

```

# Yannique Hecht
# Harvardx: PH125.1x - (1) Data Science: R Basics
# SECTION 3: INDEXING, DATA WRANGLING, PLOTS
# ASSESSMENTS

# # # ASSESSMENT 3.1: INDEXING

# # LOGICAL VECTORS

# Store the murder rate per 100,000 for each state, in `murder_rate`
[REDACTED]

# Store the `murder_rate < 1` in `low`
[REDACTED]

# # WHICH

# Store the murder rate per 100,000 for each state, in murder_rate
[REDACTED]

# Store the murder_rate < 1 in low
[REDACTED]

# Get the indices of entries that are below 1
[REDACTED]

# # ORDERING VECTORS

# Store the murder rate per 100,000 for each state, in murder_rate
[REDACTED]

# Store the murder_rate < 1 in low
[REDACTED]

# Names of states with murder rates lower than 1
[REDACTED]

# # FILTERING

# Store the murder rate per 100,000 for each state, in `murder_rate`
[REDACTED]

# Store the `murder_rate < 1` in `low`

```

```
[REDACTED]

# Create a vector ind for states in the Northeast and with murder
rates lower than 1.
[REDACTED]
```

```
# Names of states in `ind`
[REDACTED]
```

```
# Store the murder rate per 100,000 for each state, in murder_rate
[REDACTED]
```

```
# Compute the average murder rate using `mean` and store it in object
named `avg`
[REDACTED]
```

```
# How many states have murder rates below avg ? Check using sum
[REDACTED]
```

```
# # MATCH
```

```
# Store the 3 abbreviations in a vector called `abbs` (remember that
they are character vectors and need quotes)
[REDACTED]
```

```
# Match the abbs to the murders$abb and store in ind
[REDACTED]
```

```
# Print state names from ind
[REDACTED]
```

```
# # %in%
```

```
# Store the 5 abbreviations in `abbs`. (remember that they are
character vectors)
[REDACTED]
```

```
# Use the %in% command to check if the entries of abbs are
abbreviations in the the murders data frame
```

```
a [REDACTED]
```

```
# # LOGICAL OPERATOR
```

```
# Store the 5 abbreviations in `abbs`. (remember that they are
character vectors)
```

```
[REDACTED]
```

```
# Use the %in% command to check if the entries of abbs are  
abbreviations in the the murders data frame
```

```
[REDACTED]
```

```
# # # ASSESSMENT 3.2: BASIC DATA WRANGLING
```

```
# # DPLYR
```

```
# Loading data
```

```
[REDACTED]
```

```
[REDACTED]
```

```
# Loading dplyr
```

```
[REDACTED]
```

```
# Redefine murders so that it includes a column named rate with the  
per 100,000 murder rates
```

```
[REDACTED]
```

```
# # MUTATE
```

```
# Note that if you want ranks from highest to lowest you can take the  
negative and then compute the ranks
```

```
[REDACTED]
```

```
[REDACTED]
```

```
# Defining rate
```

```
[REDACTED]
```

```
# Redefine murders to include a column named rank  
# with the ranks of rate from highest to lowest
```

```
[REDACTED]
```

```
# # SELECT
```

```
# Load dplyr
```

```
[REDACTED]
```

```
# Use select to only show state names and abbreviations from murders
```

```
[REDACTED]
```

```
# # FILTER
```

```
# Add the necessary columns
```

```
[REDACTED]
```

```
[REDACTED]
```

```
# Filter to show the top 5 states with the highest murder rates
```

```
[REDACTED]
```

```
# # FILTER WITH !=
```

```
# Use filter to create a new data frame no_south
```

```
[REDACTED]
```

```
# Use nrow() to calculate the number of rows
```

```
[REDACTED]
```

```
# # FILTER WITH %IN%
```

```
# Create a new data frame called murders_nw with only the states from  
the northeast and the west
```

```
[REDACTED]
```

```
# Number of states (rows) in this category
```

```
[REDACTED]
```

```
# # FILTERING BY 2 CONDITIONS
```

```
# add the rate column
```

```
[REDACTED] [REDACTED]
```

```
[REDACTED]
```

```
# Create a table, call it my_states, that satisfies both the  
conditions
```

```
[REDACTED]
```

```
[REDACTED]
```

```
# Use select to show only the state name, the murder rate and the  
rank
```

```
[REDACTED]
```

```
# # USING THE PIPE %>%
```

```
# add the rate column
```

```
[REDACTED]
```

```
# Create a table, call it my_states, that satisfies both the conditions
```

```
[REDACTED]
```

```
# Use select to show only the state name, the murder rate and the rank
```

```
[REDACTED]
```

```
# # MUTATE, FILTER & SELECT
```

```
# Loading the libraries
```

```
[REDACTED]
```

```
# Create new data frame called my_states (with specifications in the instructions)
```

```
[REDACTED]
```

```
# # # ASSESSMENT 3.3: BASIC PLOTS
```

```
# # SCATTERPLOTS
```

```
# Load the datasets and define some variables
```

```
[REDACTED]
```

```
population_in_millions <- murders$population/10^6
```

```
[REDACTED]
```

```
plot(population_in_millions, total_gun_murders)
```

```
# Transform population using the log10 transformation and save to object log10_population
```

```
[REDACTED]
```

```
# Transform total gun murders using log10 transformation and save to  
object log10_total_gun_murders
```

```
# Create a scatterplot with the log scale transformed population and  
murders
```

```
# # HISTOGRAMS
```

```
# Store the population in millions and save to population_in_millions
```

```
# Create a histogram of this variable
```

```
# # BOXPLOTS
```

```
# Create a boxplot of state populations by region for the murders  
dataset
```

```
# # # SECTION 3 ASSESSMENT
```

```
# # Q1 First, determine the average height in this dataset. Then  
create a logical vector ind with the indices for those individuals  
who are above average height.
```

```
# # Q2 How many individuals in the dataset are above average height  
and are female?
```

```
# # Q3 If you use mean on a logical (TRUE/FALSE) vector, it returns  
the proportion of observations that are TRUE. What proportion of  
individuals in the dataset are female?
```

```
# # Q4a This question takes you through three steps to determine the  
sex of the individual with the minimum height.
```

Q4b Use the match() function to determine the index of the individual with the minimum height.

```
[REDACTED]
```

Q4c Subset the sex column of the dataset by the index in 4b to determine the individual's sex.

```
[REDACTED]
```

Q5a Determine the maximum height.

```
[REDACTED]
```

Q5b Write code to create a vector x that includes the integers between the minimum and maximum heights.

```
[REDACTED]
```

Q5c How many of the integers in x are NOT heights in the dataset?

```
[REDACTED]
```

Q6a What is the height in centimeters of the 18th individual (index 18)?

```
[REDACTED]
```

```
[REDACTED]
```

Q6b What is the mean height in centimeters?

```
[REDACTED]
```

7a How many females are in the heights2 dataset?

```
[REDACTED]
```

```
[REDACTED]
```

Q7b What is the mean height of the females in centimeters?

```
[REDACTED]
```

Q8 Plot the percent palmitic acid versus palmitoleic acid in a scatterplot. What relationship do you see?

```
[REDACTED]
```

Q9 Create a histogram of the percentage of eicosenoic acid in olive. Which of the following is true?

```
[REDACTED]
```

10 Make a boxplot of palmitic acid percentage in olive with separate distributions for each region. Which region has the most variable palmitic acid percentage? We can determine variability from

the range of values each regions's palmitic acid percentage covers:
