

Problem Set 1: Manipulation of mammalian sleep data

Your name and student ID

Today's date

```
BEGIN ASSIGNMENT
requirements: requirements.R
generate: true
files:
  - data
```

Instructions

- Solutions will be released on Wednesday, September 1st.
- This semester, problem sets are for practice only and will not be turned in for marks.

Helpful hints:

- Every function you need to use was taught during lecture! So you may need to revisit the lecture code to help you along by opening the relevant files on Datahub. Alternatively, you may wish to view the code in the condensed PDFs posted on the course website. Good luck!
- Knit your file early and often to minimize knitting errors! If you copy and paste code for the slides, you are bound to get an error that is hard to diagnose. Typing out the code is the way to smooth knitting! We recommend knitting your file each time after you write a few sentences/add a new code chunk, so you can detect the source of knitting errors more easily. This will save you and the GSIs from frustration!
- It is good practice to not allow your code to run off the page. To avoid this, have a look at your knitted PDF and ensure all the code fits in the file. If it doesn't look right, go back to your .Rmd file and add spaces (new lines) using the return or enter key so that the code runs onto the next line.

Begin by knitting this document by pushing the “Knit” button above. As you fill in code and text in the document, you can re-knit (push the button again) and see how the document changes. It is important to re-knit often, because if there is any error in your code, the file will not generate a PDF, so our advice is to knit early and often!

Using dplyr to investigate sleep times in mammals

The data file `sleep.csv` contains the sleeptimes and weights for a set of mammals. Hit the green arrow icon in the line below to execute the two lines of code in the code chunk, or execute them line by line by placing your cursor on the first line and hitting `cmd + enter` on Mac or `ctrl + enter` on PC.

```
library(dplyr)
```

```
## Warning: package 'dplyr' was built under R version 4.0.5
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##   filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##   intersect, setdiff, setequal, union
```

```
library(readr)
```

```
sleep <- read_csv("data/sleep.csv")
```

```
##
```

```
## -- Column specification -----
```

```
## cols(
```

```
##   name = col_character(),
```

```
##   genus = col_character(),
```

```
##   vore = col_character(),
```

```
##   order = col_character(),
```

```
##   conservation = col_character(),
```

```
##   sleep_total = col_double(),
```

```
##   sleep_rem = col_double(),
```

```
##   sleep_cycle = col_double(),
```

```
##   awake = col_double(),
```

```
##   brainwt = col_double(),
```

```
##   bodywt = col_double()
```

```
## )
```

- The `library` command loads the library `dplyr` into memory.
- The `readr` library contains functions to read in the dataset.
- The `dplyr` library contains functions we will use to manipulate data.

Notice that an object called `sleep` appeared in the Environment tab under “Data”.

1. [2 points] Use four useful functions discussed in lecture to examine the sleep data set:

```
# Text inside a code chunk that begins with "#" is called a comment.  
# We sometimes use comments to explain code to you in plain English.  
# Write your four functions below these comments, replacing the placeholder  
# text "<<<<YOUR CODE HERE>>>>". Remember, code does *not* begin with a "#"
```

```
"<<<<YOUR CODE HERE>>>>"
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
"<<<<YOUR CODE HERE>>>>"
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
"<<<<YOUR CODE HERE>>>>"
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
"<<<<YOUR CODE HERE>>>>"
```

```
## [1] "<<<<YOUR CODE HERE>>>>"
```

```
# Then, assign p1 to a vector of your function names, in alphabetical order.  
# For example, assigning p0 to a vector of fruits looks like this:  
# p0 <- c("apple", "banana", "orange")
```

```
p1 <- c("dim", "head", "names", "str") #SOLUTION
```

```
BEGIN QUESTION  
name: p1  
manual: false  
points: 2
```

```
## Test ##  
testthat::expect_true(length(p1)== 4,  
                        info = "p1a: Checking p1 has 4 items in a list")
```

```
## Test ##  
testthat::expect_true(p1[1] == "dim" & p1[2] == "head" & p1[3] == "names" & p1[4] == "str",  
                      info = "p1b: Checking the names of the 4 functions in alphabetical order")
```

Description of the variables found in the sleep dataset:

| Column name | Description |
|--------------|--|
| name | common name |
| genus | taxonomic rank |
| vore | carnivore, omnivore or herbivore? |
| order | taxonomic rank |
| conservation | the conservation status of the mammal |
| sleep_total | total amount of sleep, in hours |
| sleep_rem | Rapid eye movement (REM) sleep, in hours |
| sleep_cycle | length of sleep cycle, in hours |
| awake | amount of time spent awake, in hours |
| brainwt | brain weight in kilograms |
| bodywt | body weight in kilograms |

