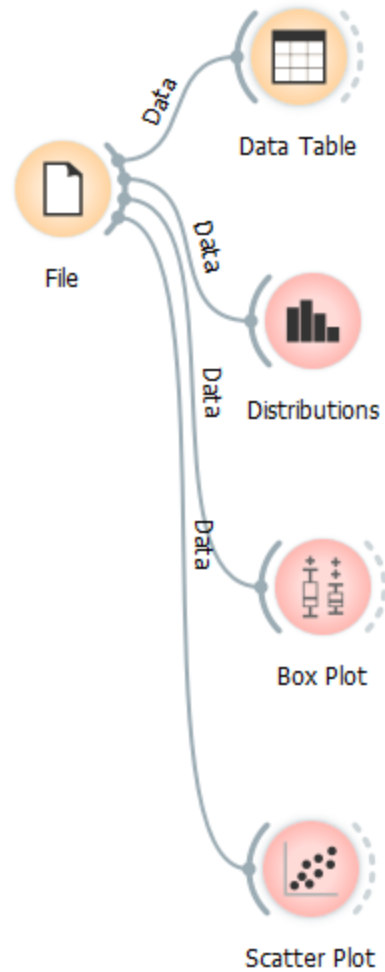
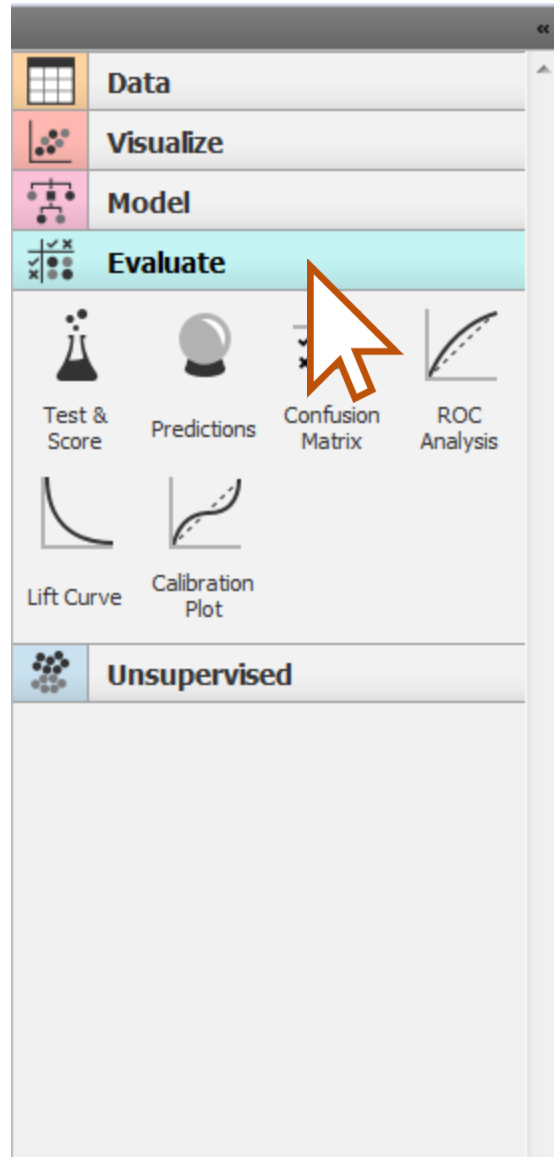


