List-columns in data.table: Nesting and unnesting data tables and vectors

Tyson S. Barrett¹

¹ Utah State University

5 Abstract

3

6

The use of *list-columns* in data frames and tibbles in the R statistical environment is well documented (e.g. Bryan, 2018), providing a cognitively efficient way to organize results of complex data (e.g. several statistical models, groupings of text, data summaries, or even graphics) with corresponding data. For example, one can store student information within classrooms, player information within teams, or even analyses within groups. This allows the data to be of variable sizes without overly complicating or adding redundancies to the structure of the data. In turn, this can improve the reliability to appropriately analyze the data.

Because of its efficiency and speed, being able to use data.table to work with list-columns would be beneficial in many data contexts (e.g. to reduce memory usage in large data sets). Herein, I demonstrate how one can create list-columns in a data table using the by argument in data.table and purrr::map(). This is done using an example data set containing information on professional basketball players in the United States. I compare the behavior of the data.table approach to the dplyr::group_nest() function, one of the several powerful tidyverse nesting functions. Results using bench::mark() show the speed and efficiency of using data.table to work with list-columns.

Keywords: data.table, dplyr, list-columns, nesting

7 Introduction

The use of *list-columns* in data frames and tibbles in the R statistical environment (R Core Team, 2018) provides a cognitively efficient way to organize complex data (e.g. several statistical models, groupings of text, data summaries, or even graphics) with corresponding

Correspondence concerning this article should be addressed to Tyson S. Barrett, 2800 Old Main, Logan, UT 84322. E-mail: tyson.barrett@usu.edu

data in a concise manner. It has become a common approach to wrangling data in the tidyverse, with functions across dplyr and tidyr providing functionality to work with list-columns (Bryan, 2018; Wickham et al., 2019; Wickham & Henry, 2019). This format is often called "nested" data, where information is, in essence, nested within a column of data.

For example, list-columns can be used to nest data regarding students within class-rooms, players within teams, measures within individuals, and text within chapters. This allows the user to do certain data manipulations within each group in a consistent, controlled manner. This can ensure that accidentally including data from other groups does not occur. Furthermore, nesting can reduce the difficulty to appropriately analyze the data stored in the list-column. Using functions like lapply() or purrr::map*() makes further analysis of the nested data more intuitive and error-free.

Because of its efficiency and speed, being able to use data.table to work with list-columns would be beneficial in many data contexts (e.g. to reduce memory usage in large data sets, increase speed of calculations). Herein, I demonstrate how one can create list-columns in a data table using the by argument in data.table (using a custom function) and the purrr::map*() functions. I further highlight the dplyr::group_nest() function and show a slightly more efficient approach when using a data table. Using bench::mark(), I assess the speed and efficiency of using data.table to work with list-columns.

This article relies on several powerful R packages, including data.table, dplyr, bench, tidyr, papaja, stringr, ggplot2, ggbeeswarm, ggrepel, performance, rvest, and lobstr (Aust & Barth, 2018; Clarke & Sherrill-Mix, 2017; Dowle & Srinivasan, 2019; Hester, 2019; Lüdecke, Makowski, & Waggoner, 2019; Slowikowski, 2018; Wickham, 2016, 2019c, 2019b, 2019a; Wickham et al., 2019; Wickham & Henry, 2019).

Example Data

Throughout much of this paper, I demonstrate the use of *list-columns* in data.table using data from NBA Stuffer. These data are downloaded, providing information on players from the 2017-2018 and 2018-2019 seasons. To do so, I first read in the HTML data, then extract the tables with player data by year (using a short custom function), add indicators, and then combine each into a single data table. Each step is shown in the code below.

```
url_2018 <- "https://www.nbastuffer.com/2017-2018-nba-player-stats/"
url_2019 <- "https://www.nbastuffer.com/2018-2019-nba-player-stats/"
players_2018 <- read_html(url_2018)
players_2019 <- read_html(url_2019)

extract_fun <- function(html){
  tabs <- html_nodes(html, "table")[2] %>%
    html_table(fill = TRUE)
  tabs[[1]]
}
```

Below is a subset of this imported data set, showing only four of the variables and the first six rows.

```
##
             full name team year mpg
                                       ppg apg
          Aaron Brooks Min 2018 5.9
  ## 1:
                                       2.3 0.6
43
  ## 2:
          Aaron Gordon Orl 2018 32.9 17.6 2.3
  ## 3: Aaron Harrison Dal 2018 25.9
                                      6.7 1.2
  ## 4:
         Aaron Jackson Hou 2018 34.5
                                      8.0 1.0
  ## 5:
           Abdel Nader Bos 2018 10.9 3.0 0.5
         Adreian Payne Orl 2018 8.5
  ## 6:
                                      4.2 0.0
```

49

50

52

53

54

Nesting with data.table

In dplyr, the group_nest() function is valuable when creating list-columns based on a grouping variable. It takes the data by group and puts it all in a list-column. Figure 1 highlights the process of taking a data frame and creating a nested data frame with a list-column. That is, all data from variables x, y, and z relating to each group is split into a distinct data frame and stored within the data column.

