

# Statistical Inference with Linear & Logistic Regression



...

GW Libraries Workshop  
Dan Kerchner ~ February 19, 2021

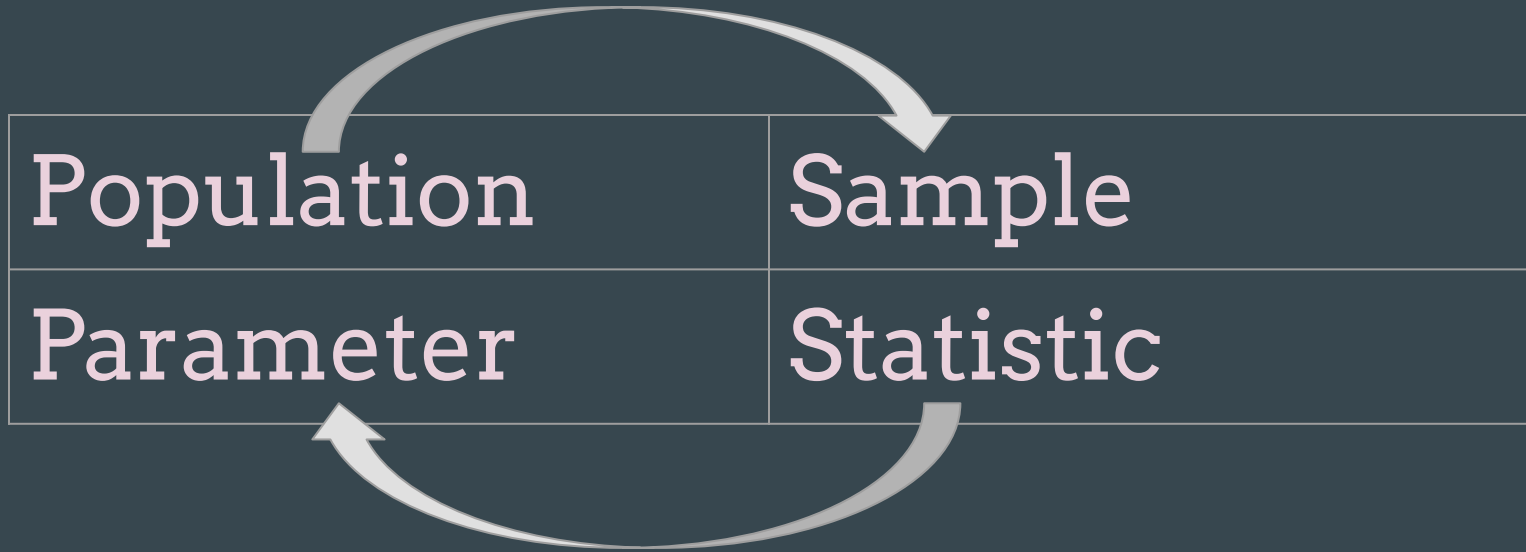
[go.gwu.edu/rstats](https://go.gwu.edu/rstats)

# Logistics

- Just speak up OR use the Zoom chat
- Kiri can provide individual help
- Plan for 1 brief ☕ break

# Super-Brief Review of Inference for Regression

# High-Level Objective



# Regression Model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

(and an interaction term might look like  $\beta_{12} X_1 X_2$ )

## Interpretation

$\beta_0$  = Mean  $Y$  when  $X_i$  values are 0

$\beta_i$  = Mean change in  $Y$  for a 1-point increase in  $X_i$ ,  
adjusting for other  $X$  variables

# Correlation between dependent & independent variables

**Pearson** correlation ( $\rho$ ) measures strength and direction (+/-) of linear association between  $X_i$  and  $Y$ . Ranges from -1 to 1.

If relationship looks non-linear (but monotonic) then **Spearman** correlation should be considered.

