

Introduction to the `data.table` package in R

Matthew Dowle

Revised: February 27, 2014

(A later revision may be available on the [homepage](#))

Introduction

This vignette is aimed at those who are already familiar with R—in particular, creating and using objects of class `data.frame`. We aim for this quick introduction to be readable in **10 minutes**, covering the main features in brief, namely: 1. Keys; 2. Fast Grouping; and 3. Fast time series join. For the context in which this document sits, please briefly check the last section, Further Resources.

`data.table` is not *automatically* better or faster. The user has to climb a short learning curve, experiment, and then use its features well. For example, this document explains the difference between a *vector scan* and a *binary search*. Although both extraction methods are available in `data.table`, if a user continues to use vector scans (as in a `data.frame`), it will ‘work’, but one will miss out on the benefits that `data.table` provides.

Creation

Recall that we create a `data.frame` using the function `data.frame()`:

```
> DF = data.frame(x=c("b", "b", "b", "a", "a"), v=rnorm(5))
> DF

  x      v
1 b  1.0426719
2 b  0.4616091
3 b  1.1801684
4 a -0.1541805
5 a -0.4361733
```

A `data.table` is created in exactly the same way:

```
> DT = data.table(x=c("b", "b", "b", "a", "a"), v=rnorm(5))
> DT

  x      v
1: b -0.4729847
2: b -0.9235720
3: b  1.0287125
4: a  0.0472751
5: a -1.2112938
```

Observe that a `data.table` prints the row numbers slightly differently. There is nothing significant about that. We can easily convert existing `data.frame` objects to `data.table`.

```
> CARS = data.table(cars)
> head(CARS)
```

```

      speed dist
1:      4     2
2:      4    10
3:      7     4
4:      7    22
5:      8    16
6:      9    10

```

We have just created two `data.tables`: `DT` and `CARS`. It is often useful to see a list of all `data.tables` in memory:

```

> tables()

      NAME NROW MB COLS      KEY
[1,] CARS   50  1  speed,dist
[2,] DT     5  1   x,v
Total: 2MB

```

The MB column is useful to quickly assess memory use and to spot if any redundant tables can be removed to free up memory. Just like `data.frames`, `data.tables` must fit inside RAM.

Some users regularly work with 20 or more tables in memory, rather like a database. The result of `tables()` is itself a `data.table`, returned silently, so that `tables()` can be used in programs. `tables()` is unrelated to the base function `table()`.

To see the column types¹ :

```

> sapply(DT, class)

      x      v
"character" "numeric"

```

You may have noticed the empty column `KEY` in the result of `tables()` earlier above. This is the subject of the next section, the first of the 3 main features of the package.

1. Keys

Let's start by considering `data.frame`, specifically `rownames` (or in English, *row names*). That is, the multiple names belonging to a single row. The multiple names belonging to the single row? That is not what we are used to in a `data.frame`. We know that each row has at most one name. A person has at least two names, a first name and a second name. That is useful to organise a telephone directory, for example, which is sorted by surname, then first name. However, each row in a `data.frame` can only have one name.

A *key* consists of one or more columns of `rownames`, which may be integer, factor, character or some other class, not simply character. Furthermore, the rows are sorted by the key. Therefore, a `data.table` can have at most one key, because it cannot be sorted in more than one way.

Uniqueness is not enforced, i.e., duplicate key values are allowed. Since the rows are sorted by the key, any duplicates in the key will appear consecutively.

