

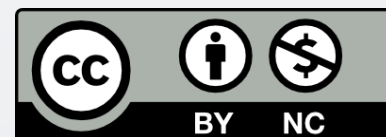


# THE #RDATATABLE PACKAGE

+ new developments in v1.9.7

Arun Srinivasan

SEP 3'16, BUDAPEST



@ARUN\_SRINIV

# WHO AM I?

- Bioinformatician / Comp. Biologist
- **data.table** user, co-developer since late 2013
- Previous: Data scientist @Open Analytics
- **Future**: Lead engineer @investment mgmt. firm

# MOST UNDERRATED PACKAGE



**Conor Nash**

@conornash



 **Follow**

Data.table is the most underrated R package. It has saved me \*days\* in waiting for analyses to complete.

August 2016

# MOST UNDERRATED PACKAGE



**Mehdi Nemlaghi**

@Mehdi\_Nemlaghi



**Follow**

@freakonometrics "setkey" function is so powerful, so innovative for #rstats. Imho, "Data.table" package is kind of underrated...

May 2015

- Homepage: <http://r-datatable.com>
- Since 2006 on CRAN, >30 releases so far
- >5500 unit tests, ~89% coverage (using covr)
- >260 packages import/depend/suggest data.table
- ~12.6 packages per month since Sep'15
- 8th most starred R package on Github (METACRAN)
- >4400 Q on StackOverflow. 3rd amongst R packages



# POWERFUL



**Alexander Flyax**

@aflyax



 **Follow**

somebody should just write a version of [#Rstat](#)'s data.table for [#python](#). end of story. nothing as powerful exists at the moment.

(With no intent on fuelling language wars)

# GREAT SADNESS



**Jim Savage**

@khakieconomist



 **Follow**

With great sadness I was forced to start using  
data.table today.



# DATA.TABLE DATA.TABLE DATA.TABLE



**Joey Reid**  
@JoeyPReid



 **Follow**

data.table  
data.table  
data.table  
data.table  
ggplot2  
rstan  
knitr

#7FavPackages

# TALK OVERVIEW

- **data.table's** philosophy
  - concise + straightforward code
  - fast + memory efficient
- New features and improvements in **v1.9.7**
  - fwrite, conditional joins, parallel sort & other optimisations

# TALK OVERVIEW

- **data.table's** philosophy
  - concise + straightforward code
  - fast + memory efficient
- New features and improvements in v1.9.7
  - fwrite, conditional joins, parallel sort & other optimisations

# DATA FRAMES

- are columnar data structures

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

2 column data.frame

# DATA FRAMES

- are columnar data structures
  - 2D — rows and columns

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

2 column data.frame

# DATA FRAMES

- are columnar data structures
  - 2D — rows and columns

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

2 column data.frame

# DATA FRAMES

- are columnar data structures
- 2D — rows and columns
- subset rows — `X[X$id != "a", ]`

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

# DATA FRAMES

- are columnar data structures
- 2D — rows and columns
- subset rows — `X[X$id != "a", ]`
- select columns — `X[, "val"]`

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6



# DATA FRAMES

- are columnar data structures
- 2D — rows and columns
- subset rows — `X[X$id != "a", ]`
- select columns — `X[, "val"]`
- subset rows & select columns —  
`X[X$id != "a", "val"]`

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

# DATA FRAMES

- are columnar data structures
- 2D — rows and columns
- subset rows — `X[X$id != "a", ]`
- select columns — `X[, "val"]`
- subset rows & select columns —  
`X[X$id != "a", "val"]`
- that's pretty much it...

X

	id	val
1	b	4
2	a	2
3	a	3
4	c	1
5	c	5
6	b	6

# 1. HOW TO COMPUTE ON COLUMNS?

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	1.5	5
4	2	apq	0.9	10
5	2	apq	0.3	13

For `code != "abd"`,  
get `sum(valA)`

1.9

# 1. HOW TO COMPUTE ON COLUMNS?

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	1.5	5
4	2	apq	0.9	10
5	2	apq	0.3	13

```
sum(DF[DF$code != "abd", "valA"])
```

1.9

# 2. GROUPED AGGREGATE

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	1.5	5
4	2	apq	0.9	10
5	2	apq	0.3	13

For `code != "abd"`,  
get `sum(valA)` and `sum(valB)`  
for each `id`

	id	valA	valB
1	1	0.7	18
2	2	1.2	23

# 2. GROUPED AGGREGATE

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	1.5	5
4	2	apq	0.9	10
5	2	apq	0.3	13

```
aggregate(cbind(valA, valB) ~ id,  
          DF[DF$code != "abd", ],  
          sum)
```

	id	valA	valB
1	1	0.7	18
2	2	1.2	23

# 3. SIMPLE UPDATE

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	1.5	5
4	2	apq	0.9	10
5	2	apq	0.3	13

