# Week 8: What Can Go Wrong and How To Fix It, Diagnostics and Solutions

Brandon Stewart<sup>1</sup>

Princeton

October 19-23, 2020

<sup>1</sup>These slides are heavily influenced by Matt Blackwell, Adam Glynn, Jens Hainmueller, Erin Hartman and Kevin Quinn.

Stewart (Princeton)

Week 8: Diagnostics and Solutions

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference

#### Thinking About Problems

#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points
- Robust Regression Methods
   Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# Pulling Back the Curtain

- There is a pedagogical tension between teaching you to be conversant in the field as it is now and preparing you for methods of the future.
- We skew a bit more towards the latter, because most people won't continue to study methods long term—I want your knowledge to stay current as long as possible!
- Towards the beginning of the course there is a lot of material that is simply foundational material you simply need to know, but we are starting to get into the practice of statistics and research which means fewer black and white answers.
- This week we will talk about how to diagnose and fix problems which is naturally a pretty context-specific activity.
- Today is setting the table for how we think about that problem. I think this is philosophically interesting and you may want to revisit at the end of the week.

## **Five Themes**

- (1) Clearly define your goal.
- (2) Examine your model.
- (3) Diagnosis through treatment.
- (4) Don't expect a free lunch.
- (5) Re-examine defaults.

# (1) Clearly Define Your Goal

- Best practices are rarely best for all applications, people generally have a distribution of applications in mind when they give advice so best to be specific about your application.
- When thinking about whether an estimator is biased, always think 'biased with respect to what?'
- Predictive generalization is one unifying way to solve questions about what is best (particularly in machine learning).
- You may not know how to precisely define your goal yet (last few weeks are about causal goals). That's okay!



# Residuals are important. Look at them.

# (2) Examine Your Model

- There are so many ways that something can go wrong, particularly when you don't know every nuance of the data.
- Examining the model helps us detect cases where assumptions have failed to hold or data is contaminated.
- A challenge is that this formally disrupts most inference infrastructure which assumes you fix models a priori and only do one test.
- Double checking that that things haven't gone horribly wrong is most likely okay, but there is always a slippery slope argument that leads to what Gelman calls the 'garden of forking paths.'
- There is always an inherent risk in looking at the data more than once. 'Double dipping' can be a serious problem.
- One way out of this we won't discuss is train-test splits.

# (3) Diagnosis Through Treatment

- When diagnosing problems it can often be hard to tell how consequential they are.
- One strategy is diagnosis through treatment, i.e. we just use a procedure that would address the problem and see if it changes our conclusions rather than diagnosing in the first place.
- We will use this strategy to replace formal diagnostics in a few different places this week.
- The downside to this approach is that sometimes you end up less clear about what the problem was in the first place.
- We will talk about successes of this approach but also some catastrophic failures.

# (4) Don't Expect a Free Lunch

- Most things I show you will have tradeoffs.
- This shouldn't be a surprise—if we had a procedure that uniformly dominated all others I wouldn't show you options, I'd just show you the one.
- Common tradeoffs include:
  - bias vs. variance
  - interpretability vs. flexibility
  - data dependence vs. assumption dependence

# (5) Re-Examine Defaults

- In many ways statistics is the science of defaults.
- Defaults are hugely important because people rarely change defaults.
- For example, arguably we all use classical standard errors because that's what lm() produces. This is way packages like estimatr use robust standard errors by default.
- But remember, different solutions are better for different settings—defaults are a function of their time and context.
- The history of 20th century social science is defined by scarcity—data was hard to find, surveys were costly to field, computing resources were expensive and difficult to access. The methods we use are a consequence of that scarcity.
- But now we have abundance—new forms of data, cheap surveys, huge computation! It is changing the methods we consider.

# Topics We Will Cover

- Non-Normal Errors
- Extreme Values
- Robust Regression
- Non-linearity
- Clustering

Regrettably we won't have time to cover two important areas: missing data and sensitivity analysis.

• Five themes for thinking about problems.

Next Time: Non-normal Errors!

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference



#### Non-Normality

#### Extreme Values

- Outliers
- Leverage Points
- Influence Points
- Robust Regression Methods
   Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# Review of the Normality assumption

• In matrix notation:

$$\mathbf{u}|\mathbf{X}\sim\mathcal{N}(\mathbf{0},\sigma_{u}^{2}\mathbf{I})$$

• Equivalent to:

$$u_i | \mathbf{x}'_i \sim \mathcal{N}(\mathbf{0}, \sigma_u^2)$$

• Fix  $\mathbf{x}'_i$  and the distribution of errors are Normal.

## Consequences of non-Normal Errors?

- In small samples:
  - Sampling distribution of  $\widehat{\beta}$  will not be Normal
  - Test statistics will not have t or F distributions
  - Probability of Type I error will not be  $\alpha$
  - ▶  $1 \alpha$  confidence interval will not have  $1 \alpha$  coverage
- In large samples:
  - Sampling distribution of  $\widehat{eta} pprox$  Normal by the CLT
  - Test statistics will be  $\approx t$  or F by the CLT
  - Probability of Type I error  $pprox \alpha$
  - $1 \alpha$  confidence interval will have  $\approx 1 \alpha$  coverage
- The sample size (*n*) needed for approximation to hold depends on how far the errors are from Normal.
- Reasonable question: if we have enough data are non-normal errors even a problem?

# Clarifying a Point of Confusion: Marginal versus Conditional

- Be careful with this assumption: distribution of the error, not the distribution of the outcome is the key assumption
- The marginal distribution of *y* can be non-Normal even if the conditional distribution is Normal!
- The plausibility depends on the X chosen by the researcher.

# Example: Is this a Violation?



# Example: Is this a Violation?



# How to Diagnose?

- Assumption is about unobserved errors  $\mathbf{u} = \mathbf{y} \mathbf{X}\boldsymbol{\beta}$
- We can only observe residuals,  $\widehat{\mathbf{u}} = \mathbf{y} \mathbf{X}\widehat{\boldsymbol{eta}}$
- If distribution of residuals  $\approx$  distribution of errors, we could check residuals as a proxy for the errors.
- Unfortunately, this is not true—the distribution of the residuals is more complicated

Solution: Carefully investigate the residuals numerically and graphically.

To understand the relationship between residuals and errors, we need to derive the distribution of the residuals (which we will do over the next few slides).

