

Data frame operations

Robert McDonald

2018-08-11

Contents

1	Basic operations on dataframes	1
2	Dataframes vs tibbles	2
2.1	Creating dataframes	2
2.2	Creating tibbles	2
2.3	Subsetting and extraction operators	3
2.4	Extracting a column	3
2.4.1	Extract one column from a dataframe	4
2.4.2	Extract one column from a tibble	4
2.4.3	Discussion	5
2.5	Subsetting vs extraction	6
2.6	Add a column	6
2.7	Delete columns	7
2.8	Add rows	7
2.9	Filter rows	7
	References	8

1 Basic operations on dataframes

When you analyze data in R, you will almost certainly use dataframes, which are rectangular objects that can store multiple types of data at once, subject to the restriction that any given column only contains data of one type (numeric, character, logical, etc.). When performing your analysis, you will undertake basic operations, which can include extracting portions of the dataframe, subsetting the data frame, deleting rows and columns, adding rows and columns, sorting, renaming, summarizing, etc.

There are generally multiple ways to undertake a given dataframe manipulation in base R. The variety can be confusing, but to make things more complicated, Hadley Wickham has introduced tibbles as an alternative to dataframes. As we will see, tibbles standardize and simplify behavior. This necessarily means that they behave differently than data frames in certain respects.

I recommend that you use tibbles. Although tibbles are sometimes annoying, they behave more consistently than data frames and there is less to remember.

Much of this document discusses the closely related concepts of *subsetting* a data frame, and *extracting* from a data frame. Here is what we mean by these terms:

- **subsetting** means creating a new object that has fewer rows or columns, but where the new object has the same class as the old object. A dataframe remains a dataframe, for example.
- **extraction** means creating a *vector* that has a class appropriate to the vector contents (numeric, character, logical, etc.) By definition, if you extract a column from a dataframe or tibble, it is no longer a dataframe or a tibble; it is a vector.

It will be critical to distinguish these two concepts. We will see that an important contribution of tibbles is that subsetting operators only subset, they never extract. This is not true with dataframes and is an important reason to prefer tibbles.

The purpose of this document is to illustrate the different ways of performing various operations, and to illustrate differences between dataframes and tibbles. Even if you prefer one construct and stick to it, the odds are that you will encounter examples that work differently and co-workers who have different preferences. You need to understand both dataframes and tibbles. The good news is that in understanding tibbles, you will come to understand dataframes more deeply.

2 Dataframes vs tibbles

2.1 Creating dataframes

We begin by defining two data frames, named `df0` and `df`, which contain recipe information:¹

```
df0 <- data.frame(what=c('crust', 'cookie', 'cake', 'pasta'),
                  flour=c(3, 3, 1, 3),
                  butter=c(2, 2, 1, 0),
                  liquid=c(1, 0, 0, 0),
                  sugar=c(0, 1, 1, 0),
                  egg=c(0, 0, 1, 2)
                  )

df <- data.frame(what=c('crust', 'cookie', 'cake', 'pasta'),
                 flour=c(3, 3, 1, 3),
                 butter=c(2, 2, 1, 0),
                 liquid=c(1, 0, 0, 0),
                 sugar=c(0, 1, 1, 0),
                 egg=c(0, 0, 1, 2),
                 stringsAsFactors=FALSE
                 )
all.equal(df0, df, check.attributes=FALSE)
[1] "Component \"what\": 'current' is not a factor"
```

By default, the `data.frame` function converts strings to factors. To change this behavior, it is common to set `stringsAsFactors=FALSE`, as when creating `df`. In the rest of this document we will assume that we do *not* want strings automatically converted to factors.

2.2 Creating tibbles

We now create a tibble, `dft`, using the same inputs. The `tibble` function does not convert strings to factors:

```
dft <- tibble(what=c('crust', 'cookie', 'cake', 'pasta'),
              flour=c(3, 3, 1, 3),
              butter=c(2, 2, 1, 0),
              liquid=c(1, 0, 0, 0),
              sugar=c(0, 1, 1, 0),
              egg=c(0, 0, 1, 2)
              )
```

¹This example was inspired by Ruhlman (2009). Ratios are by weight. Any errors related to recipe ratios are mine.

```
all.equal(df, dft, check.attributes=FALSE)
[1] TRUE
```

Although `df` and `dft` contain the same data with the same data types, the `all.equal` function will complain that `df` and `dft` have different classes unless we specify `check.attributes=FALSE`.

