



## data.table

Aprende R para Ciencia de Datos en [www.datademia.es](http://www.datademia.es)



### data.table

data.table es un paquete de R que proporciona una versión de alto rendimiento del marco de datos de base R con sintaxis y mejoras de funciones para facilitar su uso, comodidad y velocidad de programación.

Carga el paquete:

```
> library(data.table)
```

Crea un data.table DT

```
DT <- data.table(V1=c(1L, 2L),
                 V2=LETTERS[1:3],
                 V3=round(rnorm(4), 4),
                 V4=1:12)
```



Forma general: DT[i, j, by]

"En DT, selecciona filas usando **i**, luego calcula **j** y agrupa por **by**"

### Selecciona filas usando i

```
> DT[3:5, ]           Selecciona de la 3ra a la 5ta fila
> DT[3:5]           Selecciona de la 3ra a la 5ta fila
> DT[V2=="A"]       Selecciona filas que tienen el valor A en la columna V2
> DT[V2 %in% c("A", "C")] Selecciona filas que tengan el valor A o C en la columna V2
```

### Manipula columnas usando j

```
> DT[, V2]           Devuelve V2 como un vector
[1] "A" "B" "C" "A" "B" "C" ...
> DT[, .(V2, V3)]    Devuelve V2 y V3 como un data.table
> DT[, sum(V1)]       Devuelve la suma de todos los elementos de V1 en un vector
[1] 18
> DT[, .(sum(V1), sd(V3))] Devuelve la suma de todos los elementos de V1 y la desviación estándar de V3 en una tabla de datos.
V1 V2
1: 18 0.4546055

>DT[, .(Aggregate=sum(V1), Sd.V3=sd(V3))] Igual que el anterior, con nuevos nombres.
Aggregate Sd.V3
1: 18 0.4546055

> DT[, .(V1, Sd.V3=sd(V3))] Selecciona la columna V2 y calcula la desviación estándar de V3, que devuelve un valor único y se recicla

> DT[, .(print(V2), plot(V3), NULL)] Imprime la columna V2 y traza V3
```

### Usando j por grupo

```
> DT[, .(V4.Sum=sum(V4)), by=V1] Calcula la suma de V4 para cada grupo en V1
V1 V4.Sum
1: 1 36
2: 2 42

>DT[, .(V4.Sum=sum(V4)), by=. (V1, V2)] Calcula la suma de V4 para cada grupo en V1 y V2
>DT[, .(V4.Sum=sum(V4)), by=sign(V1-1)] Calcula la suma de V4 para cada grupo en el signo (V1-1)
sign V4.Sum
1: 0 36
2: 1 42

> DT[, .(V4.Sum=sum(V4)), by=. (V1.01=sign(V1-1))] Igual que el anterior, con un nuevo nombre para la variable por la que está agrupando
Calcula la suma de V4 para cada grupo en V1 después de seleccionar las primeras 5 filas
> DT[1:5, .(V4.Sum=sum(V4)), by=V1] Cuenta el número de filas para cada grupo en V1

> DT[, .N, by=V1]
```

### Añadiendo/Actualizando Columnas por referencia en j usando :=

```
> DT[, V1 := round(exp(V1), 2)] V1 se actualiza por lo que está después: =
> DT Devuelve el resultado llamando a DT
V1 V2 V3 V4
1: 2.72 A -0.1107 1
2: 7.39 B -0.1427 2
3: 2.72 C -1.8893 3
4: 7.39 A -0.3571 4
...
> DT[, c("V1", "V2") := list(round(exp(V1), 2), LETTERS[4:6])] Las columnas V1 y V2 se actualizan con lo que está después de :=
> DT[, ' := (V1=round(exp(V1), 2), V2=LETTERS[4:6]) ] ] Alternativa a la anterior. Con [, ], imprime el resultado en la pantalla
V1 V2 V3 V4
1: 15.18 D -0.1107 1
2: 1619.71 E -0.1427 2
3: 15.18 F -1.8893 3
4: 1619.71 D -0.3571 4

> DT[, V1 := NULL] Elimina V1
> DT[, c("V1", "V2") := NULL] Elimina las columnas V1 y V2
> Co1s.chosen=c("A", "B") Elimina la columna con el nombre Co1s.chosen
> DT[, Co1s.Chosen := NULL] Eliminar las columnas especificadas en la variable Co1s.chosen

> DT[, (Co1s.Chosen) := NULL]
```