Let's remind ourselves of our tables:

```

> tables()

      NAME NROW MB COLS      KEY
[1,] CARS   50  1  speed,dist
[2,] DT     5  1   x,v
Total: 2MB

> DT

```

¹As from v1.8.0, `data.table()` no longer converts `character` to `factor`.

```

      x          v
1: b -0.4729847
2: b -0.9235720
3: b  1.0287125
4: a  0.0472751
5: a -1.2112938

```

No keys have been set yet. We can use `data.frame` syntax in a `data.table`, too.

```

> DT[2,]

      x          v
1: b -0.9235720

> DT[DT$x=="b",]

      x          v
1: b -0.4729847
2: b -0.9235720
3: b  1.0287125

```

But since there are no rownames, the following does not work:

```

> cat(try(DT["b",],silent=TRUE))

```

```

Error in `[.data.table`(DT, "b", ) :

```

```

  When i is a data.table (or character vector), x must be keyed (i.e. sorted, and, marked as sorted)

```

The error message tells us we need to use `setkey()`:

```

> setkey(DT,x)
> DT

      x          v
1: a  0.0472751
2: a -1.2112938
3: b -0.4729847
4: b -0.9235720
5: b  1.0287125

```

Notice that the rows in `DT` have been re-ordered according to the values of `x`. The two "a" rows have moved to the top. We can confirm that `DT` does indeed have a key using `haskey()`, `key()`, `attributes()`, or just running `tables()`.

```

> tables()

      NAME NROW MB COLS      KEY
[1,] CARS   50 1  speed,dist
[2,] DT     5 1  x,v         x
Total: 2MB

```

Now that we are sure `DT` has a key, let's try again:

```

> DT["b",]

      x          v
1: b -0.4729847
2: b -0.9235720
3: b  1.0287125

```

By default all the rows in the group are returned². The `mult` argument (short for *multiple*) allows the first or last row of the group to be returned instead.

```
> DT["b",mult="first"]
```

```
      x          v
1: b -0.4729847
```

```
> DT["b",mult="last"]
```

```
      x          v
1: b  1.028712
```

The comma is optional.

```
> DT["b"]
```

```
      x          v
1: b -0.4729847
2: b -0.9235720
3: b  1.0287125
```

Let's now create a new `data.frame`. We will make it large enough to demonstrate the difference between a *vector scan* and a *binary search*.

```
> grpsize = ceiling(1e7/26^2) # 10 million rows, 676 groups
```

```
[1] 14793
```

```
> tt=system.time( DF <- data.frame(
+   x=rep(LETTERS,each=26*grpsize),
+   y=rep(letters,each=grpsize),
+   v=runiform(grpsize*26^2),
+   stringsAsFactors=FALSE)
+ )
```

```
      user  system elapsed
10.340   1.136  11.511
```

```
> head(DF,3)
```

```
      x y          v
1 A a  0.18132669
2 A a  0.08638986
3 A a  0.54725073
```

```
> tail(DF,3)
```

```
      x y          v
10000066 Z z  0.61219283
10000067 Z z  0.04143928
10000068 Z z  0.21274337
```

```
> dim(DF)
```

```
[1] 10000068      3
```

²In contrast to a `data.frame` where only the first rowname is returned when the rownames contain duplicates.

We might say that R has created a 3 column table and *inserted* 10,000,068 rows. It took 11.511 secs, so it inserted 868,740 rows per second. This is normal in base R. Notice that we set `stringsAsFactors=FALSE`. This makes it a little faster for a fairer comparison, but feel free to experiment.

Let's extract an arbitrary group from DF:

```
> tt=system.time(ans1 <- DF[DF$x=="R" & DF$y=="h",]) # 'vector scan'

   user  system elapsed
12.525   0.436  12.993

> head(ans1,3)

      x y      v
6642058 R h 0.4438668
6642059 R h 0.2292900
6642060 R h 0.5743599

> dim(ans1)

[1] 14793     3
```

Now convert to a `data.table` and extract the same group:

```
> DT = as.data.table(DF) # but normally use fread() or data.table() directly, originally
> system.time(setkey(DT,x,y)) # one-off cost, usually

   user  system elapsed
 0.412   0.056   0.472

> ss=system.time(ans2 <- DT[J("R","h")]) # binary search

   user  system elapsed
 0.012   0.000   0.009

> head(ans2,3)

      x y      v
1: R h 0.4438668
2: R h 0.2292900
3: R h 0.5743599

> dim(ans2)

[1] 14793     3

> identical(ans1$v, ans2$v)

[1] TRUE
```

At 0.009 seconds, this was **1443** times faster than 12.993 seconds, and produced precisely the same result. If you are thinking that a few seconds is not much to save, it's the relative speedup that's important. The vector scan is linear, but the binary search is $O(\log n)$. It scales. If a task taking 10 hours is sped up by 100 times to 6 minutes, that is significant³.

