

2021/2022

EXERCICI 2. ¿ES VENENOSA ESTA SETA?

Hemos recogido información sobre 16 variedades de setas. Para cada una de ellas conocemos el «color», el «tamaño» y la «forma», así como si es una variedad «comestible» o «venenosa». El objetivo de la actividad es determinar qué características de las setas («color», «tamaño» o «forma») aportan más información sobre su carácter comestible o venenoso. Utiliza el paquete `FSelectorRcpp` para resolver el ejercicio.

Table 2.2: Datos de una muestra de 16 setas

Color	Tamaño	Forma	Comestible
Amarillo	Pequeña	Redonda	Sí
Amarillo	Pequeña	Redonda	No
Verde	Pequeña	Irregular	Sí
Verde	Grande	Irregular	No
Amarillo	Grande	Redonda	Sí
Amarillo	Pequeña	Redonda	Sí
Amarillo	Pequeña	Redonda	Sí
Amarillo	Pequeña	Redonda	Sí
Verde	Pequeña	Redonda	No
Amarillo	Grande	Redonda	No
Amarillo	Grande	Redonda	Sí
Amarillo	Grande	Redonda	No
Amarillo	Grande	Redonda	No
Amarillo	Grande	Redonda	No
Amarillo	Grande	Redonda	No
Amarillo	Pequeña	Irregular	Sí
Amarillo	Grande	Irregular	Sí

Es posible realizar esta misma actividad con un dataset más completo que contiene información sobre 8124 variedades de setas disponible en: <http://archive.ics.uci.edu/ml/datasets/Mushroom>. Cada seta está representada en una fila con datos en 23 columnas. La primera columna indica si es comestible (e, edible) o venenosa (p, poisonous). Las 22 variables restantes indican características como la forma, la textura, el color, el olor, etc. Todas las variables son categóricas discretas (adoptan un único valor). Por ejemplo, la segunda columna indica la forma del sombrero (cap shape) y puede adoptar seis valores (bell=b, conical=c, convex=x, flat=f, knobbed=k, sunken=s). El objetivo de la actividad

2021/2022

es determinar cuál de las 22 características es la que aporta más información sobre el carácter comestible o venenoso de una seta.

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Contains R code to load the 'Rcpp' and 'fSelectorRcpp' libraries, read the 'bolets2.txt' file, and display the first 10 rows of the 'bolets' dataset.
- Console:** Shows the execution of the code, displaying the first 10 rows of the dataset:

	Color	Tamaño	Forma	Comestible
1	Amarillo	Pequena	Redonda	Si
2	Amarillo	Pequena	Redonda	No
3	Verde	Pequena	Irregular	Si
4	Verde	Grande	Irregular	No
5	Amarillo	Grande	Redonda	Si
6	Amarillo	Pequena	Redonda	Si
7	Amarillo	Pequena	Redonda	Si
8	Amarillo	Pequena	Redonda	Si
9	Verde	Pequena	Redonda	No
10	Amarillo	Grande	Redonda	No

- Environment:** Shows the 'bolets' data frame with 16 observations and 4 variables.
- Files:** Lists files in the project directory, including '.Rhistory', 'bolets.Rproj', 'bolets.txt', and 'bolets2.txt'.

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Contains the same R code as the first screenshot, but with an additional line to calculate the information gain for the 'Comestible' variable.
- Console:** Shows the execution of the code, displaying the first 10 rows of the dataset and the result of the information gain calculation:

	Color	Tamaño	Forma	Comestible
1	Amarillo	Pequena	Redonda	Si
2	Amarillo	Pequena	Redonda	No
3	Verde	Pequena	Irregular	Si
4	Verde	Grande	Irregular	No
5	Amarillo	Grande	Redonda	Si
6	Amarillo	Pequena	Redonda	Si
7	Amarillo	Pequena	Redonda	Si
8	Amarillo	Pequena	Redonda	Si
9	Verde	Pequena	Redonda	No
10	Amarillo	Grande	Redonda	No

attributes importance

	Color	Tamaño	Forma
1	0.02461657		
2		0.07336502	
3			0.02487004

- Environment:** Shows the 'bolets' data frame with 16 observations and 4 variables.
- Files:** Lists files in the project directory, including '.Rhistory', 'bolets.Rproj', 'bolets.txt', and 'bolets2.txt'.

2021/2022

La variable tamaño/medida de las setas, aporta más información sobre el carácter comestible o venenoso de esta, con el valor más alto del 0.07336502, seguido de otras variables que no aportan tanta información como el color y la forma que tienen valores más pequeños (0.02).