# 36-617: Applied Linear Models

Introduction Brian Junker 132E Baker Hall brian@stat.cmu.edu

#### **Classes Will Be Recorded**

- I plan to record and share these classes with Zoom with you for later viewing/review.
  - They will be available at cmu.canvas.edu
  - I'm sure the first few lectures will be pretty choppy, since I am just getting used to the tech in this classroom!
- Watching recorded lectures does not replace being in class.
  - There is a "participation" component in your grade
  - If at any time you cannot be physically in class (travel problems, illness, etc.) I can give you a zoom link to join class remotely.

# Outline

- Introduction
- Course Schedule & Syllabus Stuff
- Valid vs Useful
- Quick Review of Univariate Regression
- R!
- Reading:
  - Read/skim Ch's 1-3 of Sheather
  - Look at Ch 5 of Sheather more closely
  - Supplemental:
    - G&H CH 3
    - ISLR 3.1,3.2, 3.3.1,3.3.2

## Introduction – About Us

Instructor
 Brian Junker
 brian@stat.cmu.edu

<u>TA</u> Lorenzo Tomaselli Itomasel@stat.cmu.edu

- Office hours:
  - MW Noon-1pm
  - 132E Baker Hall
  - ...or by appt(in person or Zoom)
- Office hours:

  - □ ...or by appt

# Introduction – About The Course

- <u>Technical material</u>: The machinery of linear regression and its generalizations
- Disposition: When is a model adequate?
- <u>Translation:</u> ABA<sup>-1</sup>:
  - A: Translate from real world to quantitative question
  - B: Answer quantitative question using Statistics
  - A<sup>-1</sup>: Translate back to real word
  - Lather, rinse, repeat...
- Communication:
  - IMRaD -> IDMRaD
  - Clear sentences, paragraphs, sections.

# Introduction – Course Materials

#### Technical material:

- Sheather (2009). A Modern Approach to Regression with R. NY: Springer \*
- James et al. (2013). Introduction to Statistical Learning with R. NY: Springer \*

#### Supplementary texts:

- Gelman, A. & Hill, J. (2007). Data Analysis Using Regression and Multilevel/Hierarchical Models. NY: Cambridge Univ Press.
- Lynch, S. M. (2007). Introduction to applied Bayesian statistics and estimation for social scientists. NY: Springer \*

#### Software:

- R, Rstudio, Stan (more on the next slide)
- □ LaTeX (If you are new to TeX/LaTeX, try overleaf.com–free to CMU)
- <u>Help:</u> I will post some links, but get used to googling!

# Introduction - Computing

- We'll mostly be working in R, RStudio and extensions of R (primarily stan).
  - □ R: <u>https://cran.r-project.org/</u>
  - RStudio: <u>https://www.rstudio.com/</u>
    - Rmarkdown: (just try it in RStudio, use google for help)
  - Stan: <u>https://mc-stan.org/users/interfaces/rstan</u>
- In class I'll use "raw R" more than RStudio, but you will find RStudio and Rmarkdown very convenient for hw, etc.
- It is possible to do some of the work in Python, but the lectures, technical work and project are strongly geared to facilities in R.

# Introduction – Online Resources

- Online resources:
  - I will record lectures and make them available on Canvas.
  - Canvas (canvas.cmu.edu)
    - All course materials
    - Your grades
    - Monday quizzes
    - Submit and peer review project papers
  - Gradescope (within Canvas)
    - Submit weekly homework
  - Piazza (within Canvas)
    - Great for asking (and answering) questions outside of class
    - The TA and I will also monitor Piazza

