

Manipulating and analyzing data with dplyr; Exporting data

Data Carpentry contributors

```
## Warning: package 'tidyverse' was built under R version 3.4.2
## Warning: package 'tidyr' was built under R version 3.4.2
## Warning: package 'purrr' was built under R version 3.4.2
## Warning: package 'dplyr' was built under R version 3.4.2
```

Manipulating and analyzing data with dplyr

Learning Objectives

- Describe the purpose of the **dplyr** and **tidyr** packages.
 - Select certain columns in a data frame with the **dplyr** function **select**.
 - Select certain rows in a data frame according to filtering conditions with the **dplyr** function **filter**.
 - Link the output of one **dplyr** function to the input of another function with the ‘pipe’ operator **%>%**.
 - Add new columns to a data frame that are functions of existing columns with **mutate**.
 - Use the split-apply-combine concept for data analysis.
 - Use **summarize**, **group_by**, and **tally** to split a data frame into groups of observations, apply a summary statistics for each group, and then combine the results.
 - Describe the concept of a wide and a long table format and for which purpose those formats are useful.
 - Describe what key-value pairs are.
 - Reshape a data frame from long to wide format and back with the **spread** and **gather** commands from the **tidyr** package.
 - Export a data frame to a .csv file.
-

Data Manipulation using dplyr and tidyr

Bracket subsetting is handy, but it can be cumbersome and difficult to read, especially for complicated operations. Enter **dplyr**. **dplyr** is a package for making tabular data manipulation easier. It pairs nicely with **tidyr** which enables you to swiftly convert between different data formats for plotting and analysis.

Packages in R are basically sets of additional functions that let you do more stuff. The functions we’ve been using so far, like **str()** or **data.frame()**, come built into R; packages give you access to more of them. Before you use a package for the first time you need to install it on your machine, and then you should import it in every subsequent R session when you need it. You should already have installed the **tidyverse** package. This is an “umbrella-package” that installs several packages useful for data analysis which work together well such as **tidyr**, **dplyr**, **ggplot2**, **tibble**, etc.

The **tidyverse** package tries to address 3 major problems with some of base R functions: 1. The results from a base R function sometimes depend on the type of data. 2. Using R expressions in a non standard way, which can be confusing for new learners. 3. Hidden arguments, having default operations that new learners are not aware of.

We have seen in our previous lesson that when building or importing a data frame, the columns that contain characters (i.e., text) are coerced (=converted) into the factor data type. We had to set **stringsAsFactors** to **FALSE** to avoid this hidden argument to convert our data type.

This time will use the **tidyverse** package to read the data and avoid having to set **stringsAsFactors** to **FALSE**

To load the package type:

```
library("tidyverse")    ## load the tidyverse packages, incl. dplyr
```

What are dplyr and tidyr?

The package **dplyr** provides easy tools for the most common data manipulation tasks. It is built to work directly with data frames, with many common tasks optimized by being written in a compiled language (C++). An additional feature is the ability to work directly with data stored in an external database. The benefits of doing this are that the data can be managed natively in a relational database, queries can be conducted on that database, and only the results of the query are returned.

This addresses a common problem with R in that all operations are conducted in-memory and thus the amount of data you can work with is limited by available memory. The database connections essentially remove that limitation in that you can connect to a database of many hundreds of GB, conduct queries on it directly, and pull back into R only what you need for analysis.

The package **tidyr** addresses the common problem of wanting to reshape your data for plotting and use by different R functions. Sometimes we want data sets where we have one row per measurement. Sometimes we want a data frame where each measurement type has its own column, and rows are instead more aggregated groups - like plots or aquaria. Moving back and forth between these formats is nontrivial, and **tidyr** gives you tools for this and more sophisticated data manipulation.

To learn more about **dplyr** and **tidyr** after the workshop, you may want to check out this handy data transformation with **dplyr** cheatsheet and this one about **tidyr**.

We'll read in our data using the **read_csv()** function, from the tidyverse package **readr**, instead of **read.csv()**.

```
surveys <- read_csv("data/portal_data_joined.csv")
```

```
#> Parsed with column specification:
#> cols(
#>   record_id = col_integer(),
#>   month = col_integer(),
#>   day = col_integer(),
#>   year = col_integer(),
#>   plot_id = col_integer(),
#>   species_id = col_character(),
#>   sex = col_character(),
#>   hindfoot_length = col_integer(),
#>   weight = col_integer(),
#>   genus = col_character(),
#>   species = col_character(),
#>   taxa = col_character(),
#>   plot_type = col_character()
```

```
#> )  
## inspect the data  
str(surveys)
```

Notice that the class of the data is now `tbl_df`. This is referred to as a “tibble”. Tibbles are almost identical to R’s standard data frames, but they tweak some of the old behaviors of data frames. For our purposes the only differences between data frames and tibbles are that:

1. When you print a tibble, R displays the data type of each column under its name; it prints only the first few rows of data and only as many columns as fit on one screen.
2. Columns of class `character` are never automatically converted into factors.