Overall, this function is efficient and fast but—by relying on data.table—it can be somewhat faster. This will be shown using the following function, which relies solely on the syntax of data.table, using the j and by arguments as shown below.

```
group_nest_dt <- function(dt, ..., .key = "data"){
   stopifnot(is.data.table(dt))

  by <- substitute(list(...))

  dt <- dt[, list(list(.SD)), by = eval(by)]
   setnames(dt, old = "V1", new = .key)
   dt
}</pre>
```

First thing to note, is that in the data table, we create a list within a list containing
the .SD special object. This object is all the data in the data table except for the variables
that are in the by argument. The by argument, before being evaluated within the data
table, first becomes an unevaluated list of bare variable names and then evaluate it within
the data.table syntax. In essence, this function takes a data table, then creates a list of
the data table per group specified in the by argument.

```
head(group_nest_dt(players, team))
```

```
##
                        data
         team
   ## 1:
          Min <data.table>
      2:
          Orl <data.table>
   ##
      3:
          Dal <data.table>
67
   ## 4:
          Hou <data.table>
   ## 5:
          Bos <data.table>
69
          Ind <data.table>
   ## 6:
70
```

The syntax and output are nearly identical to the dplyr::group_nest() function but has data tables in the list-column instead of tibbles.

```
head(group_nest(players, team))
```

```
## # A tibble: 6 x 2
73
   ##
        team
               data
74
   ##
        <chr> <list>
75
               <tibble [44 x 30]>
   ## 1 Atl
76
   ## 2 Bos
               <tibble [37 x 30]>
77
               <tibble [41 x 30]>
   ## 3 Bro
   ## 4 Cha
               <tibble [34 x 30]>
   ## 5 Chi
               <tibble [43 x 30]>
   ## 6 Cle
               <tibble [49 x 30]>
```

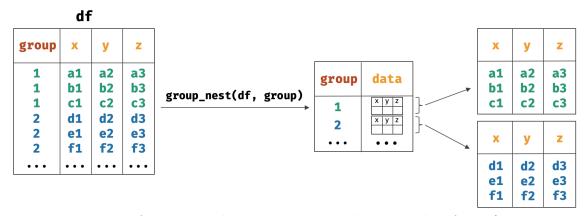


Figure 1. Diagram of one approach to creating a list-column in a data frame (i.e. creating a nested data frame).

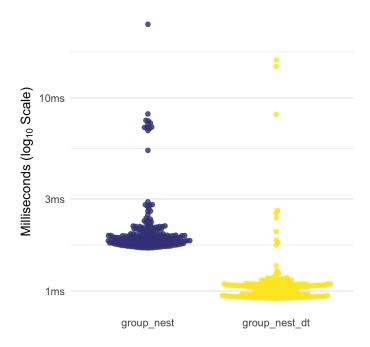


Figure 2. Speed comparisons for each nesting approach. Note the scale of the y-axis is \log_{10} .

Given both perform very similar data manipulations, it is of interest to see if there are differences in memory and speed performance. Figure 2 presents the timings from bench::mark() across the two approaches, showing group_nest_dt() is often faster, although differences for this size of data set are not meaningful. The memory allocated is also very similar, with group_nest_dt() allocating 463KB and group_nest() allocating 335KB.

This nesting approach can be used with multiple grouping variables too. For example, I show how a user could nest by both team and year, as is done below.

head(group_nest_dt(players, team, year))

```
##
         team year
90
          Min 2018 <data.table>
   ##
      1:
91
      2:
          Orl 2018 <data.table>
   ##
      3:
          Dal 2018 <data.table>
93
      4:
          Hou 2018 <data.table>
          Bos 2018 <data.table>
   ##
      5:
          Ind 2018 <data.table>
```

82

83

84

85

86

88

98

99

100

Analyses within the Nested Data

Often, the nested data can provide an intuitive format to run several analyses to understand key features of the data within the groups. Below, the relationship between points-per-game and assists-per-game for each team and year is modeled and then the R^2

116

117

118

119

of the models are extracted. Since performance::r2() provides two versions of R^2 , I then grab only the first of the two types.

```
##
103
          team year
                            data ppg_apg
                                               r2_list r2_ppg_apg
   ## 1:
          Min 2018 <data.table>
                                     <lm> <r2_generic>
                                                        0.4662060
104
          Orl 2018 <data.table>
                                     <lm> <r2_generic> 0.4357684
105
          Dal 2018 <data.table>
                                     <lm> <r2_generic>
   ## 3:
                                                        0.4305347
106
   ## 4:
          Hou 2018 <data.table>
                                     <lm> <r2_generic>
                                                        0.6967150
107
          Bos 2018 <data.table>
                                     <lm> <r2 generic>
                                                        0.6043402
108
   ## 5:
          Ind 2018 <data.table>
                                     <lm> <r2_generic>
   ## 6:
                                                        0.6060465
```

This produces two list-columns (ppg_apg and r2_list) and a numeric vector (r2_ppg_apg) all organized by team and year. This information is then readily available to plot. For example, one can look at how related points-per-game and assists-per-game are by team and year—in essence, showing which teams have players who both score and assist. The example plot is shown in Figure 3.

