# PH142 Review Session: Week 1

PH142 GSI team
July 6th, 2023

# Announcement

- Deadlines
  - Lecture quizzes on Gradescope
  - Lab 1: due 07/07 on Datahub
  - Lab 2: due 07/10 on Datahub
  - Homework 1&2: not turned in - for you to use as practice
- Midterm I: released 07/14 - available until 5pm 7/15
- Data Project
  - Check out instructions on https://ph142-ucb.github.io/su23/data-proj/.
  - Part I due on 07/17, 10pm PST

# Objectives

- Summarize key course technologies, resources, and policies.
- Review materials from lectures 1-2.
  - PPDAC Approach
  - Categorical Data Visualization
  - Intro to R

# Key Technologies & Resources

Course website: https://ph142-ucb.github.io/su23/

# Key Technologies & Resources

Accessing slides and recordings:

# Key Technologies & Resources

Course Calendar:

# Key Technologies & Resources

Ed discussion: edstem.org

# Key Technologies & Resources

Gradescope: gradescope.com

# Key Technologies & Resources

Datahub:

# PPDAC Approach

A clear statement of what we are trying to achieve

Problem

Conclusions are drawn about what has been learned about answering the Problem

Conclusion

The procedures we use to carry out the study

Plan

Data

Analysis

Summarization and analysis of the data to answer the questions posed by the Problem

The data collected according to the Plan

# Visualization of Categorical Data: "ggplot2"

1. Install and load the "ggplot2" package
   a. install.packages("ggplot2")
   b. library(ggplot2)
2. Specify your data, and what to have on the x and y axes
3. Create a plot: geom_ functions (many options)
   a. geom_bar, geom_histogram, geom_point, geom_line
   b. Tip: picture how you want to visualize your data in your head first, then pick the function that helps you achieve your goal :)
4. Change the style of your plot
   a. labs() function: update your main title, axis names, caption
   b. theme() function: change the size of your title, font, and position

# Visualization of Categorical Data: "ggplot2"

e.g. https://ggplot2.tidyverse.org/reference/ggplot.html

```
ggplot(df, aes(gp, y)) +
  geom_point() +
  geom_point(data = ds, aes(y = mean), colour = 'red', size = 3)
```

# Types of Variables

- Categorical: a variable that has grouping levels
  - Nominal: no underlying order or rank, e.g. blood type, zip code
  - Ordinal: with an underlying order or rank, e.g. blood pressure level (low, normal, high)
- Quantitative: a numeric variable which you can perform mathematical operations on
  - Discrete: can be counted, e.g. the number of cookies in the bag you got from a bakery
  - Continuous: can be measured precisely, with a rule or scale, e.g. 5.34 grams of cornstarch

# Intro to R

- Library
  - A library is a package of functions, and you can load this package of functions by running library(ggplot2), library(dplyr)
  - Make sure you have them first, otherwise you need to do install.packages() first
- Read your data: e.g. how to read a .csv file
  - `library(readr)`
  - `mydata <- read_csv('my_data.csv')`
- Some functions to get a quick look of your data
  - head(mydata): shows the first 6 rows of the dataset
  - dim(mydata): shows the total number of rows by the total number of columns
  - names(mydata) or colnames(mydata): shows all the variable names (column names) of the dataset
  - str(mydata): summarizes the information above and more

# Data manipulation: functions in library(dplyr)

First, do `library(dplyr)` to have the package in your environment.

- rename() → renames variables (columns)
  - `new_dataset <- old_dataset %>% rename(new_name = old_name)`
  - or: `new_dataset <- rename(old_dataset, new_name = old_name)`

- select() → subsets variables (columns)
  - `smaller_data <- old_data %>% select(variable1, variable2, variable3)`
  - `smaller_data <- select(old_data, variable1, variable2, variable3)`
  - `smaller_data <- select(old_data, variable1:variable3)`
  - To keep all variables other than variable1: `smaller_data <- old_data %>% select(- variable1)`

- arrange() → orders observations (rows) by a certain variable (column) or variables (columns)
  - `lake_data %>% arrange(ph)`
  - `lake_data %>% arrange(age_data, ph)`