Selecting columns and filtering rows

We’re going to learn some of the most common `dplyr` functions: `select()`, `filter()`, `mutate()`, `group_by()`, and `summarize()`. To select columns of a data frame, use `select()`. The first argument to this function is the data frame (`surveys`), and the subsequent arguments are the columns to keep.

```
select(surveys, plot_id, species_id, weight)
```

To choose rows based on a specific criteria, use `filter()`:

```
filter(surveys, year == 1995)
```

Pipes

What if you want to select and filter at the same time? There are three ways to do this: use intermediate steps, nested functions, or pipes.

With intermediate steps, you create a temporary data frame and use that as input to the next function, like this:

```
surveys2 <- filter(surveys, weight < 5)  
surveys_sml <- select(surveys2, species_id, sex, weight)
```

This is readable, but can clutter up your workspace with lots of objects that you have to name individually. With multiple steps, that can be hard to keep track of.

You can also nest functions (i.e. one function inside of another), like this:

```
surveys_sml <- select(filter(surveys, weight < 5), species_id, sex, weight)
```

This is handy, but can be difficult to read if too many functions are nested, as R evaluates the expression from the inside out (in this case, filtering, then selecting).

The last option, *pipes*, are a recent addition to R. Pipes let you take the output of one function and send it directly to the next, which is useful when you need to do many things to the same dataset. Pipes in R look like `%>%` and are made available via the `magrittr` package, installed automatically with `dplyr`. If you use RStudio, you can type the pipe with `Ctrl + Shift + M` if you have a PC or `Cmd + Shift + M` if you have a Mac.

```
surveys %>%  
  filter(weight < 5) %>%  
  select(species_id, sex, weight)
```

In the above code, we use the pipe to send the `surveys` dataset first through `filter()` to keep rows where `weight` is less than 5, then through `select()` to keep only the `species_id`, `sex`, and `weight` columns. Since

%>% takes the object on its left and passes it as the first argument to the function on its right, we don't need to explicitly include the data frame as an argument to the `filter()` and `select()` functions any more.

Some may find it helpful to read the pipe like the word “then”. For instance, in the above example, we took the data frame `surveys`, *then* we *filtered* for rows with `weight < 5`, *then* we *selected* columns `species_id`, `sex`, and `weight`. The `dplyr` functions by themselves are somewhat simple, but by combining them into linear workflows with the pipe, we can accomplish more complex manipulations of data frames.

If we want to create a new object with this smaller version of the data, we can assign it a new name:

```
surveys_sml <- surveys %>%  
  filter(weight < 5) %>%  
  select(species_id, sex, weight)  
  
surveys_sml
```

Note that the final data frame is the leftmost part of this expression.

Challenge

Using pipes, subset the `surveys` data to include individuals collected before 1995 and retain only the columns `year`, `sex`, and `weight`.

Mutate

Frequently you'll want to create new columns based on the values in existing columns, for example to do unit conversions, or to find the ratio of values in two columns. For this we'll use `mutate()`.

To create a new column of weight in kg:

```
surveys %>%  
  mutate(weight_kg = weight / 1000)
```

You can also create a second new column based on the first new column within the same call of `mutate()`:

```
surveys %>%  
  mutate(weight_kg = weight / 1000,  
         weight_kg2 = weight_kg * 2)
```

If this runs off your screen and you just want to see the first few rows, you can use a pipe to view the `head()` of the data. (Pipes work with non-`dplyr` functions, too, as long as the `dplyr` or `magrittr` package is loaded).

```
surveys %>%  
  mutate(weight_kg = weight / 1000) %>%  
  head()
```

The first few rows of the output are full of NAs, so if we wanted to remove those we could insert a `filter()` in the chain:

```
surveys %>%  
  filter(!is.na(weight)) %>%  
  mutate(weight_kg = weight / 1000) %>%  
  head()
```

`is.na()` is a function that determines whether something is an NA. The `!` symbol negates the result, so we're asking for every row where `weight` *is not* an NA.

Challenge

Create a new data frame from the `surveys` data that meets the following criteria: contains only the `species_id` column and a new column called `hindfoot_half` containing values that are half the `hindfoot_length` values. In this `hindfoot_half` column, there are no NAs and all values are less than 30.

Hint: think about how the commands should be ordered to produce this data frame!

Split-apply-combine data analysis and the `summarize()` function

Many data analysis tasks can be approached using the *split-apply-combine* paradigm: split the data into groups, apply some analysis to each group, and then combine the results. `dplyr` makes this very easy through the use of the `group_by()` function.