# Defining the Hat Matrix

- We want to figure out how the distribution of errors  $\boldsymbol{u}$  relates to the distribution of residuals  $\widehat{\boldsymbol{u}}.$
- To get there let's write  $\hat{\mathbf{u}}$  in terms of y, then we will be able to replace  $\mathbf{y}$  with  $\mathbf{X}\boldsymbol{\beta} + \mathbf{u}$ .

$$\begin{aligned} \widehat{\mathbf{u}} &= \mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}} \\ &= \mathbf{y} - \mathbf{X} \left( \mathbf{X}' \mathbf{X} \right)^{-1} \mathbf{X}' \mathbf{y} \\ &\equiv \mathbf{y} - \mathbf{H} \mathbf{y} \\ &= (\mathbf{I} - \mathbf{H}) \mathbf{y} \end{aligned}$$

•  $\mathbf{H} = \mathbf{X} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'$  is the hat matrix because it puts the "hat" on y:

$$\widehat{\mathbf{y}} = \mathbf{H}\mathbf{y}$$

- it has a few nice properties:
  - **H** is an  $n \times n$  symmetric matrix
  - **H** is idempotent: HH = H

## Relating the Residuals to the Errors

With the hat matrix, we are ready to relate the residuals to the errors.

$$\begin{aligned} \widehat{\mathbf{u}} &= (\mathbf{I} - \mathbf{H})(\mathbf{y}) \\ &= (\mathbf{I} - \mathbf{H})(\mathbf{X}\boldsymbol{\beta} + \mathbf{u}) \\ &= (\mathbf{I} - \mathbf{H})\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \mathbf{H})\mathbf{u} \\ &= \mathbf{I}\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\left(\mathbf{X}'\mathbf{X}\right)^{-1}\mathbf{X}'\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \mathbf{H})\mathbf{u} \\ &= \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \mathbf{H})\mathbf{u} \\ &= (\mathbf{I} - \mathbf{H})\mathbf{u} \end{aligned}$$

- $\bullet$  Residuals  $\widehat{u}$  are a linear function of the errors, u.
- For instance,

$$\widehat{u}_1 = (1 - h_{11})u_1 - \sum_{i=2}^n h_{1i}u_i$$

• Note that each residual is a function of all of the errors.

# Characterizing the Distribution of the Residuals

What can we say about the distribution of the residuals now that we have the expression:  $\hat{u} = (I - H)u$ .

$$E[\hat{\mathbf{u}}] = (\mathbf{I} - \mathbf{H})E[\mathbf{u}] = \mathbf{0}$$
$$V[\hat{\mathbf{u}}] = \sigma_u^2(\mathbf{I} - \mathbf{H})$$

The variance of the *i*th residual  $\hat{u}_i$  is  $V[\hat{u}_i] = \sigma^2(1 - h_{ii})$ , where  $h_{ii}$  is the *i*th diagonal element of the matrix **H** (called the hat value).

# Properties of the Distribution of Residuals

Notice in contrast to the unobserved errors, the estimated residuals have some different properties. They

- are not independent (because they must satisfy the two constraints  $\sum_{i=1}^{n} \hat{u}_i = 0$  and  $\sum_{i=1}^{n} \hat{u}_i x_i = 0$ )
- do not have the same variance. The variance of the residuals varies across data points V[û<sub>i</sub>] = σ<sup>2</sup>(1 - h<sub>ii</sub>), even though the unobserved errors all have the same variance σ<sup>2</sup>

These properties make it hard to learn about the errors (which our assumptions are about and we don't have access to) from our residuals (which we have estimated and can examine). This is inconvenient for diagnostics.

What if we could transform the residuals to address the two issues above?

### Standardized Residuals

Let's address the second problem (unequal variances) by standardizing  $\hat{u}_i$ , i.e., dividing by unit *i*'s estimated standard deviation.

This produces standardized (or "internally studentized") residuals:

$$\hat{u}'_i = rac{\hat{u}_i}{\hat{\sigma}\sqrt{1-h_{ii}}}$$

where  $\hat{\sigma} = \sqrt{\hat{\sigma}^2}$  and  $\hat{\sigma}^2 = \frac{\hat{u}'\hat{u}}{n-(k+1)}$  is our usual estimate of the error variance.

The standardized residuals are still not ideal, since the numerator and denominator of  $\hat{u}'_i$  are not independent. This makes the distribution of  $\hat{u}'_i$  nonstandard. If the distribution is non-standard, we can't easily check for violations.

## Studentized Residuals

If we remove observation *i* from the estimation of  $\sigma^2$ , then we can eliminate the dependence and the result will have a standard distribution.

• estimate error variance without residual *i*:

$$\widehat{\sigma}_{-i}^2 = rac{\widehat{\mathbf{u}}'\widehat{\mathbf{u}} - \widehat{u}_i^2/(1-h_{ii})}{n-k-2}$$

• Use this *i*-free estimate to standardize, which creates the studentized residuals:

$$\widehat{u}_i^* = \frac{\widehat{u}_i}{\widehat{\sigma}_{-i}\sqrt{1-h_{ii}}}$$

- If the errors are Normal, the studentized residuals,  $\hat{u}_i^*$ , follow a t distribution with (n k 2) degrees of freedom.
- Deviations from this *t* distribution of the residuals imply violation of Normality in the errors.

## Example: Buchanan votes in Florida, 2000

Wand et al. show that the ballot caused 2,000 Democratic voters to vote by mistake for Buchanan, a number more than enough to have tipped the vote in FL from Bush to Gore, thus giving him FL's 25 electoral votes and the presidency.



FIGURE 1. The Palm Beach County Bufferfly Ballot

### Example: Buchanan votes in Florida, 2000



## Example: Buchanan Votes in Florida

- Now that our studentized residuals follow a known standard distribution, we can proceed with diagnostic analysis for the nonnormal errors.
- We examine data from the 2000 presidential election in Florida used in Wand et al. (2001).
- Our analysis takes place at the county level and we will regress the number of Buchanan votes in each county on the total number of votes in each county.

# Buchanan Votes and Total Votes

```
____ R Code _____
> mod1 <- lm(buchanan00~TotalVotes00,data=dta)</pre>
> summary(mod1)
Residuals:
   Min 10 Median 30 Max
-947.05 -41.74 -19.47 20.20 2350.54
Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 5.423e+01 4.914e+01 1.104 0.274
TotalVotes00 2.323e-03 3.104e-04 7.483 2.42e-10 ***
Residual standard error: 332.7 on 65 degrees of freedom
Multiple R-squared: 0.4628, Adjusted R-squared: 0.4545
F-statistic: 56 on 1 and 65 DF, p-value: 2.417e-10
> residuals <- resid(mod1)</pre>
> standardized residuals <- rstandard(mod1)</pre>
> studentized residuals <- rstudent(mod1)</pre>
> dotchart(residuals,dta$name,cex=.7,xlab="Residuals")
```

## Plotting the residuals



### Plotting the residuals



Histogram of student.resids



student.resids

Stewart (Princeton)

### Plotting the residuals



October 19-23, 2020

31 / 127

# Quantile-Quantile plots

- How can we easily compare our actual distribution of residuals to the theoretical distribution?
- Quantile-quantile plot or QQ-plot is useful for comparing distributions
- Plots the quantiles of one distribution against those of another distribution
- For example, one point is the  $(m_x, m_y)$  where  $m_x$  is the median of the x distribution and  $m_y$  is the median for the y distribution
- If distributions are equal  $\implies$  45 degree line

# Good QQ-plot



t quantiles
# Buchanan QQ-plot



# How Can we Deal with non-Normal Errors?

- Drop or change problematic observations (could be a bad idea unless you have some reason to believe the data are wrong or corrupted)
- Add variables to **X** (remember that the errors are defined in terms of explanatory variables)
- Use transformations (this may work, but a transformation affects all the assumptions of the model)
- Use estimators other than OLS that are robust to nonnormality (two videos from now!)

