# Working with real data

## 11/2/2022

## Getting Started

**Some more R Packages**

In this lab, we will use the following R packages:

- The suite of **tidyverse** packages: for data wrangling and data visualization
- **openintro**: for data and custom functions with the OpenIntro resources

If these packages are not already available in your R environment, install them by typing the following three lines of code into the console of your RStudio session, pressing the enter/return key after each one.

Note that you can check to see which packages (and which versions) are installed by inspecting the *Packages* tab in the lower right panel of RStudio.

```
install.packages("tidyverse")
install.packages("openintro")
```

Next, you need to load these packages in your working environment. We do this with the `library` function. Run the following three lines in your console.

```
library(tidyverse)
library(openintro)
library(infer)
```

**The data**

Every two years, the Centers for Disease Control and Prevention conduct the Youth Risk Behavior Surveillance System (YRBSS) survey, where it takes data from high schoolers (9th through 12th grade), to analyze health patterns. You will work with a selected group of variables from a random sample of observations during one of the years the YRBSS was conducted.

Load the `yrbss` data set into your workspace.

```
data('yrbss', package='openintro')
```

This command instructs R to load some data. You should see that the environment area in the upper righthand corner of the RStudio window now lists a data set called `yrbss` that has 13583 observations on 13 variables (some numerical, some categorical). The meaning of each variable can be found by bringing up the help file:

```
?yrbss
```

Printing the whole dataset in the console is not that useful. One advantage of RStudio is that it comes with a built-in data viewer. Click on the name `yrbss` in the *Environment* pane (upper right window) that lists the objects in your environment. This will bring up an alternative display of the data set in the *Data Viewer* (upper left window). You can close the data viewer by clicking on the `x` in the upper left hand corner.

Note that the row numbers in the first column are not part of yrbss data. R adds them as part of its printout to help you make visual comparisons. You can think of them as the index that you see on the left side of a spreadsheet. In fact, the comparison to a spreadsheet will generally be helpful. R has stored yrbss data in a kind of spreadsheet or table called a *data frame.*

You can see the dimensions of this data frame as well as the names of the variables and the first few observations by typing:

```
glimpse(yrbss)
```

```
## Rows: 13,583
## Columns: 13
## $ age                    <int> 14, 14, 15, 15, 15, 15, 15, 14, 15, 15, 15, 1~
## $ gender                 <chr> "female", "female", "female", "female", "fema~
## $ grade                  <chr> "9", "9", "9", "9", "9", "9", "9", "9", "9", ~
## $ hispanic               <chr> "not", "not", "hispanic", "not", "not", "not"~
## $ race                   <chr> "Black or African American", "Black or Africa~
## $ height                 <dbl> NA, NA, 1.73, 1.60, 1.50, 1.57, 1.65, 1.88, 1~
## $ weight                 <dbl> NA, NA, 84.37, 55.79, 46.72, 67.13, 131.54, 7~
## $ helmet_12m             <chr> "never", "never", "never", "never", "did not ~
## $ text_while_driving_30d <chr> "0", NA, "30", "0", "did not drive", "did not~
## $ physically_active_7d   <int> 4, 2, 7, 0, 2, 1, 4, 4, 5, 0, 0, 0, 4, 7, 7, ~
## $ hours_tv_per_school_day <chr> "5+", "5+", "5+", "2", "3", "5+", "5+", "5+",~
## $ strength_training_7d   <int> 0, 0, 0, 0, 1, 0, 2, 0, 3, 0, 3, 0, 0, 7, 7, ~
## $ school_night_hours_sleep <chr> "8", "6", "<5", "6", "9", "8", "9", "6", "<5"~
```

## Exploratory data analysis

Let us start with analyzing the weight of the participants in kilograms: `weight`.

Using visualization and summary statistics, describe the distribution of weights. The `summary` function can be useful.

```
summary(yrbss$weight)
```

```
##    Min. 1st Qu.  Median    Mean 3rd Qu.    Max.    NA's
##   29.94   56.25   64.41   67.91   76.20  180.99    1004
```

Next, consider the possible relationship between a high schooler's weight and their physical activity.

To do so, let's create a new variable `physical_3plus`, which will be coded as either "yes" if they are physically active for at least 3 days a week, and "no" if not.

```
yrbss <- yrbss %>%
  mutate(physical_3plus = ifelse(yrbss$physically_active_7d > 2, "yes", "no"))
```

The `%>%` operator is called the **piping** operator. It takes the output of the previous expression and pipes it into the first argument of the function in the following one. In other words, `x %>% f(y)` is equivalent to `f(x, y)`.

**A note on piping:**  Note that we can read these two lines of code as the following:

*"Take the **yrbss** dataset and **pipe** it into the **mutate** function. Mutate the **yrbss** data set by creating a new variable called **physical_3plus** that is True if and only if if the person is physically active at least 3 days a week. Then assign the resulting dataset to the object called **yrbss**, i.e. overwrite the old **yrbss** dataset with the new one containing the new variable."*

We can compare the means of the distributions using the following to first group the data by the `physical_3plus` variable, and then calculate the mean `weight` in these groups using the `mean` function while ignoring missing values by setting the `na.rm` argument to `TRUE`.

```
yrbss %>%
  group_by(physical_3plus) %>%
  summarise(mean_weight = mean(weight, na.rm = TRUE))
```

```
## # A tibble: 3 x 2
##   physical_3plus mean_weight
##   <chr>                <dbl>
## 1 no                    66.7
## 2 yes                   68.4
## 3 <NA>                  69.9
```

There is an observed difference, but is this difference statistically significant? In order to answer this question we should conduct a hypothesis test...

---