We can do vector scans in `data.table`, too. In other words we can use `data.table` *badly*.

```
> system.time(ans1 <- DT[x=="R" & y=="h",]) # works but is using data.table badly
```

³We wonder how many people are deploying parallel techniques to code that is vector scanning

```

      user system elapsed
5.473   0.176   5.661

> system.time(ans2 <- DF[DF$x=="R" & DF$y=="h",]) # the data.frame way

      user system elapsed
11.941   0.408  12.403

> mapply(identical, ans1, ans2)

      x     y     v
TRUE TRUE TRUE

```

If the phone book analogy helped, the **1443** times speedup should not be surprising. We use the key to take advantage of the fact that the table is sorted and use binary search to find the matching rows. We didn't vector scan; we didn't use `==`.

When we used `DT$x=="R"` we *scanned* the entire column `x`, testing each and every value to see if it equalled "R". We did it again in the `y` column, testing for "h". Then `&` combined the two logical results to create a single logical vector which was passed to the `[]` method, which in turn searched it for `TRUE` and returned those rows. These were *vectorized* operations. They occurred internally in R and were very fast, but they were scans. *We* did those scans because *we* wrote that R code.

When `i` is itself a `data.table`, we say that we are *joining* the two `data.tables`. In this case, we are joining `DT` to the 1 row, 2 column table returned by `data.table("R", "h")`. Since we do this a lot, there is an alias for `data.tables` called `J()`, short for join.

```

> identical( DT[J("R", "h"),],
+           DT[data.table("R", "h"),])

[1] TRUE

```

Both vector scanning and binary search are available in `data.table`, but one way of using `data.table` is much better than the other.

The join syntax is a short, fast to write and easy to maintain. Passing a `data.table` into a `data.table` subset is analogous to `A[B]` syntax in base R where `A` is a matrix and `B` is a 2-column matrix⁴. In fact, the `A[B]` syntax in base R inspired the `data.table` package. There are other types of join and further arguments which are beyond the scope of this quick introduction.

The merge method of `data.table` is very similar to `X[Y]`, but there are some differences. See FAQ 1.12.

This first section has been about the first argument to `[]`, namely `i`. The next section has to do with the 2nd argument `j`.

2. Fast grouping

The second argument to `[]` is `j`, which may consist of one or more expressions whose arguments are (unquoted) column names, as if the column names were variables.

```

> DT[, sum(v)]

[1] 4999454

```

When we supply a `j` expression and a 'by' list of expressions, the `j` expression is repeated for each 'by' group:

```

> DT[, sum(v), by=x]

```

⁴Subsetting a keyed `data.table` by a `n`-column `data.table` is consistent with subsetting a `n`-dimension array by a `n`-column matrix in base R

```

      x      V1
1: A 192338.7
2: B 192303.1
3: C 192269.5
4: D 192501.1
5: E 192349.4
6: F 192518.3
7: G 192784.4
8: H 192376.0
9: I 192151.5
10: J 191802.7
11: K 192134.2
12: L 192088.9
13: M 192379.2
14: N 192064.3
15: O 192356.5
16: P 192200.6
17: Q 192138.1
18: R 192247.5
19: S 192209.3
20: T 192163.6
21: U 192354.6
22: V 192328.3
23: W 192695.7
24: X 192110.2
25: Y 192251.7
26: Z 192336.6
      x      V1

