

## factors and tables in R

Factors are data type in R that resembles that of nominal (categorical) variables in statistics. Factors are nonnumeric in nature, even if coded to numerical values in practice. Factors are viewed in R as vectors with some additional information attached; the extra information being a record of distinct values in a vector, referred to as **levels**. The illustration below displays the results of a factor coercion:

```
1 > x <- c(5,12,13,12)
2 > xf <- factor(x)
3 > xf
4 [1] 5 12 13 12
5 Levels: 5 12 13
```

The **levels** in the above illustration represent the distinct values in vector **xf**—**5, 12, 13**.

```
1 > str(xf) #note the nonnumeric "5" representation
2 Factor w/ 3 levels "5","12","13": 1 2 3 2
3
4 > unclass(xf) #note that the core of xf is the level representations of 1, 2, and 3
5 [1] 1 2 3 2
6 attr(,"levels")
7 [1] "5" "12" "13"
```

In addition to the behavior in the above illustration, the **length()** of a factor is still represented as the length of the data, opposed to a count of levels, etc.

```
1 > length(xf)
2 [1] 4
```

Anticipated levels can be added into the factor's construction. For example, the current data represents 3 separate levels (categories, labels, classes, etc.) and a fourth is anticipated to exist in a new dataset:

```
1 > x #four-element vector x
2 [1] 5 12 13 12
3
4 > xf #factor xf with 3 levels (5, 12, 13)
5 [1] 5 12 13 12
6 Levels: 5 12 13
7
8 > xff <- factor(x, levels = c(5,12,13,88)) #factor xf with 4 levels (5, 12, 13, 88)
9 > xff
10 [1] 5 12 13 12
11 Levels: 5 12 13 88
12
13 > xff[2] <- 88 #example assignment of additional factor
14 > xff
15 [1] 5 88 13 12
16 Levels: 5 12 13 88
```

Conversely, undefined levels within a factor cannot be added if the level within the factor does not exist:

```
1 > xff[2] <- 28
2 Warning message:
3 In `[<-`factor`(`*tmp*`, 2, value = 28) :
4 invalid factor level, NA generated
```

## factors and functions $f(\text{factor})$ in R

In the context of factors, the **tapply()** function is leveraged to apply functions against factor elements.

```
1 tapply(X, INDEX< FUN = NULL, ..., simplify = TRUE)
```

The call **tapply(x, f, g)** has vector **x**, factor (or list of multiple factors) **f**, and function **g**.

The motivating example illustration is vector **x** of ages of voters and factor **f** of voter party affiliation. Each factor in **f** must consist of equal lengths to **x**. If a component of **f** is a vector, it is coerced into a factor by application of function **as.factor()** to the vector in **f**.

The operation performed by **tapply()** temporarily splits **x** into groups, each corresponding to a level (or combination of levels) in factor **f**, finally applying **g()** to the resulting subvectors of **x** (illustrated below):

```
1 > ages <- c(25,26,55,37,21,42)
2 > affils <- c("R","D","D","R","U","D")
3 > tapply(ages, affils, mean)
4   D  R  U
5  41 31 21
```

The example illustration is expanded to break down groups by both gender and age. When 2 or more factors are applied, each produces a set of groups computed against each other under **AND** logic. The **mean()** function will then be applied to return the each of the four subgroups inclusive of the operation:

1) Male < 25; 2) Male > 25; 3) Female < 25; 4) Female > 25

```
1 > d <- data.frame(list(gender = c("M","M","F","M","F","F"),
2 +                       age = c(47,59,21,32,33,24),
3 +                       income = c(55000,8800,32450,76500,123000,45650)))
4 > d
5   gender age income
6 1      M  47  55000
7 2      M  59   8800
8 3      F  21  32450
9 4      M  32  76500
10 5      F  33 123000
11 6      F  24  45650
12
13 > d$over25 <- ifelse(d$age > 25, 1 ,0)
14 > d
15   gender age income over25
16 1      M  47  55000      1
17 2      M  59   8800      1
18 3      F  21  32450      0
19 4      M  32  76500      1
20 5      F  33 123000      1
21 6      F  24  45650      0
22
23 > tapply(d$income, list(d$gender, d$over25),mean)
24      0      1
25 F 39050 123000.00
26 M   NA    46766.67
```

The two specified factors (Age < OR > 25 and Gender) each contain two levels. Therefore, **tapply()** split (partitioned) the data into four groups; each group representative of each possible permutation of gender and age. Lastly, the **mean()** function was applied to each group individually.

Using the **split()** function simply forms the vector into groups, opposed to splitting a vector into groups and **then** apply a function against each group; seen in the **tapply()** function.

```
1 split(x, f, drop = FALSE, ...)
```

The vector **x** and factor **f** serve similar roles in the **split()** function as they did in the **tapply()** function. Note that in the **split()** function, **x** is allowed to be a dataframe; the latter is not true with **tapply()**.

```
1 > d
2   gender age income over25
3 1      M  47  55000      1
4 2      M  59   8800      1
5 3      F  21  32450      0
6 4      M  32  76500      1
7 5      F  33 123000      1
8 6      F  24  45650      0
9
10 > split(d$income, list(d$gender, d$over25)) #list output denoted by $
11 $F.0
12 [1] 32450 45650
13
14 $M.0
15 numeric(0)
16
17 $F.1
18 [1] 123000
19
20 $M.1
21 [1] 55000 8800 76500
```

Note that the output is a list with components denoted by **\$**; the vectors are represented by the combination of factors that make up the results (e.g. **\$F.0** = Female < 25, **\$M.1** = Male > 25).