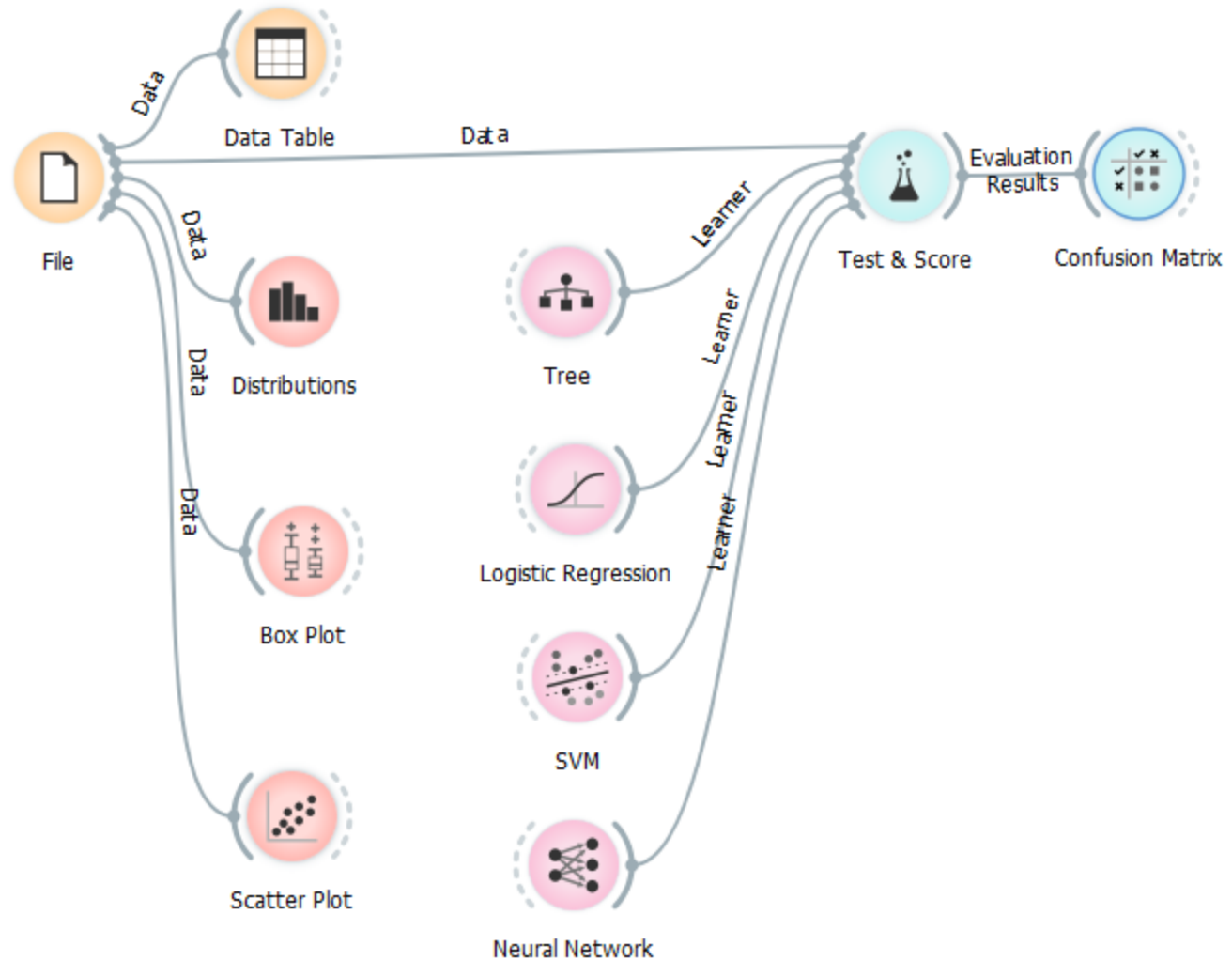
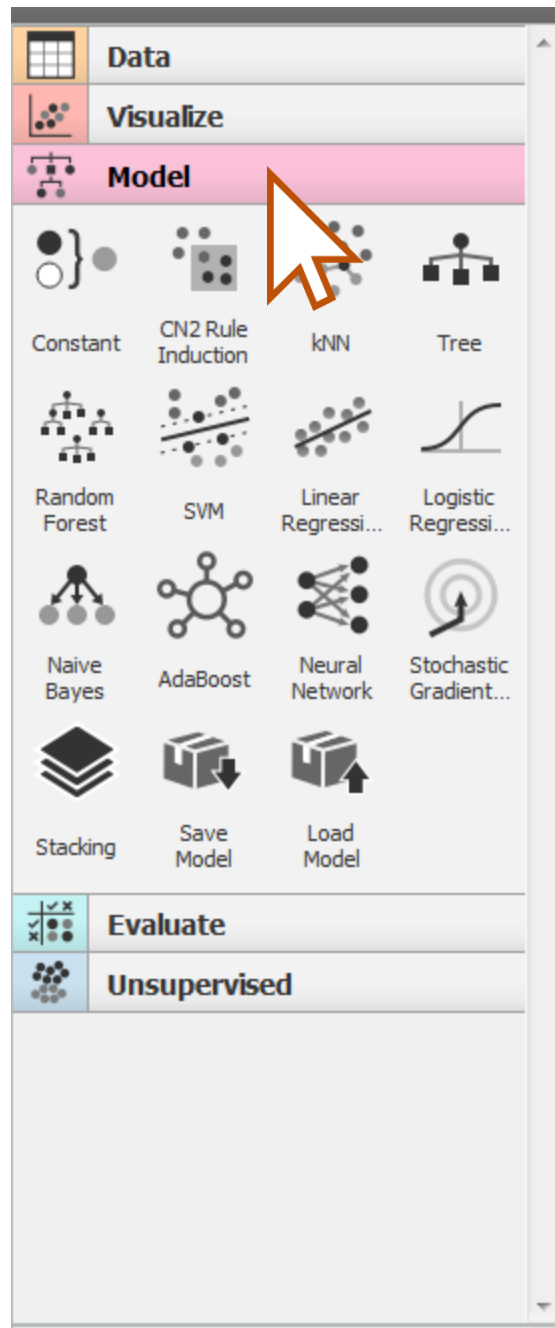


		Predicted			
		Iris-setosa	Iris-versicolor	Iris-virginica	Σ
Actual	Iris-setosa	50	0	0	50
	Iris-versicolor	0	49	1	50
	Iris-virginica	0	5	45	50
Σ		50	54	46	150