2. [2 points] Write code to select a set of columns. Specifically select the awake, brainwt, and bodywt columns. Assign this smaller dataset to a data frame called sleep_small

```
sleep_small <- select(sleep, awake, brainwt, bodywt) #SOLUTION
```

BEGIN QUESTION

name: p2
manual: false
points: 2

Test

```
testthat::expect_true(is.data.frame(sleep_small),  
  info = "p2a: Checking sleep_small is a dataframe")
```

Test

```
testthat::expect_true(ncol(sleep_small) == 3,  
  info = "p2b: Checking sleep_small has 3 columns")
```

Test

```
testthat::expect_true(all(names(sleep_small) == c("awake", "brainwt", "bodywt")),  
  info = "p2c: Checking sleep_small has 'awake', 'brainwt', and 'bodywt'")
```

3. [1 point] To select a range of columns by name, use the ‘:’ (colon) operator. Redo the selection for question 1, but use the colon operator. Assign this to `sleep_small_colon`. Note that this returns the same data frame as the previous problem, but is not recommended in practice because it depends on the ordering of the columns and isn’t explicit in the columns that are selected, whereas selecting columns by name offers much higher readability for someone else looking at your code later on.

```
sleep_small_colon <- sleep %>% select(awake:bodywt) #SOLUTION
```

BEGIN QUESTION

```
name: p3
manual: false
points: 1
```

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_small_colon),
  info = "p3a: Checking sleep_small_colon is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_small_colon) == 3,
  info = "p3b: Checking sleep_small_colon has 3 columns")
```

```
## Test ##
```

```
testthat::expect_true(all(names(sleep_small_colon) == c("awake", "brainwt", "bodywt")),
  info = "p3c: Checking sleep_small_colon has 'awake', 'brainwt', and 'bodywt'")
```

4. [1 point] From the original dataset `sleep` select all the columns except for the `vore` variable. Assign this to `sleep_no_vore`.

```
sleep_no_vore <- sleep %>% select(-vore) #SOLUTION
```

BEGIN QUESTION

name: p4

manual: false

points: 1

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_no_vore),  
  info = "p4a: Checking sleep_small_colon is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_no_vore) == 10,  
  info = "p4b: Checking sleep_small_vore has 10 columns")
```

```
## Test ##
```

```
testthat::expect_true(!("vore" %in% names(sleep_no_vore)),  
  info = "p4c: Checking sleep_no_vore has no columns with 'vore'")
```

5. [1 point] Run the following chunk of code.

```
select(sleep, starts_with("sl"))
```

```
## # A tibble: 83 x 3
##   sleep_total sleep_rem sleep_cycle
##   <dbl>      <dbl>      <dbl>
## 1      12.1      NA         NA
## 2       17       1.8        NA
## 3      14.4       2.4        NA
## 4      14.9       2.3        0.133
## 5        4       0.7        0.667
## 6      14.4       2.2        0.767
## 7       8.7       1.4        0.383
## 8        7      NA         NA
## 9      10.1       2.9        0.333
## 10       3      NA         NA
## # ... with 73 more rows
```

What does it return? Copy your choice and assign it to p5

```
# p5 <- "returns the number of columns that start with sl"
# p5 <- "returns all columns that start with sl"
# p5 <- "returns all rows that start with sl"
# p5 <- "returns all animals whose names start with sl"
```

```
p5 <- "returns all columns that start with sl" #SOLUTION
```

BEGIN QUESTION

```
name: p5
manual: false
points: 1
```

```
## Test ##
```

```
testthat::expect_true(p5 == "returns all columns that start with sl",
  info = "Checking response...")
```

```
select(sleep, starts_with("sl"))
```

```
## # A tibble: 83 x 3
##   sleep_total sleep_rem sleep_cycle
##   <dbl>      <dbl>      <dbl>
## 1      12.1      NA         NA
## 2       17       1.8        NA
## 3      14.4       2.4        NA
## 4      14.9       2.3        0.133
## 5       4        0.7        0.667
## 6      14.4       2.2        0.767
## 7       8.7       1.4        0.383
## 8       7        NA         NA
## 9      10.1       2.9        0.333
## 10      3        NA         NA
## # ... with 73 more rows
```

6. [1 point] Rewrite the previous chunk of code using the pipe operator. Assign this to `sleep_sl`.

```
sleep_sl <- sleep %>% select(starts_with("sl")) #SOLUTION
```