# Data manipulation: functions in library(dplyr)

- filter() → selects a subset of rows by certain conditions
  - If we want condition A AND condition B to be satisfied, use **,** or **&**
  - If we want condition A OR condition B to be satisfied, use **|** or **%in%**
  - `lake_data %>% filter(age_data == "recent")`
  - `lake_data %>% filter(lakes %in% c("Alligator", "Blue Cypress"))`
  - `lake_data %>% filter(ph > 6 | chlorophyll > 30)`

- mutate() → creates new variables
  - `lake_data_new <- lake_data %>% mutate(actual_fish_sample = number_fish_sampled * 100)`

- group_by() → groups the data by a categorical variable
  - `lake_data %>% group_by(age_data)%>% summarize(mean_ph = mean(ph))`

- summarize() → applies summary functions to calculate statistics
  - `lake_data %>% summarize(mean_ph = mean(ph), sd_ph = sd(ph))`

# Measure of Central Tendency

- Mean and median are **approximately equal** when…
  - Distribution is symmetric
  - Data has one peak
  - There are no outliers

- Outliers: large effect on the mean

- Skewed data: mean ≠ median
  - Skewed right:
    ## mean > median
  - Skewed left:
    ## mean < median

Lump of data follows the tail

Lump of data at the start

Long tail at the start

Long tail follows the lump of data

**Negative skew (left-skewed distribution)**

**Positive skew (right-skewed distribution)**

# Measures of Spread

- Range = max - min
- IQR = Q3 - Q1
  - Five number summary in R!
    ```
    CS_dat %>% summarize(min = min(cs_rate),
    Q1 = quantile(cs_rate, 0.25), median = median(cs_rate),
    Q3 = quantile(cs_rate, 0.75), max = max(cs_rate))
    ```
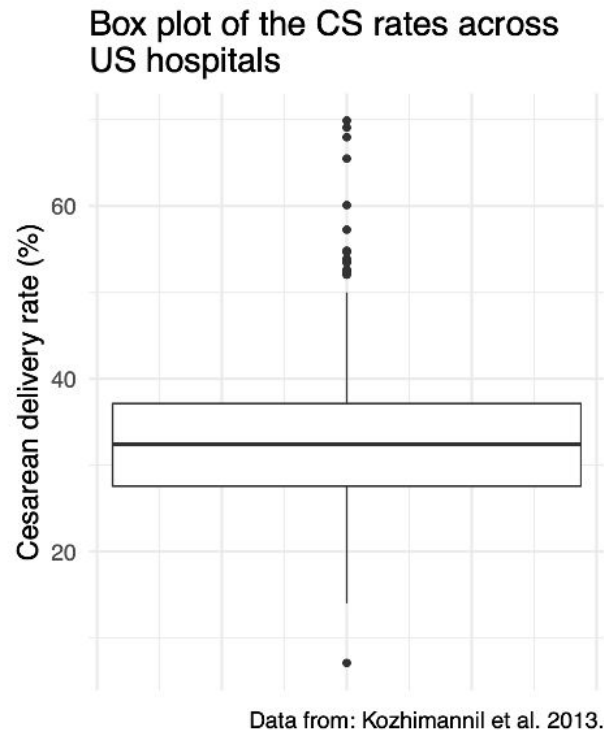- Sample variance (s^2)
- Sample standard deviation (s)
  ```
  CS_dat %>% summarize(cs_sd = sd(cs_rate), cs_var = var(cs_rate))
  ```

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2}$$

# Box plot

- Center line → median
- Top of box → Q3
- Bottom of box → Q1
- Top of top whisker → max value or highest point that is below Q3 + 1.5*IQR
- Bottom of bottom whisker → min value or lowest point that is above Q1 - 1.5*IQR
- Data points above and below whiskers → outliers

```
ggplot(CS_dat, aes(y = cs_rate)) +

geom_boxplot() +

ylab("Cesarean delivery rate (%)") +

labs(title = "Box plot of the CS rates across US hospitals",
    caption = "Data from: Kozhimannil et al. 2013.") +

theme_minimal(base_size = 15) +

scale_x_continuous(labels = NULL) # removes the labels from the x axis
```

Box plot of the CS rates across US hospitals



Data from: Kozhimannil et al. 2013.

# Common Errors

- Two different code chunks are named the same thing.
- The same variable names that are listed in the instructions are used in your work.
- If your data isn't running, try reloading your past code chunks first.
- If you want to see the output of your data, just retype the name of your variable in a new line within the same code chunk and run again.

# Questions?