



## CURSOS DE VERANO 2014

**APROXIMACIÓN PRÁCTICA A LA CIENCIA DE DATOS  
Y BIG DATA: HERRAMIENTAS KNIME, R, HADOOP Y  
MAHOUT**

**Introducción a KNIME**

María José del Jesus



# KNIME

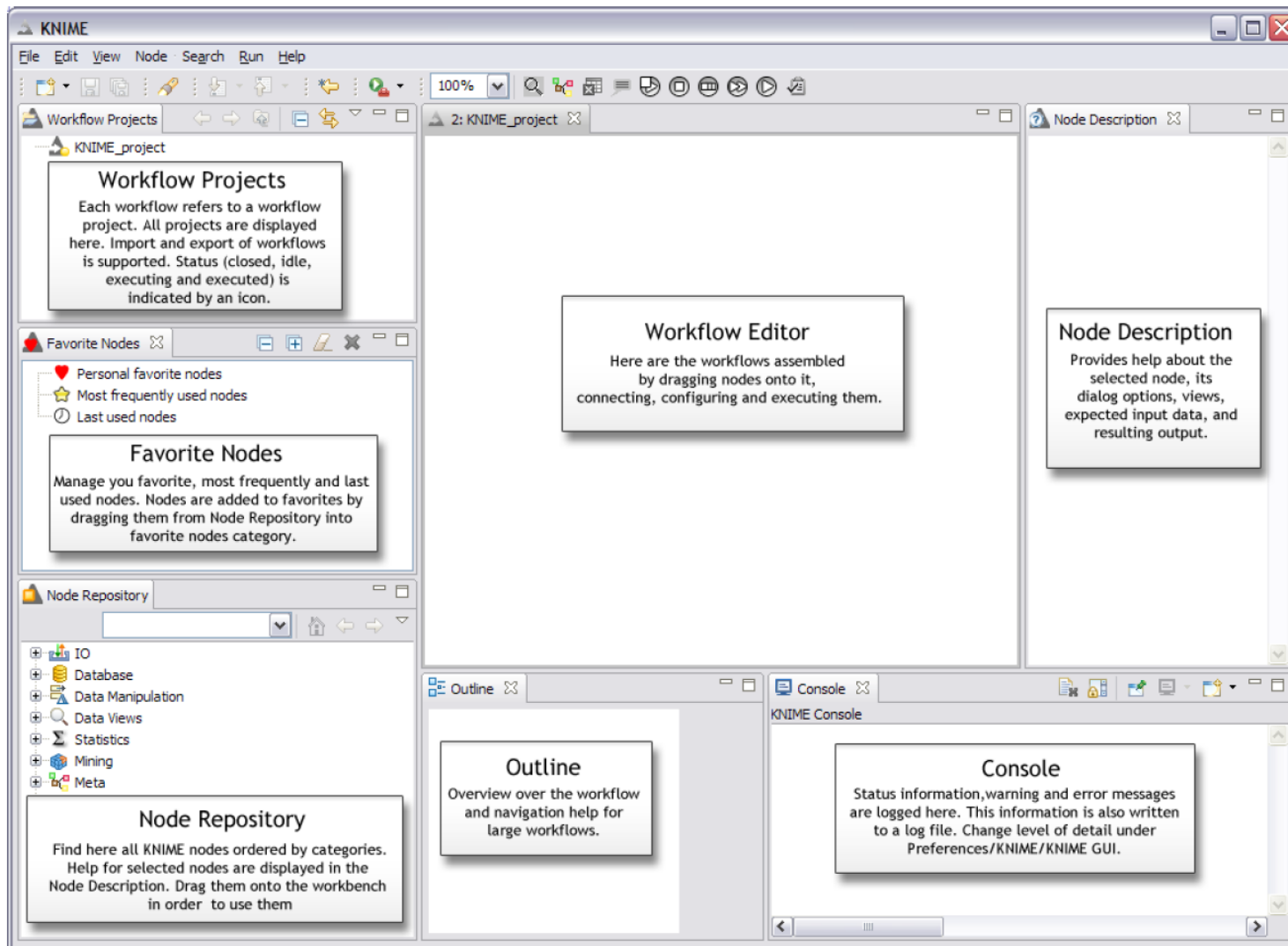
- ¿Por qué aprender KNIME?
- El entorno de trabajo de KNIME
- Ejemplos



## ¿Por qué aprender KNIME?

- KNIME es una herramienta Open Source
- Ofrece el ciclo completo de KDD
  - Visualización de datos
  - Pre-procesado de datos
  - Extracción de modelos mediante algoritmos de Minería de Datos
  - Comparación de modelos
  - Análisis de resultados
- Integración con Weka y R
- Información sobre instalación y uso:
  - <https://www.knime.org>

# El entorno de trabajo en KNIME



## El entorno de trabajo en KNIME

### Conceptos básicos

- Un proceso de análisis de datos se representa mediante un **flujo de trabajo**
- **Construcción de un flujo de trabajo:**
  - arrastrar nodos del Almacén de nodos y pegar en el editor de flujo de trabajo
  - Conectar nodos
- Estado de un **nodo**
  - Rojo → hay que configurar antes de ejecutar
  - Amarillo → el nodo está preparado para ejecutar
  - Verde → el nodo se ha ejecutado
- **Puertos** (de entrada o salida)
  - Solo se pueden conectar puertos del mismo tipo
    - Datos (triángulo amarillo): transfieren tablas de datos entre nodos
    - Bases de datos (cuadrado marrón)
    - PMML: transfieren modelos ya aprendidos
    - Otros puertos



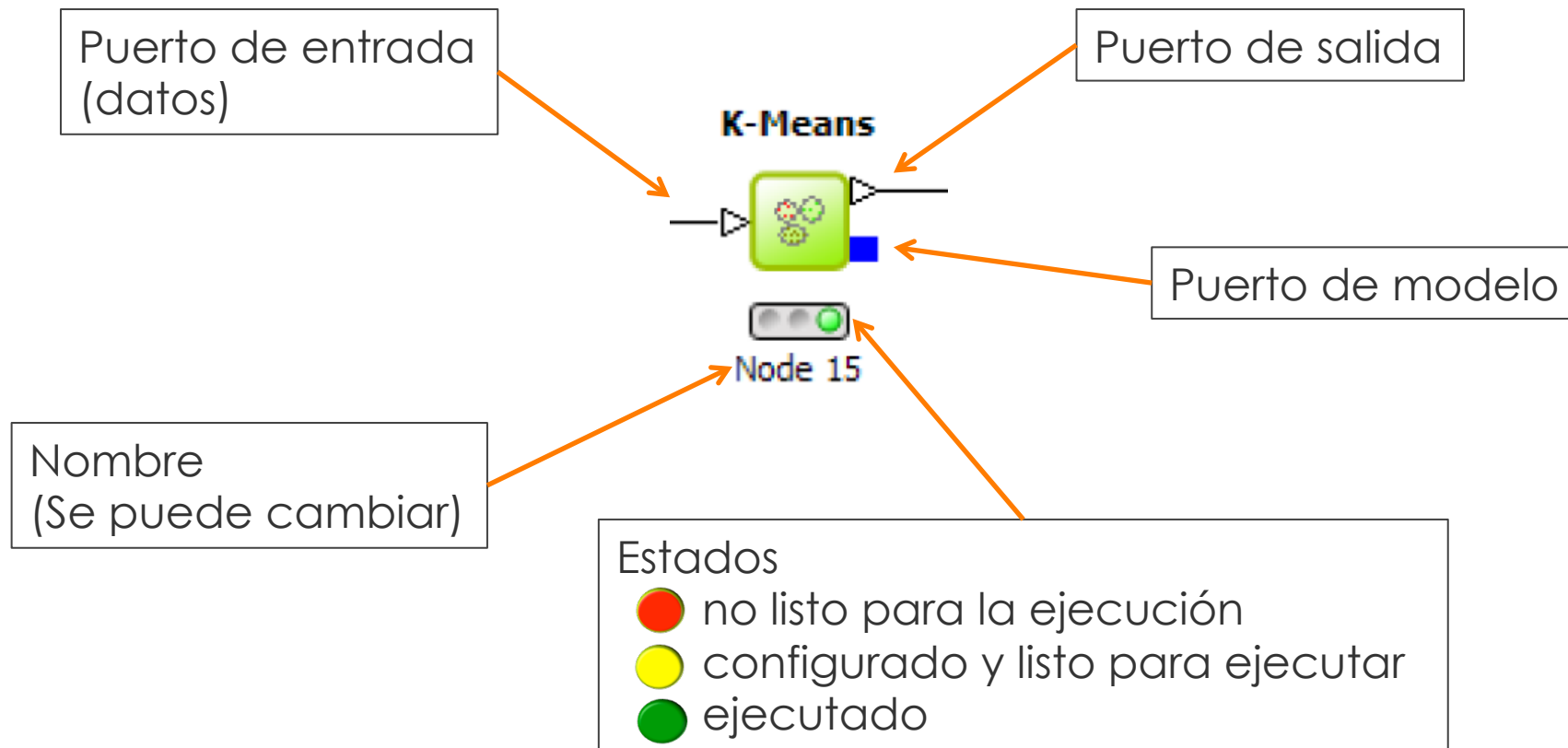
## El entorno de trabajo en KNIME

Acciones básicas:

- ❑ Crear un proyecto
- ❑ Utilizar nodos
- ❑ Construir un flujo de datos
- ❑ Nodo color manager
- ❑ Configuración de nodos
- ❑ Ejecución de flujo de datos
- ❑ Resultados HiLiting

## El entorno de trabajo en KNIME ► nodos

- Son unidades de procesamiento de un workflow



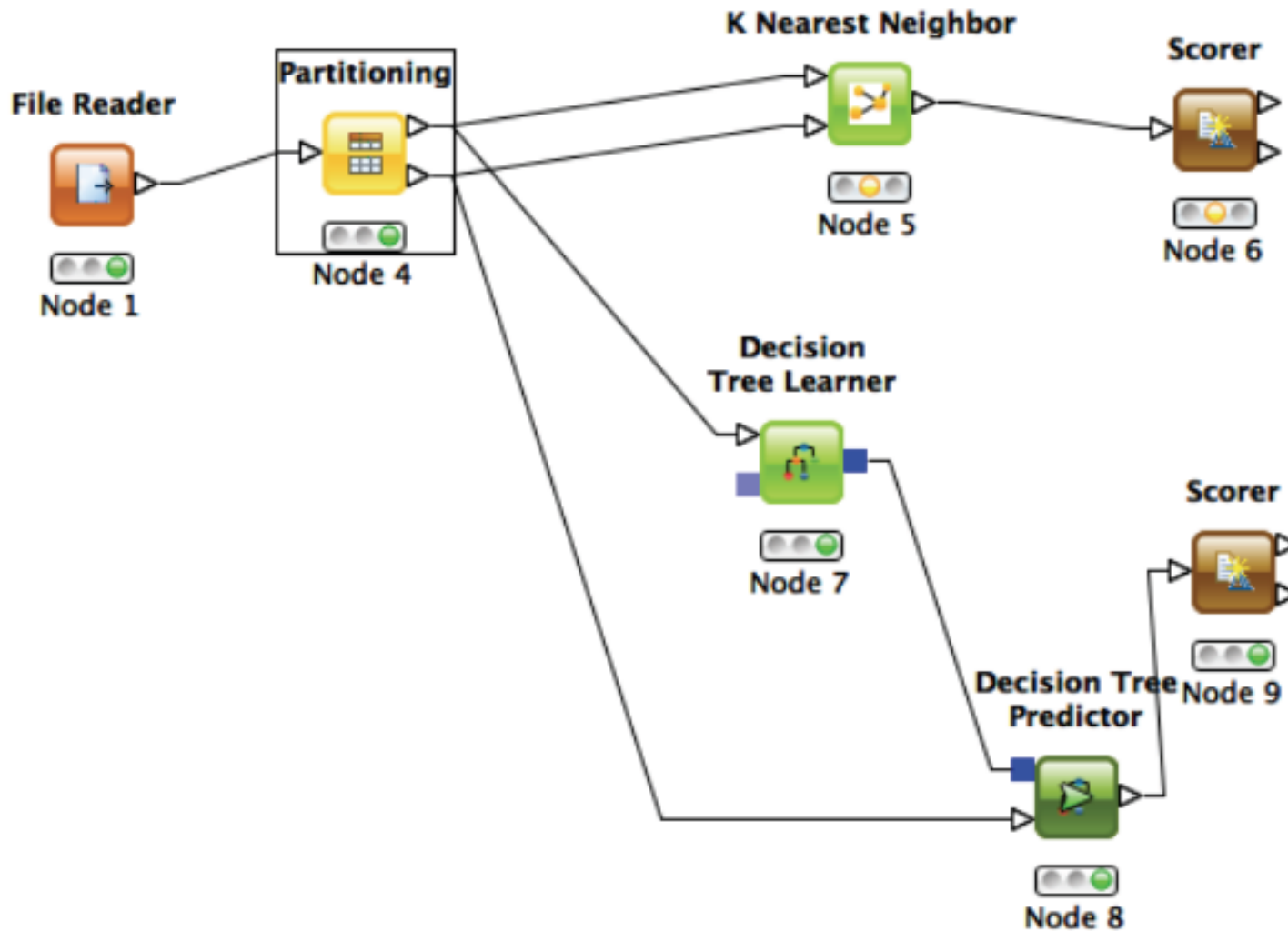


## El entorno de trabajo en KNIME ► Construir un flujo de datos

- Se construye un flujo arrastrando y soltando los nodos desde el almacén de nodos al editor de proyectos y conectándolos entre ellos
- Los datos se transportan entre nodos a través de los puertos
- Es necesario, una vez colocados los nodos en el editor, conectar la salida de cada nodo con el predecesor



# El entorno de trabajo en KNIME ► Construir un flujo de datos





## El entorno de trabajo en KNIME ► nodo Color Manager

- Permite colorear los resultados generados a partir de los datos de entrada
- El coloreo afecta a muchas vistas y ayuda a diferenciar los datos
- Si se inserta este nodo en el flujo de trabajo, los datos se codifican según los colores determinados por el Color Manager

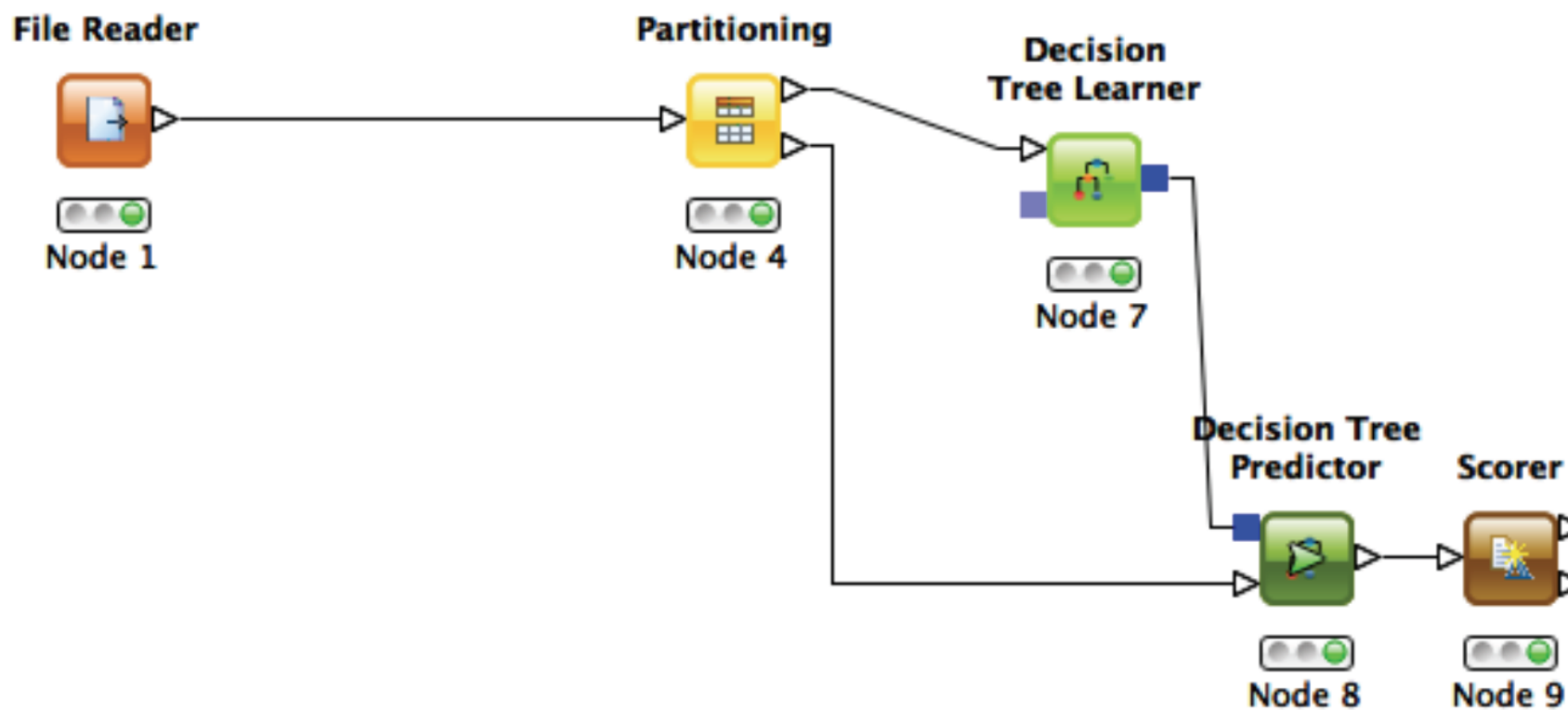
## El entorno de trabajo en KNIME ► ejecución del flujo de datos

