

## data frames in R

Intuitively, data frames are similar to matrices; having two-dimensional rows-and-columns structure. The difference exists where each column in a list can be of different classes (modes). For example, one column can be numerical values, while another contains character strings. Therefore, just as lists are heterogeneous analogs of vectors in one-dimension, data frames are heterogeneous analogs of matrices for two-dimensional data.

Concretely, a data frame is a list, with the components of that list being equal-length vectors. R does not allow components to be other types of objects, including other data frames. Therefore, the result is a heterogeneous-data analogs of arrays, not common in practice. It will thus be assumed that all components of a data frame are vectors.

## creating data frames in R

When building data frames in R, the default selection sets **stringsAsFactors = TRUE**, converting character vectors into **factors** (discussed later).

```
1 > kids <- c("Jack", "Jill")
2 > ages <- c(12,10)
3 > d <- data.frame(kids,ages, stringsAsFactors = FALSE)
4 > d
5   kids ages
6 1 Jack  12
7 2 Jill  10
```

## accessing data frames in R

Considering that **d** is a list, it can be accessed via component index values or component names:

```
1 > d[[1]]           #accessing dataframe d with components indices
2 [1] "Jack" "Jill"
3
4 > d$kids           #accessing dataframe d with component names
5 [1] "Jack" "Jill"
```

Additionally, data frame **d** can be treated as a matrix to view components (also illustrated in its structure):

```
1 > d[,1]            #print column one of dataframe d
2 [1] "Jack" "Jill"
3
4 > str(d)           #examine matrix-like structure of dataframe d
5 'data.frame':  2 obs. of  2 variables:
6  $ kids: chr  "Jack"
```

The structure of dataframe **d** illustrates two observations—two rows—that stores the data on two variables—two columns. It is safer and more clear to extra data from matrices using name references of the components with the **\$** operator. However, in the construction of general code (packages and functions), it is necessary to use matrix-like notation **[ ]** and is also useful for extracting sub dataframes.

## extracting sub data frames in R

Matrix operations also apply to data frames where appropriate. Most noted is often **filtering** of dataframes to extract various sub dataframes relevant to analysis.

```
1  > exams                                     #print dataframe exams
2    Exam 1 Exam 2 Quiz
3    1     2.0    3.3 4.0
4    2     3.3    2.0 3.7
5    3     4.0    4.0 4.0
6    4     2.3    0.0 3.3
7    5     2.3    1.0 3.3
8    6     3.3    3.7 4.0
9
10 > exams[2:5,]                             #subset rows 2 through 5 of dataframe exams
11    Exam 1 Exam 2 Quiz
12    2     3.3     2 3.7
13    3     4.0     4 4.0
14    4     2.3     0 3.3
15    5     2.3     1 3.3
16
17 > exams[2:5,2]                             #subset rows 2-5 of column 2 in exams (returned as vector)
18 [1] 2 4 0 1
19
20 > class(exams[2:5,2])                      #print class of dataframe exams subset above
21 [1] "numeric"
22
23 > exams[2:5,2,drop = FALSE]                #specify the drop argument to maintain matrix integrity
24    Exam 2
25    2     2
26    3     4
27    4     0
28    5     1
29
30 > class(exams[2:5,2,drop = FALSE])          #print the class of the subset with drop = FALSE
31 [1] "data.frame"
```

Additionally, dataframes can be filtered through sub dataframe extraction and specification of constraints.

```
1  > exams[exams$`Exam 1` >= 2.5,]
2    Exam 1 Exam 2 Quiz
3    2     3.3    2.0 3.7
4    3     4.0    4.0 4.0
5    6     3.3    3.7 4.0
```

## data frames with NA values in R

Assuming that NA values exist within the rows/columns of a dataframe, certain operations cannot be properly performed. In these cases, R should be specifically assigned to remove NA values from operations.

```
1  x <- c(2, NA, 4)
2  > mean(x)
3  [1] NA
4
5  > mean(x, na.rm = TRUE)
6  [1] 3
```

When dealing with NA values in larger datasets, it might be necessary to remove all observations where NA values are present. This is achieved through the **complete.cases()** function as illustrated on the proceeding page:

```

1  > exams                                     #print the exams dataframe with NA value in last row
2      Exam 1 Exam 2 Quiz
3  1      2.0      3.3  4.0
4  2      3.3      2.0  3.7
5  3      4.0      4.0  4.0
6  4      2.3      0.0  3.3
7  5      2.3      1.0  3.3
8  6      3.3      3.7  4.0
9  7      2.0      NA  4.0
10
11 > mean(exams)                             #attempt to return the mean of the exam dataframe
12 [1] NA
13 Warning message:
14 In mean.default(exams) : argument is not numeric or logical: returning NA
15
16 > complete.cases(exams)                   #print the result of rows that exclude NA values
17 [1] TRUE TRUE TRUE TRUE TRUE TRUE FALSE
18
19 > cmlExams <- exams [complete.cases(exams),] #subset the dataframe from non-NA obs.
20 > cmlExams
21      Exam 1 Exam 2 Quiz
22 1      2.0      3.3  4.0
23 2      3.3      2.0  3.7
24 3      4.0      4.0  4.0
25 4      2.3      0.0  3.3
26 5      2.3      1.0  3.3
27 6      3.3      3.7  4.0