Unnesting with data.table

After performing the manipulations or analyses within the nest, it can often be necessary to unnest to finalize analyses. Again, like with group_nest_dt(), the unnest_dt() function below relies solely on the syntax of data.table, using the j and by arguments as shown below.

```
unnest_dt <- function(dt, col, id){
  stopifnot(is.data.table(dt))

  by <- substitute(id)
  col <- substitute(unlist(col, recursive = FALSE))

  dt[, eval(col), by = eval(by)]
}</pre>
```

121

122

123

136

137

138

139

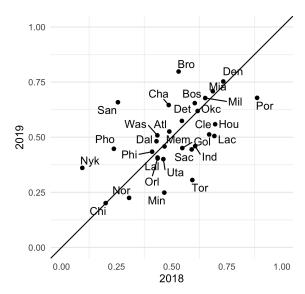


Figure 3. Example analysis performed using nested data to provide information for each team and year.

This function can be used to unnest a data table, like the players_nested data table from earlier, where the nested column can be a data table, data frame, or tibble. Below, the data column in the table is unnested by team and year and then a few of the variables are selected for demonstration purposes.

```
##
              team year
                                   full_name pos
124
                                                     age gp
                                                              mpg
   ##
               Min 2018
                                Aaron Brooks
                                               PG 33.00 32
          1:
                                                              5.9
125
               Min 2018
                              Andrew Wiggins
                                               SF 22.00 82 36.3
   ##
          2:
          3:
               Min 2018
                               Anthony Brown
                                                SF 25.00
                                                              3.7
   ##
127
                                Cole Aldrich
                                                 C 29.00 21
   ##
          4:
               Min 2018
                                                              2.3
128
   ##
          5:
               Min 2018
                                Derrick Rose
                                               PG 29.00
                                                             12.4
129
   ##
130
   ## 1223:
               Det 2019
                              Svi Mykhailiuk
                                                 G 21.84
                                                           3
                                                              6.6
131
   ## 1224:
               Det 2019
                               Zaza Pachulia
                                                 C 35.16 68 12.9
132
               Det 2019 Glenn Robinson III G-F 25.26 47 13.0
    ## 1225:
133
    ## 1226:
               Det 2019
                                   Ish Smith
                                                 G 30.77 56 22.3
134
    ## 1227:
               Det 2019
                                Khyri Thomas
                                                 G 22.92 26 7.5
135
```

Again, this function is quick and efficient. Figure 4 presents the timings from bench::mark() across the two unnesting approaches, showing the data.table approach is much faster. The memory allocated is about half for the data.table approach here, with unnest_dt() allocating 912KB and tidyr::unnest() allocating 1.83MB.

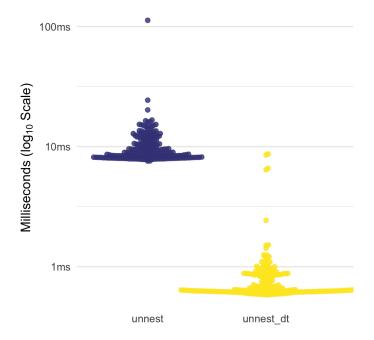


Figure 4. Speed comparisons for each unnesting approach. Note the scale of the y-axis is \log_{10} .

Unnesting Vectors with data.table

A slight variation of this function can be used for list-columns with atomic vectors instead of data tables. A function like the following works well.

In players_nested, the r2_list column is a list of numeric vectors. This can be unnested as shown below, providing the two measures of \mathbb{R}^2 per team per year.

```
##
             team year
                                r2
145
    ##
              Min 2018
                         0.466206
146
              Min 2018 0.4280779
    ##
147
    ##
              Orl 2018 0.4357684
148
              Orl 2018 0.4025783
    ##
149
         5:
              Dal 2018 0.4305347
150
   ##
151
    ## 116:
              Lac 2019 0.4808586
152
              Phi 2019 0.4342685
    ## 117:
153
              Phi 2019 0.4106964
    ## 118:
154
              Det 2019 0.5740963
    ## 119:
155
    ## 120:
              Det 2019 0.550435
156
```

158

159

160

161

162

163

164

165

168

169

Memory Usage of List-Columns

Last item to demonstrate herein is the computer memory usage of different formats of data tables with the same data. We can use the following large data sets in wide format, nested wide format, long format, and nested wide format to make brief comparisons.