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

Valid if joint distribution of  $X, Y$  is bivariate normal

# Linear Regression Model

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

where  $Y$  is a continuous outcome

## Interpretation

$\beta_0$  = Mean  $Y$  when  $X_i$  values are 0

$\beta_i$  = Mean change in  $Y$  for a 1-point increase in  $X_i$ ,  
adjusting for other  $X$  variables

# Linear Regression Assumptions

- Observations are independent
- Linearity
- Homoscedasticity
- Normality



# GLMs - Generalized Linear Models

$$g(\mu) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

where  $\mu = E(Y)$

## Common link functions

$$g(\mu) = \mu$$

$$g(\mu) = \log(\mu)$$

$$g(\mu) = \log(\mu / (1 - \mu))$$

# GLMs: Logistic Regression Model

$$\log\left(\overset{\text{odds}}{\frac{p}{1-p}}\right) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

where  $p$  = Probability of  $Y = 1$

$1-p$  = Probability of  $Y = 0$

## Interpretation

$\beta_i$  = log OR (Odds Ratio) for having  $Y = 1$  for a 1-point increase in  $X_i$ ,  
*adjusting for other predictors*

$e^{\beta_i}$  = OR for having  $Y = 1$  for a 1-point increase in  $X_i$

# Logistic Regression Assumptions

- Binary outcome (0, 1)
- Each predictor is linearly related to the log odds of the outcome

# Inference for Regression Modeling

## Confidence Interval

95% CI for  $\beta = (0.44, 0.49)$

## Hypothesis Testing

$H_0: \beta = \beta_0 \leftarrow$  Null Hypothesis

$H_A: \beta \neq \beta_0 \leftarrow$  Alternative Hypothesis

p-value: Chance that we are rejecting  $H_0$  when we should not be

# Goals

A photograph of a soccer game on a dirt field. A goalpost is visible on the right side of the field. Several players are on the field, including one in the foreground who is jumping or falling. The background shows a line of trees under a blue sky with some clouds. The word "Goals" is overlaid in the center of the image.

# Today's Goal

- Learn to use R to read in data and conduct regression analysis and associated inference tests
  - Checking assumptions
  - Visualizing
  - Computing p-values, regression coefficients, confidence intervals, and odds ratios (for logistic models)

# Today: 2 Scenarios

- Linear Regression (continuous outcome)
  - Single variable (1 continuous)
  - Multivariable (1 continuous, 1 categorical)
- Logistic Regression (categorical outcome)
  - Multivariable (continuous and categorical)

# Today's Data Set: Framingham Heart Study

- [framinghamheartstudy.org](https://www.framinghamheartstudy.org)
- Long-term prospective study of the etiology of cardiovascular disease among a population of subjects in Framingham, MA
- Began in 1948 with 5,209 subjects
- Is the source of the term "risk factor"
- Over 3,000 peer-reviewed papers published based on this study
- Participants were each followed for a total of 24 years for cardiovascular events (heart attack, stroke, death, etc.)



## Some Handy R Links

# Tutorials

- RStudio R paths: [education.rstudio.com/learn/](https://education.rstudio.com/learn/)
- Data Carpentry & Software Carpentry:
  - [datacarpentry.org](https://datacarpentry.org) and [software-carpentry.org](https://software-carpentry.org)
- LinkedIn Learning @ GW: [go.gwu.edu/linkedinlearning](https://go.gwu.edu/linkedinlearning)
- [r-tutor.com/r-introduction](https://r-tutor.com/r-introduction) & [r-tutor.com/elementary-statistics](https://r-tutor.com/elementary-statistics)
- UCLA Data Analysis Examples: [stats.idre.ucla.edu/other/dae/](https://stats.idre.ucla.edu/other/dae/)
- Visualizing regression results:  
[worldbank.github.io/r-econ-visual-library/RegressionCoef.html](https://worldbank.github.io/r-econ-visual-library/RegressionCoef.html)
- R Graph Gallery (w/code): [r-graph-gallery.com](https://r-graph-gallery.com)