The `summarize()` function

`group_by()` is often used together with `summarize()`, which collapses each group into a single-row summary of that group. `group_by()` takes as arguments the column names that contain the **categorical** variables for which you want to calculate the summary statistics. So to compute the mean `weight` by sex:

```
surveys %>%
  group_by(sex) %>%
  summarize(mean_weight = mean(weight, na.rm = TRUE))
```

You may also have noticed that the output from these calls doesn't run off the screen anymore. It's one of the advantages of `tbl_df` over data frame.

You can also group by multiple columns:

```
surveys %>%
  group_by(sex, species_id) %>%
  summarize(mean_weight = mean(weight, na.rm = TRUE))
```

When grouping both by `sex` and `species_id`, the first rows are for individuals that escaped before their sex could be determined and weighted. You may notice that the last column does not contain NA but NaN (which refers to "Not a Number"). To avoid this, we can remove the missing values for weight before we attempt to calculate the summary statistics on weight. Because the missing values are removed first, we can omit `na.rm = TRUE` when computing the mean:

```
surveys %>%
  filter(!is.na(weight)) %>%
  group_by(sex, species_id) %>%
  summarize(mean_weight = mean(weight))
```

Here, again, the output from these calls doesn't run off the screen anymore. If you want to display more data, you can use the `print()` function at the end of your chain with the argument `n` specifying the number of rows to display:

```
surveys %>%
  filter(!is.na(weight)) %>%
  group_by(sex, species_id) %>%
  summarize(mean_weight = mean(weight)) %>%
  print(n = 15)
```

Once the data are grouped, you can also summarize multiple variables at the same time (and not necessarily on the same variable). For instance, we could add a column indicating the minimum weight for each species for each sex:

```
surveys %>%
  filter(!is.na(weight)) %>%
  group_by(sex, species_id) %>%
  summarize(mean_weight = mean(weight),
            min_weight = min(weight))
```

Tallying

When working with data, we often want to know the number of observations found for each factor or combination of factors. For this task, **dplyr** provides `tally()`. For example, if we wanted to group by sex and find the number of rows of data for each sex, we would do:

```
surveys %>%
  group_by(sex) %>%
  tally()
```

Here, `tally()` is the action applied to the groups created by `group_by()` and counts the total number of records for each category.

Challenge

1. How many individuals were caught in each `plot_type` surveyed?
2. Use `group_by()` and `summarize()` to find the mean, min, and max hindfoot length for each species (using `species_id`).
3. What was the heaviest animal measured in each year? Return the columns `year`, `genus`, `species_id`, and `weight`.
4. You saw above how to count the number of individuals of each `sex` using a combination of `group_by()` and `tally()`. How could you get the same result using `group_by()` and `summarize()`? Hint: see `?n`.

Reshaping with gather and spread

In the spreadsheet lesson we discussed how to structure our data leading to the four rules defining a tidy dataset:

1. Each variable has its own column
2. Each observation has its own row
3. Each value must have its own cell
4. Each type of observational unit forms a table

Here we examine the fourth rule: Each type of observational unit forms a table.

In `surveys`, the rows of `surveys` contain the values of variables associated with each record (the unit), values such the weight or sex of each animal associated with each record. What if instead of comparing records, we wanted to compare the different mean weight of each species between plots? (Ignoring `plot_type` for simplicity).

We'd need to create a new table where each row (the unit) is comprised of values of variables associated with each plot. In practical terms this means the values of the species in `genus` would become the names of column variables and the cells would contain the values of the mean weight observed on each plot.

Having created a new table, it is therefore straightforward to explore the relationship between the weight of different species within, and between, the plots. The key point here is that we are still following a tidy data structure, but we have **reshaped** the data according to the observations of interest: average species weight per plot instead of recordings per date.

The opposite transformation would be to transform column names into values of a variable.

We can do both these of transformations with two `tidyr` functions, `spread()` and `gather()`.

Spreading

`spread()` takes three principal arguments:

1. the data
2. the *key* column variable whose values will become new column names.
3. the *value* column variable whose values will fill the new column variables.

Further arguments include `fill` which, if set, fills in missing values with the value provided.

Let's use `spread()` to transform surveys to find the mean weight of each species in each plot over the entire survey period. We use `filter()`, `group_by()` and `summarise()` to filter our observations and variables of interest, and create a new variable for the `mean_weight`. We use the pipe as before too.

```
surveys_gw <- surveys %>%  
  filter(!is.na(weight)) %>%  
  group_by(genus, plot_id) %>%  
  summarize(mean_weight = mean(weight))  
  
str(surveys_gw)
```

This yields `surveys_gw` where the observations for each plot are spread across multiple rows, 196 observations of 13 variables. Using `spread()` to key on `genus` with values from `mean_weight` this becomes 24 observations of 11 variables, one row for each plot. We again use pipes:

```
surveys_spread <- surveys_gw %>%  
  spread(key = genus, value = mean_weight)  
  
str(surveys_spread)
```

We could now plot comparisons between the weight of species in different plots, although we may wish to fill in the missing values first.

```
surveys_gw %>%  
  spread(genus, mean_weight, fill = 0) %>%  
  head
```

Gathering

The opposing situation could occur if we had been provided with data in the form of `surveys_spread`, where the genus names are column names, but we wish to treat them as values of a genus variable instead.