### **Buchanan Revisited**

Let's delete Palm Beach and also use log transformations for both variables

```
##
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.48597 0.37889 -6.561 1.09e-08 ***
## log(edaytotal) 0.70311 0.03621 19.417 < 2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4362 on 64 degrees of freedom
## Multiple R-squared: 0.8549, Adjusted R-squared: 0.8526
## F-statistic: 377 on 1 and 64 DF, p-value: < 2.2e-16</pre>
```

#### **Buchanan Revisited**



Histogram of resids.nopb

Histogram of stand.resids.nopb



Histogram of student.resids.nopb



### **Buchanan Revisited**



# A Note of Caution About Log Transformations

- Log transformations are a standard approach in the literature and intro regression classes
- They are extremely helpful for data that is skewed (e.g. a few very large positive values)
- Generally you want to convert these findings back to the original scale for interpretation
- Remember the complexities of log transforms from Week 5!

# We Covered

- Non-Normal Data
- Hat Matrix
- Transforming Residuals

#### Next Time: Extreme Values

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference



#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points

#### Robust Regression Methods

• Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# The Trouble with Norway

- Lange and Garrett (1985): organizational and political power of labor interact to improve economic growth
- Jackman (1987): relationship just due to North Sea Oil?
- Table guide:
  - $x_1 =$ organizational power of labor
  - $x_2 = \text{political power of labor}$
  - Parentheses contain t-statistics

	Constant	<i>x</i> <sub>1</sub>	<i>x</i> <sub>2</sub>	$x_1 \cdot x_2$
Norway Obs Included	.814	192	278	.137
	(4.7)	(2.0)	(2.4)	(2.9)
Norway Obs Excluded	.641	068	138	.054
	(4.8)	(0.9)	(1.5)	(1.3)

# Creative Curve Fitting with Norway



# The Most Important Lesson: Check Your Data

"Do not attempt to build a model on a set of poor data! In human surveys, one often finds 14-inch men, 1000-pound women, students with 'no' lungs, and so on. In manufacturing data, one can find 10,000 pounds of material in a 100 pound capacity barrel, and similar obvious errors.

All the planning, and training in the world will not eliminate these sorts of problems. In our decades of experience with 'messy data,' we have yet to find a large data set completely free of such quality problems."

Draper and Smith (1981, p. 418)

#### Carefully Examine the Data First!!

- Examine summary statistics: summary(data)
- Scatterplot matrix for densities and bivariate relationships:
   E.g. scatterplotMatrix(data) from car library.
- Further conditional plots for multivariate data: E.g. ggplot2

# Three Types of Extreme Values

- Outlier: extreme in the y direction
- 2 Leverage point: extreme in one x direction
- Influence point: extreme in both directions
- Not all of these are problematic
- If the data are truly "contaminated" (come from a different distribution), can cause inefficiency and possibly bias
- Can be a violation of iid (not identically distributed)

# **Outlier** Definition



- An **outlier** is a data point with very large regression errors,  $u_i$
- Very distant from the rest of the data in the y-dimension
- Increases standard errors (by increasing  $\widehat{\sigma}^2$ )
- No bias if typical in the x's

# **Detecting Outliers**

- Look at standardized residuals,  $\hat{u}'_i$ ?
  - but  $\widehat{\sigma}^2$  could be biased upwards by the large residual from the outlier
  - Makes detecting residuals harder
- Possible solution: use studentized residuals

$$\widehat{u}_i^* = \frac{\widehat{u}_i}{\widehat{\sigma}_{-i}\sqrt{1-h_i}}$$

•  $\widehat{\sigma}>\widehat{\sigma}_{-i}$  because we drop the large residual from the outlier, and so  $\widehat{u}_i'<\widehat{u}_i^*$ 

# Cutoff Rules for Outliers

- The studentized residuals follow a t distribution,  $u_i^* \sim t_{n-k-2}$ , when  $u_i \sim N(0, \sigma^2)$
- Rule of thumb:  $|\widehat{u}_i^*| > 2$  will be relatively rare
- Extreme outliers,  $|\widehat{u}_i^*| > 4-5$  are much less likely
- People usually adjust cutoff for multiple testing

#### **Buchanan outliers**



#### What to do about outliers

- Is the data corrupted?
  - Fix the observation (obvious data entry errors)
  - Remove the observation
  - Be transparent either way
- Is the outlier part of the data generating process?
  - Transform the dependent variable (log(y))
  - Use a method that is robust to outliers (robust regression)
- Key question is what is the goal? If you want to estimate the expectation of a distribution and a property of that distribution is extreme observations, that's just part of the story.

A Cautionary Tale: The "Discovery" of the Ozone Hole

- In the late 70s, NASA used an automated data processing program on satellite measurements of atmospheric data to track changes in atmospheric variables such as ozone.
- This data "quality control" algorithm rejected abnormally low readings of ozone over the Antarctic as unreasonable.
- This delayed the detection of the ozone hole by several years until British Antarctic Survey scientists discovered it based on analysis of their own observations (*Nature*, May 1985).
- The ozone hole was detected in satellite data only when the raw data was reprocessed. When the software was rerun without the pre-processing flags, the ozone hole was seen as far back as 1976.

# A Sociological Cautionary Tale

Comment on Herring, ASR, April 2009

# **Does Diversity Pay? A Replication of Herring (2009)**



American Sociological Review 2017, Vol. 82(4) 857–867 © American Sociological Association 2017 DOI: 10.1177/0003122417714422 journals.sagepub.com/home/asr



#### Dragana Stojmenovska,<sup>a</sup> Thijs Bol,<sup>a</sup> and Thomas Leopold<sup>a</sup>

#### Abstract

In an influential article published in the *American Sociological Review* in 2009, Herring finds that diverse workforces are beneficial for business. His analysis supports seven out of eight hypotheses on the positive effects of gender and racial diversity on sales revenue, number of customers, perceived relative market share, and perceived relative profitability. This comment points out that Herring's analysis contains two errors. First, missing codes on the outcome variables are treated as substantive codes. Second, two control variables—company size and establishment size—are highly skewed, and this skew obscures their positive associations with the predictor and outcome variables. We replicate Herring's analysis correcting for both errors. The findings support only one of the original eight hypotheses, suggesting that diversity is nonconsequential, rather than beneficial, to business success.

