Chapter 4 The tidyverse

Up to now we have been manipulating vectors by reordering and subsetting them through indexing. However, once we start more advanced analyses, the preferred unit for data storage is not the vector but the data frame. In this chapter we learn to work directly with data frames, which greatly facilitate the organization of information. We will be using data frames for the majority of this book. We will focus on a specific data format referred to as *tidy* and on specific collection of packages that are particularly helpful for working with *tidy* data referred to as the *tidyverse*.

We can load all the tidyverse packages at once by installing and loading the tidyverse package:

library(tidyverse)

We will learn how to implement the tidyverse approach throughout the book, but before delving into the details, in this chapter we introduce some of the most widely used tidyverse functionality, starting with the **dplyr** package for manipulating data frames and the **purrr** package for working with <u>functions</u>. Note that the tidyverse also includes a graphing package, **ggplot2**, which we introduce later in Chapter 7 in the Data Visualization part of the book; the **readr** package discussed in Chapter 5; and many others. In this chapter, we first introduce the concept of *tidy data* and then demonstrate how we use the tidyverse to work with data frames in this format.

4.1 Tidy data

We say that a data table is in *tidy* format if each row represents one observation and columns represent the different variables available for each of these observations. The murders dataset is an example of a tidy data frame.

#>	state	abb	region	population	total
#> 1	Alabama	AL	South	4779736	135
#> 2	Alaska	AK	West	710231	19
# > 3	Arizona	ΑZ	West	6392017	232
# > 4	Arkansas	AR	South	2915918	93
#> 5	California	CA	West	37253956	1257
#> 6	Colorado	C0	West	5029196	65

Each row represent a state with each of the five columns providing a different variable related to these states: name, abbreviation, region, population, and total murders.

To see how the same information can be provided in different formats, consider the following example:

#>	country	year	fertility	
#> 1	Germany	1960	2.41	
#> 2	South Korea	1960	6.16	This is a great example.
# > 3	Germany	1961	2.44	see below and consider the difference.
# > 4	South Korea	1961	5.99	
#> 5	Germany	1962	2.47	
<i>#</i> > 6	South Korea	1962	5.79	

This tidy dataset provides fertility rates for two countries across the years. This is a tidy dataset because each row presents one observation with the three variables being country, year, and fertility rate. However, this dataset originally came in another format and was reshaped for the **dslabs** package. Originally, the data was in the following format:

#> country 1960 1961 1962
#> 1 Germany 2.41 2.44 2.47
#> 2 South Korea 6.16 5.99 5.79

The same information is provided, but there are two important differences in the format: 1) each row includes several observations and 2) one of the variables, year, is stored in the header. For the tidyverse packages to be optimally used, data need to be reshaped into tidy format,

which you will learn to do in the Data Wrangling part of the book. Until then, we will use example datasets that are already in tidy format.

Although not immediately obvious, as you go through the book you will start to appreciate the advantages of working in a framework in which functions use tidy formats for both inputs and outputs. You will see how this permits the data analyst to focus on more important aspects of the analysis rather than the format of the data.

4.2 Exercises

- 1. Examine the built-in dataset co2. Which of the following is true:
 - a. co2 is tidy data: it has one year for each row.
 - b. co2 is not tidy: we need at least one column with a character vector.
 - c. co2 is not tidy: it is a matrix instead of a data frame.
 - d. co2 is not tidy: to be tidy we would have to wrangle it to have three columns (year, month and value), then each co2 observation would have a row.
- 2. Examine the built-in dataset ChickWeight . Which of the following is true:
 - a. ChickWeight is not tidy: each chick has more than one row.
 - b. ChickWeight is tidy: each observation (a weight) is represented by one row. The chick from which this measurement came is one of the variables.
 - c. ChickWeight is not tidy: we are missing the year column.
 - d. ChickWeight is tidy: it is stored in a data frame.
- 3. Examine the built-in dataset BOD . Which of the following is true:
 - a. BOD is not tidy: it only has six rows.
 - b. BOD is not tidy: the first column is just an index.
 - c. BOD is tidy: each row is an observation with two values (time and demand)
 - d. BOD is tidy: all small datasets are tidy by definition.
- 4. Which of the following built-in datasets is tidy (you can pick more than one):
 - a. BJsales
 - b. EuStockMarkets
 - c. DNase
 - d. Formaldehyde

- e. Orange
- f. UCBAdmissions

4.3 Manipulating data frames

The **dplyr** package from the **tidyverse** introduces functions that perform some of the most common operations when working with data frames and uses names for these functions that are relatively easy to remember. For instance, to change the data table by adding a new column, we use <u>mutate</u>. To filter the data table to a subset of rows, we use <u>filter</u>. Finally, to subset the data by selecting specific columns, we use <u>select</u>.

