

# DATA ANALYSIS THE DATA.TABLE WAY

**General form:** DT[i, j, by] → “Take DT, subset rows using i, then calculate j grouped by by”

| CREATE A DATA TABLE                 |   |       |    |         |       |
|-------------------------------------|---|-------|----|---------|-------|
| Create a data.table and call it DT. | <pre>library(data.table) set.seed(45L) DT &lt;- data.table(V1=c(1L,2L),                   V2=LETTERS[1:3],                   V3=round(rnorm(4),4),                   V4=1:12)</pre> | > DT  | V1 | V2      | V3 V4 |
|                                     |   | 1: 1  | A  | -1.1727 | 1     |
|                                     |   | 2: 2  | B  | -0.3825 | 2     |
|                                     |   | 3: 1  | C  | -1.0604 | 3     |
|                                     |   | 4: 2  | A  | 0.6651  | 4     |
|                                     |   | 5: 1  | B  | -1.1727 | 5     |
|                                     |   | 6: 2  | C  | -0.3825 | 6     |
|                                     |   | 7: 1  | A  | -1.0604 | 7     |
|                                     |   | 8: 2  | B  | 0.6651  | 8     |
|                                     |   | 9: 1  | C  | -1.1727 | 9     |
|                                     |   | 10: 2 | A  | -0.3825 | 10    |
|                                     |   | 11: 1 | B  | -1.0604 | 11    |
|                                     |   | 12: 2 | C  | 0.6651  | 12    |

## SUBSETTING ROWS USING `i`

| What?  | Example                                      | Notes  | Output   |
|--|--|--|--|
| Subsetting rows by numbers.  | DT[3:5, ] #or DT[3:5]                        | Selects third to fifth row.  | V1 V2 V3 V4<br>1: 1 C -1.0604 3<br>2: 2 A 0.6651 4<br>3: 1 B -1.1727 5   |
| Use column names to select rows in i based on a condition using fast automatic indexing. Or for selecting on multiple values:<br><code>DT[column %in% c("value1", "value2")]</code> , which selects all rows that have value1 or value2 in column. | DT[V2 == "A"]<br><br>DT[V2 %in% c("A", "C")] | Selects all rows that have value A in column V2.<br><br>Select all rows that have the value A or C in column V2. | V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 A 0.6651 4<br>3: 1 A -1.0604 7<br>4: 2 A -0.3825 10<br><br>V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 1 C -1.0604 3<br>...<br>7: 2 A -0.3825 10<br>8: 2 C 0.6651 12 |

## MANIPULATING ON COLUMNS IN `j`

| What?  | Example                                      | Notes   | Output  |
|--|--|---|---|
| Select 1 column in j.  | DT[, V2]                                     | Column V2 is returned as a vector.  | [1] "A" "B" "C" "A"<br>"B" "C" ...  |
| Select several columns in j.   | DT[, .(V2, V3)]                              | Columns V2 and V3 are returned as a data.table.   | V2 V3<br>1: A -1.1727<br>2: B -0.3825<br>3: C -1.0604<br>...<br>...                       |
| .() is an alias to list(). If .() is used, the returned value is a data.table. If .() is not used, the result is a vector. |  |   |   |
| Call functions in j.   | DT[, sum(V1)]                                | Returns the sum of all elements of column V1 in a vector.                                       | [1] 18  |
| Computing on several columns.  | DT[, .(sum(V1), sd(V3))]                     | Returns the sum of all elements of column V1 and the standard deviation of V3 in a data.table.  | V1 V2<br>1: 18 0.7634655  |
| Assigning column names to computed columns.  | DT[, .(Aggregate = sum(V1), Sd.V3 = sd(V3))] | The same as above, but with new names.  | Aggregate Sd.V3<br>1: 18 0.7634655  |
| Columns get recycled if different length.  | DT[, .(V1, Sd.V3 = sd(V3))]                  | Selects column V1, and compute std. dev. of V3, which returns a single value and gets recycled. | V1 Sd.V3<br>1: 1 0.7634655<br>2: 2 0.7634655<br>...<br>11: 1 0.7634655<br>12: 2 0.7634655 |
| Multiple expressions can be wrapped in curly braces.   | DT[, {print(V2)<br>plot(V3)<br>NULL}]        | Print column V2 and plot V3.  | [1] "A" "B" "C" "A"<br>"B" "C" ...<br>#And a plot   |