The printed data frame looks like this:

```
kable(df)
```

what	flour	butter	liquid	sugar	egg
crust	3	2	1	0	0
cookie	3	2	0	1	0
cake	1	1	0	1	1
pasta	3	0	0	0	2

The printed tibble looks identical to the printed data frame:

```
kable(dft)
```

what	flour	butter	liquid	sugar	egg
crust	3	2	1	0	0
cookie	3	2	0	1	0
cake	1	1	0	1	1
pasta	3	0	0	0	2

When dealing with tibbles it can be helpful to customize the options that govern the way a tibble appears when you print it in the console. You can get a list of available options with `?tibble::tibble-options`. I like the following settings, but you can do something different and you can change them at any time by issuing another `options` command.

```
options(digits=2 ## how many digits to print?
, tibble.width=Inf ## how many tibble columns to print?
, tibble.print_min=10 ## 10 is the default
, pillar.subtle=FALSE ## don't highlight significant digits
)
```

2.3 Subsetting and extraction operators

When dealing with dataframes and tibbles there are three data subsetting and extraction constructs we need to understand:

- “[” Single square brackets are used to *subset*. (There is an important exception to this with dataframes, for which subsetting to a single column also extracts. This does not happen with tibbles.)
- “[[” Double square brackets *extract* a single column at a time
- “\$” The dollar sign also *extracts*, serving the same purpose as double square brackets

2.4 Extracting a column

A common operation with data frames is to create a subset by extracting columns. We will illustrate ways to do this with both dataframes and tibbles.

2.4.1 Extract one column from a dataframe

Suppose we want to extract the flour column from the dataframe `df`:

```
c1 <- df[, 'flour']
c2 <- df$flour
c3 <- df[['flour']]
c4 <- df['flour'][[1]]
class(c1) ## this is a vector
[1] "numeric"
all.equal(c1, c2)
[1] TRUE
all.equal(c2, c3)
[1] TRUE
all.equal(c3, c4)
[1] TRUE
mean(c1)
[1] 2.5
mean(c2)
[1] 2.5
mean(c3)
[1] 2.5
mean(c4)
[1] 2.5
```

Note that in every case we have a numeric vector, and we can compute the mean of this vector.

2.4.2 Extract one column from a tibble

We repeat the previous syntax, only operating on a tibble rather than a dataframe.

```
c1t <- dft[, 'flour'] ## This subsets but does not extract!
c2t <- dft$flour
c3t <- dft[['flour']]
c4t <- dft['flour'][[1]]
class(c1t) ## this is a tibble, not a vector
[1] "tbl_df"      "tbl"        "data.frame"
all.equal(c1t, c2t)
[1] "Cols in y but not x: `c(3, 3, 1, 3)`."
[2] "Cols in x but not y: `flour`."
all.equal(c2t, c3t)
[1] TRUE
all.equal(c3t, c4t)
[1] TRUE
mean(c1t)
Warning in mean.default(c1t): argument is not numeric or logical: returning
NA
[1] NA
mean(c2t)
[1] 2.5
mean(c3t)
[1] 2.5
mean(c4t)
[1] 2.5
```

Because `c1t` is not a numeric vector, we cannot compute its mean. If you wish to compute the mean of a tibble column, you need to convert it to a numeric vector or else use `base::colMeans`, which is designed to operate on data frames.

```
test1 <- dft[, 'flour']
test4 <- dft %>% select(flour) %>% summarize(mean=mean(flour))
mean(test1[[1]])
[1] 2.5
mean(pull(test1))
[1] 2.5
colMeans(test1)
flour
 2.5
test4
# A tibble: 1 x 1
  mean
  <dbl>
1    2.5
test4[[1]]
[1] 2.5
```

Notice that `test4` is a tibble! The dplyr functions return a tibble, even when you have just computed a single number. We can extract the result using `[[1]]`.

2.4.3 Discussion

The last example illustrates an important difference between a dataframe and a tibble. If we extract the flour column from the tibble using the single bracket, the single column remains a tibble. If we extract the flour column from the data frame using the same syntax, the single column is converted to a vector.

Contrast this with extracting the flour column from either `df` or `dft` using double brackets or the `$` synonym, either of which creates a numeric vector. The tibble behavior is, in my opinion, appropriate, because the single bracket is a subsetting operator, not an extraction operator. You might ask why we care. The issue in this case is that the `mean` function (for example) cannot operate directly on either a dataframe or tibble column. It can only operate on a vector. Using either `[[1]]` or `pull` permits us to calculate the mean of `c1t`.² The tibble philosophy is that if you use a subset operator on a tibble, you get a tibble.

Here is another example:

```
onecol <- df[, c('flour')] ## this is a vector
twocol <- df[, c('flour', 'liquid')] ## this is a dataframe
```

You can see that you get a different structure when you extract one or two columns from a dataframe. Extracting a single column gives you a numeric vector, while extracting two columns gives you a dataframe.

Another source of confusion with dataframes is that that single brackets behave differently with a dataframe depending upon whether or not there is a row indicator in the expression:

```
c1 <- df[, 'flour'] ## this is a vector of length 4
c5 <- df['flour']   ## this is a dataframe (list) of length 1
all.equal(c1, c5) ## c1 and c5 differ
```

²There is also an `unlist` function. By default, this creates a named vector, although you can prevent the creation of names by specifying an option, as in `unlist(test1, use.names=FALSE)`. I find `[[1]]` or the `pull` function to be cleaner solutions.

```

[1] "Modes: numeric, list"
[2] "Lengths: 4, 1"
[3] "names for current but not for target"
[4] "Attributes: < target is NULL, current is list >"
[5] "target is numeric, current is data.frame"
all.equal(c1, c5[[1]]) ## the same
[1] TRUE

```

The object `c1` is a vector. The object `c5` is a subset, a new dataframe consisting of the `flour` column from the original dataframe. When you perform operations on a `tibble`, you get a new `tibble`. If you want a vector, you have to ask for a vector by using the extraction function `[[1]]` (which I recommend) or by using the `dplyr` `pull` function‘.