4.3.1 Adding a column with mutate

We want all the necessary information for our analysis to be included in the data table. So the first task is to add the murder rates to our murders data frame. The function mutate takes the data frame as a first argument and the name and values of the variable as a second argument using the convention name = values . So, to add murder rates, we use:

```
library(dslabs)
data("murders")
murders <- mutate(murders, rate = total / population * 100000)</pre>
```

Notice that here we used total and population inside the function, which are objects that are **not** defined in our workspace. But why don't we get an error?

This is one of **dplyr**'s main features. Functions in this package, such as mutate, know to look for variables in the data frame provided in the first argument. In the call to mutate above, total will have the values in murders\$total. This approach makes the code much more readable.

We can see that the new column is added:

<pre>head(murders)</pre>									
#>		state	abb	region	population	total rate			
#>	1	Alabama	AL	South	4779736	135 2.82			
#>	2	Alaska	AK	West	710231	19 2.68			
#>	3	Arizona	ΑZ	West	6392017	232 3.63			
#>	4	Arkansas	AR	South	2915918	93 3.19			
#>	5	California	СА	West	37253956	1257 3 . 37			
#>	6	Colorado	С0	West	5029196	65 1.29			

Although we have overwritten the original murders object, this does not change the object that loaded with data(murders). If we load the murders data again, the original will overwrite our mutated version.

4.3.2 Subsetting with filter

Now <u>suppose that we want to filter the data table to only show the entries for which the murder</u> <u>rate is lower than 0.71.</u> To do this we use the filter function, which takes the data table as the first argument and then the conditional statement as the second. Like mutate , we can use the unquoted variable names from murders inside the function and it will know we mean the columns and not objects in the workspace.

filter(murders, rate <= 0.71)</pre>

#>		state	abb	region	population	total	rate
#>	1	Hawaii	HI	West	1360301	7	0.515
#>	2	Iowa	IA	North Central	3046355	21	0.689
#>	3	New Hampshire	NH	Northeast	1316470	5	0.380
#>	4	North Dakota	ND	North Central	672591	4	0.595
#>	5	Vermont	VT	Northeast	625741	2	0.320

4.3.3 Selecting columns with select

Although our data table only has six columns, some data tables include hundreds. If we want to view just a few, we can use the **dplyr** select function. In the code below we select three columns, assign this to a new object and then filter the new object:

```
new_table <- select(murders, state, region, rate)</pre>
filter(new_table, rate <= 0.71)</pre>
#>
             state
                           region rate
#> 1
                             West 0.515
            Hawaii
#> 2
              Iowa North Central 0.689
#> 3 New Hampshire
                        Northeast 0.380
#> 4 North Dakota North Central 0.595
#> 5
           Vermont
                        Northeast 0.320
```

In the call to select, the first argument murders is an object, but state, region, and rate are variable names.

4.4 Exercises

1. Load the **dplyr** package and the murders dataset.

```
library(dplyr)
library(dslabs)
data(murders)
```

You can add columns using the **dplyr** function mutate. This function is aware of the column names and inside the function you can call them unquoted:

murders <- mutate(murders, population_in_millions = population / 10^6)</pre>

We can write population rather than murders\$population. The function mutate knows we are grabbing columns from murders.

Use the function mutate to add a murders column named rate with the per 100,000 murder rate as in the example code above. Make sure you redefine murders as done in the example code above (murders <- [your code]) so we can keep using this variable.

2. If rank(x) gives you the ranks of x from lowest to highest, rank(-x) gives you the ranks from highest to lowest. Use the function mutate to add a column rank containing the rank, from highest to lowest murder rate. Make sure you redefine murders so we can keep using this https://rafalab.github.io/dsbook/tidyverse.html

variable.

3. With **dplyr**, we can use select to show only certain columns. For example, with this code we would only show the states and population sizes:

```
select(murders, state, population) %>% head()
```

Use select to show the state names and abbreviations in murders. Do not redefine murders, just show the results.

4. The **dplyr** function filter is used to choose specific rows of the data frame to keep. Unlike select which is for columns, filter is for rows. For example, you can show just the New York row like this:

```
filter(murders, state == "New York")
```

You can use other logical vectors to filter rows.

Use filter to show the top 5 states with the highest murder rates. After we add murder rate and rank, do not change the murders dataset, just show the result. Remember that you can filter based on the rank column.

5. We can remove rows using the != operator. For example, to remove Florida, we would do this:

no_florida <- filter(murders, state != "Florida")</pre>

Create a new data frame called no_south that removes states from the South region. How many states are in this category? You can use the function nrow for this.