- Cuando los nodos del flujo tienen color amarillo, se puede ejecutar
- Los nodos se ejecutan de izquierda a derecha
  - Un nodo puede ejecutarse cuando todos los nodos predecesores han terminado su ejecución
- Formas de ejecución:
  - Por nodo (con la opción Execute)
  - Ejecutando el último nodo del flujo (KNIME ejecuta los predecesores)
  - Seleccionando varios nodos y disparando la ejecución (KNIME determina el orden y ejecuta nodos en paralelo, si es posible)

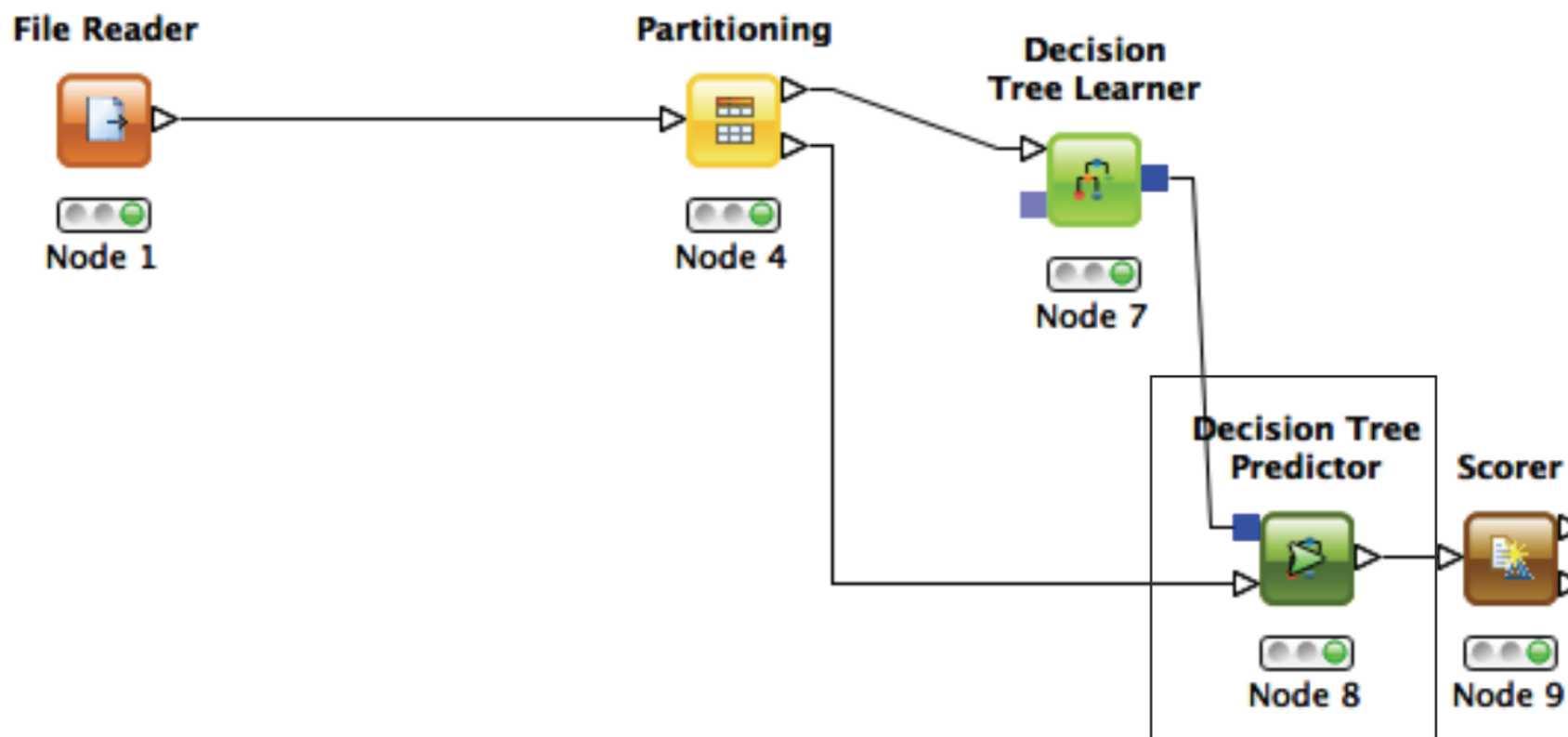
## El entorno de trabajo en KNIME ► HiLiting

- Si se seleccionan datos en una vista y se aplica “Hilite” sobre ellos, se podrá ver el efecto “Hilite” sobre los datos en el resto de vistas del flujo que soporten esta opción
- Los datos seleccionados se resaltarán en color naranja

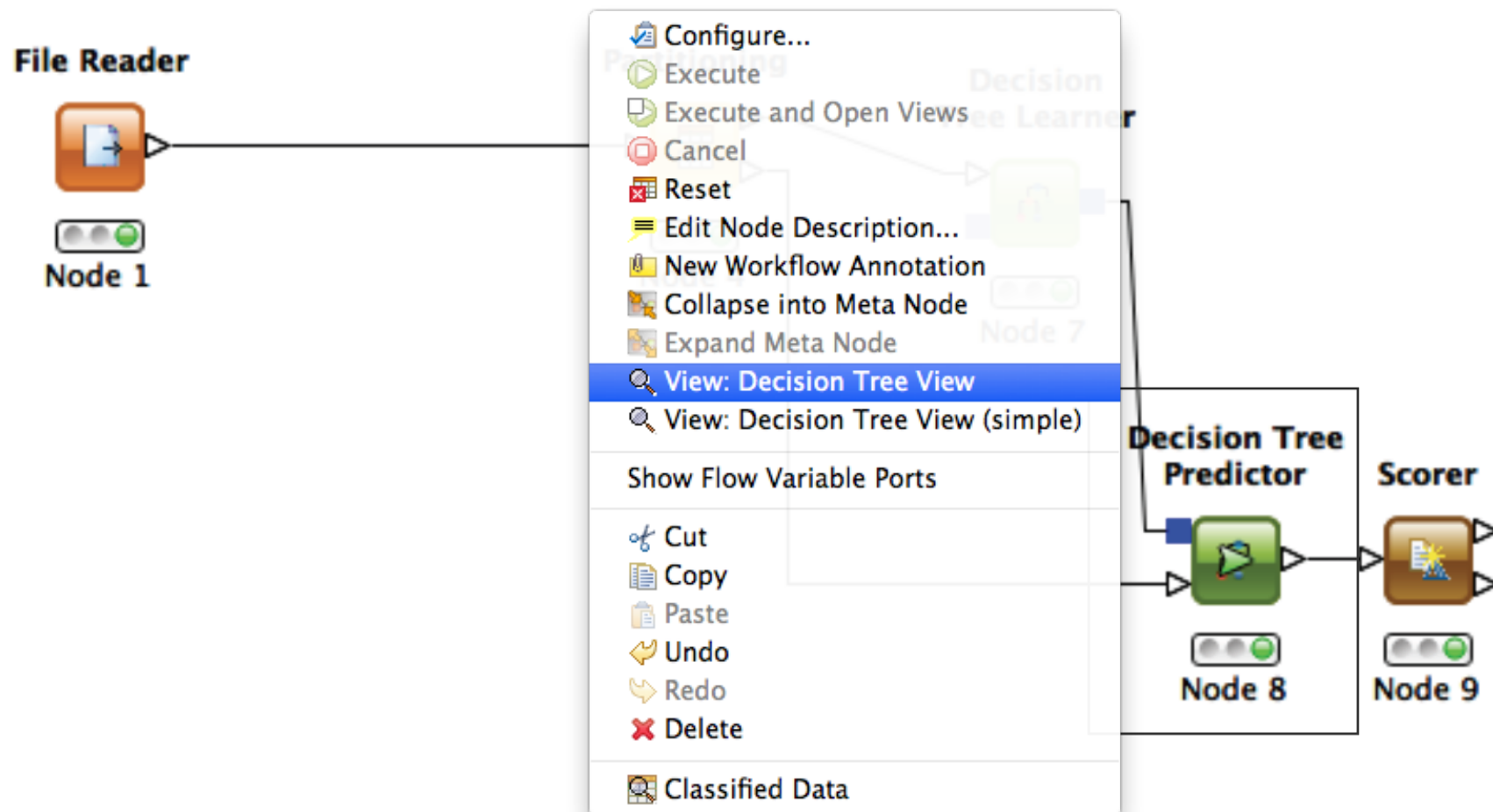
# El entorno de trabajo en KNIME ► HiLiting



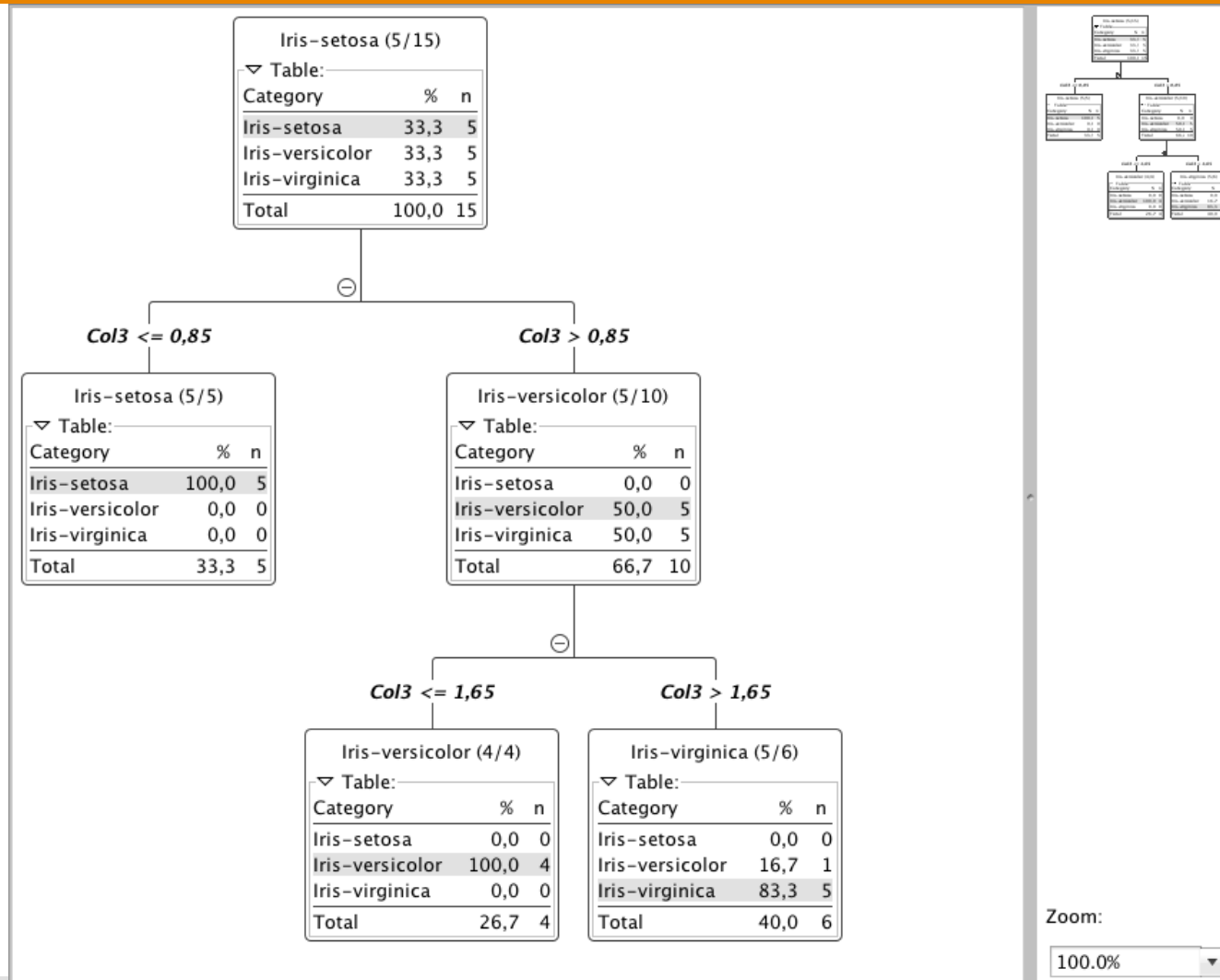
## El entorno de trabajo en KNIME ► HiLiting



## El entorno de trabajo en KNIME ► HiLiting



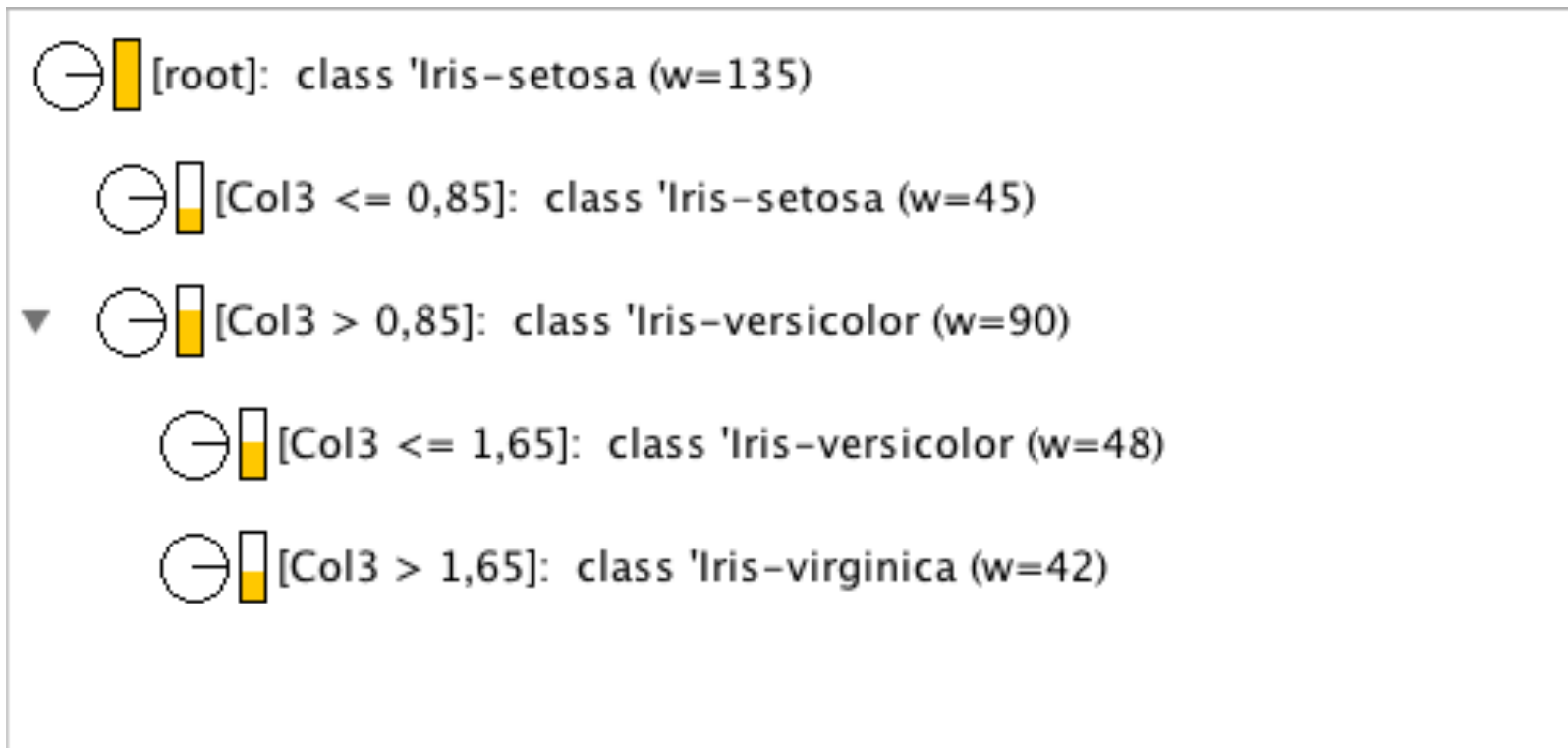
# El entorno de trabajo en KNIME ► HiLiting