#### General schedule for the semester

Week	Dates	Tentative Topics	Tentative Sources
	Aug 22–26	Getting Ready	Sheather, Ch 1, 2, 3
Week 1	Aug 29, 31	Intro, Multiple Regression	Sheather, Ch 5
			ISLR 3.1, 3.2
			G&H Ch 3
Week 2	Sep 5 (no class <sup>1</sup> ), Sep 7	Qualitative Predictors	ISLR, 3.3.1, 3.3.2, handouts
Week 3	Sep 12, 13	Diagnostics & Transformations	Sheather, Ch 6
			ISLR 3.3.3
			G&H Ch 4
Week 4	Sep 19, 21	Variable Selection	Sheather Ch 7
			ISLR Ch 6
			G&H Ch 4
Week 5	Sep 26, 28	Logistic Regression	Sheather Ch 8
		Take-Home Midterm Assigned	ISLR 4.3
			G&H Ch 5
Week 6	Oct 3, 5	Generalized Linear Models	Sheather Ch 8
		Take-Home Midterm Due	G&H Ch 6
Week 7	Oct 10, 12	Nonparametric Regression	ISLR Ch 7
	00110,12	Nonparametric Regression	Sheather Appx
			handouts
	Oct 17-21		
		FALL BREAK	
Week 8	Oct 24, 26	Causal Reasoning	G&H Ch 9, 10
Week 9	Oct 31, Nov 2	Multilevel and Mixed Effects Models	Sheather, 10.1
			Intercepts: G&H Ch 12
		Project assigned	Slopes: G&H Ch 13
Week 10	Nov 7, 9	Multilevel logistic regression & GLMs	G&H parts of Ch's 14 & 15
Week 11	Nov 14, 16	Residuals, Estimation and Model Se-	Handouts; stuff from G&H
WCCK II	1107 14, 10	lection	Handouts, stun Hom O&H
Week 12	Nov 21, 23 (no class) <sup>2</sup>	Bayes & Shrinkage	Lynch Ch 3, 4
HOUR ID	1107 21, 25 (no cidas)	Dayto to bir inkuge	G&H parts of Ch 16
Week 13	Nov 28, 30	STAN and MLM's	Lynch, Ch 9
HOUR IS	1107 20, 50		G&H parts of Ch 17
		Project due	-
Week 14	Dec 5, 7	MLM's with STAN	Handouts; maybe parts of
			G&H Ch's 18, 21

<sup>1</sup>No Class Sep 5: Labor Day (US Holiday).

<sup>2</sup>No Class Nov 23: Thanksgiving (US Holiday) Nov 24.

# Syllabus Stuff – Work & Rules

- 20%: 10-ish HW's
  - Please feel free to work with each other on hw;
     BUT you must list who you worked with.
- 10%: Monday Quizzes (on weekly reading/materials)
- 25%: Take-Home Midterm
- 25%: Final Report Project
- 10% Peer review of projects
- 10% Participation (Do I remember your name? What you did in class? In office hours? On Piazza?)
- Credit where credit is due
  - Please list any <u>person</u> or any <u>source</u> you consulted in doing your work, in a list of references at the end of hw, project, take-home
- All hw will be submitted via Canvas (Gradescope)
  - Generally we will not accept late hw or late take-homes...

# Reading & HW

- The first HW assignment will be available on Canvas later today (due next week).
- Reading:
  - You should already have read Sheather Ch's 1-3
  - For this week read Sheather Ch 5
    - ISLR 3.1 & 3.2, and G&H CH 3 are good supplemental reading
  - □ For next week: Sheather 5.3, and ISLR 3.3.1 & 3.3.2
- Pdf's for lecture notes etc. in the file area on Canvas
- In general, you should do the reading before each week's classes.

### Introduction – Level

- Hopefully you have seen calculus-based prob & stat, matrix algebra, and a little linear regression.
   We need to talk like statisticians!
- You have all different levels of experience with
   Applied regression and statistical modeling
   R
  - Writing scientific reports
- Fill in the gaps
  - Learn on your own (Google)! Help each other!
  - Ask Lorenzo and me!

Let's take a break and think about these two quotes...

- "[I]t makes sense to base inferences or conclusions only on valid models"
   S.J. Sheather (2007)
- "All models are wrong but some are useful"
   G.E.P. Box (1978)

## Linear Regression – The Model

• Let 
$$X_i = (x_{i0}, ..., x_{ip})$$
 and  $\beta = (\beta_0, ..., \beta_p)^T$ ; then  
 $y_i = \beta_0 X_{i0} + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \epsilon_i$   
 $= X_i \beta + \epsilon_i$ 

• If we also stack  $Y = (y_1, ..., y_n)^T$ ,  $X = (X_1^T, ..., X_n^T)^T$ , and  $\epsilon = (\epsilon_1, ..., \epsilon_n)^T$ , we can write