BEGIN QUESTION

```
name: p6
manual: false
points: 1
```

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_sl),
  info = "p6a: Checking sleep_sl is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_sl) == 3,
  info = "p6b: Checking sleep_sl has 3 columns")
```

```
## Test ##
```

```
testthat::expect_true(("sleep_total" %in% names(sleep_sl)) &&
  ("sleep_rem" %in% names(sleep_sl)) &&
  ("sleep_cycle" %in% names(sleep_sl)),
  info = "p6c: Checking sleep_sl has the 3 columns that start with sl")
```


7. [1 point] From the original sleep dataset, filter the rows for mammals that sleep a total of more than 16 hours. Assign this to sleep_over16.

```
sleep_over16 <- sleep %>% filter(sleep_total > 16) #SOLUTION
```

BEGIN QUESTION

name: p7

manual: false

points: 1

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_over16),  
  info = "p7a: Checking sleep_over16 is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_over16) == 11,  
  info = "p7b: Checking sleep_over16 has 11 columns")
```

```
## Test ##
```

```
testthat::expect_true(nrow(sleep_over16) == 8,  
  info = "p7c: Checking sleep_over16 has 8 rows")
```

8. [2 points] Filter the rows for mammals that sleep a total of more than 16 hours and have a body weight of greater than 1 kilogram. Assign this to `sleep_mammals`.

```
sleep_mammals <- sleep %>% filter(sleep_total > 16 & bodywt > 1) #SOLUTION
```

BEGIN QUESTION

```
name: p8
manual: false
points: 2
```

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_mammals),
  info = "p8a: Checking sleep_mammals is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_mammals) == 11,
  info = "p8b: Chekcing sleep_mammals has 11 columns")
```

```
## Test ##
```

```
testthat::expect_true(nrow(sleep_mammals) == 3,
  info = "p8c: Checking sleep_mammals has 3 rows")
```

9. [1 point] Suppose you are specifically interested in the sleep of horses and giraffes. From the original sleep dataset, assign sleep_hg to a data frame for horses and giraffes only.

```
sleep_hg <- sleep %>% filter(name %in% c("Horse", "Giraffe")) #SOLUTION
```

BEGIN QUESTION

name: p9
manual: false
points: 1

Test

```
testthat::expect_true(is.data.frame(sleep_hg),  
  info = "p9a: Checking sleep_hg is a dataframe")
```

Test

```
testthat::expect_true(ncol(sleep_hg) == 11,  
  info = "p9b: Checking sleep_hg has 11 columns")
```

Test

```
testthat::expect_true(nrow(sleep_hg) == 2,  
  info = "p9c: Checking sleep_hg has 2 rows")
```

Test

```
testthat::expect_true("Horse" %in% sleep_hg$name &&  
  "Giraffe" %in% sleep_hg$name,  
  info = "p9d: Checking sleep_hg has the correct rows")
```

10. [1 point] From the original dataset, order the dataset by sleep time from shortest sleep time to longest sleep time. Assign this to sleep_time.

```
sleep_time <- sleep %>% arrange(sleep_total) #SOLUTION
```

BEGIN QUESTION

name: p10
manual: false
points: 1

Test

```
testthat::expect_true(is.data.frame(sleep_time),  
  info = "p10a: Checking sleep_time is a dataframe")
```

Test

```
testthat::expect_true(ncol(sleep_time) == 11,  
  info = "p10b: Checking sleep_time has 11 columns")
```

Test

```
testthat::expect_true(nrow(sleep_time) == 83,  
  info = "p10c: Checking sleep_time has 83 rows")
```

11. [1 point] Now order for longest sleep time to shortest sleep time. Assign this to `sleep_rev`.

```
sleep_rev <- sleep %>% arrange(-sleep_total) #SOLUTION
```

BEGIN QUESTION

```
name: p11
manual: false
points: 1
```

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_rev),
  info = "p11a: Checking sleep_rev is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_rev) == 11,
  info = "p11b: Checking sleep_rev has 11 columns")
```

```
## Test ##
```

```
testthat::expect_true(nrow(sleep_rev) == 83,
  info = "p11c: Checking sleep_rev has 83 rows")
```

12. [2 points] Suppose you are interested in the order of sleep time, but according to whether the animal is a carnivore, herbivore, or omnivore. Rewrite the above statement to order sleep time according to the type of “-vore” that then animal is. Call this “sleep_time_rev”:

```
sleep_time_rev <- sleep %>% arrange(vore, -sleep_total) #SOLUTION
```