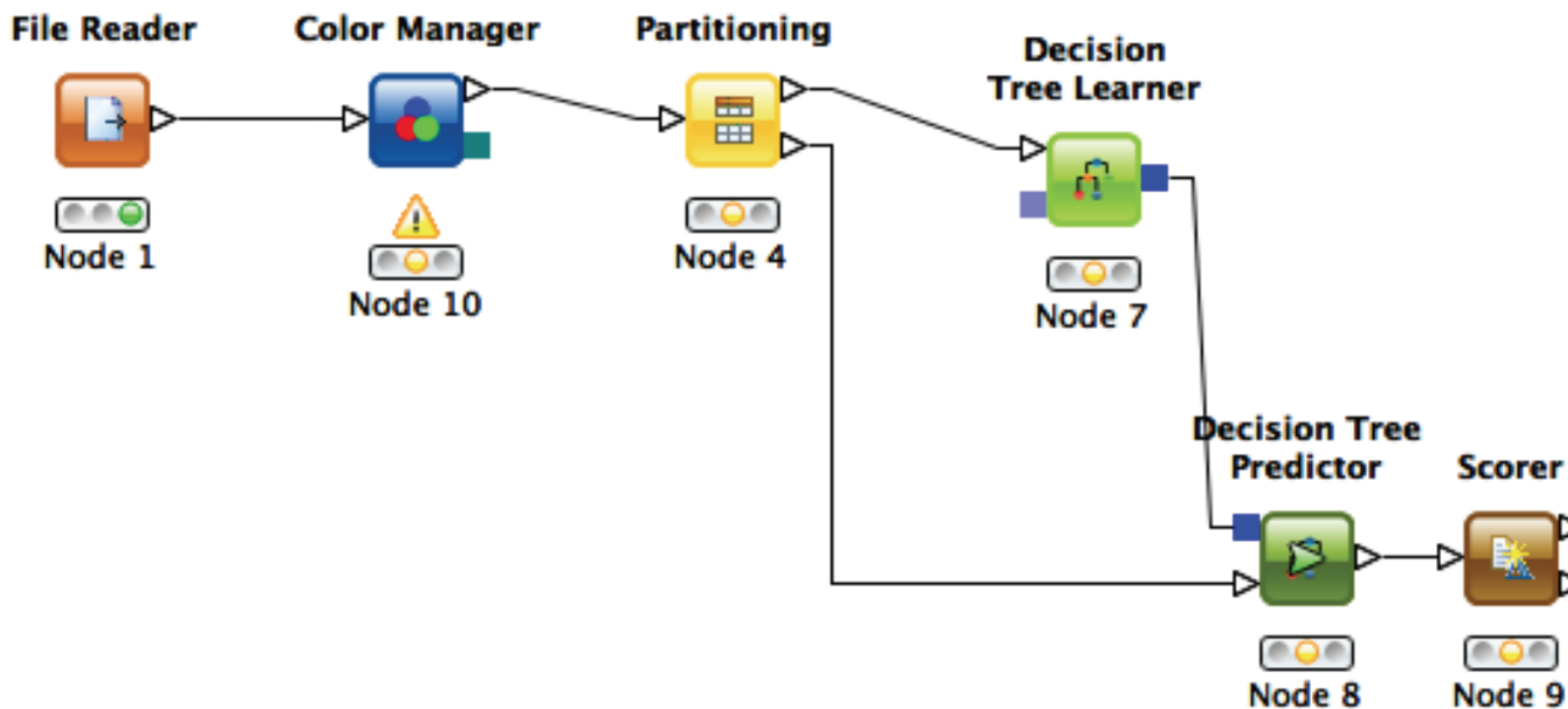


## El entorno de trabajo en KNIME ► HiLiting

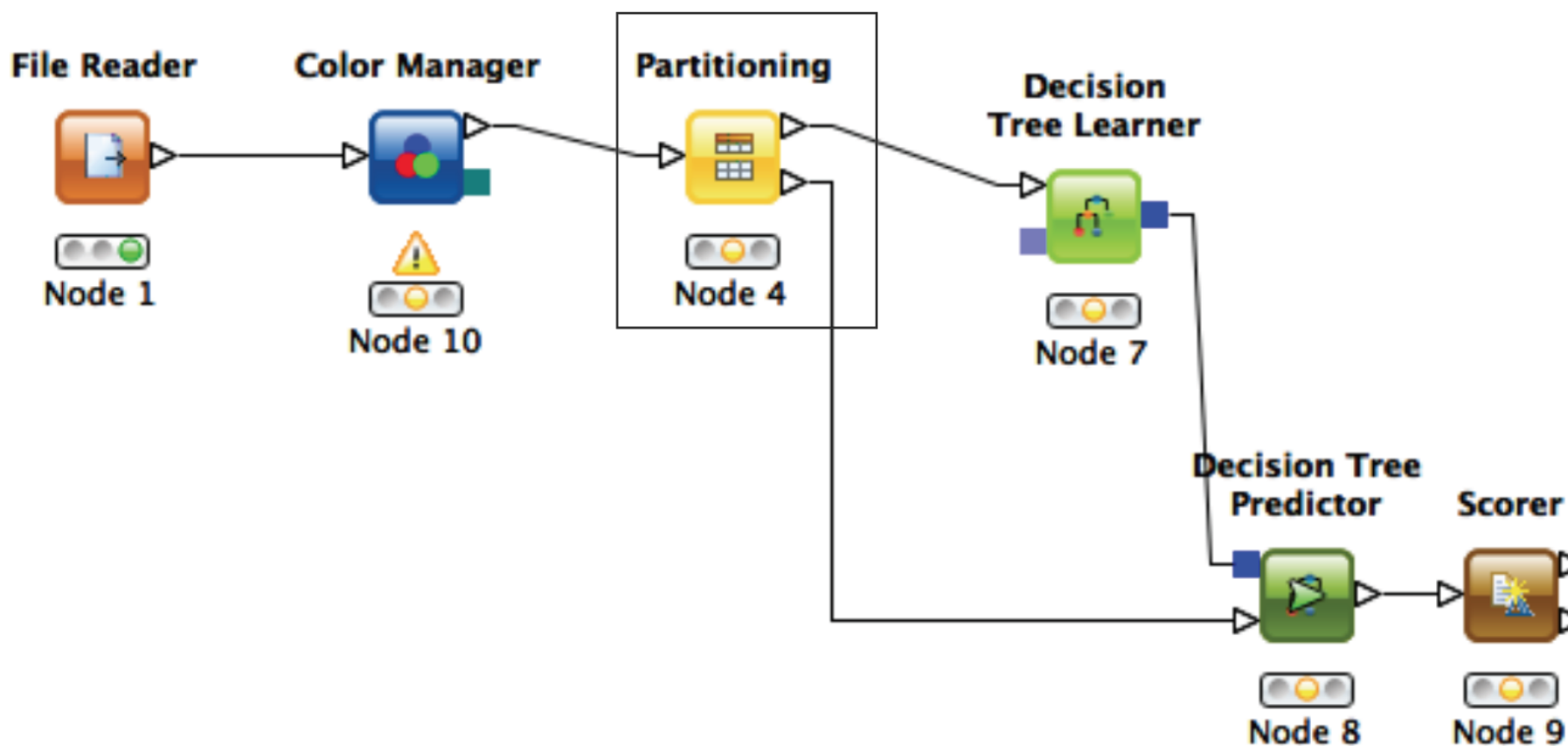
### Árbol simple



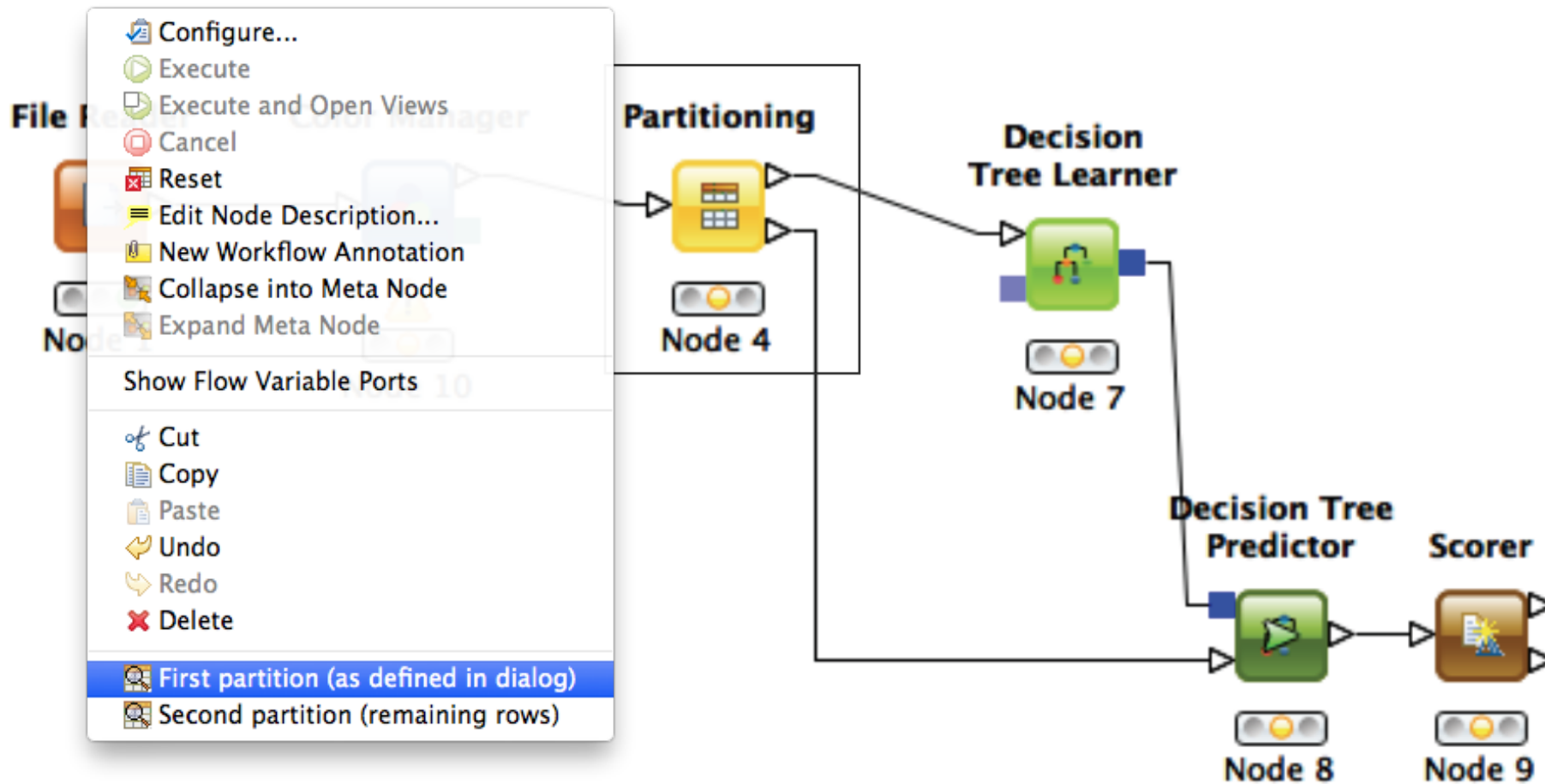
# El entorno de trabajo en KNIME ► HiLiting



# El entorno de trabajo en KNIME ► HiLiting



# El entorno de trabajo en KNIME ► HiLiting



# El entorno de trabajo en KNIME ► HiLiting

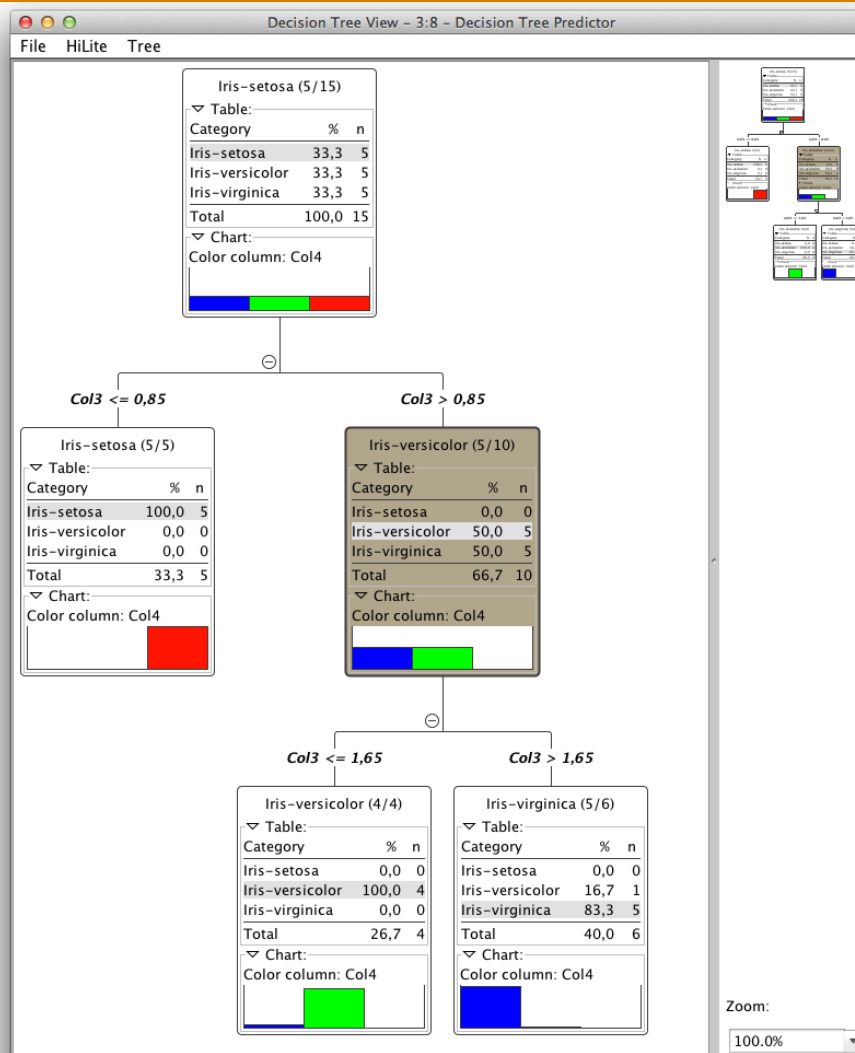
First partition (as defined in dialog) - 3:4 - Partitioning

File

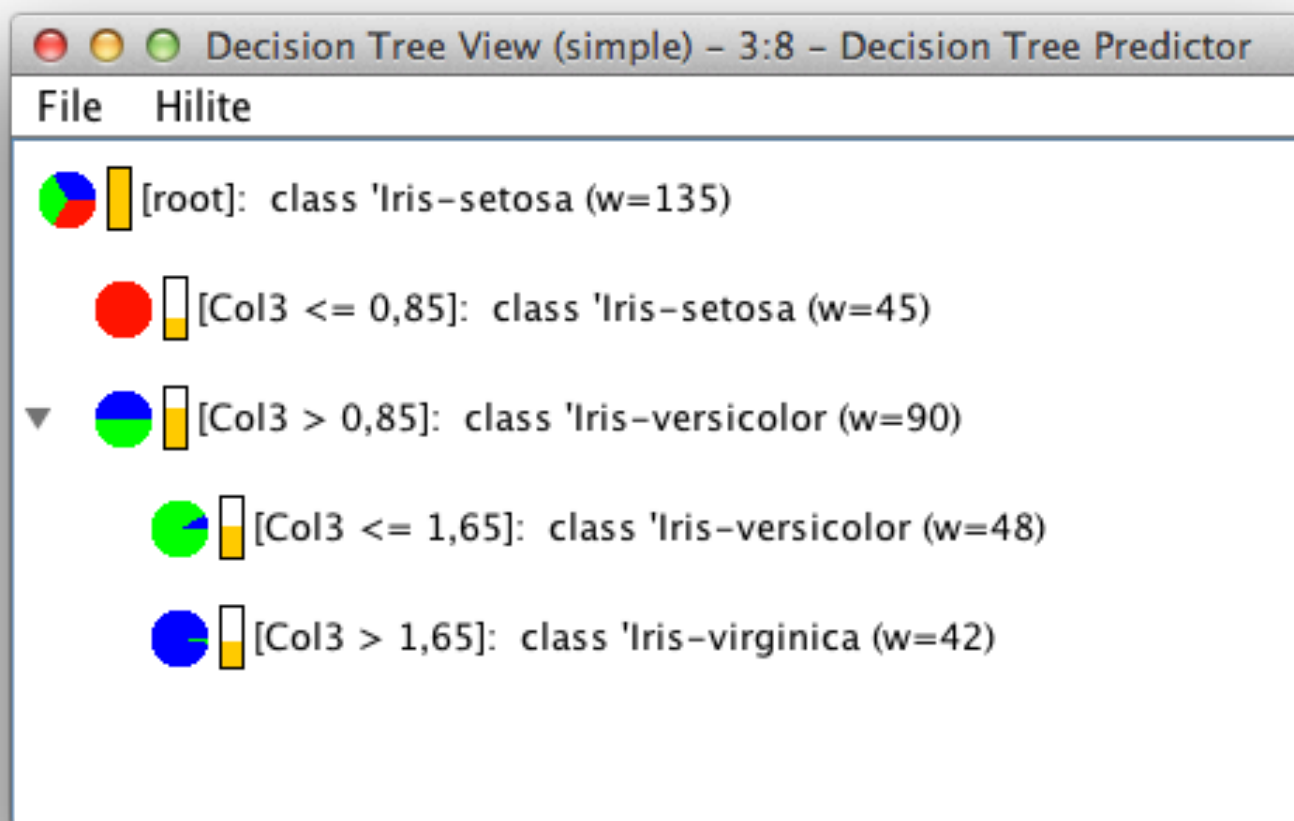
Table "default" - Rows: 15    Spec - Columns: 5    Properties ►

Row ID	D Col0	D Col1	D Col2	D Col3	S Col4
Row0	5.1	3.5	1.4	0.2	Iris-setosa
Row5	5.4	3.9	1.7	0.4	Iris-setosa
Row15	5.7	4.4	1.5	0.4	Iris-setosa
Row36	5.5	3.5	1.3	0.2	Iris-setosa
Row48	5.3	3.7	1.5	0.2	Iris-setosa
Row52	6.9	3.1	4.9	1.5	Iris-versic...
Row64	5.6	2.9	3.6	1.3	Iris-versic...
Row68	6.2	2.2	4.5	1.5	Iris-versic...
Row70	5.9	3.2	4.8	1.8	Iris-versic...
Row91	6.1	3	4.6	1.4	Iris-versic...
Row107	7.3	2.9	6.3	1.8	Iris-virginica
Row112	6.8	3	5.5	2.1	Iris-virginica
Row121	5.6	2.8	4.9	2	Iris-virginica
Row126	6.2	2.8	4.8	1.8	Iris-virginica
Row127	6.1	3	4.9	1.8	Iris-virginica