For `code == "abd"`,  
update `valA`  
with `NA`

# 3. SIMPLE UPDATE

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	NA	5
4	2	apq	0.9	10
5	2	apq	0.3	13

For `code == "abd"`,  
update `valA`  
with `NA`



# 3. SIMPLE UPDATE

DF

	id	code	valA	valB
1	1	abc	0.1	11
2	1	abc	0.6	7
3	1	abd	NA	5
4	2	apq	0.9	10
5	2	apq	0.3	13

```
DF[DF$code == "abd", "valA"] <- NA
```

# CAN WE BE MORE CONSISTENT?

```
sum(DF[DF$code != "abd", "valA"])
```

How to get *sum* of both *valA* and *valB*?  
Or *sum* of *valA* and *valB* combined?

```
aggregate(cbind(valA, valB) ~ id,  
          DF[DF$code != "abd", ],  
          sum)
```

New function. Formula interface.  
Unwanted columns are subsetting.  
How to get *sum(valA)* and *mean(valB)*?

```
DF[DF$code == "abd", "valA"] <- NA
```

Entire expression is now to the left  
of the "<-" operator

# ENHANCED DATA FRAMES

- Three main enhancements:
  1. Allow **column names** to be seen **as variables** within **[...]**
  2. Since they're variables, we can **do computations** on them **directly**, i.e, within **[...]**
  3. Additional argument **by**

# DATA TABLES

- are columnar data structures as well

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6

2 column data.table

# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6

2 column data.table

# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6

2 column data.table

# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns
- subset rows — `X[id != "a", ]`

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6

# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns
- subset rows — `X[id != "a", ]`
- select columns — `X[, val]`

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6



# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns
  - subset rows — `X[id != "a", ]`
  - select columns — `X[, val]`
  - *compute on* columns — `X[, mean(val)]`

X

	id	val	
1:	b	4	mean 3.5
2:	a	2	
3:	a	3	
4:	c	1	
5:	c	5	
6:	b	6	

# DATA TABLES

- are columnar data structures as well
  - 2D — rows and columns
  - subset rows — `X[id != "a", ]`
  - select columns — `X[, val]`
  - *compute on* columns — `X[, mean(val)]`
  - subset rows **&** select / *compute on* columns
    - `X[id != "a", mean(val)]`

X

	id	val	
1:	b	4	
2:	a	2	
3:	a	3	
4:	c	1	
5:	c	5	
6:	b	6	

*mean*  
4.0

# DATA TABLES

- are columnar data structures as well
- 2D — rows and columns
- subset rows — `X[id != "a", ]`
- select columns — `X[, val]`
- *compute on* columns — `X[, mean(val)]`
- subset rows **&** select / *compute on* columns  
— `X[id != "a", mean(val)]`
- *virtual* 3rd dimension — **group by**

X

	id	val
1:	b	4
2:	a	2
3:	a	3
4:	c	1
5:	c	5
6:	b	6

# DATA TABLES

- think in terms of basic units — **rows**, **columns** and **groups**
- data.table syntax provides *placeholder* for each of them

General form: **DT**[**i**, **j**, **by**]

On which rows

What to do?

Grouped by  
what?

# EQUIVALENT DATA TABLE CODE

```
sum(DF[DF$code != "abd", "valA"])
```

```
DT[code != "abd", sum(valA)]
```

```
aggregate(cbind(valA, valB) ~ id,  
          DF[DF$code != "abd", ],  
          sum)
```

```
DT[code != "abd",  
   .(sum(valA), sum(valB)),  
   by = id]
```

```
DF[DF$code == "abd", "valA"] <- NA
```

```
DT[code == "abd", valA := NA]
```

# TWO TABLES

A

	id	code	valA	valB
1:	1	abc	0.1	11
2:	1	abc	0.6	7
3:	1	abd	1.5	5
4:	2	apq	0.9	10
5:	2	apq	0.3	13

B

	id	code	mul
1:	1	abd	2.0
2:	2	apq	0.5
3:	3	abc	1.7

Update **valA** with **valA\*mul** while matching on **id, code**

# TWO TABLES

A

	id	code	valA	valB
1:	1	abc	0.1	11
2:	1	abc	0.6	7
3:	1	abd	1.5	5
4:	2	apq	0.9	10
5:	2	apq	0.3	13

B

	id	code	mul
1:	1	abd	2.0
2:	2	apq	0.5
3:	3	abc	1.7

Update **valA** with **valA\*mul** while matching on **id, code**

```
A[B, on = .(id, code),  
  valA := valA * mul]
```

on which rows?  
what to do?