6. We can also use %in% to filter with **dplyr**. You can therefore see the data from New York and Texas like this:

```
filter(murders, state %in% c("New York", "Texas"))
```

Create a new data frame called murders_nw with only the states from the Northeast and the West. How many states are in this category?

7. Suppose you want to live in the Northeast or West **and** want the murder rate to be less than 1. We want to see the data for the states satisfying these options. Note that you can use logical operators with filter. Here is an example in which we filter to keep only small states in the Northeast region.

filter(murders, population < 5000000 & region == "Northeast")</pre>

Make sure murders has been defined with rate and rank and still has all states. Create a table called my_states that contains rows for states satisfying both the conditions: it is in the Northeast or West and the murder rate is less than 1. Use select to show only the state name, the rate, and the rank.

4.5 The pipe: %>%

With **dplyr** we can perform a series of operations, for example select and then filter, by sending the results of one function to another using what is called the **pipe operator**: %>%. Some details are included below.

We wrote code above to show three variables (state, region, rate) for states that have murder rates below 0.71. To do this, we defined the intermediate object <code>new_table</code> . In **dplyr** we can write code that looks more like a description of what we want to do without intermediate objects:

original data \rightarrow select \rightarrow filter

For such an operation, we can use the pipe %>%. The code looks like this:

```
murders %>% select(state, region, rate) %>% filter(rate <= 0.71)</pre>
#>
             state
                           region rate
#> 1
            Hawaii
                             West 0.515
#> 2
              Iowa North Central 0.689
#> 3 New Hampshire
                       Northeast 0.380
      North Dakota North Central 0,595
#> 4
                       Northeast 0,320
#> 5
           Vermont
```

This line of code is equivalent to the two lines of code above. What is going on here?

In general, the pipe **sends** the result of the left side of the pipe to be the first argument of the function on the right side of the pipe. Here is a very simple example:

```
16 %>% sqrt()
#> [1] 4
```

We can continue to pipe values along:

```
16 %>% sqrt() %>% log2()
#> [1] 2
```

The above statement is equivalent to log2(sqrt(16)).

Remember that the pipe sends values to the first argument, so we can define other arguments as if the first argument is already defined:

```
16 %>% sqrt() %>% log(base = 2)
#> [1] 2
```

Therefore, when using the pipe with data frames and **dplyr**, we no longer need to specify the required first argument since the **dplyr** functions we have described all take the data as the first argument. In the code we wrote:

```
murders %>% select(state, region, rate) %>% filter(rate <= 0.71)</pre>
```

murders is the first argument of the select function, and the new data frame (formerly new_table) is the first argument of the filter function.

Note that the pipe works well with functions where the first argument is the input data. Functions in **tidyverse** packages like **dplyr** have this format and can be used easily with the pipe.

4.6 Exercises

1. The pipe %>% can be used to perform operations sequentially without having to define intermediate objects. Start by redefining murder to include rate and rank.

In the solution to the previous exercise, we did the following:

```
my_states <- filter(murders, region %in% c("Northeast", "West") &
    rate < 1)</pre>
```

```
select(my_states, state, rate, rank)
```

The pipe %>% permits us to perform both operations sequentially without having to define an intermediate variable my_states. We therefore could have mutated and selected in the same line like this:

Notice that select no longer has a data frame as the first argument. The first argument is assumed to be the result of the operation conducted right before the %>%.

Repeat the previous exercise, but now instead of creating a new object, show the result and only include the state, rate, and rank columns. Use a pipe %>% to do this in just one line.

2. Reset murders to the original table by using data(murders). Use a pipe to create a new data frame called my_states that considers only states in the Northeast or West which have a murder rate lower than 1, and contains only the state, rate and rank columns. The pipe should also have four components separated by three %>%. The code should look something like this:

```
my_states <- murders %>%
mutate SOMETHING %>%
filter SOMETHING %>%
select SOMETHING
```

4.7 Summarizing data

An important part of exploratory data analysis is summarizing data. The average and standard deviation are two examples of widely used summary statistics. More informative summaries can often be achieved by first splitting data into groups. In this section, we cover two new **dplyr** verbs that make these computations easier: <u>summarize</u> and <u>group_by</u>. We learn to access resulting values using the pull function.

4.7.1 summarize

The summarize function in **dplyr** provides a way to compute summary statistics with intuitive and readable code. We start with a simple example based on heights. The heights dataset includes heights and sex reported by students in an in-class survey.

```
library(dplyr)
library(dslabs)
data(heights)
```

The following code computes the average and standard deviation for females:

```
s <- heights %>%
filter(sex == "Female") %>%
summarize(average = mean(height), standard_deviation = sd(height))
s
#> average standard_deviation
#> 1 64.9 3.76
```

This takes our original data table as input, filters it to keep only females, and then produces a new summarized table with just the average and the standard deviation of heights. We get to choose the names of the columns of the resulting table. For example, above we decided to use average and standard_deviation, but we could have used other names just the same.