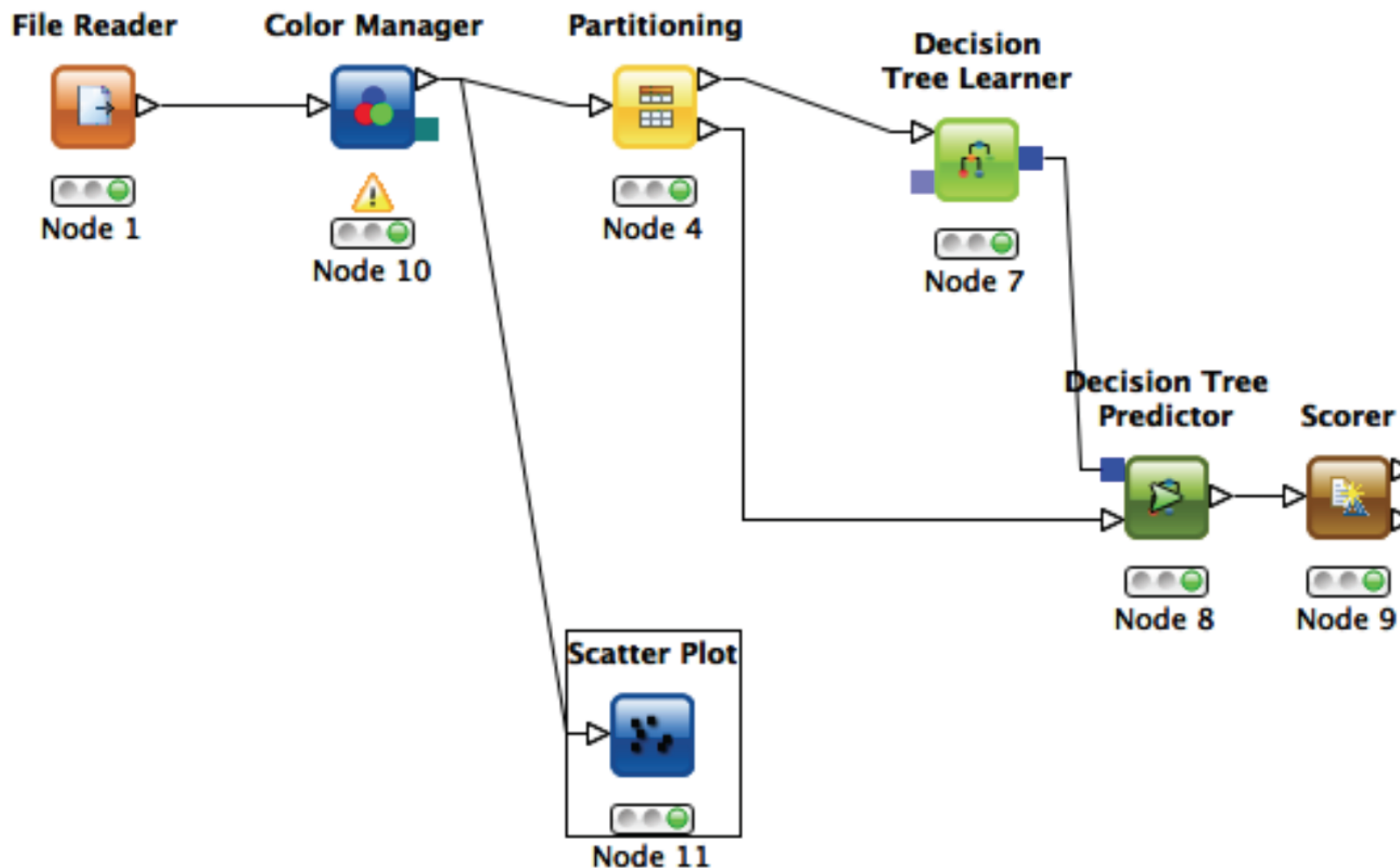
# El entorno de trabajo en KNIME ► HiLiting



## El entorno de trabajo en KNIME ► HiLiting

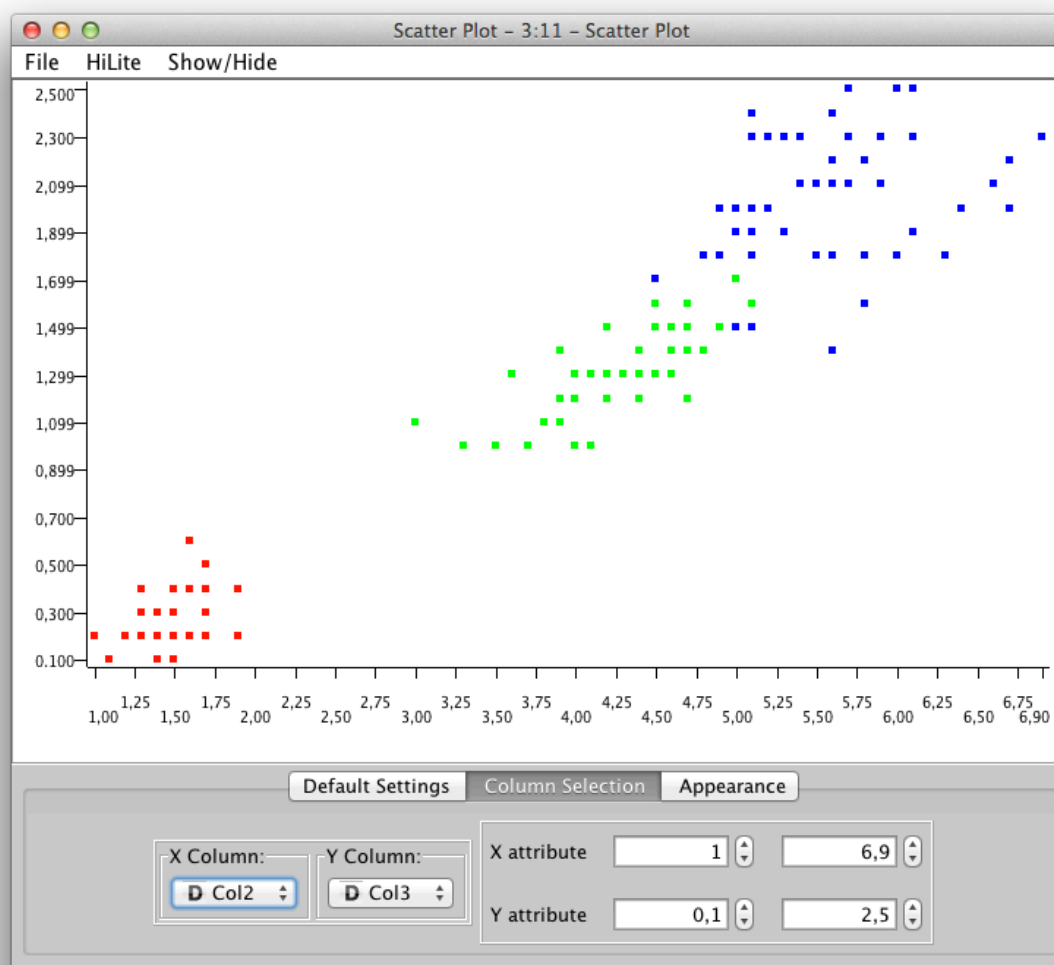


# El entorno de trabajo en KNIME ► HiLiting

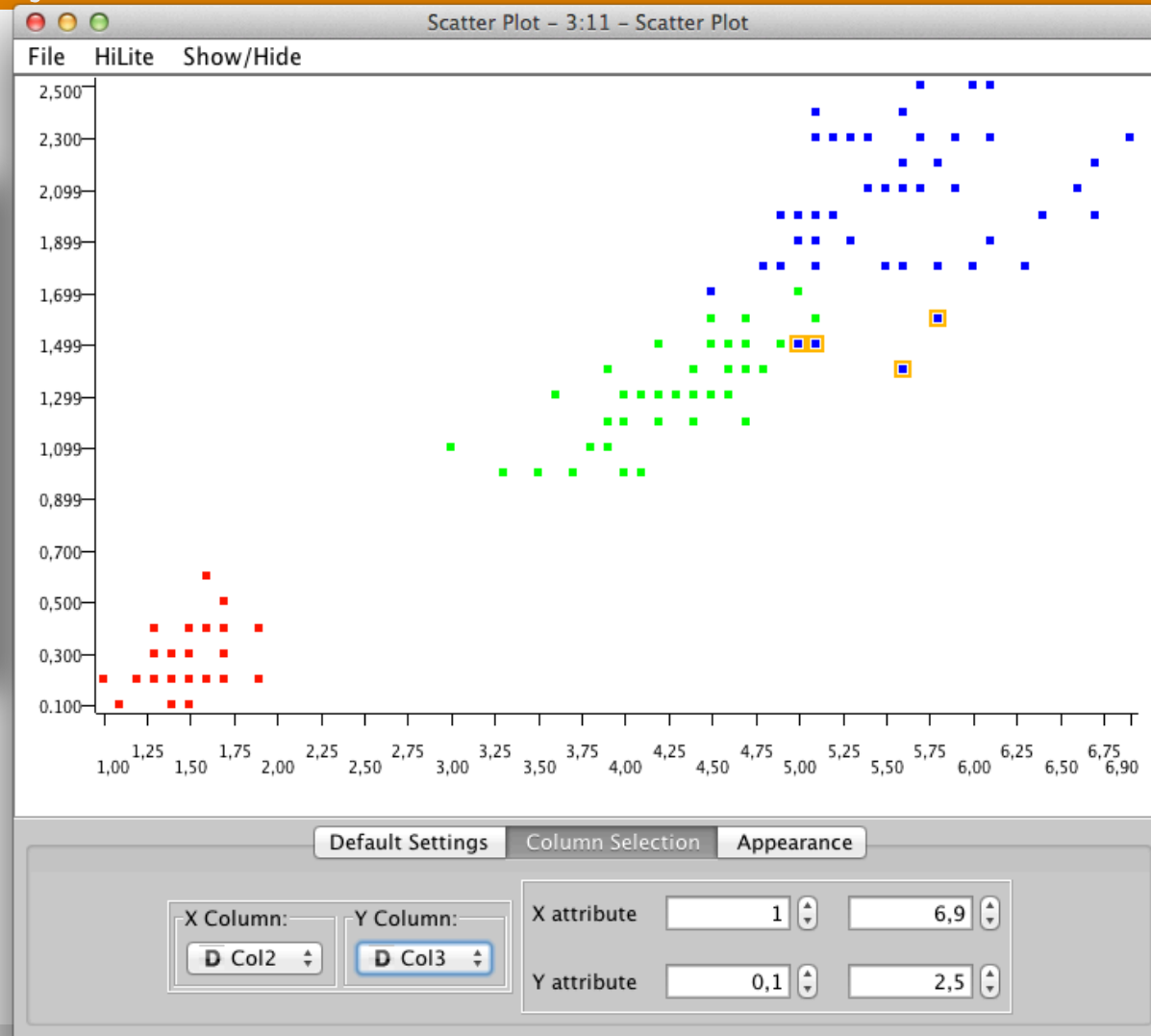
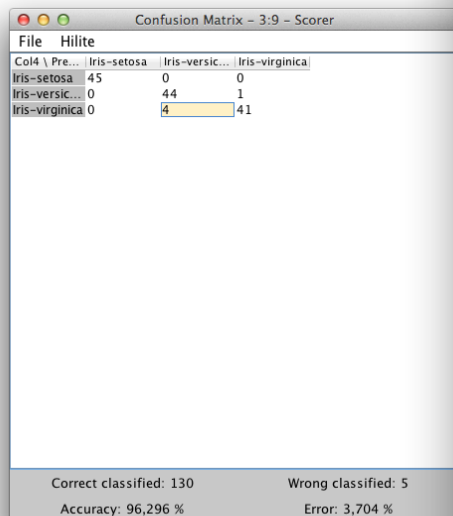




# El entorno de trabajo en KNIME ► HiLiting



# El entorno de trabajo en KNIME ► HiLiting



## El entorno de trabajo en KNIME ► Hotkeys

Task	Hotkey	Description
Node Configuration	F6	opens the configuration dialog of a node
	F7	executes selected nodes
Node Execution	Shift + F7	executes all configured nodes
	Shift + F10	executes configured nodes and opens all views
	F9	Cancels selected running nodes
	Shift + F9	Cancels all running nodes
Move Nodes and Annotations	Ctrl + Shift + Arrow	moves a selected node in the workflow editor
	Ctrl + Shift + PgUp/PgDown	Moves the selected up or down in z order
	F8	resets selected nodes
Workflow Operations	Ctrl + S	Saves the workflow
	Ctrl + Shift + S	Saves all open workflows
	Ctrl + Shift + W	Closes all open workflows
Meta-node	Shift + F12	Opens meta-node wizard



## Ejemplo: Iris

- ▣ Carga de datos
- ▣ Visualización
- ▣ Análisis predictivo
- ▣ Análisis descriptivo



<https://archive.ics.uci.edu/ml/datasets/Iris>

- ▣ N° ejemplos: 150
- ▣ N° variables: 4
- ▣ N° clases: 3 (50/50/50)



## Ejemplo: Iris



1. Cargar el fichero de ejemplos iris.dat
2. Obtener medidas estadísticas
3. Asignar a cada ejemplo un color en función de la clase a la que pertenece
4. Visualizar el conjunto de ejemplos en base a pares de variables
  - ▣ Determinar el par de variables “más relevantes”
5. Realizar una partición con hold-out al 60% estratificada
6. Visualizar el conjunto de test en base a las dos variables seleccionadas en el paso 4

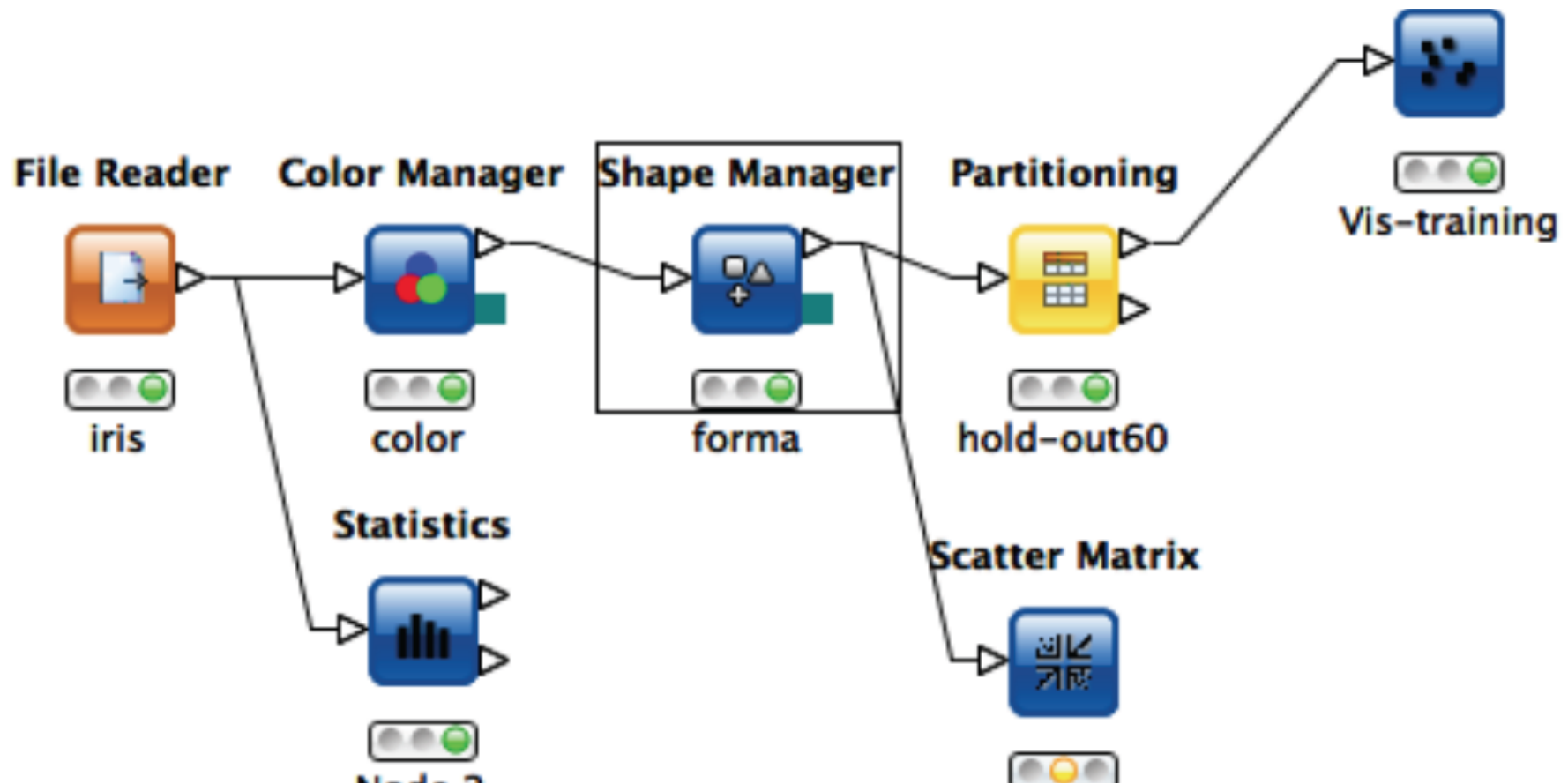


## Ejemplo: Iris

1. File Reader: IO → Read
2. Statistics: Statistics
  - ▣ Calcula y muestra estadísticas
3. Color Manager: Data Views → Property
  - ▣ Asigna colores a las clases
4. Scatter Matrix : Data Views → Utility
  - ▣ Visualiza los ejemplos según pares de variables (scatter plots)
5. Partitioning: Data Manipulation → Row → Transform
  - ▣ Hold-out
6. Scatter Plot: Data Views → Utility



# Ejemplo 1: Iris ▶ Carga de datos y visualización



# Ejemplo de visualizar



- Configure...
- Execute
- Execute and Open Views
- Cancel
- Reset
- Edit Node Description...
- New Workflow Annotation
- Collapse into Meta Node
- Expand Meta Node
- View: Statistics View**
- Show Flow Variable Ports
- Cut
- Copy
- Paste
- Undo
- Redo
- Delete
- Statistics Table
- Occurrences Table