#### A Sociological Cautionary Tale

In our correspondence with Herring, he did not offer a definitive explanation for these discrepancies, but indicated that he may have treated all codes other than "not applicable" (-999) as substantive codes. Given (1) the large difference between his sample size and the number of valid observations in the NOS, and (2) the large number of missing values due to reasons other than "not applicable"in particular for sales revenue and number of customers-this coding error appears likely to account for much of the discrepancies. This means, for example, that 206 business organizations in which the sales revenue was unknown were treated as if they had sales of 88,888,888,888 US Dollars. Yet, even when we replicated this error (i.e., keeping all organizations with missing values other than -999 in our sample), we were unable to recover Herring's sample sizes, although the differences were smaller.

Week 8: Diagnostics and Solutions

# Leverage Point Definition



- Values that are extreme in the x direction
- That is, values far from the center of the covariate distribution
- Decrease SEs (more X variation)
- No bias if typical in y dimension

#### Leverage Points: Hat values

To measure leverage in multivariate data we will go back to the hat matrix H:

$$\hat{\mathbf{y}} = \mathbf{X}\hat{oldsymbol{eta}} = \mathbf{X}\left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}$$

**H** is  $n \times n$ , symmetric, and idempotent. It generates fitted values as follows:

$$\hat{y}_i = \mathbf{h}'_i \mathbf{y} = \begin{bmatrix} h_{i,1} & h_{i,2} & \cdots & h_{i,n} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{j=1}^n h_{i,j} y_j$$

Therefore,

- $h_{ij}$  dictates how important  $y_j$  is for the fitted value  $\hat{y}_i$  (regardless of the actual value of  $y_j$ , since **H** depends only on **X**)
- The diagonal entries  $h_{ii} = \sum_{j=1}^{n} h_{ij}^2$ , so they summarize how important  $y_i$  is for all the fitted values. We call them the hat values or leverages and a single subscript notation is used:  $h_i = h_{ii}$
- Intuitively, the hat values measure how far a unit's vector of characteristics x<sub>i</sub> is from the vector of means of X
- Rule of thumb: examine hat values greater than 2(k+1)/n

Stewart (Princeton)

Appendix: Facts about Hat Values

- $\sum_{i=1}^{n} h_i = k+1$
- $1/n \ge h_i \ge 1$  for all i
- $Var[\widehat{u}_i] = (1 h_i)\sigma^2$
- With a simple linear regression, we have

$$h_i = \frac{1}{n} + \frac{(X_i - \overline{X})^2}{\sum_{j=1}^n (X_j - \overline{X})^2}$$

#### Buchanan hats



### Influence points



- An influence point is one that is both an outlier (extreme in Y) and a leverage point (extreme in X).
- Causes the regression line to move toward it.

# Detecting Influence Points/Bad Leverage Points

#### • Influence Points:

Influence on coefficients = Leverage  $\times$  Outlyingness

• More formally: Measure the change that occurs in the slope estimates when an observation is removed from the data set. Let

$$D_{ij} = \hat{\beta}_j - \hat{\beta}_{j(-i)}, \quad i = 1, \dots, n, \quad j = 0, \dots, k$$

where  $\hat{\beta}_{j(-i)}$  is the estimate of the *j*th coefficient from the same regression once observation *i* has been removed from the data set.

•  $D_{ij}$  is called the DFbeta, which measures the influence of observation *i* on the estimated coefficient for the *j*th explanatory variable.

#### Standardized Influence

To make comparisons across coefficients, it is helpful to scale  $D_{ij}$  by the estimated standard error of the coefficients:

$$D_{ij}^* = \frac{\hat{\beta}_j - \hat{\beta}_{j(-i)}}{\hat{SE}_{-i}(\hat{\beta}_j)}$$

where  $D_{ii}^*$  is called **DFbetaS**.

- D<sup>\*</sup><sub>ij</sub> > 0 implies that removing observation *i* decreases the estimate of β<sub>j</sub> → obs *i* has a positive influence on β<sub>j</sub>.
- $D_{ij}^* < 0$  implies that removing observation *i* increases the estimate of  $\beta_j \rightarrow \text{obs } i$  has a negative influence on  $\beta_j$ .
- Values of  $|D_{ij}^*| > 2/\sqrt{n}$  are an indication of high influence.
- In R: dfbetas(model)

### Buchanan influence

```
##
##
  Coefficients:
##
                 Estimate Std. Error t value Pr(>|t|)
## (Intercept) -2.935e+01 5.520e+01 -0.532 0.59686
## edaytotal 1.100e-03 4.797e-04 2.293 0.02529 *
## absnbuchanan 6.895e+00 2.129e+00 3.238 0.00195 **
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 317.2 on 61 degrees of freedom
##
     (3 observations deleted due to missingness)
## Multiple R-squared: 0.5361, Adjusted R-squared: 0.5209
## F-statistic: 35.24 on 2 and 61 DF, p-value: 6.711e-11
```

### Buchanan influence

##		(Intercept)	edaytotal	absnbuchanan
##	1	0.3454475146	0.4050504921	-0.7505222758
##	2	-0.0234266617	-0.0241000045	-0.0131672181
##	3	0.0650795039	-0.7319311820	0.3401669862
##	4	-0.0333980968	0.0133802934	-0.0087505576
##	5	-0.0397626659	-0.0073746223	0.0096551713
##	6	-0.0009277798	0.0001505476	0.0002210247

# Buchanan influence



• Palm Beach county moves each of the coefficients by more than 3 standard errors!

Stewart (Princeton)

Week 8: Diagnostics and Solutions

October 19-23, 2020 62 / 127

# Summarizing Influence across All Coefficients

- Leverage tells us how much one data point affects a single coefficient.
- A number of summary measures exist for influence of data points across all coefficients, all involving both leverage and outlyingness.
- A popular measure is Cook's distance:

$$D_i = \underbrace{\frac{\hat{u}_i'^2}{k+1}}_{\text{outlyingness}} \times \underbrace{\frac{h_i}{1-h_i}}_{\text{leverage}}$$

where  $\hat{u}'_i$  is the standardized residual and  $h_i$  is the hat value.

- It can be shown that D<sub>i</sub> is a weighted sum of k + 1 DFbetaS's for observation i
- In R, cooks.distance(model)
- D > 4/(n-k-1) is commonly considered large
- The influence plot: the studentized residuals plotted against the hat values, size of points proportional to Cook's distance.