## DOING BY GROUP

| What?  | Example   | Notes   | Output  |
|--|---|---|---|
| Doing <b>j</b> by group.   | DT[, .(V4.Sum = sum(V4)), by=V1]                    | Calculates the sum of <b>V4</b> , for every group in <b>V1</b> .  | V1 V4.Sum<br>1: 1 36  |
| Doing <b>j</b> by several groups using <b>.()</b> .                    | DT[, .(V4.Sum = sum(V4)), by=.(V1, V2)]             | The same as above, but for every group in <b>V1</b> and <b>V2</b> .                                       | V1 V2 V4.Sum<br>1: 1 A 8<br>2: 2 B 10<br>3: 1 C 12<br>4: 2 A 14<br>5: 1 B 16<br>6: 2 C 18 |
| Call functions in <b>by</b> .  | DT[, .(V4.Sum = sum(V4)), by=sign(V1-1)]            | Calculates the sum of <b>V4</b> , for every group in <b>sign(V1-1)</b> .                                  | sign V4.Sum<br>1: 0 36<br>2: 1 42   |
| Assigning new column names in <b>by</b> .                              | DT[, .(V4.Sum = sum(V4)), by=.(V1.01 = sign(V1-1))] | Same as above, but with a new name for the variable we are grouping by.                                   | V1.01 V4.Sum<br>1: 0 36<br>2: 1 42  |
| Grouping only on a subset by specifying <b>i</b> .                     | DT[1:5, .(V4.Sum = sum(V4)), by=V1]                 | Calculates the sum of <b>V4</b> , for every group in <b>V1</b> , after subsetting on the first five rows. | V1 V4.Sum<br>1: 1 9<br>2: 2 6   |
| Using <b>.N</b> to get the total number of observations of each group. | DT[, .N, by=V1]                                     | Count the number of rows for every group in <b>V1</b> .   | V1 N<br>1: 1 6<br>2: 2 6  |

## ADDING/UPDATING COLUMNS BY REFERENCE IN USING `:=`

| What?  | Example  | Notes   | Output   |
|--|--|---|--|
| Adding/updating a column by reference using <code>:=</code> in one line.<br><b>Watch out:</b> extra assignment ( <code>DT &lt;- DT[...]</code> ) is redundant. | DT[, V1 := round(exp(V1), 2)]                                | Column <b>V1</b> is updated by what is after <code>:=</code> .<br><b>Watch out:</b> extra assignment ( <code>DT &lt;- DT[...]</code> ) is redundant.                      | Returns the result invisibly.<br>Column <b>V1</b> went from: [1] 1 2 1<br>2 ... to [1] 2.72 7.39 2.72<br>7.39 ...  |
| Adding/updating several columns by reference using <code>:=</code> .   | DT[, c("V1", "V2") := list(round(exp(V1), 2), LETTERS[4:6])] | Column <b>V1</b> and <b>V2</b> are updated by what is after <code>:=</code> .   | Returns the result invisibly.<br>Column <b>V1</b> changed as above.<br>Column <b>V2</b> went from: [1] "A"<br>"B" "C" "A" "B" "C" ... to: [1]<br>"D" "E" "F" "D" "E" "F" ... |
| Using functional <code>:=</code> .   | DT[, ' :=' (V1 = round(exp(V1), 2), V2 = LETTERS[4:6])][]    | Another way to write the same line as above this one, but easier to write comments side-by-side. Also, when <code>[]</code> is added the result is printed to the screen. | Same changes as line above this one, but the result is printed to the screen because of the <code>[]</code> at the end of the statement.                                     |
| Remove a column instantly using <code>:=</code> .  | DT[, V1 := NULL]   | Removes column <b>V1</b> .  | Returns the result invisibly.<br>Column <b>V1</b> became <code>NULL</code> .   |
| Remove several columns instantly using <code>:=</code> .   | DT[, c("V1", "V2") := NULL]                                  | Removes columns <b>V1</b> and <b>V2</b> .   | Returns the result invisibly. Column <b>V1</b> and <b>V2</b> became <code>NULL</code> .  |
| Wrap the name of a variable which contains column names in parenthesis to pass the contents of that variable to be deleted.                                    | cols.chosen = c("A", "B")<br>DT[, cols.chosen := NULL]       | <b>Watch out:</b> this deletes the column with column name <code>cols.chosen</code> .   | Returns the result invisibly.<br>Column with name <code>cols.chosen</code> became <code>NULL</code> .  |
|  | DT[, (cols.chosen) := NULL]                                  | Deletes the columns specified in the variable <code>cols.chosen</code> ( <b>V1</b> and <b>V2</b> ).   | Returns the result invisibly.<br>Columns <b>V1</b> and <b>V2</b> became <code>NULL</code> .  |