Statistics View - 3:2 - Statistics

File

Numeric columns    Nominal columns

Row ID	D Col0	D Col1	D Col2	D Col3
Minimum	4.3	2	1.1	0.1
Maximum	7.9	4.2	6.9	2.5
Mean	5.853	3.055	3.765	1.193
Std. deviation	0.855	0.422	1.785	0.771
Variance	0.732	0.178	3.185	0.595
Overall sum	702.4	366.6	451.8	143.2
No. missings	0	0	0	0
Median	?	?	?	?
Row count	120	120	120	120

Numeric columns    Nominal columns

Col0	Col1	Col2	Col3	Col4
No. missings: 0	No. missings: 0	No. missings: 0	No. missings: 0	No. missings: 0
<b>Top 20:</b> 5.0 : 8 6.3 : 7 5.1 : 6 4.9 : 6 5.4 : 6 5.7 : 6 5.5 : 6 4.8 : 5 5.8 : 5 6.5 : 5 6.0 : 5	<b>Top 20:</b> 3.0 : 25 2.8 : 11 3.1 : 10 3.4 : 10 2.9 : 10 3.2 : 8 2.5 : 6 3.5 : 5 3.8 : 5 2.7 : 5 3.3 : 4	<b>Top 20:</b> 1.4 : 11 1.5 : 9 5.1 : 7 1.3 : 6 1.6 : 6 4.5 : 6 5.6 : 5 4.7 : 4 4.0 : 4 1.7 : 3 4.9 : 3	<b>Top 20:</b> 0.2 : 23 1.3 : 11 1.5 : 9 1.8 : 8 1.4 : 7 1.0 : 7 2.3 : 7 0.4 : 6 0.1 : 6 2.0 : 6 0.3 : 5	<b>Top 20:</b> Iris-setosa : 40 Iris-versicolor : 40 Iris-virginica : 40

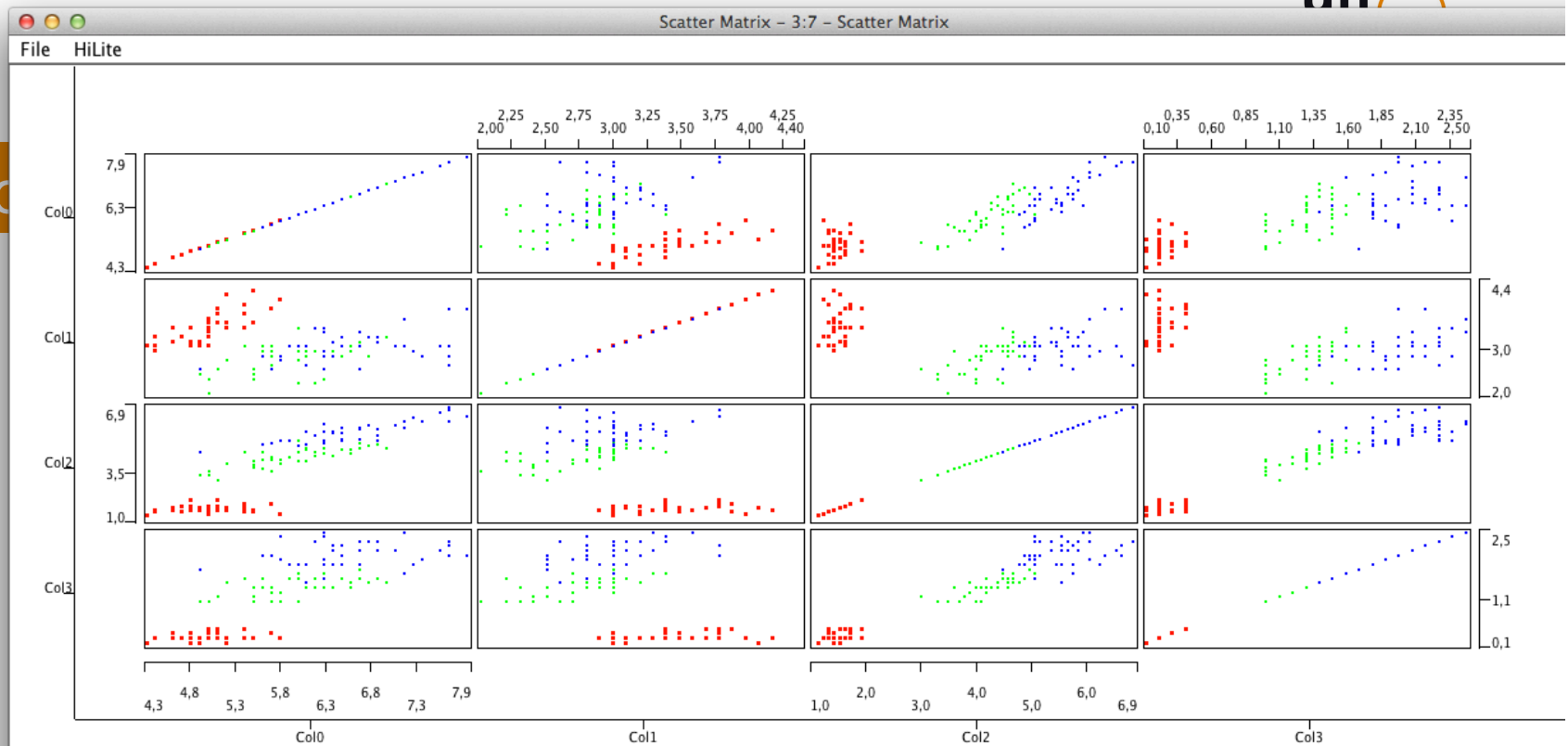


# Ejemplo

Scatter Matrix



Node 7



Default Settings Column Selection Appearance

Exclude

Column(s):  Search

Select all search hits

S Col4

Select

add >>

add all >>

<< remove

<< remove all

Include

Column(s):  Search

Select all search hits

D Col0

D Col1

D Col2

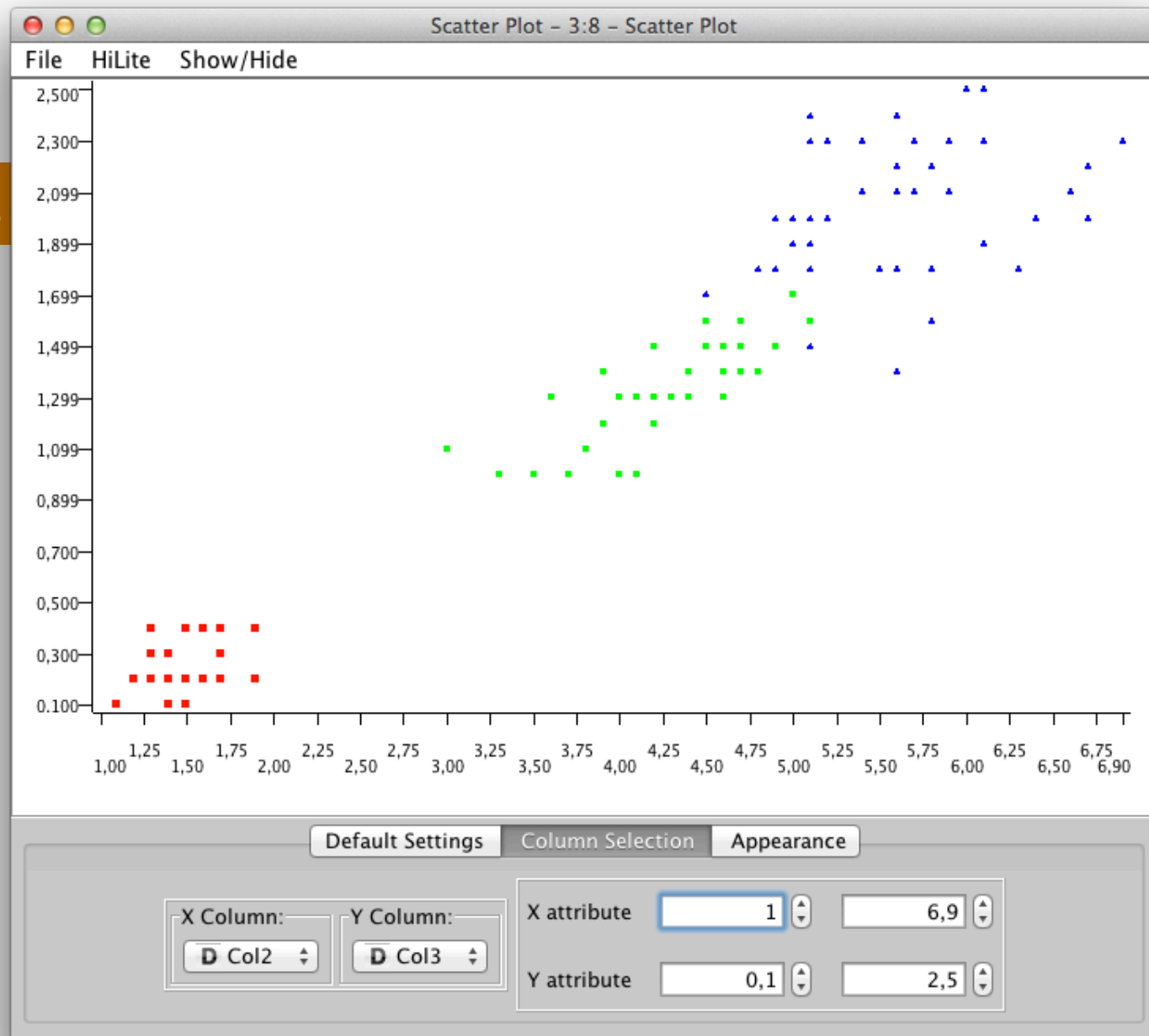
D Col3

# Ejemplo 1: Iris

Scatter Plot

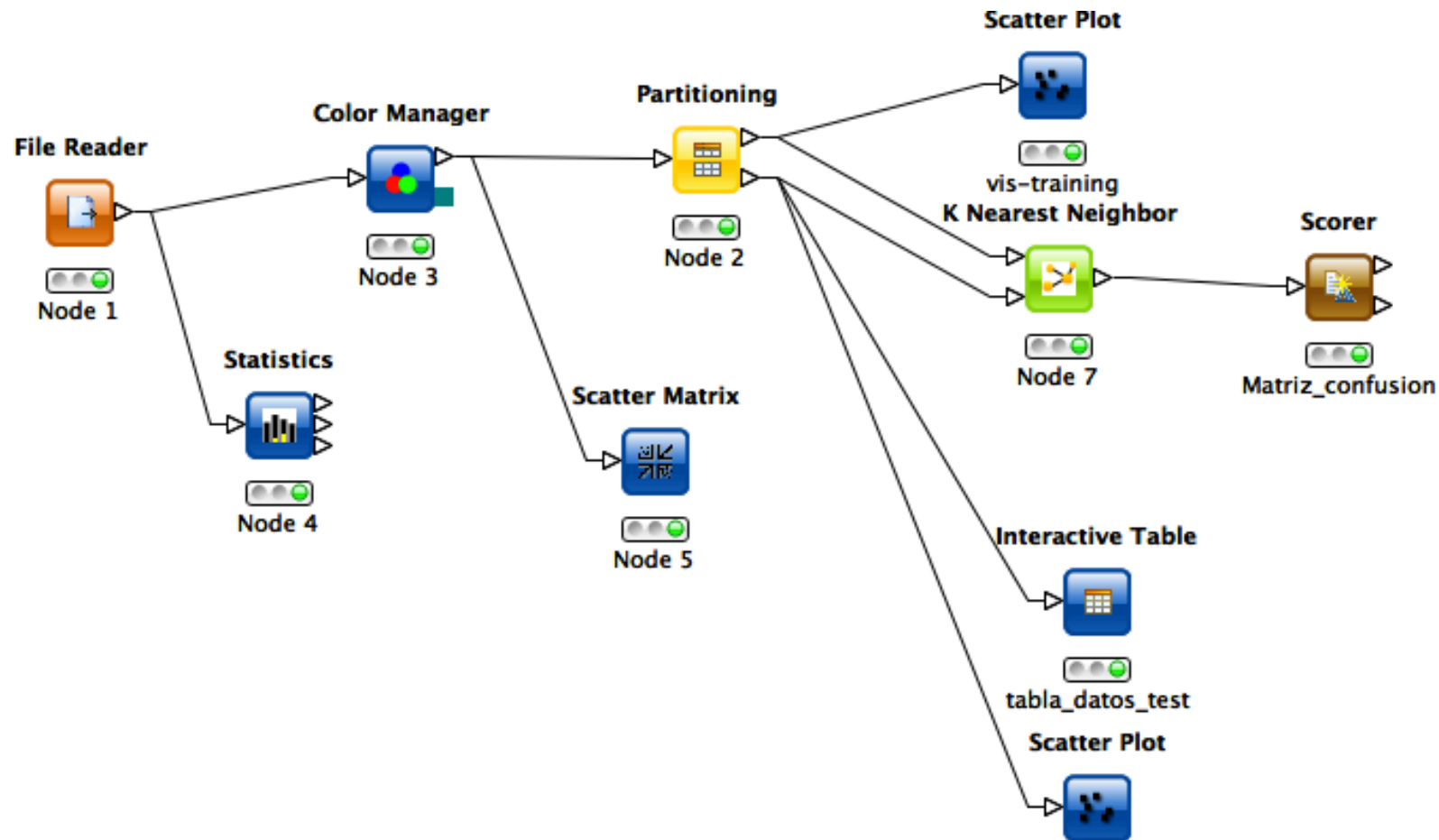


Node 8





# Ejemplo 1: Iris ► Análisis predictivo ► kNN



# Ejemplo 1: Iris

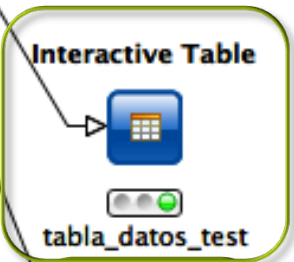
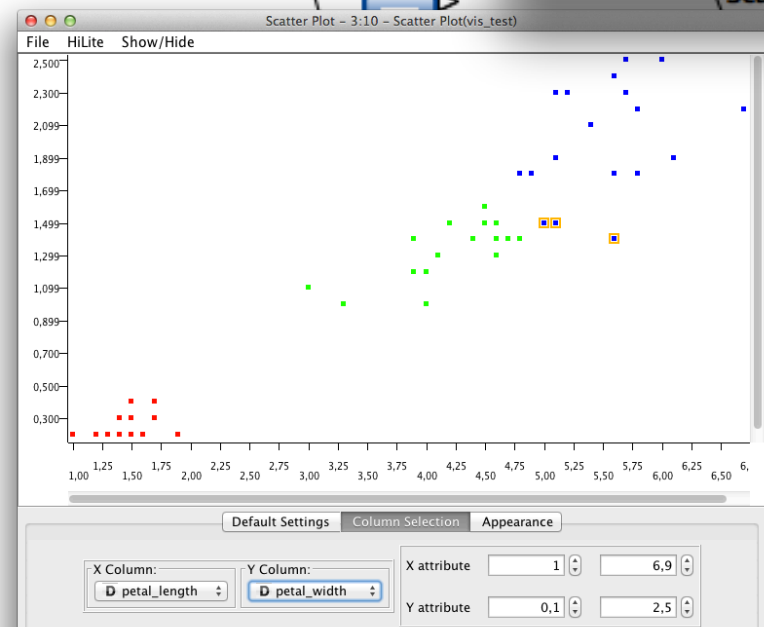
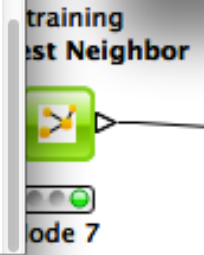
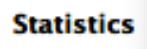
Table View - 3:9 - Interactive Table(tabla\_datos\_test)

Row ID	D sepal...	D sepal...	D petal...	D petal...	S class
Row91	6.1	5.4	4.0	1.4	Iris-versic...
Row92	5.8	2.6	4	1.2	Iris-versic...
Row93	5	2.3	3.3	1	Iris-versic...
Row98	5.1	2.5	3	1.1	Iris-versic...
Row99	5.7	2.8	4.1	1.3	Iris-versic...
Row100	6.3	3.3	6	2.5	Iris-virginica
Row103	6.3	2.9	5.6	1.8	Iris-virginica
Row104	6.5	3	5.8	2.2	Iris-virginica
Row108	6.7	2.5	5.8	1.8	Iris-virginica
Row117	7.7	3.8	6.7	2.2	Iris-virginica
Row119	6	2.2	5	1.5	Iris-virginica
Row120	6.9	3.2	5.7	2.3	Iris-virginica
Row123	6.3	2.7	4.9	1.8	Iris-virginica
Row126	6.2	2.8	4.8	1.8	Iris-virginica
Row130	7.4	2.8	6.1	1.9	Iris-virginica
Row133	6.3	2.8	5.1	1.5	Iris-virginica
Row134	6.1	2.6	5.6	1.4	Iris-virginica
Row136	6.3	3.4	5.6	2.4	Iris-virginica
Row138	6	3	4.8	1.8	Iris-virginica
Row139	6.9	3.1	5.4	2.1	Iris-virginica
Row140	6.7	3.1	5.6	2.4	Iris-virginica
Row141	6.9	3.1	5.1	2.3	Iris-virginica
Row142	5.8	2.7	5.1	1.9	Iris-virginica
Row144	6.7	3.3	5.7	2.5	Iris-virginica
Row145	6.7	3	5.2	2.3	Iris-virginica