2.5 Subsetting vs extraction

The important thing to remember about extraction is that *you can’t extract multiple columns at once*. The reason is that you can’t know in advance the class of the extracted columns and they may be inconsistent. If column 1 is numeric and column 2 is character, there is no common data structure to which you can extract them, other than a list. Dataframe and tibble are both lists.

Of course you can *subset* multiple columns at once, because you are retaining the data frame or tibble from which subset.

2.6 Add a column

We will add “nuts” and “yeast” columns:

```

df1 <- df2 <- df3 <- df ## create copies of the original data frame
df1$nuts <- c(0, 0.5, 0, 0)
df1$yeast <- 0 ## the recycling rule in action
df2[, "nuts"] <- c(0, 0.5, 0, 0)
df2[, "yeast"] <- 0
df3 <- df %>% mutate(nuts=c(0, 0.5, 0, 0), yeast=0)
all.equal(df1, df2)
[1] TRUE
all.equal(df2, df3)
[1] TRUE
kable(df1)

```

what	flour	butter	liquid	sugar	egg	nuts	yeast
crust	3	2	1	0	0	0.0	0
cookie	3	2	0	1	0	0.5	0
cake	1	1	0	1	1	0.0	0
pasta	3	0	0	0	2	0.0	0

Very important: Notice that when using `dplyr` to construct `df3`, no quotes are needed. When using base R, quotes are needed for column names. The absence of quotation within the `dplyr` universe is a helpful simplification, but it is easy to become confused if you mix `dplyr` and base R. My advice is to stick with one or the other as much as possible.

2.7 Delete columns

Now we decide we don't need the columns we just added:

```
df4 <- df1;
df4[c('nuts', 'yeast')] <- NULL ## assigning to `NULL` deletes an object
df5 <- df1[-c(7, 8)]
df6 <- df1[~which(names(df1) %in% c('nuts', 'yeast'))]
df7 <- df1[, ~which(names(df1) %in% c('nuts', 'yeast'))]
df8 <- df1 %>% select(-nuts, -yeast)
all.equal(df4, df5)
[1] TRUE
all.equal(df4, df6)
[1] TRUE
all.equal(df4, df7)
[1] TRUE
all.equal(df4, df8)
[1] TRUE
```

2.8 Add rows

We will use the dataframe including nut and yeast columns. Notice that we get the same result whether we define the new row as a `list` or as a `data.frame`. This occurs because a `data.frame` is a `list`!

```
newrow1 <- list('bread', 1, 0, .67, 0, 0, 0, .02)
newrow2 <- list(what='bread', flour=1, butter=0,
               liquid=.67, sugar=0, egg=0, nuts=0, yeast=.02)
newrow3 <- data.frame(what='bread', flour=1, butter=0,
                     liquid=.67, sugar=0, egg=0, nuts=0, yeast=.02)
df.ar1 <- rbind(df1, newrow1)
df.ar2 <- rbind(df1, newrow2)
df.ar3 <- bind_rows(df1, newrow2) ## dplyr
df.ar4 <- bind_rows(df1, newrow3)
Warning in bind_rows(x, .id): binding character and factor vector,
coercing into character vector
all.equal(df.ar1, df.ar2)
[1] TRUE
all.equal(df.ar1, df.ar3)
[1] TRUE
all.equal(df.ar1, df.ar4)
[1] TRUE
```

2.9 Filter rows

We can choose to keep or delete rows that meet specific criteria. We will use `df.ar1` as the base dataframe.

Suppose we want only items that use butter

```
df.butter1 <- df.ar1[df.ar1$butter > 0, ]
df.butter2 <- df.ar1[df.ar1$"butter" > 0, ]
df.butter3 <- subset(df.ar1, df.ar1$butter > 0)
df.butter4 <- filter(df.ar1, butter > 0)
```

```
all.equal(df.butter1, df.butter2)
[1] TRUE
all.equal(df.butter1, df.butter3)
[1] TRUE
all.equal(df.butter1, df.butter4)
[1] TRUE
```

Suppose we want items that use butter and liquid. Any expression evaluating to a logical will work, so compound conditions will work. When using dplyr we can separate compound conditions either with a comma or by using an ampersand:

```
df.butterliquid1 <- df.ar1[df.ar1$butter & df.ar1$liquid > 0, ]
df.butterliquid2 <- filter(df.ar1, butter > 0, liquid > 0)
df.butterliquid3 <- filter(df.ar1, butter > 0 & liquid > 0)
all.equal(df.butterliquid1, df.butterliquid2)
[1] TRUE
all.equal(df.butterliquid1, df.butterliquid3)
[1] TRUE
```

References

Ruhlman, Michael. 2009. *Ratio: The Simple Codes Behind the Craft of Everyday Cooking*. Scribner.