# Influence Plot Buchanan

Influence Plot



Courtesy of Erin Hartman

Stewart (Princeton)

#### Code for Influence Plot

```
ggplot(fl_lm, aes(x = .hat, y = rstudent(fl_lm),
size = .cooksd.
col = .cooksd > 4/(nrow(fl_data) - 1 - 1),
label = fl_data$county)) +
geom_point() + geom_text(vjust = 2) +
xlab("Hat Values") + ylab("Studentized Residuals") +
geom_vline(xintercept = 2 * (fl_lm$rank - 1 + 1)/nrow(fl_data)
, linetype = 2) +
geom_hline(yintercept = c(-4, 4), linetype = 2) +
scale_color_manual("High Influence",
values = c("TRUE" = ucla_gold,
"FALSE" = ucla_blue)) +
scale_size("Cook's Distance") + theme_bw() +
theme(legend.position = c(0.9, 0.5)) + ylim(c(-7, 20)) +
xlim(c(0, 0.4)) + ggtitle("Influence Plot")
```

# A Quick Function for Standard Diagnostic Plots

- > par(mfrow=c(2,2))
- > plot(mod1)



# The Improved Model

R Code \_

- > par(mfrow=c(2,2))
- > plot(mod2)


# 'Fun With Outliers'! (via FiveThirtyEight)

#### NOV. 9, 2018, AT 12:20 PM

Something Looks Weird In Broward County. Here's What We Know About A Possible Florida Recount.

By <u>Nathaniel Rakich</u> Filed under 2018 Election





ILLUSTRATION BY FIVETHIRTYEIGHT

# 'Fun With Outliers'! (via FiveThirtyEight)



The percentage difference between votes cast for governor and votes cast for U.S. Senate in every Florida county in the 2018 midterm election, as of 8:15 a.m. on Nov. 9



Stewart (Princeton)

Week 8: Diagnostics and Solutions

October 19-23, 2020 68 / 127

# 'Fun With Outliers'! (via FiveThirtyEight)

Ballot Style 58		Seq:058	
Official General Election Ballot November 6, 2018 Broward County, Florida	Boleta Oficial De La Elecció 6 De Noviembre Del 2 Condado de Broward, F	n General Ofisyèl Jeneral E 018 6 Novann lorida Konte Browa	leksyon Bilte n 2018 rd, Florida
Ballot Instructions:	Governor and Lieutenant Governor Gobernador v Teniente Gobernador	Fourth District Court of Appeal	
completely next to your choice. Use only the marking	Gouvènè Åk Lyetnan Gouvènè (Vote for One/Vote por Uno/Vote pou Youn	Tribunal De Apelaciones Del Cuarto Distrito	
device provided or a black pen.	Ron DeSantis REP	Katriyèm Distrik Lakou Dapèl	
<ul> <li>If you make a mistake, ask for a new ballot. Do not cross out or your vote may not</li> </ul>	Jeanette Nuflez	Shall Judge Burton C. Conner of the	1
count. • To vote for a write-in candidate, fill in the oval ● and print the name dearly on	Andrew Gillum DEM Chris King	Fourth District Court of Appeal be retained in office?	
the blank line provided for the write-in candidate.	<ul> <li>Darcy G. Richardson REF Nancy Argenziano</li> </ul>	¿Deberá retenerse en su cargo al Juez Burton G. Conner del Tribunal del Cuarto Distrito de Apelaciones?	~
<ul> <li>Para votar, llene completamente el ovalo          <ul> <li>Para votar, llene</li> <li>Completamente el ovalo              </li> </ul> </li> </ul>	O Kele "KC" Gibson NPA	Éske se pou jis Burton C. Conner nan	
sólo un lápiz de punta negra o una pluma de tinta negra nara marcar la holata	Ellen Wilds	karryem osstik takou dapel rete nan pos 8 a?	
<ul> <li>Si se equivoca, pida una nueva boleta. Si borra algo o hace mancas, es posible que</li> </ul>	<ul> <li>Ryan Christopher Foley NPA John Tutton Jr</li> </ul>	O Yes/SWI	
su voto no se cuente. Para Votar por un candidato	O Bran Dimlar		
cuyo nombre no está impreso en la boleta, llene el	Ryan Howard McJury	Shall Judge Jeffrey T. Kuntz of the Fourth	
óvalo		office?	
línea en blanco provista para un candidato agregado.	O Without Facebook lake	Datastastastastas as a summer of here	
Enfómasyon Sou Bilten Vót		Jeffrey T. Kuntz del Tribunal del Cuarto	
<ul> <li>Pou vote, byen kolore tout andan oval           ki akote</li> </ul>	Attorney General Fiscal General	Distrib de Apelaciónes /	
respons ou chwazi a. Sélman sévi ak yon plim nwa	Pwokinè Jeneral Wote for One/Vote por Uno/Vote pou Youn	Éske se pou jis Jeffrey T. Kuntz nan	
oubyen ak yon kreyon pou ekri sou bilten vót la.	O Ashiey Moody REP	abiyen disak lakod dapenete nan pos	
<ul> <li>Si w fé yon erê, mande yo ba w yon nouvo bilten vô. Si w</li> </ul>	O See Share DEM	O Yes5Wi	
efase oubyen fe novuo mak, I ap posib pou vôt ou pa	C Infrar Mars Datied 10	O Nationa	
valab ankó. • Pou vote pou yon kandida ki	Chief Example Officer	O Norworkon	
pa gen non I enprime sou bilten vöt la, kolore ti oval la	Controlador Estatal	Shall Judge Carole Y. Taylor of the	
<ul> <li>epi ekri non kandida a sou liy vid la rezève pou ekri</li> </ul>	(Vote for One/Vote por Uno/Vote pou Youn	Fourth District Court of Appeal be retained in office?	
non yon kandida.	O Jimmy Patronis REP		
United States Senator Senador De Los Estados Unidos Senaté Etazini	O Jeremy Ring DEM	Debera retenense en su cargo al Juaz Carole Y. Taylor del Tribunal del Cuarto Distrito de Apelaciones?	
vote for Unevote por Unovote pour Tour	Write-in/Escribin/A lekri	Éska sa novi is Carola V. Tavlor nan	
O Bill Nelson DEM	Commissioner of Agriculture	katriyêm distrik lakou dapêl rete nan pôs li a?	
°	Komisyonè Agrikiti (Vote for One/Vote por Uno/Vote pou Youn	○ Yes/SiWi	
Write-In/Escribit/A lekri	O Matt Caldwell REP	O Nañaña	
Representative in Congress	<ul> <li>Nicole "Nikki" Fried</li> <li>DEM</li> </ul>	Circuit Judge, 17th Judicial Circuit,	
Representante En El Congreso	Justice of the Supreme Court	Juez De Circuito, Circuito 17Mo, Gruno 38	
Reprezantan Nan Kongrè	Magistrado en el Tribunal Supremo	Jij Itineran nan 17èm Sikui Gercup 38	
Vote for OneVote por UnoVote pou Your	Jistis Nan Lakou Siprèm	(Vote for One/Vote por Uno/Vote pou Youn)	
O Alcee L Hastings DEM		O Jason Allen-Rosner	1
•	Stat Justice Alan Lawson of the Supreme Court be retained in office?	<ul> <li>Stefanie Camille Moon</li> </ul>	
Write-it/Escribit/A lekri	Deberá referense en el caron al	Circuit Judge, 17th Judicial Circuit	1