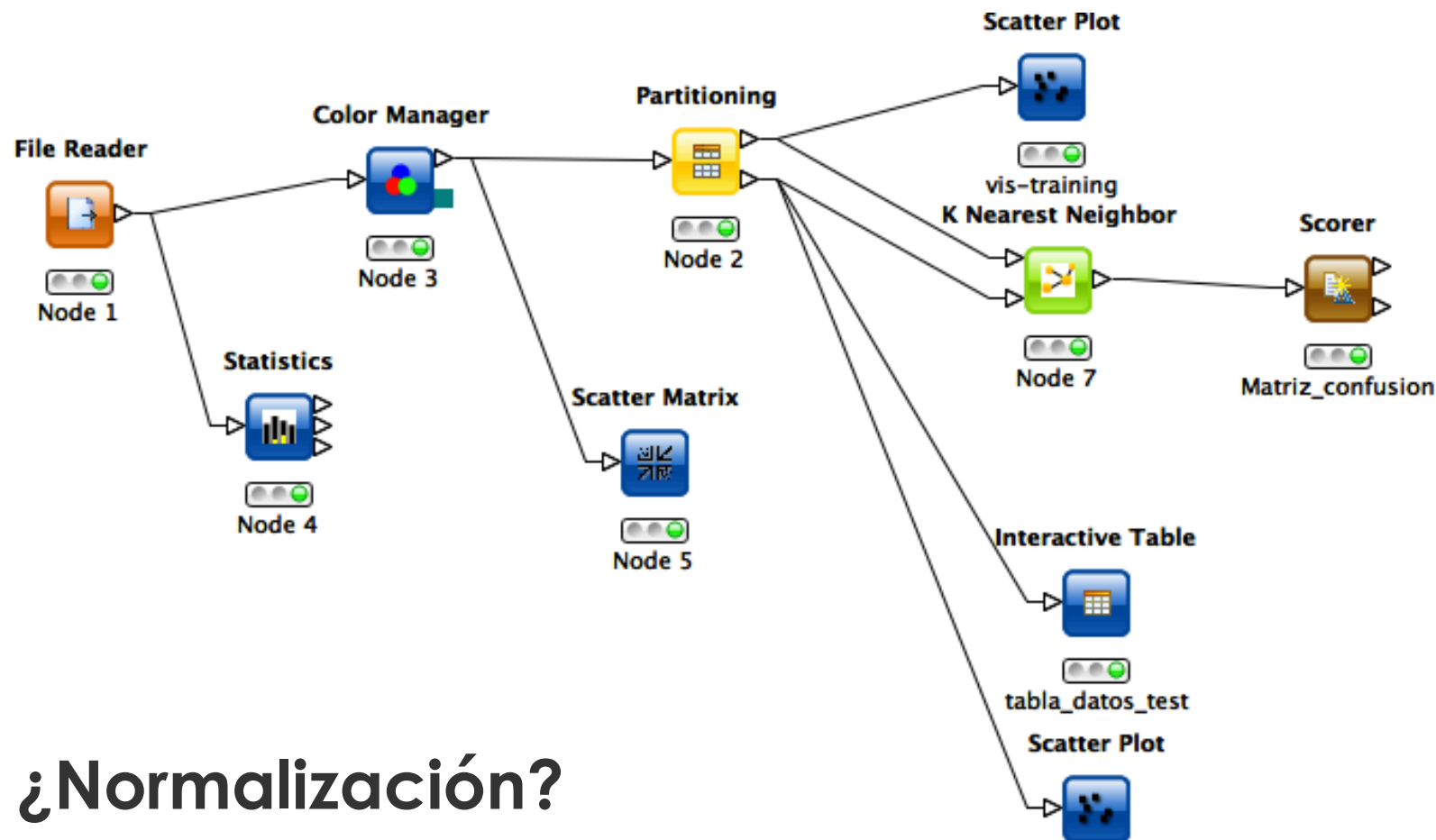
Confusion Matrix - 3:8 - Scorer(Matrix\_...)

class \ Cla...	Iris-setosa	Iris-versic...	Iris-virginica
Iris-setosa	20	0	0
Iris-versic...	0	19	1
Iris-virginica	0	3	17

Correct classified: 56      Wrong classified: 4  
 Accuracy: 93,333 %      Error: 6,667 %  
 Cohen's kappa ( $\kappa$ ) 0,9



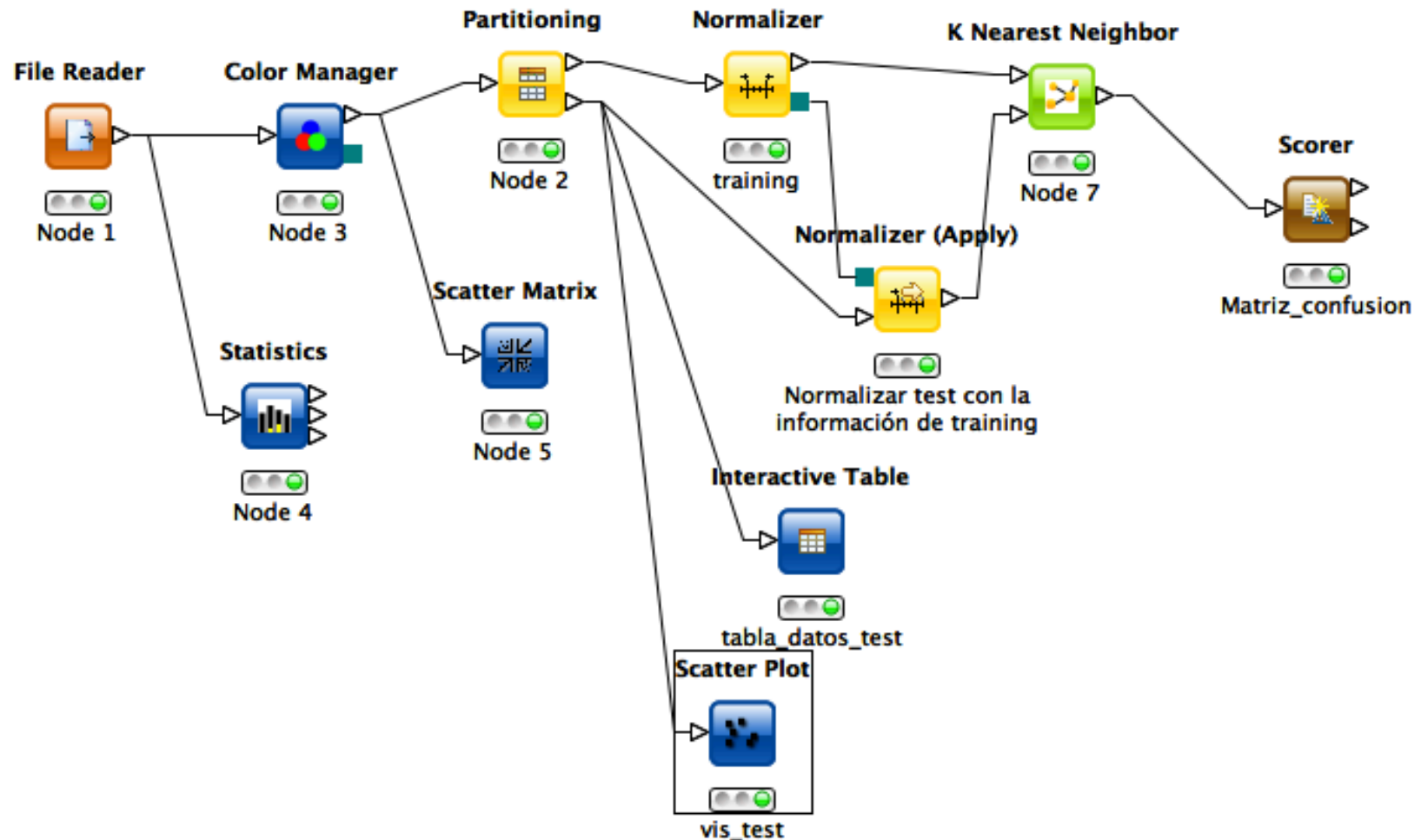
## Ejemplo 1: Iris ► Análisis predictivo ► kNN



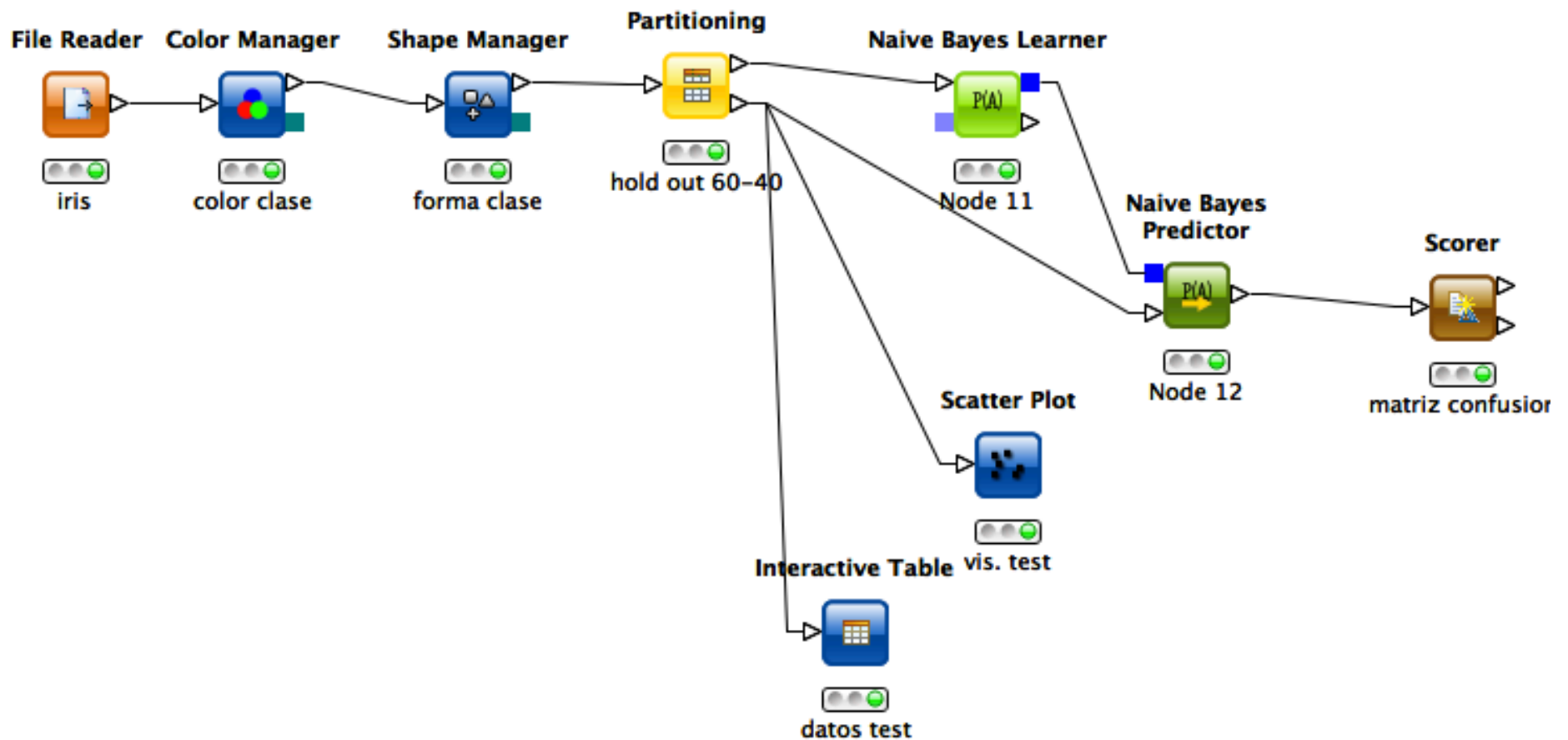
# ¿Normalización?

Normalizer: Data Manipulation → Column → Transform

# Ejemplo 1: Iris ▶ Análisis predictivo ▶ kNN

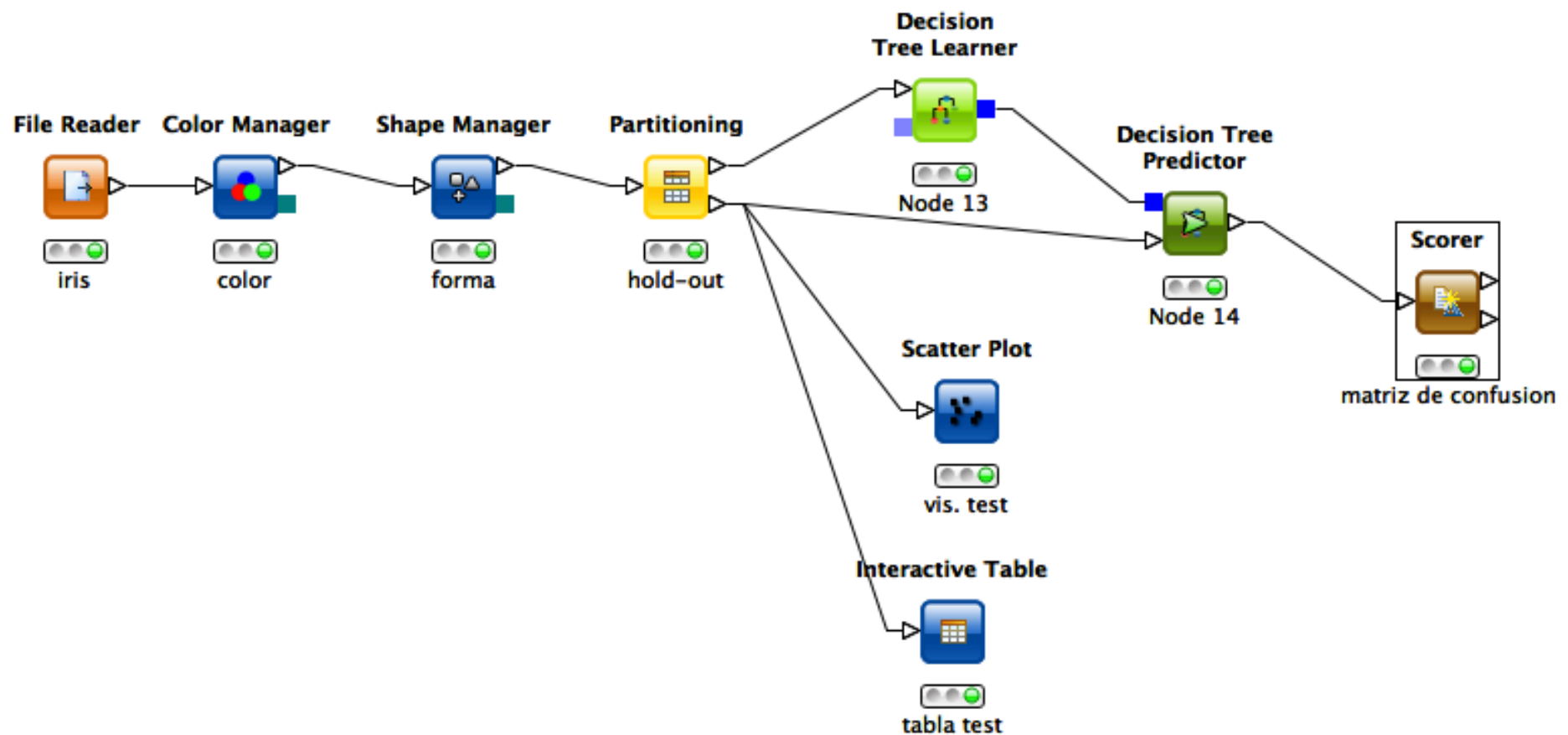


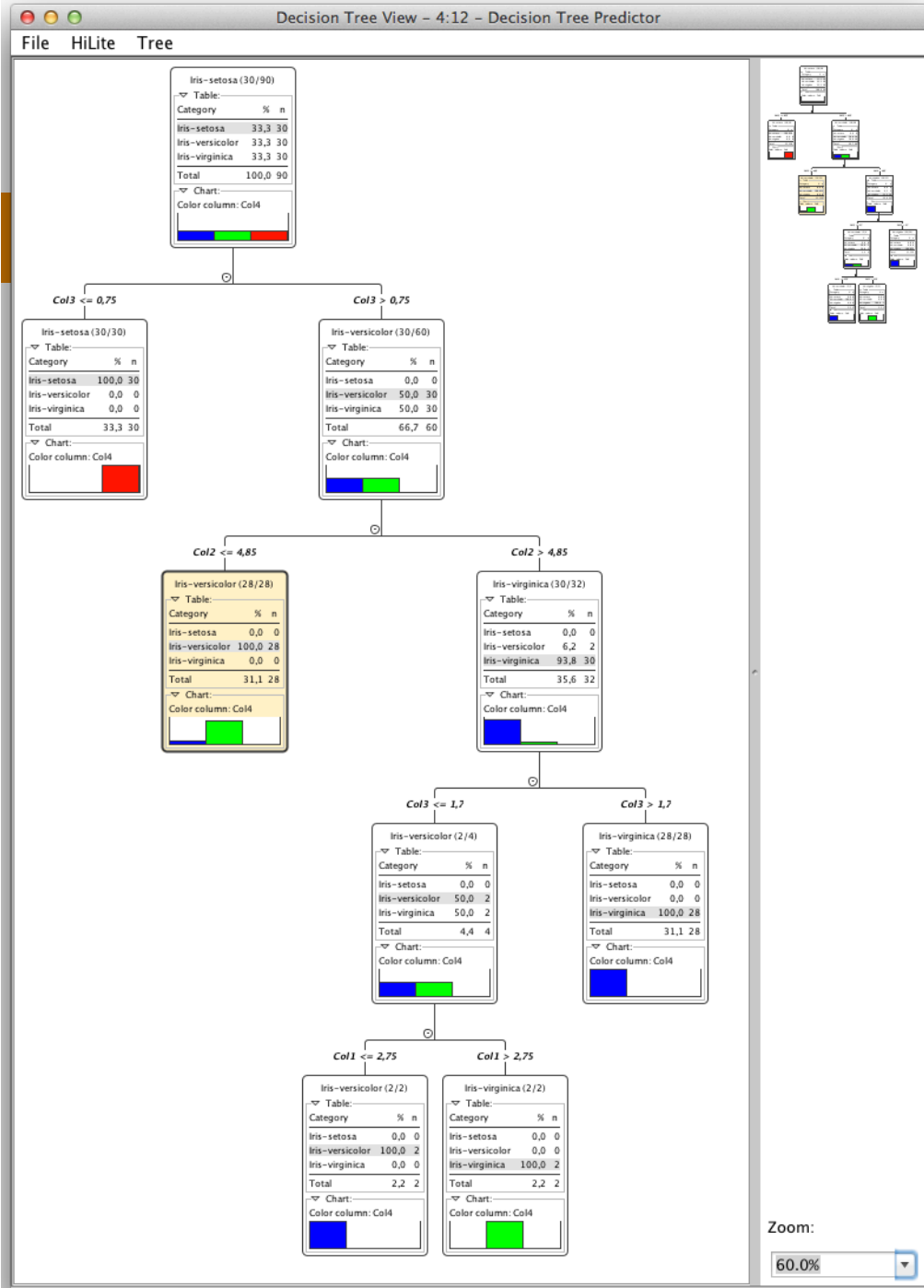
# Ejemplo 1: Iris ► Análisis predictivo ► Naïve Bayes