Because the resulting table stored in s is a data frame, we can access the components with the accessor \$:

```
s$average
#> [1] 64.9
s$standard_deviation
#> [1] 3.76
```

As with most other **dplyr** functions, summarize is aware of the variable names and we can use them directly. So when inside the call to the summarize function we write mean(height), the function is accessing the column with the name "height" and then computing the average of the resulting numeric vector. We can compute any other summary that operates on vectors and returns a single value. For example, we can add the median, minimum, and maximum heights like this:

We can obtain these three values with just one line using the quantile function: for example, quantile(x, c(0,0.5,1)) returns the min (0th percentile), median (50th percentile), and max (100th percentile) of the vector x. However, if we attempt to use a function like this that returns two or more values inside summarize :

```
heights %>%
filter(sex == "Female") %>%
summarize(range = quantile(height, c(0, 0.5, 1)))
```

we will receive an error: Error: expecting result of length one, got : 2. With the function summarize, we can only call functions that return a single value. In Section 4.12, we will learn how to deal with functions that return more than one value.

For another example of how we can use the summarize function, let's compute the average murder rate for the United States. Remember our data table includes total murders and population size for each state and we have already used **dplyr** to add a murder rate column:

```
murders <- murders %>% mutate(rate = total/population*100000)
```

Remember that the US murder rate is **not** the average of the state murder rates:

```
summarize(murders, mean(rate))
#> mean(rate)
#> 1 2.78
```

This is because in the computation above the small states are given the same weight as the large ones. The US murder rate is the total number of murders in the US divided by the total US population. So the correct computation is:

```
us_murder_rate <- murders %>%
  summarize(rate = sum(total) / sum(population) * 100000)
us_murder_rate
#> rate
#> 1 3.03
```

This computation counts larger states proportionally to their size which results in a larger value.

4.7.2 pull

The us_murder_rate object defined above represents just one number. Yet we are storing it in a data frame:

```
class(us_murder_rate)
#> [1] "data.frame"
```

since, as most dplyr functions, summarize always returns a data frame.

This might be problematic if we want to use this result with functions that require a numeric value. Here we show a useful trick for accessing values stored in data when using pipes: when a data object is piped that object and its columns can be accessed using the pull function. To understand what we mean take a look at this line of code:

```
us_murder_rate %>% pull(rate)
#> [1] 3.03
```

This returns the value in the rate column of us_murder_rate making it equivalent to us_murder_rate\$rate .

To get a number from the original data table with one line of code we can type:

```
us_murder_rate <- murders %>%
    summarize(rate = sum(total) / sum(population) * 100000) %>%
    pull(rate)
    us_murder_rate
    #> [1] 3.03
which is now a numeric:
```

class(us_murder_rate)
#> [1] "numeric"

4.7.3 Group then summarize with group_by

A common operation in data exploration is to first split data into groups and then compute summaries for each group. For example, we may want to compute the average and standard deviation for men's and women's heights separately. The <u>group by</u> function helps us do this.

If we type this:

heights %>% group_by(sex) #> # A tibble: 1,050 x 2 #> # Groups: sex [2] height #> sex <fct> <dbl> #> #> 1 Male 75 #> 2 Male 70 #> 3 Male 68 #> 4 Male 74 #> 5 Male 61 #> # ... with 1,045 more rows

The result does not look very different from heights, except we see Groups: sex [2] when we print the object. Although not immediately obvious from its appearance, this is now a special data frame called a *grouped data frame*, and *dplyr* functions, in particular summarize, will behave differently when acting on this object. Conceptually, you can think of this table as many tables, with the same columns but not necessarily the same number of rows, stacked together in one object. When we summarize the data after grouping, this is what happens:

```
heights %>%
  group_by(sex) %>%
  summarize(average = mean(height), standard deviation = sd(height))
#> `summarise()` ungrouping output (override with `.groups` argument)
#> # A tibble: 2 x 3
            average standard deviation
     sex
#>
     <fct>
              <dbl>
#>
                                  <dbl>
#> 1 Female
               64.9
                                  3.76
#> 2 Male
               69.3
                                   3.61
```

The summarize function applies the summarization to each group separately.

