# WELCOME TO NDACAN MONTHLY OFFICE HOURS!

#### NATIONAL DATA ARCHIVE ON CHILD ABUSE AND NEGLECT DUKE UNIVERSITY, CORNELL UNIVERSITY, & UNIVERSITY OF CALIFORNIA: SAN FRANCISCO





- The session will begin at 11am EST
  - 11:00 11:30am LeaRn with NDACAN (Introduction to R)
  - 11:30 12:00pm Office hours breakout sessions
- Please submit LeaRn questions to the Q&A box
- This session is being recorded.
- See ZOOM Help Center for connection issues: <u>https://support.zoom.us/hc/en-us</u>
  - If issues persist and solutions cannot be found through Zoom, contact Andres Arroyo at aa 17@cornell.edu.

# LEARN WITH NDACAN

Presented by Frank Edwards

2

## MATERIALS FOR THIS COURSE

- Course Box folder (<u>https://cornell.box.com/v/LeaRn-with-R-NDACAN-2024-2025</u>) contains
  - Data (will be released as used in the lessons)
    - Census state-level data, 2015-2019
    - AFCARS state-aggregate data, 2015-2019
    - AFCARS (FAKE) individual-level data, 2016-2019
    - NYTD (FAKE) individual-level data, 2017 Cohort
  - Documentation/codebooks for the provided datasets
  - Slides used in each week's lesson
  - Exercises as that correspond to each week's lesson
  - An .R file that will have example, usable R code for each lesson will be updated and appended with code from each lesson

# WEEK 8: INTRODUCTION TO ADVANCED MODELING

May 16, 2025



### DATA USED IN THIS WEEK'S EXAMPLE CODE

- AFCARS fake aggregated data ./data/afcars\_aggreg\_suppressed.csv
- AFCARS fake individual-level data ./data/afcars\_2016\_2019\_indv\_fake.csv
  - Simulated foster care data following the AFCARS structure
  - Can order full data from NDACAN:
    - https://www.ndacan.acf.hhs.gov/datasets/request-dataset.cfm
- Census data ./data/
  - Full data available from NIH / NCI SEER

# GENERALIZED LINEAR MODELS

### FROM LINEAR REGRESSION

We define the expected value of y in a linear regression model as

$$E[y_i] = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

### TO GENERALIZED LINEAR MODELS

We can define a generalized linear regression for the expectation of y as

$$g(E[y_i]) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

For a wide range of link functions g

### TWO COMMONLY USED LINK FUNCTIONS

Logistic regression

$$logit(E[y_i]) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

Poisson regression

$$log(E[y_i]) = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}$$

#### CHANGES TO THE ERROR STRUCTURE

Linear regression

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \varepsilon_i \\ \varepsilon_i \sim N(0,\sigma^2) \end{aligned}$$

Logistic regression

$$p = logit^{-1}(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i})$$
  
$$y_i \sim Bernoulli(p)$$

Poisson regression

$$\begin{aligned} \lambda_i &= e^{\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i}} \\ & y \sim Poisson(\lambda) \end{aligned}$$

### ESTIMATING GLMS IN R

We specify the distributional 'family' to be used with the glm() function

Logistic regression:  $glm(y \sim x, data = mydata, family = "binomial")$ Poisson regression:  $glm(y \sim x, data = mydata, family = "poisson")$ 

# MULTILEVEL MODELS

### FROM LINEAR REGRESSION

We can define a multilevel model with varying intercepts for cluster j as

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \delta_j + \varepsilon_i \\ & \varepsilon \sim N(0, \sigma^2) \\ & \delta \sim N(0, \sigma^2) \end{aligned}$$

### ESTIMATION IN R

The Ime4 package provides a flexible set of multilevel models

Linear, multilevel slopes:

 $lmer(y \sim x + (|j|), data = mydata)$ 

Logistic, multilevel slopes and intercepts:

 $glmer(y \sim x + (i|j), data = mydata, family = "binomial")$ 

Cheat sheet for Ime4 formulas:

https://stats.stackexchange.com/questions/13166/rs-Imer-cheat-sheet

# OVER TO RSTUDIO