BEGIN QUESTION

name: p12
manual: false
points: 2

Test

```
testthat::expect_true(is.data.frame(sleep_time_rev),  
  info = "p12a: Checking sleep_time_rev is a dataframe")
```

Test

```
testthat::expect_true(ncol(sleep_time_rev) == 11,  
  info = "p12b: Checking sleep_time_rev has 11 columns")
```

Test

```
testthat::expect_true(nrow(sleep_time_rev) == 83,  
  info = "p12c: Checking sleep_time_rev has 83 rows")
```

13. [1 point] Create a new column called `rem_proportion` which is the ratio of rem sleep to total amount of sleep. Assign this new data frame to `sleep_ratio` from `sleep` data.

```
sleep_ratio <- sleep %>% mutate(rem_proportion = sleep_rem/sleep_total) #SOLUTION
```

BEGIN QUESTION

name: p13
manual: false
points: 1

Test

```
testthat::expect_true(is.data.frame(sleep_ratio),  
  info = "p13a: Checking sleep_ratio is a dataframe")
```

Test

```
testthat::expect_true(ncol(sleep_ratio) == 12,  
  info = "p13b: Checking sleep_time_rev has 12 collumns")
```

Test

```
testthat::expect_true(nrow(sleep_ratio) == 83,  
  info = "p13c: Checking sleep_ratio has 83 rows")
```

14. [1 point] Add a second column called `bodywt_grams` which is the `bodywt` column in grams.

```
sleep_r_bw <- sleep %>% mutate(rem_proportion = sleep_rem/sleep_total, bodywt_grams = bodywt * 1000) #S
```

BEGIN QUESTION

name: p14
manual: false
points: 1

```
## Test ##
```

```
testthat::expect_true(is.data.frame(sleep_r_bw),  
  info = "p14a: Checking sleep_r_bw is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(sleep_r_bw) == 13,  
  info = "p14b: Checking sleep_r_bw has 13 columns")
```

```
## Test ##
```

```
testthat::expect_true(nrow(sleep_r_bw) == 83,  
  info = "p14c: Checking sleep_r_bw has 83 rows")
```


15. [1 point] Calculate the average sleep time across all the animals in the dataset using a `dplyr` function and assign it to the variable `avg_sleep_time`. Your answer should be a data frame of 1 observation and 1 variable called `sleep_avg`

```
avg_sleep_time <- sleep %>% summarize(sleep_avg = mean(sleep_total)) #SOLUTION
```

BEGIN QUESTION

name: p15
manual: false
points: 1

Test

```
testthat::expect_true(is.data.frame(avg_sleep_time),  
  info = "p15a: Checking avg_sleep_time is a dataframe")
```

Test

```
testthat::expect_true(ncol(avg_sleep_time) == 1 &&  
  nrow(avg_sleep_time) == 1,  
  info = "p15b: Checking avg_sleep_time has 1 row and 1 column")
```

Test

```
testthat::expect_true(is.numeric(avg_sleep_time$sleep_avg),  
  info = "p15c: Checking sleep_avg column is numeric")
```

Test

```
testthat::expect_true(all.equal(avg_sleep_time$sleep_avg, 10.43373, tol = 0.01),  
  info = "p15d: Checking sleep avg column is 10.4337")
```

16. [2 points] Calculate the average sleep time for each type of “-vore”. Hint: you’ll need to use two dplyr functions! The column names should be vore and sleep_avg. Call this dataframe avg_by_vore

```
. = " # BEGIN PROMPT
avg_by_vore <- NULL # YOUR CODE HERE
" # END PROMPT

# BEGIN SOLUTION NO PROMPT
avg_by_vore <- sleep %>%
  group_by(vore) %>%
  summarize(sleep_avg = mean(sleep_total))
# END SOLUTION
```

```
BEGIN QUESTION
name: p16
manual: false
points: 1
```

```
## Test ##
```

```
testthat::expect_true(is.data.frame(avg_by_vore),
  info = "p16a: Checking avg_by_vore is a dataframe")
```

```
## Test ##
```

```
testthat::expect_true(ncol(avg_by_vore) == 2 &&
  nrow(avg_by_vore) == 5,
  info = "p16b: Checking avg_by_vore has 5 rows and 2 columns")
```

```
## Test ##
```

```
testthat::expect_true(identical(names(avg_by_vore), c("vore", "sleep_avg")),
  info = "p16c: Checking column names are vore and sleep_avg")
```

END