For another example, let's compute the median murder rate in the four regions of the country:

mu	murders %>%							
9	group_by(region) %>%							
:	<pre>summarize(median_rate = median(rate))</pre>							
#>	`s	summarise()` u	ngrouping output	(override with `	.groups` argument)			
#>	#	A tibble: 4 x	2					
#>		region	median_rate					
#>		<fct></fct>	<db1></db1>					
#>	1	Northeast	1.80					
#>	2	South	3.40					
#>	3	North Central	1.97					
#>	4	West	1.29					

4.8 Sorting data frames

When examining a dataset, it is often convenient to sort the table by the different columns. We know about the order and sort function, but for ordering entire tables, the **dplyr** function arrange is useful. For example, here we order the states by population size:

```
murders %>%
  arrange(population) %>%
  head()
#>
                    state abb
                                     region population total
                                                                rate
#> 1
                  Wyoming WY
                                       West
                                                563626
                                                            5 0.887
#> 2 District of Columbia
                                                           99 16.453
                           DC
                                      South
                                                601723
#> 3
                  Vermont VT
                                  Northeast
                                                625741
                                                            2
                                                               0.320
             North Dakota ND North Central
                                                              0.595
#> 4
                                                672591
                                                            4
#> 5
                   Alaska AK
                                       West
                                                710231
                                                              2.675
                                                           19
#> 6
             South Dakota SD North Central
                                                814180
                                                            8
                                                               0,983
```

With arrange we get to decide which column to sort by. To see the states by murder rate, from lowest to highest, we arrange by rate instead:

```
murders %>%
  arrange(rate) %>%
  head()
#>
             state abb
                              region population total rate
#> 1
           Vermont
                           Northeast
                                                    2 0.320
                   VT
                                         625741
#> 2 New Hampshire
                           Northeast
                                        1316470
                                                    5 0.380
                    NH
                                                    7 0.515
#> 3
            Hawaii HI
                                        1360301
                                West
#> 4
     North Dakota ND North Central
                                         672591
                                                    4 0.595
#> 5
             Iowa IA North Central
                                                   21 0.689
                                        3046355
                                                   12 0.766
#> 6
             Idaho ID
                                        1567582
                                West
```

Note that the default behavior is to order in ascending order. In **dplyr**, the function desc transforms a vector so that it is in descending order. To sort the table in descending order, we can type:

```
weird syntax!
Seems like it should be
arrange(rate, increasing=F) or
desc(arrange(rate))
but this is just how
desc() works I guess.
```

4.8.1 Nested sorting

If we are ordering by a column with ties, we can use a second column to break the tie. Similarly, a third column can be used to break ties between first and second and so on. Here we order by region, then within region we order by murder rate:

```
murders %>%
  arrange(region, rate) %>%
  head()
#>
             state abb
                          region population total rate
#> 1
           Vermont VT Northeast
                                                2 0.320
                                     625741
#> 2 New Hampshire NH Northeast
                                    1316470
                                                5 0.380
             Maine ME Northeast
#> 3
                                    1328361
                                               11 0.828
#> 4 Rhode Island RI Northeast
                                    1052567
                                               16 1.520
#> 5 Massachusetts MA Northeast
                                    6547629
                                              118 1.802
          New York NY Northeast
                                   19378102
                                              517 2.668
#> 6
```

4.8.2 The top n

In the code above, we have used the function head to avoid having the page fill up with the entire dataset. If we want to see a larger proportion, we can use the top_n function. This function takes a data frame as it's first argument, the number of rows to show in the second, and the variable to filter by in the third. Here is an example of how to see the top 5 rows:

mu	<pre>murders %>% top_n(5, rate) or top_n(murders, 5, rate)</pre>								
#>	state	abb		region	population	total	rate		
#>	1 District of Columbia	DC		South	601723	99	16.45		
#>	2 Louisiana	LA		South	4533372	351	7.74		
#>	3 Maryland	MD		South	5773552	293	5.07		
#>	4 Missouri	МО	North	Central	5988927	321	5.36		
#>	5 South Carolina	SC		South	4625364	207	4.48		

Note that rows are not sorted by rate, only filtered. If we want to sort, we need to use arrange. Note that if the third argument is left blank, top_n filters by the last column.

4.9 Exercises

For these exercises, we will be using the data from the survey collected by the United States National Center for Health Statistics (NCHS). This center has conducted a series of health and nutrition surveys since the 1960's. Starting in 1999, about 5,000 individuals of all ages have been interviewed every year and they complete the health examination component of the survey. Part of the data is made available via the **NHANES** package. Once you install the **NHANES** package, you can load the data like this:

```
library(NHANES)
data(NHANES)
```

The **NHANES** data has many missing values. The mean and sd functions in R will return NA if any of the entries of the input vector is an NA. Here is an example:

```
library(dslabs)
data(na_example)
mean(na_example)
#> [1] NA
sd(na_example)
#> [1] NA
```

To ignore the NA s we can use the na.rm argument:

```
mean(na_example, na.rm = TRUE)
#> [1] 2.3
sd(na_example, na.rm = TRUE)
#> [1] 1.22
```

Let's now explore the NHANES data.