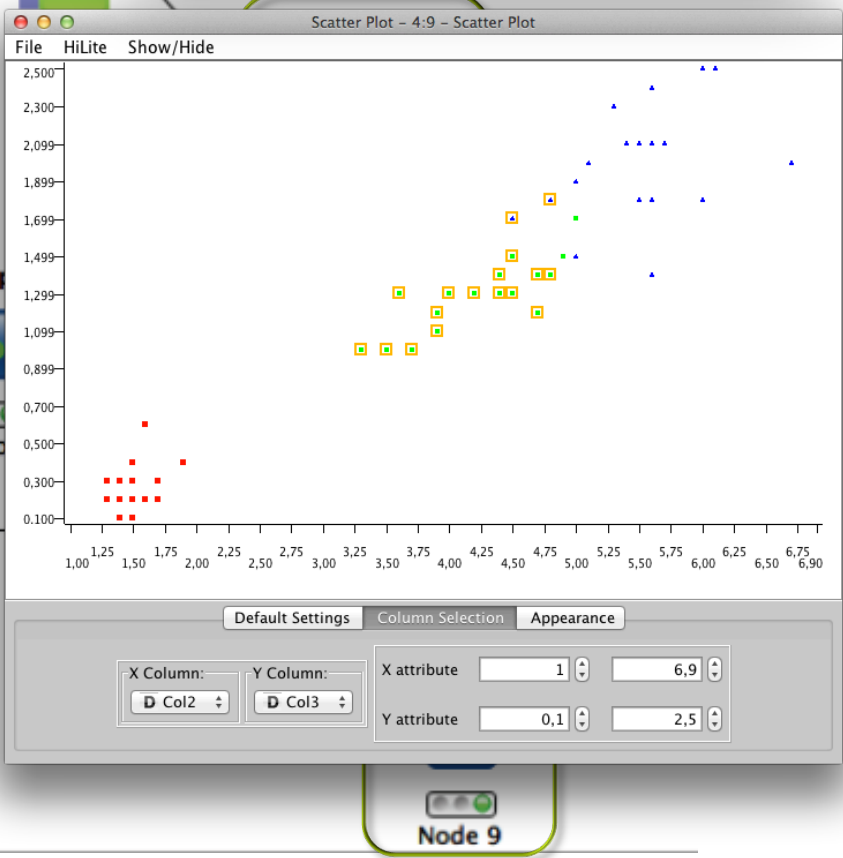
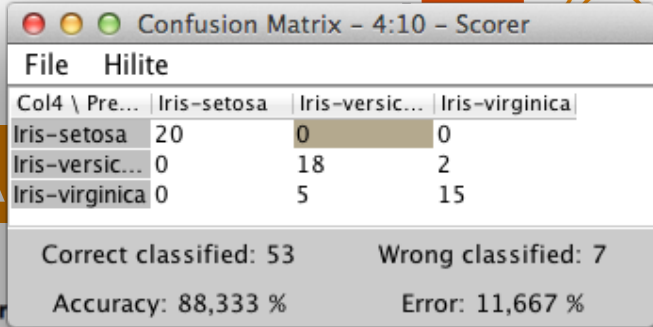
# Ejemplo 1: Iris ► Análisis predictivo ► Árboles de decisión





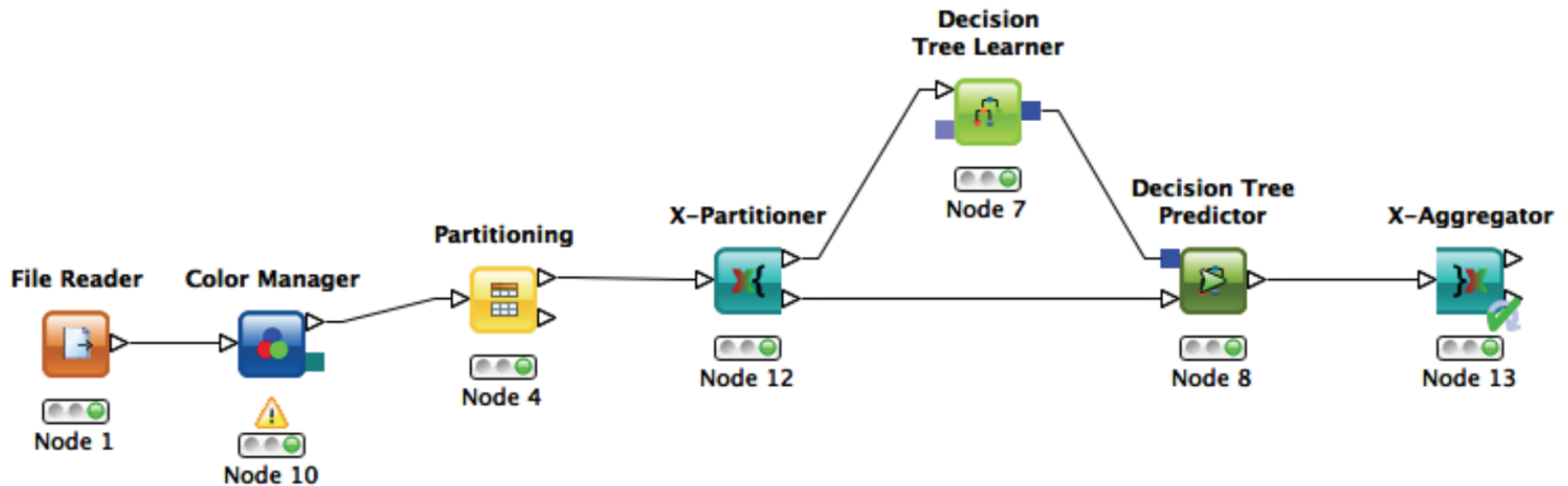
vo ▶ Á

Decision Tree Learner

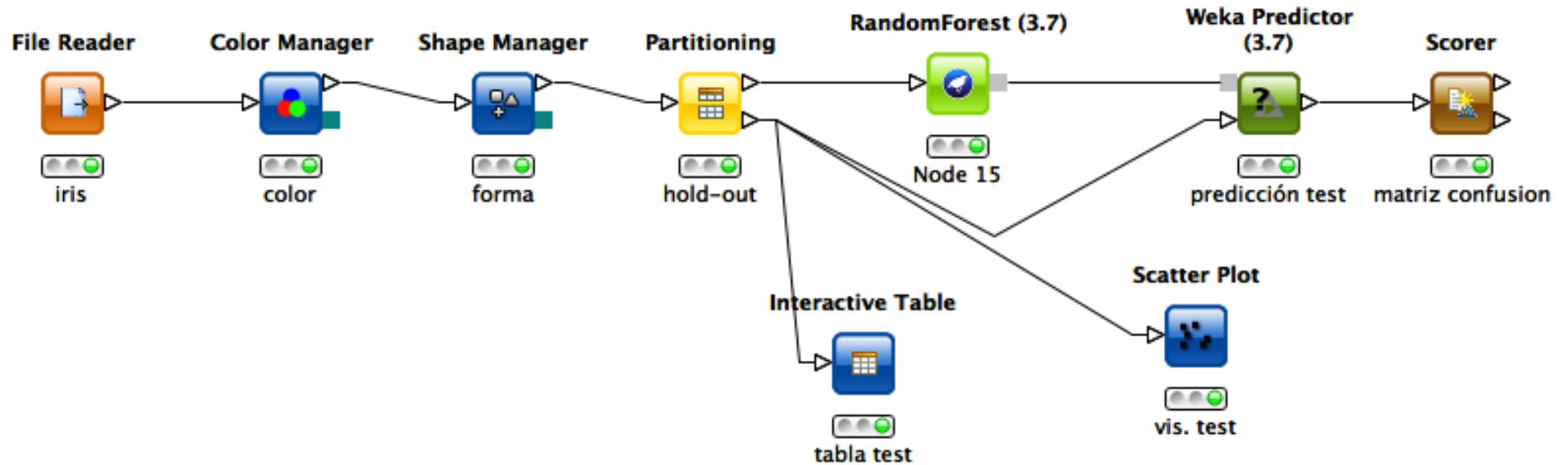


## Ejemplo 1: Iris ► Análisis predictivo

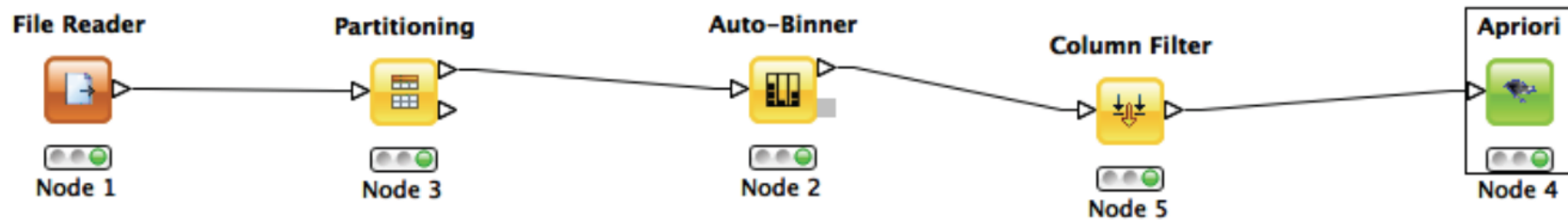
Los ciclos de validación cruzada se verán en los casos de estudio



# Ejemplo 1: Iris ► Análisis predictivo ► Random Forest



# Ejemplo 1: Iris ▶ Análisis descriptivo ▶ A priori Weka



File Reader



Node 1

```

Weka Node View - 8:6 - Apriori
File
Col3Binned=Bin1 25
0 25
Col3Binned=Bin2 21
1 16
Col3Binned=Bin3 25
1 14
Col3Binned=Bin4 19
2 19

Size of set of large itemsets L(2): 7

Large Itemsets L(2):
Col0Binned=Bin1 Col2Binned=Bin1 18
0 18
Col0Binned=Bin1 Col3Binned=Bin1 19
0 19
Col1Binned=Bin4 Col2Binned=Bin1 14
0 14
Col1Binned=Bin4 Col3Binned=Bin1 14
0 14
Col2Binned=Bin1 Col3Binned=Bin1 23
0 23
Col2Binned=Bin2 Col3Binned=Bin2 17
1 14
Col2Binned=Bin4 Col3Binned=Bin4 14
2 14

Size of set of large itemsets L(3): 1

Large Itemsets L(3):
Col0Binned=Bin1 Col2Binned=Bin1 Col3Binned=Bin1 18
0 18

Best rules found:

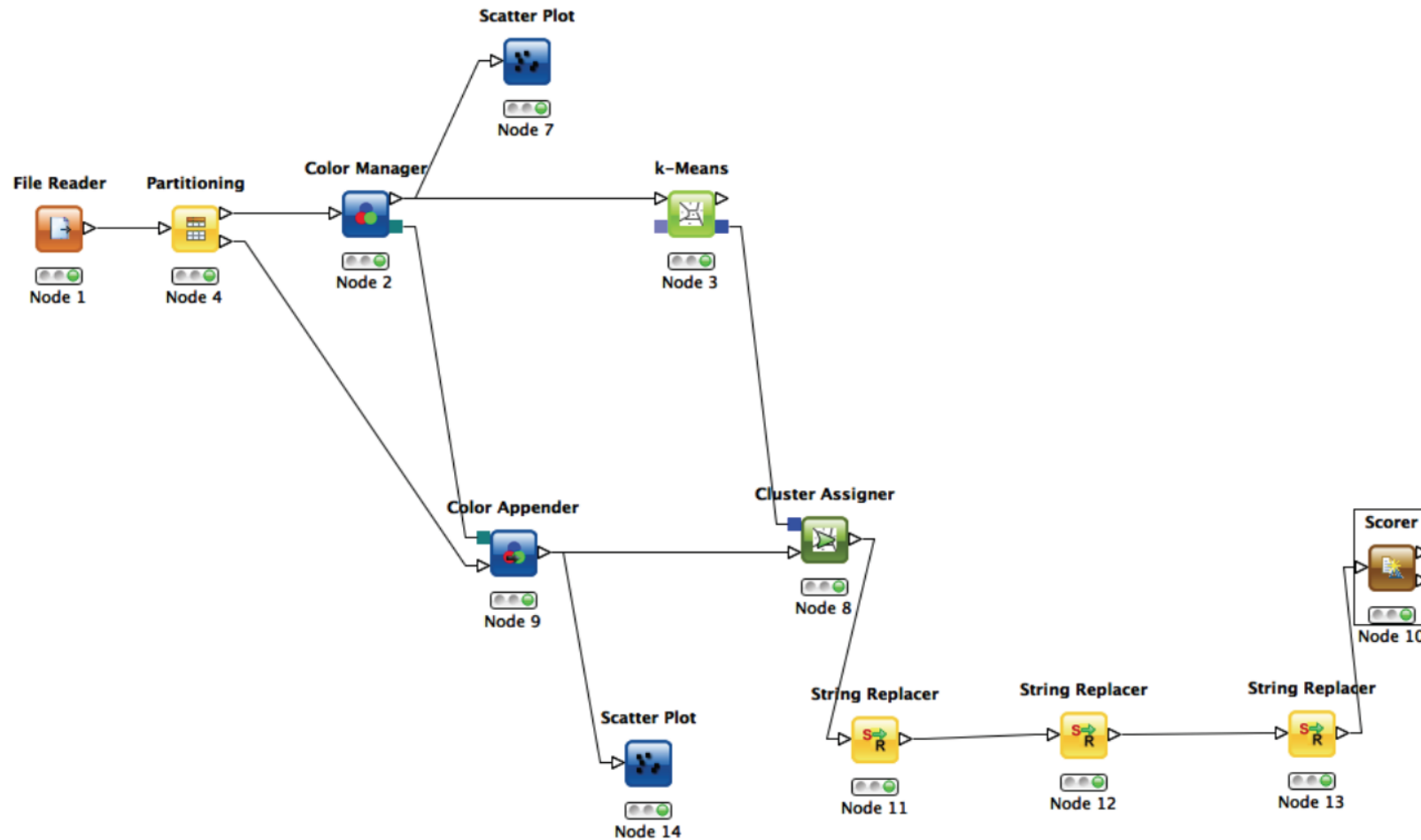
1. Col2Binned=Bin1 25 ==> Col4=Iris-setosa 25  conf:(1)
2. Col3Binned=Bin1 25 ==> Col4=Iris-setosa 25  conf:(1)
3. Col2Binned=Bin1 Col3Binned=Bin1 23 ==> Col4=Iris-setosa 23  conf:(1)
4. Col3Binned=Bin4 19 ==> Col4=Iris-virginica 19  conf:(1)
5. Col0Binned=Bin1 Col3Binned=Bin1 19 ==> Col4=Iris-setosa 19  conf:(1)
6. Col2Binned=Bin4 18 ==> Col4=Iris-virginica 18  conf:(1)
7. Col0Binned=Bin1 Col2Binned=Bin1 18 ==> Col4=Iris-setosa 18  conf:(1)
8. Col0Binned=Bin1 Col2Binned=Bin1 Col3Binned=Bin1 18 ==> Col4=Iris-setosa 18  conf:(1)
9. Col1Binned=Bin4 Col2Binned=Bin1 14 ==> Col4=Iris-setosa 14  conf:(1)
10. Col1Binned=Bin4 Col3Binned=Bin1 14 ==> Col4=Iris-setosa 14  conf:(1)