I use the lobstr package to assess the object size of each format of the same data, shown in Table 1. Not surprising, the memory usage of nested data is lower than for its none nested corresponding data. This is directly related to the reduction in redundancies in the data otherwise there. That is, the nested data has far fewer rows containing the group variable. That, alone, in this large data saves memory. For example, the size of a single column of the group variable in wide format is 4 MB; and in long format it is 16 MB By reducing a single variable in this case, we save several megabytes of memory.

Discussion

List-columns are a useful approach to structuring data into a format that can be safely cleaned, manipulated, and analyzed by groups. It also provides for a more cognitively

Table 1
Memory usage for each format of the same data

| Format | Memory (MB) |
|--------------------|-------------|
| Wide Format | 40.0 |
| Nested Wide Format | 36.0 |
| Long Format | 80.0 |
| Nested Long Format | 64.0 |

efficient way for a user to understand their data, allowing large data to be represented more concisely within groups.

The tidyverse provides several functions to work with nested data, which are relatively quick and efficient. For most data situations, these functions will do all that a user will need. However, in some situations, data.table can perform needed manipulations and analyses that cannot otherwise be done or that would take too long to complete. In these situations, and for users that prefer to use data.table, this tutorial can help provide direction in using list-columns.

Furthermore, as expected, the memory usage of nested data is lower than for its none nested corresponding data. This is due to the reduction in the redundancies present in wide and long format. This suggests that it is not only the cognitive benefits to the user that makes this format more efficient.

Limitations

There are some notable limitations to list-columns in general, and in data.table specifically. First, the three custom functions built on data.table presented herein are not well-tested and are certainly not expected to work in each case where dplyr::group_nest(), tidyr::unnest(), and other tidy functions would work. Rather, they were presented to show how a user can leverage the speed and efficiency of data.table to create, and work with, list-columns.

Second, it is important to realize that nested data can remove the ability to use vectorization across groups. Depending on the analyses being conducted, this may slow down the computation to the point that nested data actually is a hindrance to performance.

Finally, when using list-columns in tibbles, the print method provides the dimensions of each corresponding nested tibble. This method is helpful in understanding the nested data without any need to extract it. This could be a minor, but valuable, update to the print method in data.table.

Conclusions

The use of list-columns in data.table is very similar to that in the tidyverse. It provides speed and efficiency in both nesting and unnesting the data, and can be used with the purrr::map*() and other powerful functions.

201 References

- Aust, F., & Barth, M. (2018). papaja: Create APA manuscripts with R Markdown. Retrieved from https://github.com/crsh/papaja
- Bryan, J. (2018). List columns (as part of "purrr tutorial"). Retrieved from https://jennybc. github.io/purrr-tutorial/ls13_list-columns.html
- Clarke, E., & Sherrill-Mix, S. (2017). *Ggbeeswarm: Categorical scatter (violin point) plots*.

 Retrieved from https://github.com/eclarke/ggbeeswarm
- Dowle, M., & Srinivasan, A. (2019). *Data.table: Extension of 'data.frame'*. Retrieved from https://CRAN.R-project.org/package=data.table
- Hester, J. (2019). Bench: High precision timing of r expressions. Retrieved from https: //CRAN.R-project.org/package=bench
- Lüdecke, D., Makowski, D., & Waggoner, P. (2019). Performance: Assessment of regression models performance. Retrieved from https://easystats.github.io/performance/
- R Core Team. (2018). R: A language and environment for statistical computing. Vienna,
 Austria: R Foundation for Statistical Computing. Retrieved from https://www.
 R-project.org/
- Slowikowski, K. (2018). Ggrepel: Automatically position non-overlapping text labels with 'ggplot2'. Retrieved from https://CRAN.R-project.org/package=ggrepel
- Wickham, H. (2016). *Ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York. Retrieved from https://ggplot2.tidyverse.org
- Wickham, H. (2019a). Lobstr: Visualize r data structures with trees. Retrieved from https://CRAN.R-project.org/package=lobstr
- Wickham, H. (2019b). Rvest: Easily harvest (scrape) web pages. Retrieved from https: //CRAN.R-project.org/package=rvest
- Wickham, H. (2019c). Stringr: Simple, consistent wrappers for common string operations.

 Retrieved from https://CRAN.R-project.org/package=stringr
- Wickham, H., François, R., Henry, L., & Müller, K. (2019). *Dplyr: A grammar of data manipulation*. Retrieved from https://CRAN.R-project.org/package=dplyr
- Wickham, H., & Henry, L. (2019). *Tidyr: Tidy messy data*. Retrieved from https://CRAN.
 R-project.org/package=tidyr