In this situation we are gathering the column names and turning them into a pair of new variables. One variable represents the column names as values, and the other variable contains the values previously associated with the column names.

`gather()` takes four principal arguments:

1. the data
2. the *key* column variable we wish to create from column names.

3. the *values* column variable we wish to create and fill with values associated with the key.
4. the names of the columns we use to fill the key variable (or to drop).

To recreate `surveys_gw` from `surveys_spread` we would create a key called `genus` and value called `mean_weight` and use all columns except `plot_id` for the key variable. Here we drop `plot_id` column with a minus sign.

```
surveys_gather <- surveys_spread %>%
  gather(key = genus, value = mean_weight, -plot_id)

str(surveys_gather)
```

Note that now the NA genera are included in the re-gathered format. Spreading and then gathering can be a useful way to balance out a dataset so every replicate has the same composition.

We could also have used a specification for what columns to include. This can be useful if you have a large number of identifying columns, and it's easier to specify what to gather than what to leave alone. And if the columns are in a row, we don't even need to list them all out - just use the `:` operator!

```
surveys_spread %>%
  gather(key = genus, value = mean_weight, Baiomys:Spermophilus) %>%
  head
```

Challenge

1. Spread the `surveys` data frame with `year` as columns, `plot_id` as rows, and the number of genera per plot as the values. You will need to summarize before reshaping, and use the function `n_distinct()` to get the number of unique genera within a particular chunk of data. It's a powerful function! See `?n_distinct` for more.
2. Now take that data frame and `gather()` it again, so each row is a unique `plot_id` by `year` combination.
3. The `surveys` data set has two measurement columns: `hindfoot_length` and `weight`. This makes it difficult to do things like look at the relationship between mean values of each measurement per year in different plot types. Let's walk through a common solution for this type of problem. First, use `gather()` to create a dataset where we have a key column called `measurement` and a `value` column that takes on the value of either `hindfoot_length` or `weight`. *Hint:* You'll need to specify which columns are being gathered.
4. With this new data set, calculate the average of each `measurement` in each `year` for each different `plot_type`. Then `spread()` them into a data set with a column for `hindfoot_length` and `weight`. *Hint:* You only need to specify the key and value columns for `spread()`.

Exporting data

Now that you have learned how to use `dplyr` to extract information from or summarize your raw data, you may want to export these new data sets to share them with your collaborators or for archival.

Similar to the `read_csv()` function used for reading CSV files into R, the tidyverse includes a `write_csv()` function that generates CSV files from data frames.

Before using `write_csv()`, we are going to create a new folder, `data_output`, in our working directory that will store this generated data set. We don't want to write generated data sets in the same directory as our raw data. It's good practice to keep them separate. The `data` folder should only contain the raw, unaltered data, and should be left alone to make sure we don't delete or modify it. In contrast, our script will generate the

contents of the `data_output` directory, so even if the files it contains are deleted, we can always re-generate them.

In preparation for our next lesson on plotting, we are going to prepare a cleaned up version of the data set that doesn't include any missing data.

Let's start by removing observations for which the `species_id` is missing. In this data set, the missing species are represented by an empty string and not an NA. Let's also remove observations for which `weight` and the `hindfoot_length` are missing. This data set should also only contain observations of animals for which the sex has been determined:

```
surveys_complete <- surveys %>%  
  filter(species_id != "",           # remove missing species_id  
         !is.na(weight),           # remove missing weight  
         !is.na(hindfoot_length),  # remove missing hindfoot_length  
         sex != "")                # remove missing sex
```

Because we are interested in plotting how species abundances have changed through time, we are also going to remove observations for rare species (i.e., that have been observed less than 50 times). We will do this in two steps: first we are going to create a data set that counts how often each species has been observed, and filter out the rare species; then, we will extract only the observations for these more common species:

```
## Extract the most common species_id  
species_counts <- surveys_complete %>%  
  group_by(species_id) %>%  
  tally() %>%  
  filter(n >= 50)  
  
## Only keep the most common species  
surveys_complete <- surveys_complete %>%  
  filter(species_id %in% species_counts$species_id)
```

To make sure that everyone has the same data set, check that `surveys_complete` has 30463 rows and 13 columns by typing `dim(surveys_complete)`.

Now that our data set is ready, we can save it as a CSV file in our `data_output` folder.

```
write_csv(surveys_complete, path = "data_output/surveys_complete.csv")
```