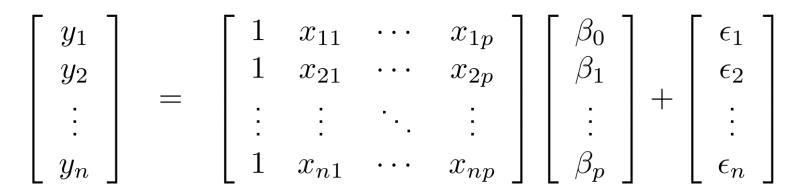
$$Y_{n \times 1} = X_{n \times k} \beta_{k \times 1} + \epsilon_{n \times 1} \quad (k = p+1)$$

#### Linear Regression – The Model

$$Y = X\beta + \epsilon$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{10} & x_{11} & \cdots & x_{1p} \\ x_{20} & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n0} & x_{n1} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

#### Usually $x_{i0} \equiv 1$ , so we get



## Linear Regression – The Model

In the model

$$y_i = X_i\beta + \epsilon_i, i = 1, \dots, n$$

it is usual to assume  $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$ 

Recall

• Y ~ N(0,1) iff 
$$f(y) = \frac{1}{\sqrt{2\pi}} e^{-y^2/2}$$

$$\hfill Y \, {}^{\sim}\, {\rm N}(\mu, \sigma^2)$$
 iff  $\frac{y-\mu}{\sigma} \sim N(0,1)$ 

$$f(y) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

**Digression – Multivariate Normal**  
• 
$$\mathbf{Y} = (\mathbf{Y}_1, \dots, \mathbf{Y}_n)^{\mathsf{T}} \sim \mathsf{N}(\mathsf{O},\mathsf{I}) = N \begin{pmatrix} \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \cdots & 1 \end{bmatrix} \end{pmatrix}$$

$$f(y_1, \dots, y_n) = \prod_{i=1}^n f(y_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} e^{-y_i^2/2}$$

• Y ~ N( $\mu$ ,  $\Sigma$ ) iff  $\Sigma^{-1/2}(Y - \mu) \sim N(0, I)$ ...and some ugly formula for f(y<sub>1</sub>, ..., y<sub>n</sub>)...

### Digression – Multivariate Normal

• When Y ~ N( $\mu$ , $\Sigma$ ), then

$$\mu = (\mu_1, \mu_2, \dots, \mu_n)^T = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix}$$

$$\mathsf{Var}(Y) = \Sigma = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1n} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1} & \sigma_{n2} & \cdots & \sigma_n^2 \end{bmatrix}$$

is the *variance-covariance matrix*:

is the *mean vector*.

$$E[y_i] = \mu_i, \ i = 1, \dots, n$$

$$Var(y_i) = \sigma_i^2, i = 1, \dots, n$$
$$Cov(y_i, y_j) = \sigma_{ij}, i, j = 1, \dots, n$$

### Regresssion – ML/LS Estimates

- Y = X $\beta$  +  $\epsilon$ ,  $\epsilon \sim N(0, \sigma^2 I)$
- So Var(Y<sub>i</sub>) = E[(Y<sub>i</sub> X<sub>i</sub>  $\beta$ )<sup>2</sup>] = E[ $\epsilon_i^2$ ] = Var( $\epsilon_i$ ) =  $\sigma^2$
- Then we can estimate (MoM!):

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - X_i \beta)^2$$

- Fitting the model is basically just finding values  $\beta$ to minimize  $\frac{1}{n}RSS$ , i.e., minimize
  - $\frac{1}{n} \sum_{i=1}^{n} (y_i X_i \beta)^2 = \frac{1}{n} (Y X \beta)^T (Y X \beta)$

It turns out that 
$$\hat{\beta} = (X^T X)^{-1} X^T y$$

### Regression – ML/LS Estimates

$$y = X\beta + \epsilon, \ \epsilon \sim N(0, \sigma^2 I)$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$
  

$$\hat{y} = X \hat{\beta} = X (X^T X)^{-1} X^T y = Hy$$
  

$$\hat{e} = y - \hat{y} = (I - H)y$$
  
"Hat matrix"