1. We will provide some basic facts about blood pressure. First let's select a group to set the standard. We will use 20-to-29-year-old females. AgeDecade is a categorical variable with these ages. Note that the category is coded like " 20-29", with a space in front! What is the average and standard deviation of systolic blood pressure as saved in the BPSysAve variable? Save it to a variable called ref.

Hint: Use filter and summarize and use the na.rm = TRUE argument when computing the average and standard deviation. You can also filter the NA values using filter.

2. Using a pipe, assign the average to a numeric variable ref_avg . Hint: Use the code similar to above and then pull .

3. Now report the min and max values for the same group.

4. Compute the average and standard deviation for females, but for each age group separately rather than a selected decade as in question 1. Note that the age groups are defined by AgeDecade . Hint: rather than filtering by age and gender, filter by Gender and then use group_by .

5. Repeat exercise 4 for males.

6. We can actually combine both summaries for exercises 4 and 5 into one line of code. This is because group_by permits us to group by more than one variable. Obtain one big summary table using group_by(AgeDecade, Gender).

7. For males between the ages of 40-49, compare systolic blood pressure across race as reported in the Race1 variable. Order the resulting table from lowest to highest average systolic blood pressure.

4.10 Tibbles

Tidy data must be stored in data frames. We introduced the data frame in Section 2.4.1 and have been using the murders data frame throughout the book. In Section 4.7.3 we introduced the group_by function, which permits stratifying data before computing summary statistics. But where is the group information stored in the data frame?

```
murders %>% group_by(region)
#> # A tibble: 51 x 6
#> # Groups:
               region [4]
#>
     state
                abb
                      region population total rate
                <chr> <fct>
     <chr>
                                  <dbl> <dbl> <dbl>
#>
#> 1 Alabama
                     South
                                          135 2.82
               AL
                                4779736
#> 2 Alaska
               AK
                     West
                                 710231
                                           19 2.68
#> 3 Arizona
               ΑZ
                     West
                                6392017
                                          232 3.63
#> 4 Arkansas
               AR
                     South
                                2915918
                                           93 3,19
#> 5 California CA
                     West
                               37253956 1257 3.37
#> # ... with 46 more rows
```

Notice that there are no columns with this information. But, if you look closely at the output above, you see the line A tibble followd by dimensions. We can learn the class of the returned object using:

murders %>% group_by(region) %>% class()
#> [1] "grouped_df" "tbl_df" "tbl" "data.frame"

The tbl, pronounced tibble, is a special kind of data frame. The functions group_by and summarize always return this type of data frame. The group_by function returns a special kind of tbl, the grouped_df. We will say more about these later. For consistency, the **dplyr** manipulation verbs (select, filter, mutate, and arrange) preserve the class of the input: if they receive a regular data frame they return a regular data frame, while if they receive a tibble they return a tibble. But tibbles are the preferred format in the tidyverse and as a result tidyverse functions that produce a data frame from scratch return a tibble. For example, in Chapter 5 we will see that tidyverse functions used to import data create tibbles.

Tibbles are very similar to data frames. In fact, you can think of them as a modern version of data frames. Nonetheless there are three important differences which we describe next.

4.10.1 Tibbles display better

The print method for tibbles is more readable than that of a data frame. To see this, compare the outputs of typing murders and the output of murders if we convert it to a tibble. We can do this using as_tibble(murders). If using RStudio, output for a tibble adjusts to your window size. To see this, change the width of your R console and notice how more/less columns are shown.

4.10.2 Subsets of tibbles are tibbles

If you subset the columns of a data frame, you may get back an object that is not a data frame, such as a vector or scalar. For example:

```
class(murders[,4])
#> [1] "numeric"
```

is not a data frame. With tibbles this does not happen:

```
class(as_tibble(murders)[,4])
                                  "data.frame"
#> [1] "tbl df"
                    "tbl"
```

This is useful in the tidyverse since functions require data frames as input.

