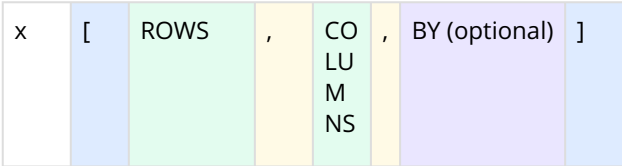# Data manipulation and filtering with data.table

By using data.table to structure our cytometry data, we can manipulate and filter the datasets extremely quickly. Subsetting rows or columns of a data.table can be achieved using the following structure:

| x | [ | ROWS | , | COLUMNS | , | BY (optional) | ] |
|---|---|------|---|---------|---|---------------|---|

Filtering approaches:

- Rows can be selected by the row number or values in a specified column
- Columns can be specified as the column numbers or the column names
- BY is used when grouping data together, but this is not demonstrated on this page

Here is a quick guide to some common operations, but a more comprehensive guide can be found on the data.table home page.

| Task | Command (with minimum variable inputs) |
|------|----------------------------------------|
| Access a 'demo' dataset stored within Spectre, and assign it to 'cell.dat'. | ```r
cell.dat <- Spectre::demo.clustered
``` |
| Subsetting data using **column** numbers | ```r
## Select the first column of the data.table
sub <- cell.dat[ ,1]
sub

## Select columns 1, 3, and 5
sub <- cell.dat[ ,c(1,3,5)]
sub

## Select columns 1 to 10, and 15
sub <- cell.dat[ ,c(1:10,15)]
sub
``` |
| Subsetting data using **column** names | ```r
## View a list of column names
as.matrix(names(cell.dat))
```

```
         [,1]
 [1,] "FileName"
 [2,] "NK11"
 [3,] "CD3"
 [4,] "CD45"
 [5,] "Ly6G"
 [6,] "CD11b"
 [7,] "B220"
 [8,] "CD8a"
 [9,] "Ly6C"
[10,] "CD4"
[11,] "NK11_asinh"
[12,] "CD3_asinh"
[13,] "CD45_asinh"
``` |

| Task | Command (with minimum variable inputs) |
|------|----------------------------------------|
| | ```[14,] "Ly6G_asinh"```<br>```[15,] "CD11b_asinh"```<br>```[16,] "B220_asinh"```<br>```[17,] "CD8a_asinh"```<br>```[18,] "Ly6C_asinh"```<br>```[19,] "CD4_asinh"```<br>```[20,] "Sample"```<br>```[21,] "Group"```<br>```[22,] "Batch"```<br>```[23,] "FlowSOM_cluster"```<br>```[24,] "FlowSOM_metacluster"```<br>```[25,] "Population"```<br>```[26,] "UMAP_X"```<br>```[27,] "UMAP_Y"```<br><br>```## Select columns names 11 to 19```<br>```cols <- names(cell.dat)[c(11:19)]```<br>```cols```<br><br>```[1] "NK11_asinh"  "CD3_asinh"   "CD45_asinh"  "Ly6G_asinh"  "CD11b_asinh"```<br>```"B220_asinh"  "CD8a_asinh"  "Ly6C_asinh"  "CD4_asinh"```<br><br>To select columns based on **column name**, either '..' needs to go before the vector of column names, or ', with = FALSE' needs to be added to the end of the data.table filtering arguments.<br><br>```## OPTION 1 - Select columns using '..'```<br>```sub <- cell.dat[ ,..cols]```<br>```sub```<br><br>```## OPTION 2 - Select columns using 'with = FALSE'```<br>```sub <- cell.dat[ ,cols, with = FALSE]```<br>```sub``` |
| Subsetting data using **row** numbers | ```## Select the first row of the data.table```<br>```sub <- cell.dat[1, ]```<br>```sub```<br><br>```## Select rows 1, 3, and 5```<br>```sub <- cell.dat[c(1,3,5), ]```<br>```sub```<br><br>```## Select rows 1 to 10, and 15```<br>```sub <- cell.dat[c(1:10,5), ]```<br>```sub``` |
| Subsetting data using **row** values in a selected column | Subsetting rows using data.table can be performed with the following structure:<br><br>| cell.dat | [ | CONDITIONS | , | | ] |<br>|---|---|---|---|---|---|<br><br>The **conditions** is essentially a filtering operation to determine which rows have a value in a specific column that is equal to, higher, or lower than a specified value. This can be performed using something like this: |

| Task | Command (with minimum variable inputs) |
|------|----------------------------------------|
| | ```## Creates a TRUE/FALSE results for which rows contain "Ly6C_asinh" values above 2
cell.dat[["Ly6C_asinh"]] > 2```<br><br>We can use this conditional TRUE/FALSE results to select which rows to include from our data.table.<br><br>```## Select rows (cells) where 'Ly6C_asinh' is above 2
sub <- cell.dat[cell.dat[["Ly6C_asinh"]] > 2, ]
sub

## Select rows (cells) where 'CD45_asinh' is below 3
sub <- cell.dat[cell.dat[["CD45_asinh"]] < 3, ]
sub

## Select rows (cells) where 'FlowSOM_metacluster' is above 2
sub <- cell.dat[cell.dat[["FlowSOM_metacluster"]] == 5,]
sub

## Select rows (cells) where 'Population' is 'Infil Macrophages'
sub <- cell.dat[cell.dat[["Population"]] == 'Infil Macrophages',]
sub``` |
| Subsetting data using multiple **row** values in multiple columns | Multiple 'or' arguments an be added by using '\|', for example: **cell.dat[A \| B \| C,].** This is an 'OR' operation, so all cells that are Ly6C_asinh > 2, in addition to all cells that are 'Infil Macrophages' will be included, rather than only cells satisfying both conditions (which would be an 'AND' operation).<br><br>```## Select rows (cells) where
    # 'Ly6C_asinh' is above 2 OR
    # 'Population' is 'Infil Macrophages'

sub <- cell.dat[cell.dat[["Ly6C_asinh"]] > 2 |
                cell.dat[["Population"]] == 'Infil Macrophages'
                ,]
sub``` |

# Extracting text before or after a symbol

Let's say we have a character vector, with a symbol dividing it into two:

```
x <- c("everything before|everything after")
```

We can use the following to extract everything **before** '|'.

```
before <- sub("\\|.*", "", x)
before
```

```
"everything before"
```

And we can use the following to extract everything **after** '|'.

```r
after <- sub(".*\\|", "", x)
after
```

```
"everything after"
```