```

## rbind() and cbind() with data frames ➡ in R

The **rbind()** and **cbind()** functions apply to dataframes equally, given compatible dataframe sizes; additional columns and rows are required to have the same length as the target dataframe. It is typical where an added row is in the form of another dataframe or list:

```

1  > d
2      kids ages
3  1 Jack   12
4  2 Jill  10
5
6  > rbind(d, list("laura", 19))
7      kids ages
8  1 Jack   12
9  2 Jill  10
10 3 laura  19

```

Additionally, new columns can be created from existing columns:

```

1  > eq <- cbind(exams, exams[,2] - exams[,1])
2  > eq
3      Exam 1 Exam 2 Quiz exams[, 2] - exams[, 1]
4  1      2.0      3.3  4.0                  1.3
5  2      3.3      2.0  3.7                 -1.3
6  3      4.0      4.0  4.0                  0.0
7  4      2.3      0.0  3.3                 -2.3
8  5      2.3      1.0  3.3                 -1.3
9  6      3.3      3.7  4.0                  0.4

```

To avoid extraneous names being assigned to calculated columns, variables can be specified with the **names()** function. Alternatively, the properties of lists (equal-length vectors) can be applied and calculated in a vectorized fashion; the names can be assigned in the midst of this optimization (proceeding page):

```

1 > exams$ExamDiff <- exams$`Exam 2` - exams$`Exam 1`
2 > exams
3   Exam 1 Exam 2 Quiz ExamDiff
4 1    2.0    3.3  4.0     1.3
5 2    3.3    2.0  3.7    -1.3
6 3    4.0    4.0  4.0     0.0
7 4    2.3    0.0  3.3    -2.3
8 5    2.3    1.0  3.3    -1.3
9 6    3.3    3.7  4.0     0.4

```

The above illustrates the addition of components to an existing list. The property of **recycling** can equally be exploited in the creation of components in existing lists.

```

1 > d
2   kids ages
3 1 Jack  12
4 2 Jill  10
5
6 > d$one <- 1
7 > d
8   kids ages one
9 1 Jack  12  1
10 2 Jill  10  1

```

## apply( ) with data frames in R

The **apply( )** function can be used against dataframes, given that the columns are of the same type.

```

1 > apply(exams,1,max) #finding the maximum grade each student received
2 [1] 4.0 3.7 4.0 3.3 3.3 4.0

```

## merging data frames in R

In relational database theory, joining tables according to common variable values is an essential task. In R, the latter is accomplished through the **merge( )** function.

```

1 > merge(x, y) #the simplest form of a join in R

```

The above illustration assumes that the two data frames have one or more columns with names in common:

```

1 > d1
2   kids states
3 1   Jack   CA
4 2   Jill   MA
5 3 Jillian  MA
6 4   John   HI
7
8 > d2
9   ages kids
10 1  10  Jill
11 2   7 Lillian
12 3  12  Jack
13
14 > d <- merge(d1, d2) #merging d1 and d2 only returns records with common properties
15 > d
16   kids states ages
17 1 Jack   CA  12
18 2 Jill  MA  10

```

Utilized the **by.x** and **by.y** arguments in the **merge( )** function handles cases where variables contain the same information for joining, but lack consistency across naming conventions (illustrated on proceeding page):

```

1 > d3
2   ages      pals
3 1  12      Jack
4 2  10      Jill
5 3   7 Lillian
6
7 > merge(d1, d3, by.x = "kids", by.y = "pals")
8   kids states ages
9 1 Jack      CA  12
10 2 Jill     MA  10

```

It is possible that duplicate matches will appear with the full results of a **merge()** application.

```

1 > d1
2   kids states
3 1  Jack      CA
4 2  Jill      MA
5 3 Jillian    MA
6 4  John      HI
7
8 > d2a <- rbind(d2, list(15, "Jill"))
9 > d2a
10  ages      kids
11 1  10      Jill
12 2   7 Lillian
13 3  12      Jack
14 4  15      Jill
15
16 > merge(d1, d2a)      #One of the Jill's has unknown residence; erroneously assigned MA
17   kids states ages
18 1 Jack      CA  12
19 2 Jill      MA  10
20 3 Jill      MA  15

```

The above illustrations displays the unintended consequences of applying certain **merge()** functions to dataframes in R. It is imperative to be aware of the data and how it is being applied in R functions.

## applying functions to data frames $f(\text{grid})$ in R

Similar to lists, the **lapply()** and **sapply()** functions can be applied to dataframes in R; remembering that dataframes are special cases of lists (list components consisting of the dataframe's columns). Therefore, calling **lapply()** on a dataframe with specified function **f()**, will be called on each of the frame's columns, with the returned valued populated in a corresponding list.

```

1 > d
2   kids ages
3 1 Jack  12
4 2 Jill  10
5
6 > d1 <- lapply(d, sort) #create list d1 with two vectors, sorted kids and sorted ages
7 > d1
8 $kids
9 [1] Jack Jill
10 $ages
11 [1] 10 12

```

Noting that **d1** is a list; dataframe coercion to a is accomplished through the **as.data.frame()** function:

```

1 > as.data.frame(d1) #coercion to a dataframe has lost column references (Jack ≠ 10)
2   kids ages
3 1 Jack  10
4 2 Jill  12

```