- The "residual SD" is the square root of  $s^2 = \frac{1}{n-k} \sum_{i=1}^n (y_i - X_i \hat{\beta})^2 = \frac{1}{n-k} (y - X \hat{\beta})^T (y - X \hat{\beta})$
- With a little more matrix algebra, Recall that k=p+1  $Var(\hat{\beta}) = (X^T X)^{-1} \sigma^2$   $Var(\hat{y}) = X(X^T X)^{-1} X^T \sigma^2 = H \sigma^2$  $Var(\hat{e}) = (I - H) \sigma^2$

### Regression – Example

#### Demographic factors and income...

> heights <- read.dta("heights.dta")</pre>

> str(heights)

'data.frame': 2029 obs. of 9 variables: \$ earn : num NA NA 50000 60000 30000 NA 50000 NA 51000 ... \$ height1: int 556555555... \$ height2: num 6 4 2 6 4 5 3 8 3 4 ... : int 2 1 2 2 2 2 2 2 2 . . . \$ sex 1 : int 1 2 1 1 1 1 3 2 1 1 ... \$ race \$ hisp : int 2 2 2 2 2 2 2 2 2 2 . . . \$ ed 12 12 16 16 16 17 16 18 17 15 ... : num \$ yearbn : num 53 50 45 32 61 33 99 36 51 64 ... \$ height : num 66 64 74 66 64 65 63 68 63 64 ...

Regression – Example $\hat{\beta} = (X^T X)^{-1} X^T y$								
> summary(lm(earn ~ height + ed, data=heights)) $SE(\hat{\beta}) = \sqrt{\operatorname{diag}(\widehat{\operatorname{Var}}(\hat{\beta}))} = \sqrt{\operatorname{diag}((X^TX)^{-1}s^2)}$								
Coefficients:								
	Estimate St	d. Error t	t value B	r(> t )				
(Intercept)	-106263.5	8564.7	-12.41	<2e-16 ***				
height	1376.7	126.7	10.87 🗲	$\frac{2e-16}{16} + \frac{1}{16} t_{j} = \hat{\beta}_{j} / SE$	$(\hat{\beta})$			
ed	2590.8	197.7	13.11	<2e-16 *** $l_j = \rho_j / \beta L$	$(\rho_j)$			

Residual standard error: 17780 on 1376 degrees of freedom
 (650 observations deleted due to missingness)
Multiple R-squared: 0.1915, Adjusted R-squared: 0.1904
F-statistic: 163 on 2 and 1376 DF, p-value: < 2.2e-16</pre>

earn  $\approx \hat{\beta}_0 + \hat{\beta}_1(\text{height}) + \hat{\beta}_2(\text{ed})$  = -106263.5 + 1376.7(height) + 2590.8(ed)  $\frac{1}{n-k}RSS = s^2 = (17780)^2 = 316128400$   $R^2 = \frac{\widehat{\text{Var}}(\hat{y})}{\widehat{\text{Var}}(y)} = \frac{SS_{reg}}{SSY} = 1 - \frac{RSS}{SSY} = 0.1915$  k = p+1 = 3 ; n-k = 1376n = 1379

#### Regression - Example

- Continuing in R...
- heights.r in the week01 folder under "Files" on Canvas...

# R!

 Some people learning R for the first time; others have done extensive data analysis projects in R.

#### Some references on R:

- http://www.cookbook-r.com
- Quick-R: <u>https://www.statmethods.net/</u>
- Online course: <u>https://www.datacamp.com/courses/free-introduction-to-r</u>
- Lately I really like <a href="https://kieranhealy.org/publications/dataviz/">https://kieranhealy.org/publications/dataviz/</a>
- If you have not used R before...
  - http://www.cs.cmu.edu/~10702/R2/Rintro.pdf provides a good start
    - If you are new to R, type into R all the commands and examples in rintro.pdf
    - If you have worked with R before, read through rintro.pdf and try to predict what would happen with each command. If you are not sure, type in that command/example.

# Summary

- Introduction
- Course Schedule & Syllabus Stuff
- Valid vs Useful
- Quick Review of Univariate Regression
- R!
- Reading:
  - Read/skim Ch's 1-3 of Sheather
  - Look at Ch 5 of Sheather more closely
  - Supplemental:
    - G&H CH 3
      - ISLR 3.1,3.2, 3.3.1,3.3.2