```

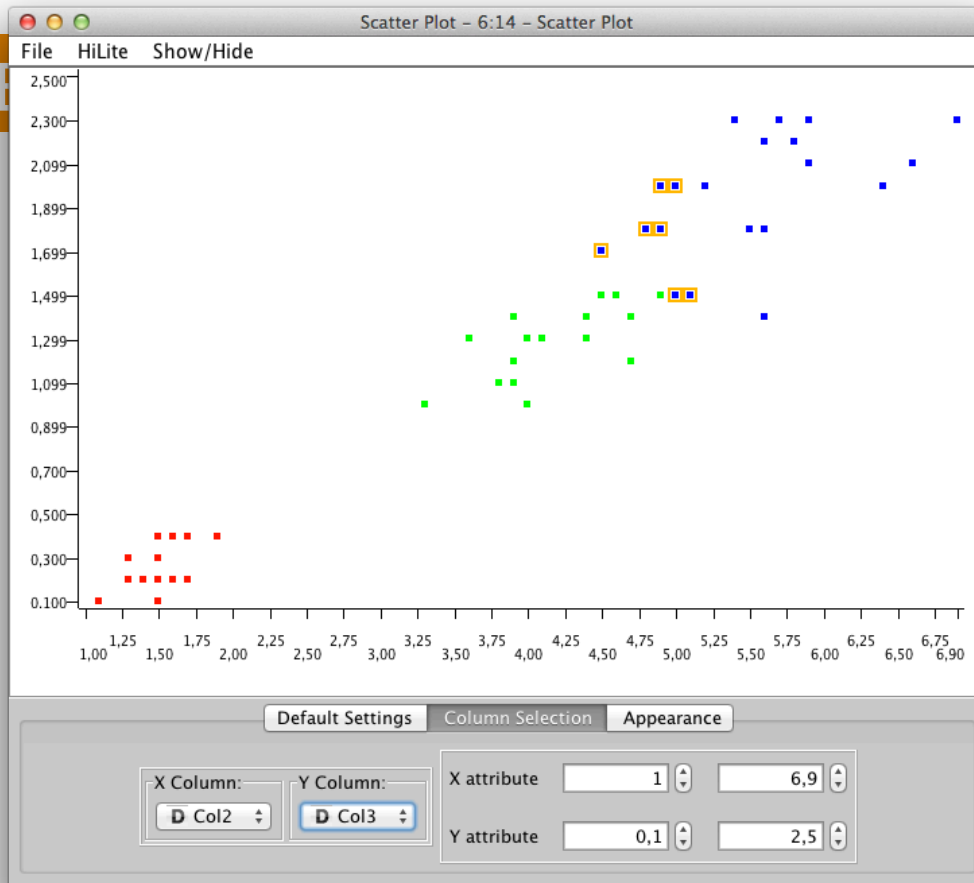
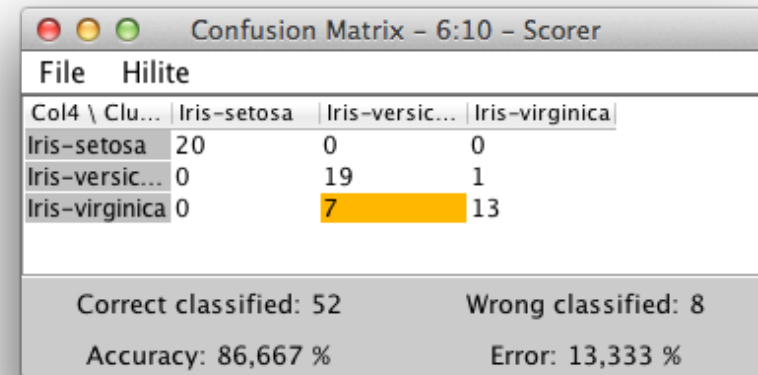


Node 4

# Ejemplo 1: Iris ▶ Análisis descriptivo ▶ k means

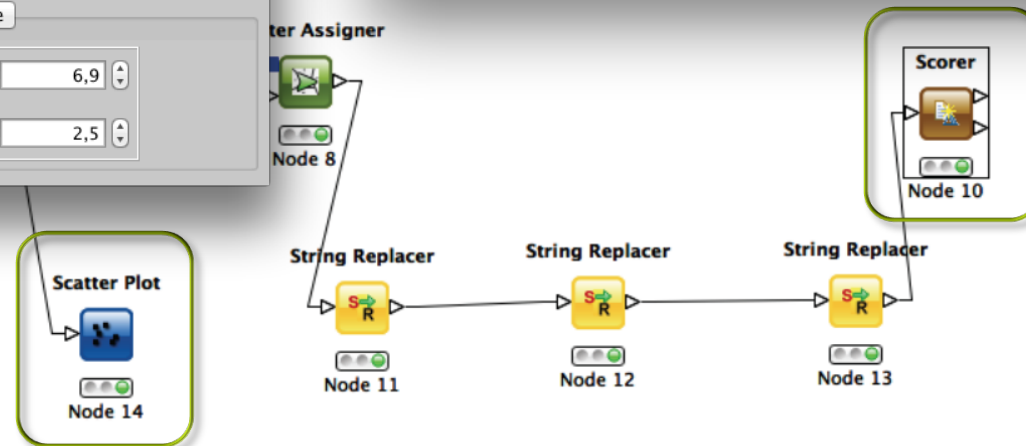


# Praktikum ▶ k means

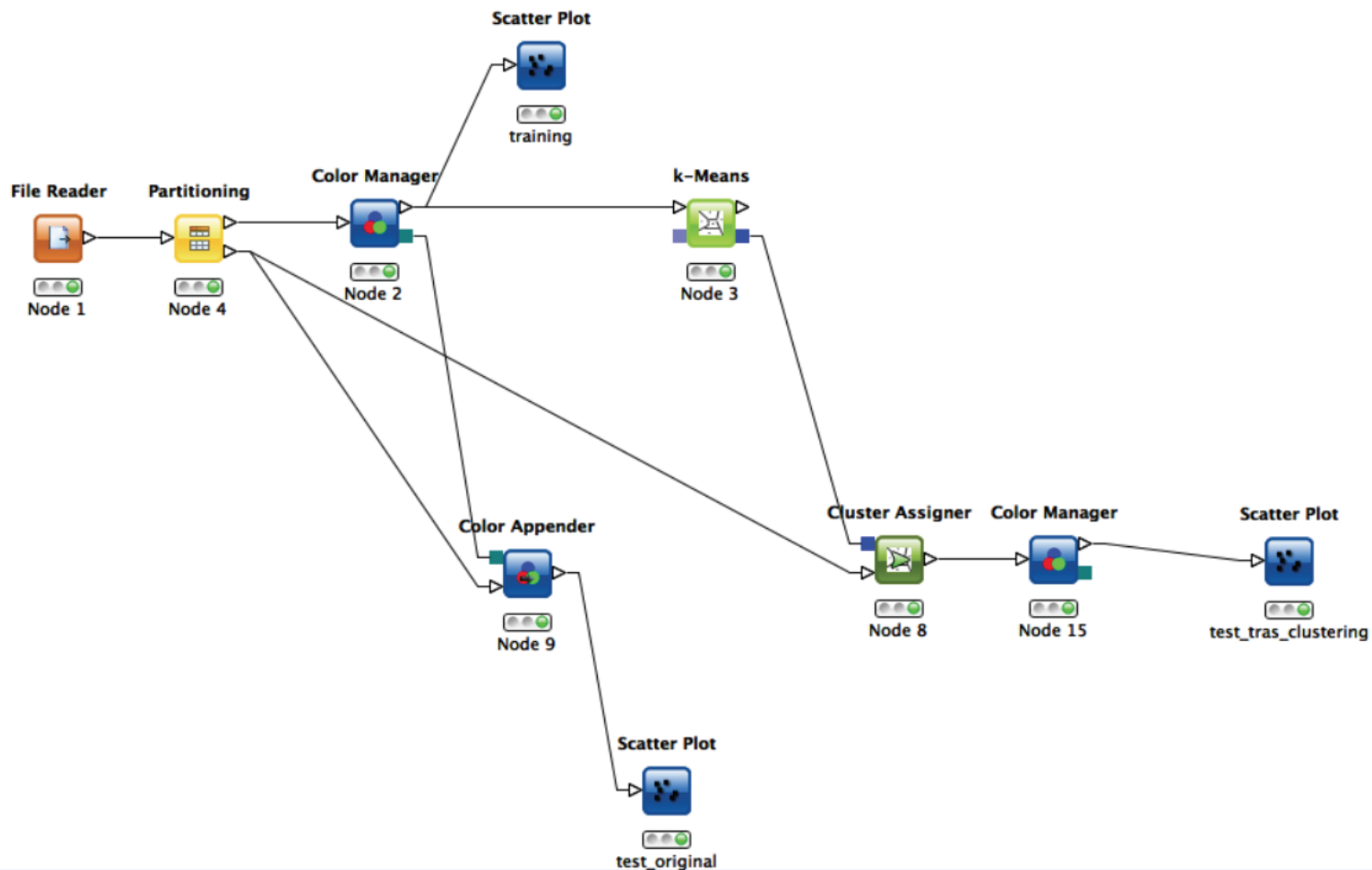
File	Hilite		
Col4 \ Clu...	Iris-setosa	Iris-versic...	Iris-virginica
Iris-setosa	20	0	0
Iris-versic...	0	19	1
Iris-virginica	0	7	13

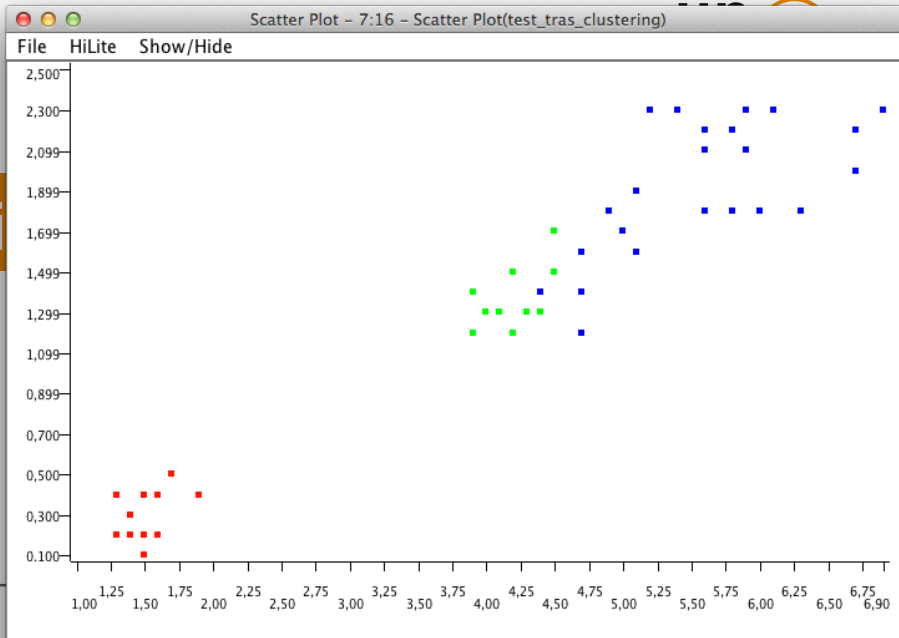
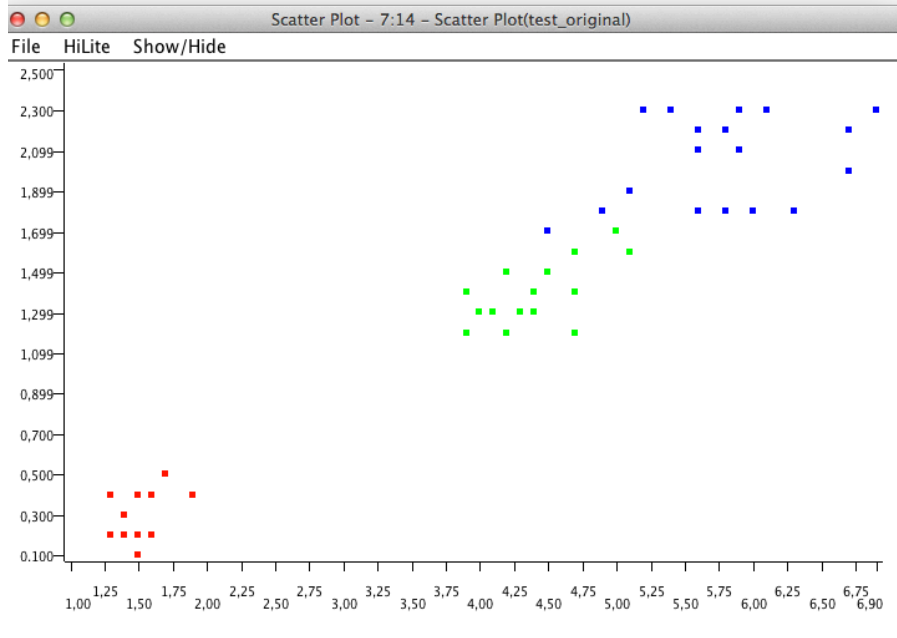
Correct classified: 52      Wrong classified: 8  
Accuracy: 86,667 %      Error: 13,333 %





# Ejemplo 1: Iris ▶ Análisis descriptivo ▶ k means





Default Settings Column Selection Appearance

X Column:  Y Column:

X attribute:

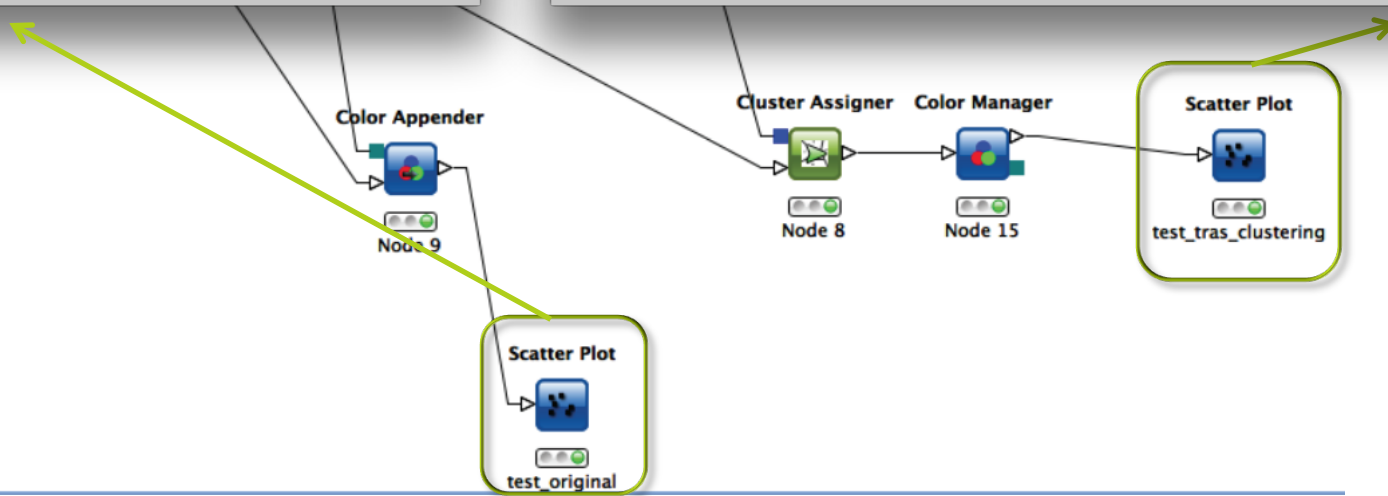
Y attribute:

Default Settings Column Selection Appearance

X Column:  Y Column:

X attribute:

Y attribute:



## Ejemplo 2: Pima

- Carga de datos
- Visualización
- Análisis predictivo
- Análisis descriptivo



<https://archive.ics.uci.edu/ml/datasets/Pima+Indians+Diabetes>

- 8 variables predictoras (continuas)
- 768 ejemplos
- 2 clases (500/268)
  
- ¿Valores perdidos?



## Ejemplo 3: Wine

- Carga de datos
- Visualización
- Análisis predictivo
- Análisis descriptivo



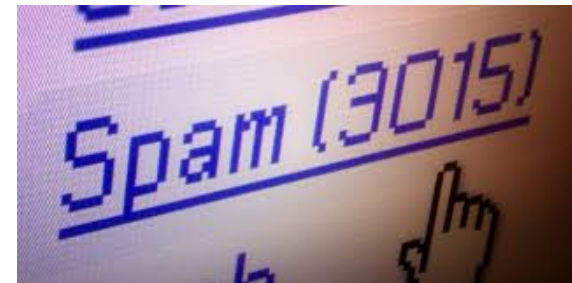
<https://archive.ics.uci.edu/ml/datasets/Wine>

- 12 variables predictoras (continuas)
- 178 ejemplos
- 3 clases (59/71/48)



## Ejemplo 4: Spam

- Carga de datos
- Visualización
- Preprocesamiento
- Análisis predictivo
- Análisis descriptivo



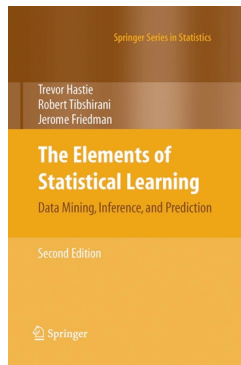
<https://archive.ics.uci.edu/ml/datasets/Spambase>

- 57 variables predictoras (55 continuas, 2 enteras)
- 4601 ejemplos
- 2 clases (2788/1813)

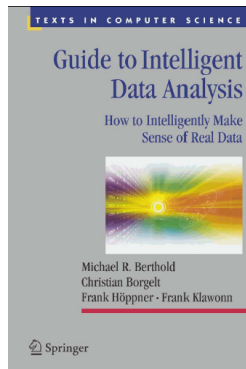




## Bibliografía



- The Elements of Statistical Learning: Data Mining, Inference, and Prediction, by Trevor Hastie, Robert Tibshirani and Jerome Friedman. 2009



- Guide to Intelligent Data Analysis. How to Intelligently Make Sense of Real Data, by M. R. Berthold, C. Borgelt, F. Höppner, F. Klawonn. Springer 2010



## Más información

- KNIME pages ([www.knime.org](http://www.knime.org))
- KNIME Tech ([tech.knime.org](http://tech.knime.org))
  - KNIME Quickstart Guide  
[https://tech.knime.org/files/KNIME\\_quickstart.pdf](https://tech.knime.org/files/KNIME_quickstart.pdf)
- KNIME TV Channel on YouTube
- 100 best KNIME Videos  
<http://meta-guide.com/videography/100-best-knime-videos>



## CURSOS DE VERANO 2014

**APROXIMACIÓN PRÁCTICA A LA CIENCIA DE DATOS  
Y BIG DATA: HERRAMIENTAS KNIME, R, HADOOP Y  
MAHOUT**

**Introducción a KNIME**

María José del Jesus