Stewart (Princeton)

Week 8: Diagnostics and Solutions

- Outliers, Leverage and Influence Points
- Always check your data!
- Don't let regression be a magic black box for you- understand what is in your data that is leading to the findings.

Next Time: Robust Regression

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference



#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points

# Robust Regression Methods Appendix: Robustness

#### Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# Limitations of the Standard Tools



- What happens when there are two influence points?
- Red line drops the red influence point
- Blue line drops the blue influence point
- Neither of the "leave-one-out" approaches helps recover the line

### The Idea of Robustness

- We have and will cover a few ideas in robust statistics (much of which is due directly or indirectly to Peter Huber).
- Robust methods are procedures that are designed to continue to provide 'reasonable' answers in the presence of violation of some assumptions.
- A lot of social scientists use robust standard errors but far fewer use robust regression tools.
- These methods used to be computationally prohibitive but haven't been for the last 10-15 years

# But What About Gauss-Markov and BLUE?

- One argument here is that even without normality, we know that Gauss-Markov is the Best Linear Unbiased Estimator (BLUE)
- How comforting should this be? Not very.
- The Linear point is an artificial restriction. It means the estimator has to be of the form  $\hat{\beta} = \mathbf{W}y$  but why only use those?
- With normality assumption we get Best Unbiased Estimator (BUE) which is quite comforting when  $n \gg p$  (number of observations much larger than number of variables).

## This Point is Not Obvious

This flies in the face of most conventional wisdom in textbooks.

"We need not look for another linear unbiased estimator, for we will not find such an estimator whose variance is smaller than the OLS estimator" - Gujarati (2004)

Quotes from Rainey and Baissa (2015) presentation

# Robustly Estimating a Location

- Let's simplify- what if we want to estimate the center of a symmetric distribution.
- Two options (of many): mean and median
- Characteristics to consider: efficiency when assumptions hold, sensitivity to assumption violation.
- For normal data  $y_i \sim \mathcal{N}(\mu, \sigma^2)$ , median is less efficient:

• 
$$V(\hat{\mu}_{\text{mean}}) = \frac{\sigma^2}{n}$$

• 
$$V(\hat{\mu}_{\text{median}}) = \frac{\pi \sigma^2}{2n}$$

- Median is  $\frac{\pi}{2}$  times larger (i.e. less efficient)
- We can measure sensitivity with the influence function which measures change in estimator based on corruption in one datapoint.

### Influence Function

- Imagine that we had a sample Y from a standard normal: -0.068, -1.282, 0.013, 0.141, -0.980, 1.63.  $\bar{Y} = -1.52$
- Now imagine we add a contaminated 7th observation which could range from -10 to +10. How would the estimator change for the median and mean?



#### Example from Fox

Stewart (Princeton)

# Breakdown Point

- The influence function showed us how one aberrant point can change the resulting estimate.
- We also want to characterize the breakdown point which is the fraction of arbitrarily bad data that the estimator can tolerate without being affected to an arbitrarily large extent
- The breakdown point of the mean is 0 because (as we have seen) a single bad data point can change things a lot.
- The median has a breakdown point of 50% because half the data can be bad without causing the median to become completely unstuck.

### *M*-estimators

- We can phrase this more generally than the mean or the median which will allow us to extend the ideas to regression via *M*-estimation
- *M*-estimators minimize a sum over an objective function  $\sum_{i=1}^{n} \rho(E)$  where *E* is  $Y_i \hat{\mu}$ 
  - The mean has  $\sum_i \rho(E) = \sum_i (Y_i \hat{\mu})^2$
  - The median has  $\sum_i \rho(E) = \sum_i |(Y_i \hat{\mu})|$
- The shape of the influence function is determined by the derivative of the objective function with respect to *E*.
- Other objectives include the Huber objective and Tukey's biweight objective which have different properties.
- Calculating robust *M* estimators often requires an iterative procedure and a careful initialization.

# M-estimation for Regression

- We can apply this to regression fairly straightforwardly. In robust M-estimators we choose  $\rho()$  so that observations with large residuals get less weight.
- Can be very robust to outliers in the Y space (less so in the X space usually)
- Some options:
  - ► Least Median Squares: choose  $\hat{\beta}$  to minimize median  $\left\{ (y_i \mathbf{x}'_i \hat{\beta}_{LMS})^2 \right\}_{i=1}^n$ . Very high breakdown point, but very inefficient.
  - Least Trimmed Squares: choose β̂ to minimize the sum of the p smallest elements of {(y<sub>i</sub> − x'<sub>i</sub>β̂<sub>LTS</sub>)<sup>2</sup>}<sup>n</sup><sub>i=1</sub>. High breakdown point and more efficient, still not as efficient as some.
  - MM-estimator: what I recommend in practice (more in appendix)
- You can find an asymptotic covariance matrix for *M*-estimators but I would bootstrap it if possible as the asymptotics kick in slowly.

```
library(MASS)
set.seed(588)
n <- 50
x < - rnorm(n)
y <- 10 - 2 x + rnorm(n)
x[1:5] <- rnorm(5, mean=5)
y[1:5] <- 10 + rnorm(5)
ols.out <- lm(y^x)
m.out <- rlm(y~x, method="M")</pre>
lms.out <- lqs(y~x, method="lms")</pre>
lts.out <- lqs(y~x, method="lts")</pre>
s.out <- lqs(y~x, method="S")</pre>
mm.out <- rlm(y~x, method="MM")</pre>
```

### Simulation Results















# Thoughts on Robust Estimators

- Robust estimators aren't commonly seen in applied social science work but perhaps they should be.
- Even though Gauss-Markov does not require normality, the L in BLUE is a fairly restrictive condition.
- In most cases I personally would start with OLS, do diagnostics and then consider a robust alternative. If I don't have time for diagnostics, maybe robust is better from the outset.
- See Baissa and Rainey (2018) "When BLUE is Not Best: Non-Normal Errors and the Linear Model" in *Political Science Research & Methods* for more on this topic.
- The Fox textbook Chapter 19 is also quite good on this and points out to the key references

# We Covered

- Robust Regression
- Appendix after these slides with some more formality on *M*-estimators.