## INDEXING AND KEYS

| What?  | Example   | Notes  | Output   |
|--|---|--|--|
| Use <code>setkey()</code> to set a key on a DT. The data is sorted on the column we specified by reference.  | <code>setkey(DT, V2)</code>                     | A key is set on column <b>V2</b> .   | Returns results invisibly.   |
| Use keys like supercharged rownames to select rows.  | <code>DT["A"]</code>                            | Returns all the rows where the key column (set to column <b>V2</b> in the line above) has the value <b>A</b> .   | V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 A 0.6651 4<br>3: 1 A -1.0604 7<br>4: 2 A -0.3825 10                  |
|  | <code>DT[c("A", "C")]</code>                    | Returns all the rows where the key column ( <b>V2</b> ) has the value <b>A</b> or <b>C</b> .   | V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 A 0.6651 4<br>...<br>7: 1 C -1.1727 9<br>8: 2 C 0.6651 12            |
| The <code>mult</code> argument is used to control which row that <code>i</code> matches to is returned, default is all.  | <code>DT["A", mult = "first"]</code>            | Returns first row of all rows that match the value <b>A</b> in the key column ( <b>V2</b> ).   | V1 V2 V3 V4<br>1: 1 A -1.1727 1  |
|  | <code>DT["A", mult = "last"]</code>             | Returns last row of all rows that match the value <b>A</b> in the key column ( <b>V2</b> ).  | V1 V2 V3 V4<br>1: 2 A -0.3825 10   |
| The <code>nomatch</code> argument is used to control what happens when a value specified in <code>i</code> has no match in the rows of the DT. Default is <code>NA</code> , but can be changed to 0. 0 means no rows will be returned for that non-matched row of <code>i</code> . | <code>DT[c("A", "D")], nomatch = 0</code>       | Returns all the rows where the key column ( <b>V2</b> ) has the value <b>A</b> or <b>D</b> . <b>A</b> is found, <b>D</b> is not so <code>NA</code> is returned for <b>D</b> .    | V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 A 0.6651 4<br>3: 1 A -1.0604 7<br>4: 2 A -0.3825 10<br>5: NA D NA NA |
| <code>by=.EACH</code> allows to group by each subset of known groups in <code>i</code> . A key needs to be set to use <code>by=.EACH</code> .  | <code>DT[c("A", "C")], sum(V4)</code>           | Returns one total sum of column <b>V4</b> , for the rows of the key column ( <b>V2</b> ) that have values <b>A</b> or <b>C</b> .   | [1] 52   |
|  | <code>DT[c("A", "C")], sum(V4), by=.EACH</code> | Returns one sum of column <b>V4</b> for the rows of column <b>V2</b> that have value <b>A</b> , and another sum for the rows of column <b>V2</b> that have value <b>C</b> .      | V2 V1<br>1: A 22<br>2: C 30  |
| Any number of columns can be set as key using <code>setkey()</code> . This way rows can be selected on 2 keys which is an equijoin.  | <code>setkey(DT, V1, V2)</code>                 | Sorts by column <b>V1</b> and then by column <b>V2</b> within each group of column <b>V1</b> .   | Returns results invisibly.   |
|  | <code>DT[.(2, "C")]</code>                      | Selects the rows that have the value <b>2</b> for the first key (column <b>V1</b> ) and the value <b>C</b> for the second key (column <b>V2</b> ).                               | V1 V2 V3 V4<br>1: 2 C -0.3825 6<br>2: 2 C 0.6651 12  |
|  | <code>DT[.(2, c("A", "C"))]</code>              | Selects the rows that have the value <b>2</b> for the first key (column <b>V1</b> ) and within those rows the value <b>A</b> or <b>C</b> for the second key (column <b>V2</b> ). | V1 V2 V3 V4<br>1: 2 A 0.6651 4<br>2: 2 A -0.3825 10<br>3: 2 C -0.3825 6<br>4: 2 C 0.6651 12                  |