```

The `by` in `data.table` is fast. Let's compare it to `tapply`.

```

> ttt=system.time(tt <- tapply(DT$v,DT$x,sum)); ttt

      user system elapsed
24.593   1.652  26.337

> sss=system.time(ss <- DT[,sum(v),by=x]); sss

      user system elapsed
 0.924   0.172   1.105

> head(tt)

      A      B      C      D      E      F
192338.7 192303.1 192269.5 192501.1 192349.4 192518.3

> head(ss)

      x      V1
1: A 192338.7
2: B 192303.1
3: C 192269.5
4: D 192501.1
5: E 192349.4
6: F 192518.3

> identical(as.vector(tt), ss$V1)

[1] TRUE

```

At 1.105 sec, this was **23** times faster than 26.337 sec, and produced precisely the same result. Next, let's group by two columns:

```
> ttt=system.time(tt <- tapply(DT$v,list(DT$x,DT$y),sum)); ttt
  user  system elapsed
26.978   2.524  29.621
> sss=system.time(ss <- DT[,sum(v),by="x,y"]); sss
  user  system elapsed
 1.128   0.268   1.401
> tt[1:5,1:5]
      a      b      c      d      e
A 7374.683 7408.510 7411.224 7464.679 7371.076
B 7400.611 7408.320 7402.351 7392.034 7351.796
C 7400.865 7403.321 7423.999 7424.538 7357.930
D 7378.675 7421.867 7417.467 7394.636 7400.371
E 7445.237 7417.741 7394.770 7359.348 7460.721
> head(ss)
  x y      V1
1: A a 7374.683
2: A b 7408.510
3: A c 7411.224
4: A d 7464.679
5: A e 7371.076
6: A f 7416.746
> identical(as.vector(t(tt)), ss$V1)
[1] TRUE
```

This was **21** times faster, and the syntax is a little simpler and easier to read.

The following features are mentioned only briefly here; further examples are in `?data.table` and the [FAQ vignette](#).

- To return several expressions, pass a `list()` to `j`.
- Each item of the list is recycled to match the length of the longest item.
- You can pass a `list()` of *expressions* of column names to `by` e.g.
`DT[,sum(v),by=list(month(dateCol),region)]`
where calling `month()` on `dateCol` is what we mean by expressions of column names.
- Any R functions from any package can be used in `j` and `by`.

3. Fast time series join

This is also known as last observation carried forward (LOCF) or a *rolling join*.

Recall that `x[i]` is a join between `data.table x` and `data.table i`. If `i` has 2 columns, the first column is matched to the first column of the key of `x`, and the 2nd column to the 2nd. An equi-join is performed, meaning that the values must be equal.

The syntax for fast rolling join is

```
x[i,roll=TRUE]
```

As before the first column of `i` is matched to `x` where the values are equal. The last column of `i` though, the 2nd one in this example, is treated specially. If no match is found, then the row before is returned, provided the first column still matches.

For examples type `example(data.table)` and follow the output at the prompt.

Other resources

This was a quick start guide. Further resources include :

- The help page describes each and every argument: `?data.table`
- The FAQs deal with distinct topics: `vignette("datatable-faq")`
- The performance tests contain more examples: `vignette("datatable-timings")`
- `test.data.table()` contains over 250 low level tests of the features
- Website: <http://datatable.r-forge.r-project.org/>
- Presentations:
 - <http://files.meetup.com/1406240/Data%20munging%20with%20SQL%20and%20R.pdf>
 - <http://www.londonr.org/LondonR-20090331/data.table.LondonR.pdf>
- YouTube Demo: <http://www.youtube.com/watch?v=rvT8XThGA8o>
- R-Forge commit logs: <http://lists.r-forge.r-project.org/pipermail/datatable-commits/>
- Mailing list : datatable-help@lists.r-forge.r-project.org
- User reviews : <http://crantastic.org/packages/data-table>