Next Time: Nonlinearity

# Appendix: Characterizing Estimator Robustness (formally)

#### Definition (Breakdown Point)

The breakdown point of an estimator is the smallest fraction of the data that can be changed an arbitrary amount to produce an arbitrarily large change in the estimate (Seber and Lee 2003, pg 82)

#### Definition (Influence Function)

Let  $F_p = (1 - p)F + p\delta_{z_0}$  where F is a probability measure,  $\delta_{z_0}$  is the point mass at  $\mathbf{z}_0 \in \mathbb{R}^k$ , and  $p \in (0, 1)$ .

Let  $T(\cdot)$  be a statistical functional. The influence function of T is

$$IF(\mathbf{z}_0; T, F) = \lim_{p \downarrow 0} \frac{T(F_p) - T(F)}{p}$$

The influence function is a function of  $\mathbf{z}_0$  given T and F. It describes how T changes with small amounts of contamination at  $z_0$  (Hampel, Rousseeuw, Ronchetti, and Stahel, (1986), p. 84).

# Appendix: S Estimators

To talk about MM-estimators we need a type of estimator called an S-estimator.

<u>S-estimators</u> work somewhat differently in that the goal is to minimize the scale estimate subject to a constraint.

An S-estimator for the regression model is defined as the values of  $\hat{\beta}_{S}$  and s that minimize s subject to the constraint:

$$\frac{1}{n}\sum_{i=1}^{n}\rho\left(\frac{y_{i}-\mathbf{x}_{i}^{\prime}\hat{\boldsymbol{\beta}}_{S}}{s}\right)\geq K$$

where K is user-defined constant (typically set to 0.5) and  $\rho : \mathbb{R} \to [0, 1]$  is a function with the following properties (Davies, 1990, p. 1653):

**1**  $\rho(0) = 1$ 

2 
$$\rho(u) = \rho(-u), u \in \mathbb{R}$$

) for some 
$$c>$$
 0,  $ho(u)>$  0 if  $|u|< c$  and  $ho(u)=$  0 if  $|u|> c$ 

<u>MM-estimators</u> are, in some sense, the best of both worlds– very high breakdown point and good efficiency.

The work by first calculating S-estimates of the scale and coefficients and then using these as starting values for a particular M-estimator.

Good properties, but costly to compute (usually impossible to compute exactly).

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference



#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points

# Robust Regression Methods Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# Nonlinearity

- We know that linear regression asymptotically gets us the best linear approximation to the conditional expectation function.
- We can always add transformations of variables (like polynomials) to expand **X** in a way that makes non-linear shapes in the original variable.
- What happens when we don't know the shape of the non-linearity?
- In Week 5 we talked about nonparametric regressions for settings with one independent variable.
- Many forms of machine learning are best thought of as nonparametric regressions in higher dimensions.
- We can often see poor fits of the conditional expectation function in the residuals, but let's instead just do diagnosis by treatment and look at some of these other approaches to modeling.



Thanks XKCD for having a comic for everything!

Stewart (Princeton)



#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points

# Robust Regression Methods Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Clustering

# **Bias-Variance Tradeoff**



#### Model Complexity

Stewart I	(Princeton)
SLEWAIL	FILLELOIL

### Example Synthetic Problem

$$y = \sin(1 + x^2) + \epsilon$$



This section adapted from slides by Radford Neal.

Stewart (Princeton)

Week 8: Diagnostics and Solutions

### Linear Basis Function Models

- We talked before about polynomials  $x^2, x^3, x^4$  for modeling non-linearities, this is a linear basis function model.
- In general the idea is to do a linear regression of y on  $\phi_1(x), \phi_2(x), \ldots, \phi_{m-1}(x)$  where  $\phi_j$  are basis functions.
- The model is now:

$$y = f(x, \beta) + \epsilon$$
$$f(x, \beta) = \beta_0 + \sum_{j=1}^{m-1} \beta_j \phi_j(x) = \beta^T \phi(x)$$

# Polynomial Basis Functions

We have already seen some basis functions. Here are OLS fits with polynomial basis functions of increasing order.



It appears that the last model is too complex and is overfitting a bit.

### Local Basis Functions

Polynomials are global basis functions, each affecting the prediction over the whole input space. Often local basis functions are more appropriate.

One choice is a Gaussian basis function

$$\phi_j(x) = \exp(-(x-\mu_j)^2)/2s^2)$$



### Gaussian Basis Fits



# Regularization

- We've seen that flexible models can lead to overfitting
- Two ways to address: limit model flexibility or use a flexible model and regularize
- Regularization is a way of expressing a preference for smoothness in our function by adding a penalty term to our optimization function.
- Here we will consider a penalty of the form  $\lambda \sum_{j=1}^{m-1} \beta_j^2$  where  $\lambda$  controls the strength of the penalty.
- The penalty trades off some bias for an improvement in variance
- The trick in general is how to set  $\lambda$

### Results

Here are the results with  $\lambda = 0.01$ :



100 / 127

### Results

Here are the results with  $\lambda = 0.1$ :



100 / 127
#### Results

Here are the results with  $\lambda = 1$ :



Stewart (Princeton)

Week 8: Diagnostics and Solutions

October 19-23, 2020

100 / 127

#### Results

Here are the results with  $\lambda = 10$ :



## Conclusions from This Example

- we can control overfitting by modifying the width of the basis function *s* or with penalty
- we will need some way in general to tune these
- we will also need some way to handle multivariate functions.

## Generalized Additive Models (GAM)

Recall the linear model,

$$y_i = \beta_0 + x_{1i}\beta_1 + x_{2i}\beta_2 + x_{3i}\beta_3 + u_i$$

For GAMs, we maintain additivity, but instead of imposing termwise linearity we allow flexible functional forms for each explanatory variable, where  $s_1(\cdot), s_2(\cdot)$ , and  $s_3(\cdot)$  are smooth functions that are estimated from the data:

$$y_i = \beta_0 + s_1(x_{1i}) + s_2(x_{2i}) + s_3(x_{3i}) + u_i$$

Generalized Additive Models (GAM)

 $y_i = \beta_0 + s_1(x_{1i}) + s_2(x_{2i}) + s_3(x_{3i}) + u_i$ 

- GAMS are semi-parametric, they strike a compromise between nonparametric methods and parametric regression
- s<sub>j</sub>(·) are usually estimated with locally weighted regression smoothers or cubic smoothing splines (but many approaches are possible)
- They do NOT give you a set of regression parameters  $\hat{\beta}$ . Instead one obtains a graphical summary of how  $E[Y|X, X_2, ..., X_k]$  varies with  $X_1$  (estimates of  $s_j(\cdot)$  at every value of  $X_{i,j}$ )
- Theory and estimation are somewhat involved, but they are easy to use:
  - gam.out <- gam(y~s(x1)+s(x2)+x3)
    plot(gam.out)</pre>
  - Multiple functions but I recommend mgcv package