## CHAINING HELPS TACK EXPRESSIONS TOGETHER AND AVOID (UNNECESSARY) INTERMEDIATE ASSIGNMENTS

| What?  | Example  | Notes   | Output                          |
|--|--|---|---------------------------------|
| Do 2 (or more) sets of statements at once by chaining them in one statement. This corresponds to <i>having</i> in SQL. | <code>DT&lt;-DT[, .(V4.Sum = sum(V4)), by=V1]</code><br><code>DT[V4.Sum &gt; 40] #no chaining</code> | First calculates sum of <b>V4</b> , grouped by <b>V1</b> . Then selects that group of which the sum is > 40 without chaining. | V1 V4.Sum<br>1: 1 36<br>2: 2 42 |
|  | <code>DT[, .(V4.Sum = sum(V4)), by=V1][V4.Sum &gt; 40 ]</code>                                       | Same as above, but with chaining.   | V1 V4.Sum<br>1: 2 42            |
| Order the results by chaining.   | <code>DT[, .(V4.Sum = sum(V4)), by=V1][order(-V1)]</code>  | Calculates sum of <b>V4</b> , grouped by <b>V1</b> , and then orders the result on <b>V1</b> .                                | V1 V4.Sum<br>1: 2 42<br>2: 1 36 |

## ADVANCED DATA TABLE OPERATIONS

| What?   | Example  | Notes  | Output  |
|---|--|--|---|
| .N contains the number of rows or the last row.   | Usable in <code>i: DT[.N-1]</code><br>Usable in <code>j: DT[, .N]</code>   | Returns the penultimate row of the data.table.<br>Returns the number of rows.  | V1 V2 V3 V4<br>1: 1 B -1.0604 11<br>[1] 12  |
| .() is an alias to <code>list()</code> and means the same. The <code>.()</code> notation is not needed when there is only one item in <code>by</code> or <code>j</code> .   | Usable in <code>j: DT[, .(V2,V3)] #or DT[,list(V2,V3)]</code><br>Usable in <code>by: DT[, mean(V3), by=.(V1,V2)]</code>                          | Columns <b>V2</b> and <b>V3</b> are returned as a data.table.<br>Returns the result of <code>j</code> , grouped by all possible combinations of groups specified in <code>by</code> .      | V2 V3<br>1: A -1.1727<br>2: B -0.3825<br>3: C -1.0604<br>...<br>V1 V2 V1<br>1: 1 A -1.11655<br>2: 2 B 0.14130<br>3: 1 C -1.11655<br>4: 2 A 0.14130<br>5: 1 B -1.11655<br>6: 2 C 0.14130   |
| .SD is a data.table and holds all the values of all columns, except the one specified in <code>by</code> . It reduces programming time but keeps readability. <code>.SD</code> is only accessible in <code>j</code> . | <code>DT[, print(.SD), by=V2]</code><br><code>DT[, .SD[c(1,.N)], by=V2]</code><br><br><code>DT[, lapply(.SD, sum), by=V2]</code>                 | To look at what <code>.SD</code> contains.<br>Selects the first and last row grouped by column <b>V2</b> .<br>Calculates the sum of all columns in <code>.SD</code> grouped by <b>V2</b> . | #All of <code>.SD</code> (output too long to display here)<br>V2 V1 V3 V4<br>1: A 1 -1.1727 1<br>2: A 2 -0.3825 10<br>3: B 2 -0.3825 2<br>4: B 1 -1.0604 11<br>5: C 1 -1.0604 3<br>6: C 2 0.6651 12<br>V2 V1 V3 V4<br>1: A 6 -1.9505 22<br>2: B 6 -1.9505 26<br>3: C 6 -1.9505 30 |
| <code>.SDcols</code> is used together with <code>.SD</code> , to specify a subset of the columns of <code>.SD</code> to be used in <code>j</code> .<br><br><code>.SDcols</code> can be the result of a function call. | <code>DT[, lapply(.SD, sum), by=V2, .SDcols = c("V3","V4")]</code><br><br><code>DT[, lapply(.SD, sum), by=V2, .SDcols = paste0("V", 3:4)]</code> | Same as above, but only for columns <b>V3</b> and <b>V4</b> of <code>.SD</code> .<br>Same result as the line above.  | V2 V3 V4<br>1: A -1.9505 22<br>2: B -1.9505 26<br>3: C -1.9505 30<br>V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 B -0.3825 2<br>3: NA C -1.0604 3<br>4: NA A 0.6651 4<br>5: 1 NA -1.1727 5<br>6: 2 NA -0.3825 6<br>7: 1 A -1.0604 7<br>8: 2 B 0.6651 8                                |