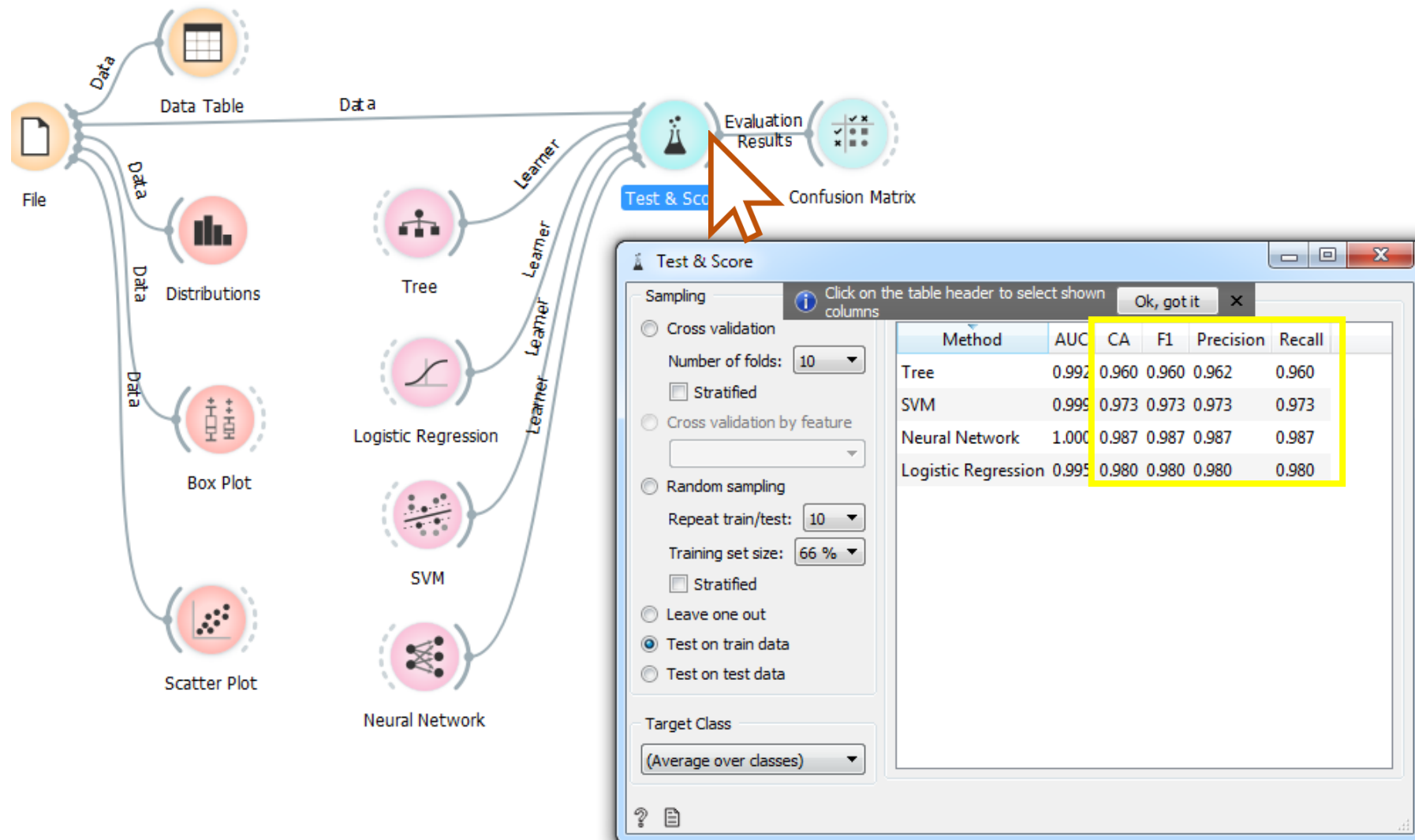
***“96% de
clasificaciones correctas”***

4to paso: Evaluar múltiples algoritmos





Métricas de desempeño facilitan comparar modelos generados por varios algoritmos

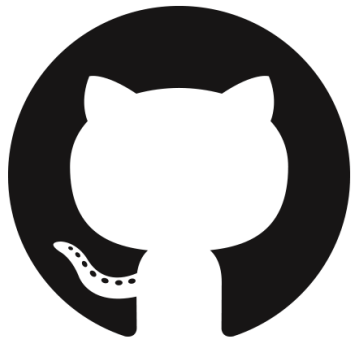


Programas de fuente abierta



Existen un sinnúmero de recursos y oportunidades para el aprender sobre este *Aprendizaje Automático*!

codecademy



fast.ai



deeplearning.ai

Udemy

MIT OPEN COURSEWARE
MASSACHUSETTS INSTITUTE OF TECHNOLOGY

edX

YouTube



kaggle