# TWO TABLES

A

	id	code	valA	valB
1:	1	abc	0.1	11
2:	1	abc	0.6	7
3:	1	abd	3.0	5
4:	2	apq	0.45	10
5:	2	apq	0.15	13

B

	id	code	mul
1:	1	abd	2.0
2:	2	apq	0.5
3:	3	abc	1.7

Update **valA** with **valA\*mul** while matching on **id, code**

```
A[B, on = .(id, code),  
  valA := valA * mul]
```

on which rows?  
what to do?



# TALK OVERVIEW

- data.table's philosophy
  - concise + straightforward code
  - fast + memory efficient
- **New features and improvements in v1.9.7**
  - **fwrite, conditional joins, parallel sort & other optimisations**

# FWRITE - PARALLEL FILE WRITER

		Laptop SSD		Server		
		4core/16gb		32core/256gb		
		10m rows		100m rows		
		=====		=====		
		Time	Size	RamDisk	HDD	Size
		Sec	GB	Time	Time	GB
<code>fwrite(DT, "fwrite.csv")</code>	csv	2	0.8	9	61	7.5
<code>write_feather(DT, "feather.bin")</code>	bin	5	1.0	27	75	9.1
<code>save(DT, file="save1.Rdata", compress=F)</code>	bin	11	1.2	90	137	12.0
<code>save(DT, file="save2.Rdata", compress=T)</code>	bin	70	0.4	647	679	2.8
<code>write.csv(DT, "write.csv.csv", **)</code>	csv	63	0.8	749	824	7.3
<code>readr::write_csv(DT, "write_csv.csv")</code>	csv	132	0.8	1997	1571	7.3

[\*\*] row.names=F, quote=F

**SOURCE:** <http://blog.h2o.ai/2016/04/fast-csv-writing-for-r/>

# FSORT - PARALLEL SORT

length	size in RAM	threads	base R	v1.9.7
500m	3.8GB	8	65s	3.9s
1b	7.6GB	32	140m	3.5s
10b	76GB	32	25m	48s

**SOURCE:** <https://www.r-project.org/dsc/2016/slides/ParallelSort.pdf>

# PARALLEL ROW SUBSETS

DT[sample(.N, .N/2)] 200e6 rows, 4 cols~4.6GB	
v1.9.6	20.0s
v1.9.7 (C, parallelised)	3.6s (16 threads)
	run time

# %BETWEEN%

x %between% c(2000, 20000) length(x) = 500e6, int, ~1.9GB		
v1.9.6	15.7s	7.2GB
v1.9.7 (C, parallelised)	1.1s (4 threads)	3.8GB
	run time	peak memory

# MEDIAN

1e6 rows, 61 columns, ~460MB  
10,000 unique groups

```
7 # v1.9.6, CRAN version
8 ans1 <- dt[, lapply(.SD, median), by=x]
9 #   user  system elapsed
10 # 19.610   0.229  19.917
11
12 # v1.9.7, devel
13 # median is internally optimised|
14 ans2 <- dt[, lapply(.SD, median), by=x]
15 #   user  system elapsed
16 #  1.195   0.007   1.207
17
18 # = 16.5x speedup.
```

# CONDITIONAL OPERATIONS

A

	id	code	valA	valB
1:	1	abc	0.1	11
2:	1	abc	0.6	7
3:	1	abd	1.5	5
4:	2	apq	0.9	10
5:	2	apq	0.3	13

B

	id	begin	end
1:	1	0.1	0.9
2:	2	0.6	0.8

Update **valB** with **NA** while matching on **id**, **valA > begin**, **valA < end**

# CONDITIONAL OPERATIONS

A

	id	code	valA	valB
1:	1	abc	0.1	11
2:	1	abc	0.6	NA
3:	1	abd	1.5	5
4:	2	apq	0.9	10
5:	2	apq	0.3	13

B

	id	begin	end
1:	1	0.1	0.9
2:	2	0.6	0.8

Update **valB** with **NA** while matching on **id**,  $\text{valA} > \text{begin}$ ,  $\text{valA} < \text{end}$

```
A[B, on = .(id, valA > begin, valA < end),  
  valB := NA]
```



# SUMMARY AND FUTURE DIRECTIONS

- data.table allows for *concise* and *straightforward* code, and is *fast* and *memory efficient*
- More efforts towards parallelisation in future
- File backed data.tables would be great feature to have soon, [#1336](#)
- Give data.table a go :-)

**Thank you for  
your attention!**

Questions?