With tibbles, if you want to access the vector that defines a column, and not get back a data frame, you need to use the accessor \$: https://rafalab.github.io/dsbook/tidyverse.html

```
class(as_tibble(murders)$population)
#> [1] "numeric"
```

A related feature is that tibbles will give you a warning if you try to access a column that does not exist. If we accidentally write Population instead of population this:

```
murders$Population
#> NULL
```

returns a NULL with no warning, which can make it harder to debug. In contrast, if we try this with a tibble we get an informative warning:

```
as_tibble(murders)$Population
#> Warning: Unknown or uninitialised column: `Population`.
#> NULL
```

4.10.3 Tibbles can have complex entries

While data frame columns need to be vectors of numbers, strings, or logical values, tibbles can have more complex objects, such as lists or functions. Also, we can create tibbles with functions:

tibble(id = c(1, 2, 3), func = c(mean, median, sd))
#> # A tibble: 3 x 2
#> id func
#> <dbl> <list>
#> 1 1 <fn>
#> 2 2 <fn>
#> 3 3 <fn>

4.10.4 Tibbles can be grouped

The function group_by returns a special kind of tibble: a grouped tibble. This class stores information that lets you know which rows are in which groups. The tidyverse functions, in particular the summarize function, are aware of the group information.

4.10.5 Create a tibble using tibble instead of data.frame

It is sometimes useful for us to create our own data frames. To create a data frame in the tibble format, you can do this by using the tibble function.

```
grades <- tibble(names = c("John", "Juan", "Jean", "Yao"),
        exam_1 = c(95, 80, 90, 85),
        exam_2 = c(90, 85, 85, 90))
```

Note that base R (without packages loaded) has a function with a very similar name,

data.frame, that can be used to create a regular data frame rather than a tibble. One other important difference is that by default data.frame coerces characters into factors without providing a warning or message:

To avoid this, we use the rather cumbersome argument stringsAsFactors :

To convert a regular data frame to a tibble, you can use the as_tibble function.

```
as_tibble(grades) %>% class()
#> [1] "tbl_df" "tbl" "data.frame"
```

4.11 The dot operator

One of the advantages of using the pipe %>% is that we do not have to keep naming new objects as we manipulate the data frame. As a quick reminder, if we want to compute the median murder rate for states in the southern states, instead of typing:

We can avoid defining any new intermediate objects by instead typing:

```
filter(murders, region == "South") %>%
  mutate(rate = total / population * 10^5) %>%
  summarize(median = median(rate)) %>%
  pull(median)
#> [1] 3.4
```

We can do this because each of these functions takes a data frame as the first argument. But what if we want to access a component of the data frame. For example, what if the pull function was not available and we wanted to access tab_2\$rate ? What data frame name would we use? The answer is the dot operator.

For example to access the rate vector without the pull function we could use

```
rates <- filter(murders, region == "South") %>%
  mutate(rate = total / population * 10^5) %>%
    .$rate Cool! This is equivalent to
    tab_2$rate
#> [1] 3.4
```

In the next section, we will see other instances in which using the . is useful.

4.12 do

The tidyverse functions know how to interpret grouped tibbles. Furthermore, to facilitate stringing commands through the pipe %>%, tidyverse functions consistently return data frames, since this assures that the output of a function is accepted as the input of another. But most R functions do not recognize grouped tibbles nor do they return data frames. The quantile function is an example we described in Section 4.7.1. The do function serves as a bridge between R functions such as quantile and the tidyverse. The do function understands grouped tibbles and always returns a data frame.

In Section 4.7.1, we noted that if we attempt to use quantile to obtain the min, median and max in one call, we will receive an error: Error: expecting result of length one, got : 2.

```
data(heights)
heights %>%
filter(sex == "Female") %>%
summarize(range = quantile(height, c(0, 0.5, 1)))
```

We can use the do function to fix this.

First we have to write a function that fits into the tidyverse approach: that is, it receives a data frame and returns a data frame.

```
my_summary <- function(dat){
    x <- quantile(dat$height, c(0, 0.5, 1))
    tibble(min = x[1], median = x[2], max = x[3])
}</pre>
```

We can now apply the function to the heights dataset to obtain the summaries:

```
heights %>%
group_by(sex) %>%
my_summary
#> # A tibble: 1 x 3
#> min median max
#> <dbl> <dbl> <dbl>
#> 1 50 68.5 82.7
```

But this is not what we want. We want a summary for each sex and the code returned just one summary. This is because my_summary is not part of the tidyverse and does not know how to handled grouped tibbles. do makes this connection:

```
heights %>%
  group_by(sex) %>%
 do(my summary(.))
#> # A tibble: 2 x 4
#> # Groups: sex [2]
             min median
#>
     sex
                           max
#>
     <fct> <dbl> <dbl> <dbl>
#> 1 Female
               51
                    65.0
                          79
#> 2 Male
               50
                    69
                          82.7
```

Note that here we need to use the dot operator. The tibble created by group_by is piped to do. Within the call to do, the name of this tibble is . and we want to send it to my_summary . If you do not use the dot, then my_summary has **no argument** and returns an error telling us that argument "dat" is missing. You can see the error by typing:

If you do not use the parenthesis, then the function is not executed and instead do tries to return the function. This gives an error because do must always return a data frame. You can see the error by typing:

```
heights %>%
group_by(sex) %>%
do(my_summary)
```