### Indices y claves

```
> setkey(DT, V2)           Establece una clave en V2
> DT["A"]                 Devuelve todas las filas donde la columna clave (establecida a V2) tiene el valor A
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
> DT[c("A", "C")]        Devuelve todas las filas donde la columna clave (V2) tiene el valor A o C
Devuelve la primera fila de todas las filas que coinciden con el valor A en la columna clave V2
Devuelve la última fila de todas las filas que coinciden con el valor A en la columna clave V2
Devuelve todas las filas donde la columna clave V2 tiene el valor A o D

> DT["A", mult="first"]
> DT["A", mult="last"]

> DT[c("A", "D")]
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
3: NA D NA NA
> DT[c("A", "D"), nomatch=0] Devuelve todas las filas donde la columna clave V2 tiene el valor A o D
V1 V2 V3 V4
1: 1 A -0.2392 1
2: 2 A -1.6148 4
> DT[c("A", "C"), sum(V4)] Devuelve la suma total de V4, para las filas de la columna clave V2 que tengan valores A o C
V2 V1
1: A 22
2: C 30

> setkey(DT, V1, V2)      Ordenar por V1 y luego por V2 dentro de cada grupo de V1 (invisible)
Selecione filas que tengan el valor 2 para la primera clave (V1) y el valor C para la segunda clave (V2)

> DT[.(2, "C")]
V1 V2 V3 V4
1: 2 C 0.3262 6
2: 2 C -1.6148 12
> DT[.(2, c("A", "C"))] Selecciona filas que tengan el valor 2 para la primera clave (V1) y dentro de esas filas el valor A o C para la segunda clave (V2)
V1 V2 V3 V4
1: 2 A -1.6148 4
2: 2 A 0.3262 10
3: 2 C 0.3262 6
```

### Familia set

```
> setnames(DT, "V2", "Rating") Cambia el nombre V2 a Rating (invisible)
> setnames(DT, c("V2", "V3"), c("V2.rating", "V3.DC")) Cambia 2 nombres de columnas (invisible)
> setcolorder(DT, c("V2", "V1", "V4", "V3")) Cambia el orden de las columnas usando el vector escogido (invisible)
```

### Operaciones avanzadas de data table

```
> DT[.N-1]             Devuelve la penúltima fila del DT
> DT[.N]               Devuelve el número de filas
> DT[, .(V2, V3)]      Devuelve V2 y V3 como data.table
> DT[, list(V2, V3)]   Devuelve V2 y V3 como data.table
> DT[, mean(V3), by=. (V1, V2)] Devuelve el resultado de j, agrupado por todas las combinaciones posibles de los grupos especificados en by
V1 V2 V3
1: 1 A 0.4053
2: 1 B 0.4053
3: 1 C 0.4053
4: 2 A -0.6443
5: 2 B -0.6443
6: 2 C -0.6443
```

### .SD & .SDcols

```
> DT[, print(.SD), by=V2] Mira lo que contiene .SD
> DT[, .SD[c(1, .N)], by=V2] Selecciona la primera y la última fila agrupadas por V2
> DT[, lapply(.SD, sum), by=V2] Calcula la suma de columnas en .SD agrupadas por V2
> DT[, lapply(.SD, sum), by=V2, .SDcols=c("V3", "V4")] Calcula la suma de V3 y V4 en .SD agrupados por V2
V2 V3 V4
1: A -0.478 22
2: B -0.478 26

Calcula la suma de V3 y V4 en .SD agrupados por V2
> DT[, lapply(.SD, sum), by=V2, .SDcols=paste0("V", 3:4)]
```

### Encadenamiento

```
> DT <- DT[, .(V4.Sum=sum(V4)), by =V1] Calcula la suma de V4, agrupado por V1
V1 V4.Sum
1: 1 36
2: 2 42

> DT[V4.Sum>40] Selecciona el grupo del cual la suma es > 40
> DT[, .(V4.Sum=sum(V4)), by=V1][V4.Sum>40] Selecciona el grupo del cual la suma es > 40 (encadenamiento)
V1 V4.Sum
1: 2 42
> DT[, .(V4.Sum=sum(V4)), by=V1][order(-V1)] Calcula la suma de V4, agrupado por V1, ordenado por V1
V1 V4.Sum
1: 2 42
2: 1 36
```