## USING THE `set()`-FAMILY

| What?   | Example  | Notes  | Output  |
|---|--|--|---|
| <code>set()</code> is used to repeatedly update rows and columns by reference. <code>Set()</code> is a loopable low overhead version of <code>:=</code> .<br><b>Watch out:</b> It can not handle grouping operations. | Syntax of <code>set(): for (i in from:to) set(DT, row, column, new value)</code> .<br><br><code>rows = list(3:4,5:6)</code><br><code>cols = 1:2</code><br><code>for (i in seq_along(rows))</code><br><code>{ set(DT,</code><br><code>      i=rows[[i]],</code><br><code>      j = cols[i],</code><br><code>      value = NA ) }</code> | Sequence along the values of <code>rows</code> , and for the values of <code>cols</code> , set the values of those elements equal to <code>NA</code> .   | Returns the result invisibly.<br>> DT<br>V1 V2 V3 V4<br>1: 1 A -1.1727 1<br>2: 2 B -0.3825 2<br>3: NA C -1.0604 3<br>4: NA A 0.6651 4<br>5: 1 NA -1.1727 5<br>6: 2 NA -0.3825 6<br>7: 1 A -1.0604 7<br>8: 2 B 0.6651 8  |
| <code>setnames()</code> is used to create or update column names by reference.  | Syntax of <code>setnames()</code> :<br><code>setnames(DT, "old", "new") []</code><br><br><code>setnames(DT, "V2", "Rating")</code><br><br><code>setnames(DT, c("V2", "V3"), c("V2.rating", "V3.DataCamp"))</code>  | Changes ( <code>set</code> ) the name of column <b>old</b> to <b>new</b> . Also, when <code>[]</code> is added at the end of any <code>set()</code> function the result is printed to the screen.<br>Sets the name of column <b>V2</b> to <b>Rating</b> . Returns the result invisibly.<br>Changes two column names. Returns the result invisibly. | Changes (set) the name of column <b>old</b> to <b>new</b> . Also, when <code>[]</code> is added at the end of any <code>set()</code> function the result is printed to the screen.<br>Sets the name of column <b>V2</b> to <b>Rating</b> . Returns the result invisibly.<br>Changes two column names. Returns the result invisibly. |
| <code>setcolorder()</code> is used to reorder columns by reference.   | <code>setcolorder(DT, "neworder")</code><br><br><code>setcolorder(DT, c("V2", "V1", "V4", "V3"))</code>  | <b>neworder</b> is a character vector of the new column name ordering.<br>Changes the column ordering to the contents of the vector.   | Returns the result invisibly. The new column order is now <code>[1] "V2" "V1" "V4" "V3"</code>  |