4.13 The purrr package

In Section 3.5 we learned about the sapply function, which permitted us to apply the same function to each element of a vector. We constructed a function and used sapply to compute the sum of the first n integers for several values of n like this:

```
compute_s_n <- function(n){
    x <- 1:n
    sum(x)
}
n <- 1:25
s_n <- sapply(n, compute_s_n)</pre>
```

This type of operation, applying the same function or procedure to elements of an object, is quite common in data analysis. The **purrr** package includes functions similar to sapply but that better interact with other tidyverse functions. The main advantage is that we can better control the output type of functions. In contrast, sapply can return several different object types; for example, we might expect a numeric result from a line of code, but sapply might convert our result to character under some circumstances. **purrr** functions will never do this: they will return objects of a specified type or return an error if this is not possible.

The first **purrr** function we will learn is <u>map</u>, which works very similar to sapply but always, without exception, returns a list:

```
library(purrr)
s_n <- map(n, compute_s_n)
class(s_n)
#> [1] "list"
```

If we want a numeric vector, we can instead use map_dbl which always returns a vector of numeric values.

```
s_n <- map_dbl(n, compute_s_n)
class(s_n)
#> [1] "numeric"
```

This produces the same results as the sapply call shown above.

A particularly useful **purrr** function for interacting with the rest of the tidyverse is <u>map_df</u>, <u>which</u> <u>always returns a tibble data frame</u>. However, the function being called needs to return a vector or a list with names. For this reason, the following code would result in a Argument 1 must have names error:

```
s_n <- map_df(n, compute_s_n)</pre>
```

We need to change the function to make this work:

```
compute_s_n <- function(n){
    x <- 1:n
    tibble(sum = sum(x)) Instead of just sum(x)
}
s_n <- map_df(n, compute_s_n)</pre>
```

The **purrr** package provides much more functionality not covered here. For more details you can consult this online resource.

4.14 Tidyverse conditionals

A typical data analysis will often involve one or more conditional operations. In Section 3.1 we described the ifelse function, which we will use extensively in this book. In this section we present two **dplyr** functions that provide further functionality for performing conditional operations.

4.14.1 case_when

The case_when function is useful for vectorizing conditional statements. It is similar to ifelse but can output any number of values, as opposed to just TRUE or FALSE. Here is an example splitting numbers into negative, positive, and 0: This isn't what he means.

He means it can condition

A common use for this function is to define categorical variables based on existing variables. For example, suppose we want to compare the murder rates in four groups of states: *New England*, *West Coast*, *South*, and *other*. For each state, we need to ask if it is in New England, if it is not we ask if it is in the West Coast, if not we ask if it is in the South, and if not we assign *other*. Here is how we use case_when to do this:

```
murders %>%
  mutate(group = case_when(
   abb %in% c("ME", "NH", "VT", "MA", "RI", "CT") ~ "New England",
                                                                              good
   abb %in% c("WA", "OR", "CA") ~ "West Coast",
                                                                            example.
    region == "South" ~ "South",
   TRUE ~ "Other")) %>%
  group_by(group) %>%
 summarize(rate = sum(total) / sum(population) * 10^5)
#> `summarise()` ungrouping output (override with `.groups` argument)
#> # A tibble: 4 x 2
#>
   group
                rate
    <chr>
                 <dbl>
#>
#> 1 New England 1.72
#> 2 Other
                 2.71
                 3.63
#> 3 South
#> 4 West Coast 2.90
```

4.14.2 between

A common operation in data analysis is to determine if a value falls inside an interval. We can check this using conditionals. For example, to check if the elements of a vector \times are between a and b we can type

x >= a & x <= b

However, this can become cumbersome, especially within the tidyverse approach. The between function performs the same operation.

between(x, a, b)

4.15 Exercises

1. Load the murders dataset. Which of the following is true?

- a. murders is in tidy format and is stored in a tibble.
- b. murders is in tidy format and is stored in a data frame.
- c. murders is not in tidy format and is stored in a tibble.
- d. murders is not in tidy format and is stored in a data frame.

2. Use as_tibble to convert the murders data table into a tibble and save it in an object called murders_tibble .

3. Use the group_by function to convert murders into a tibble that is grouped by region.

4. Write tidyverse code that is equivalent to this code:

```
exp(mean(log(murders$population)))
```

Write it using the pipe so that each function is called without arguments. Use the dot operator to access the population. Hint: The code should start with murders %>%.

5. Use the map_df to create a data frame with three columns named n, s_n, and s_n_2. The first column should contain the numbers 1 through 100. The second and third columns should each contain the sum of 1 through n with n the row number.