The example below determines the indices of which vector elements correspond to Male, Female, or Infants:

```
1 > g <- c("M","F","F","I","M","M","F")
2 > split(1:7,g)
3 $F
4 [1] 2 3 7
5
6 $I
7 [1] 4
8
9 $M
10 [1] 1 5 6
```

The **by()** function operates similar to that of **tapply()**, which is internally called within **by()**. The difference is how the **by()** function is applied to **objects**, rather than vectors. This is useful for performing certain operations against dataframes and matrices.

```
1 split(x, f, drop = FALSE, ...)
```

Calls to the **by()** function are similar to those of the **tapply()** function. The first argument **x** specifies the data, the second argument **f** specifies the grouping factor(s), and the third argument specifies the function to be applied to each factor group individually. Just as **tapply()** forms groups of indices of a vector corresponding to each factor level, **by()** finds groups of row numbers of a dataframe, creating sub dataframes for each corresponding factor level. The specified function is then called against each of the identified factor groups defined in the problem.

## working with tables in R

Given the following example:

```
1 > u <- c(22,8,33,6,8,29,-2)
2 > fl <- list(c(5,12,13,12,13,5,13),c("a","bc","a","a","bc","a","a"))
3 > tapply(u,fl,length)
4      a bc
5 5  2 NA
6 12 1  1
7 13 2  1
```

Again, **tapply()** partitions vector **u** into subvectors and applies the **length()** function to each subvector. Note that the latter operation is independent of the contents of vector **u**; the focus is strictly on the factors. The resulting **contingency table** illustrates the frequency of factors occurring in the data. Note the occurrence of **NA** representing that **5** occurred with “**bc**” in no observations from the function. The problem with the above illustration is the misrepresentation of **0** occurrences as an **NA** value; use **table()** instead:

```
1 > table(fl)      #applying the table() function to appropriately represent all cases
2      fl.2
3 fl.1 a bc
4      5  2  0
5     12 1  1
6     13 2  1
```

The first argument in the **table()** function is either a factor or a list of factors (two factors in above).

Typically, a dataframe serves as the **table()** argument. Consider the following example of polling-data:

```
1 > ct <- data.frame("Vote.for.X" = c("Yes", "Yes", "No", "Not Sure", "No"),
2 +               "Voted.for.X.Last.Time" = c("Yes", "No", "No", "Yes", "No"))
3 > ct
4   Vote.for.X Voted.for.X.Last.Time
5 1         Yes                   Yes
6 2         Yes                   No
7 3          No                   No
8 4 Not Sure                   Yes
9 5          No                   No
```

The **table()** function computes the **contingency table** for the above dataset.

```
1 > cttab <- table(ct)
2 > cttab
3               Voted.for.X.Last.Time
4 Vote.for.X No Yes
5   No      2  0
6 Not Sure  0  1
7   Yes    1  1
```

One-dimensional counts are equally available; counts on a single factor:

```
1 > table(c(5,12,13,12,8,5))
2
3      5  8 12 13
4      2  1  2  1
```

## matrix/array-like operations on tables in R

Considering that nonmathematical operators can be used on dataframes, the same applies to tables—the *cell counts component of a table object is an array*. Table cell counts can be accessed with matrix notation:

```
1 > class(cttab)
2 [1] "table"
3
4 > cttab[1,1]
5 [1] 2
6
7 > cttab[1,]
8   No Yes
9    2  0
```

Below illustrates multiplicative properties of table operations by converting cell counts to proportions:

```
1 > cttab/5
2           Voted.for.X.Last.Time
3 Vote.for.X No Yes
4   No      0.4 0.0
5   Not Sure 0.0 0.2
6   Yes      0.2 0.2
```

The **marginal values** (statistical method are the values of a variable when the variable itself is held constant while other variables are summed) are computed with the **apply()** function below:

```
1 > apply(cttab, 1, sum)
2   No Not Sure   Yes
3    2     1     2
```

The above illustration computes that marginal values of the polling data ( $1+1=2$ ,  $1+0=1$ ,  $1+1=2$ ).

The **marginal values** both dimensions can be computed simultaneously with the **addmargins()** function.

```
1 > addmargins(cttab)
2           Voted.for.X.Last.Time
3 Vote.for.X No Yes Sum
4   No      2  0  2
5   Not Sure 0  1  1
6   Yes      1  1  2
7   Sum      3  2  5
```

The dimension names of a table can be accessed through the **dimnames()** function as follows:

```
1 > dimnames(cttab)
2 $Vote.for.X
3 [1] "No"      "Not Sure" "Yes"
4
5 $Voted.for.X.Last.Time
6 [1] "No"      "Yes"
```

## other factor- and table-related functions in R

the **aggregate()** function calls **tapply()** once for each variable in a group.

```
1 aggregate(x, ...)
```

The **cut()** function is a common way to generate factors (particularly with tables) by providing data vector **x** and a set of bins defined by vector **b**. The functions determines the bins each element of **x** falls into.

```
1 cut(x, b, labels = FALSE, ...)
```