# Books you can access for free

- Free books online - Hadley Wickham:
  - R for Data Science [r4ds.had.co.nz](https://r4ds.had.co.nz)
  - Advanced R [adv-r.hadley.nz/](https://adv-r.hadley.nz/)
- Through your GW library privileges:

ADVANCED SEARCH

Search for: ☐ Catalog + Articles ☒ Catalog ☐ Articles

Subject ▼ contains ▼ R (Computer programming language)

## Reference Links

- R language (CRAN): [r-project.org](https://r-project.org)
- R search engine: [rseek.org](https://rseek.org)
- [rstudio.com](https://rstudio.com)
  - Cheat Sheets! [rstudio.com/resources/cheatsheets](https://rstudio.com/resources/cheatsheets)
- [stackoverflow.com](https://stackoverflow.com)

# Statistics+R help @ GW

R-Statistics Appointments:

[calendly.com/statistical-consulting-gw](https://calendly.com/statistical-consulting-gw)

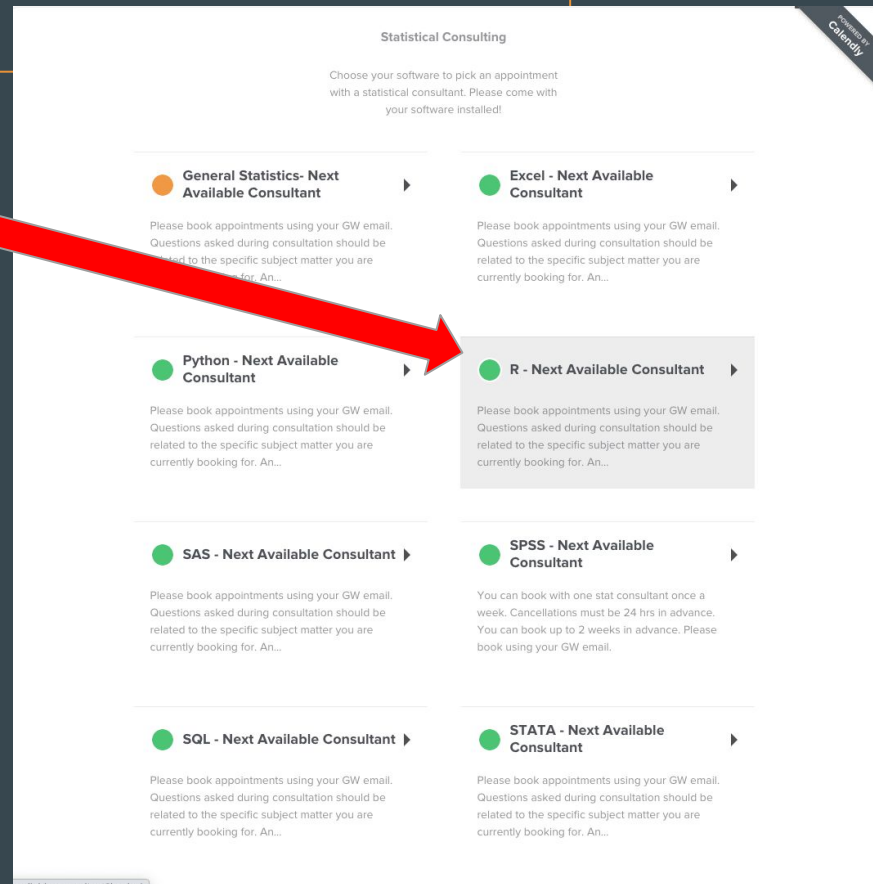
Also...

Appointments with me:

[calendly.com/kerchner](https://calendly.com/kerchner)

Coding consultations (Python, git, etc.):

[calendly.com/gwul-coding/](https://calendly.com/gwul-coding/)



# Thanks!

Dan Kerchner

[kerchner@gwu.edu](mailto:kerchner@gwu.edu)

Disovankiri Boung

[dboung@gwu.edu](mailto:dboung@gwu.edu)