# Generalized Additive Models (GAM)

The GAM approach can be extended to allow interactions  $(s_{12}(\cdot))$  between explanatory variables, but this eats up degrees of freedom so you need a lot of data.

$$y_i = \beta_0 + s_{12}(x_{1i}, x_{2i}) + s_3(x_{3i}) + u_i$$

It can also be used for hybrid models where we model some variables as parametrically and other with a flexible function:

$$y_i = \beta_0 + \beta_1 x_{1i} + s_2(x_{2i}) + s_3(x_{3i}) + u_i$$





age

106 / 127





red/green are +/- 2 s.e.



# Concluding Thoughts

- Non-linearity is pretty easy to detect and can substantially change our inferences
- GAMs are a great way to model/detect non-linearity but transformations are often simpler
- However, be wary of the global properties of transformations and polynomials
- Non-linearity concerns are most relevant for continuous covariates with a large range (age)
- NB: it is okay if you didn't follow all of this today! GAMs are tricky.

# We Covered

- Linear basis function models.
- Generalized Additive Models.

Next Time: Clustering

Where We've Been and Where We're Going...

- Last Week
  - multiple regression
- This Week
  - diagnosing problems and troubleshooting the linear model
  - $\blacktriangleright$  unusual and influential data  $\rightarrow$  robust estimation
  - $\blacktriangleright$  non-linearity  $\rightarrow$  generalized additive models
  - unusual errors  $\rightarrow$  sandwich SEs
- Next Week
  - frameworks for causal inference
- Long Run
  - $\blacktriangleright$  probability  $\rightarrow$  inference  $\rightarrow$  regression  $\rightarrow$  causal inference



#### 2 Non-Normality

#### 3 Extreme Values

- Outliers
- Leverage Points
- Influence Points
- Robust Regression Methods
   Appendix: Robustness

#### 5 Nonlinearity

- Linear Basis Function Models
- Generalized Additive Models

#### Olustering

# Clustered Dependence: Intuition

- Think back to the Gerber, Green, and Larimer (2008) social pressure mailer example.
- Their design: randomly sample households and randomly assign them to different treatment conditions
- But the measurement of turnout is at the individual level
- Violation of iid/random sampling:
  - errors of individuals within the same household are correlated
  - ► ~→ violation of homoskedasticity
- Called clustering or clustered dependence

# Clustered Dependence: notation

- Clusters: *j* = 1, . . . , *m*
- Units:  $i = 1, ..., n_j$
- $n_j$  is the number of units in cluster j
- $n = \sum_{i} n_{j}$  is the total number of units
- Units (usually) belong to a single cluster:
  - voters in households
  - individuals in states
  - students in classes
  - rulings in judges
- Especially important when outcome varies at the unit-level, y<sub>ij</sub> and the main independent variable varies at the cluster level, x<sub>j</sub>.
- Ignoring clustering is "cheating": units not independent

#### Clustered Dependence: Example Model

$$y_{ij} = \beta_0 + \beta_1 x_{ij} + \varepsilon_{ij}$$
$$= \beta_0 + \beta_1 x_{ij} + \mathbf{v}_j + \mathbf{u}_{ij}$$

- $v_j \stackrel{iid}{\sim} N(0, \rho\sigma^2)$  cluster error component
- $u_{ij} \stackrel{iid}{\sim} N(0, (1-\rho)\sigma^2)$  unit error component
- $v_j$  and  $u_{ij}$  are assumed to be independent of each other
- $ho \in (0,1)$  is called the within-cluster correlation.
- What if we ignore this structure and just use  $\varepsilon_{ij}$  as the error?
- Variance of the composite error is  $\sigma^2$ :

$$\begin{aligned} \mathsf{Var}[\varepsilon_{ij}] &= \mathsf{Var}[v_j + u_{ij}] \\ &= \mathsf{Var}[v_j] + \mathsf{Var}[u_{ij}] \\ &= \rho \sigma^2 + (1 - \rho) \sigma^2 = \sigma^2 \end{aligned}$$

#### Lack of Independence

• Covariance between two units *i* and *s* in the same cluster is  $\rho\sigma^2$ :

$$\mathsf{Cov}[\varepsilon_{ij},\varepsilon_{sj}] = \rho\sigma^2$$

• Correlation between units in the same group is just *ρ*:

$$\operatorname{Cor}[\varepsilon_{ij}, \varepsilon_{sj}] = \rho$$

• Zero covariance of two units *i* and *s* in different clusters *j* and *k*:

$$\operatorname{Cov}[\varepsilon_{ij}, \varepsilon_{sk}] = 0$$

### Example Covariance Matrix

$$\boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_{1,1} & \varepsilon_{2,1} & \varepsilon_{3,1} & \varepsilon_{4,2} & \varepsilon_{5,2} & \varepsilon_{6,2} \end{bmatrix}'$$
$$\operatorname{Var}[\boldsymbol{\varepsilon}] = \boldsymbol{\Sigma} = \begin{bmatrix} \sigma^2 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & 0 & 0 & 0 \\ \sigma^2 \cdot \rho & \sigma^2 & \sigma^2 \cdot \rho & 0 & 0 & 0 \\ \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & \sigma^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma^2 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho \\ 0 & 0 & 0 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho \\ 0 & 0 & 0 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & \sigma^2 \end{bmatrix}$$

118 / 127

# Appendix: Example 6 Units, 2 Clusters $\epsilon = [\epsilon_{1,1} \epsilon_{2,1} \epsilon_{3,1} \epsilon_{4,2} \epsilon_{5,2} \epsilon_{6,2}]'$

$$\begin{split} V[\boldsymbol{\varepsilon}] = \boldsymbol{\Sigma} = \begin{bmatrix} V[\varepsilon_{1,1}] & Cov[\varepsilon_{2,1},\varepsilon_{1,1}] & Cov[\varepsilon_{3,1},\varepsilon_{1,1}] & \cdot & \cdot & \cdot & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{2,1}] & V[\varepsilon_{2,1}] & Cov[\varepsilon_{3,1},\varepsilon_{2,1}] & \cdot & \cdot & \cdot & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{3,1}] & Cov[\varepsilon_{2,1},\varepsilon_{3,1}] & V[\varepsilon_{3,1}] & \cdot & \cdot & \cdot & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{4,2}] & Cov[\varepsilon_{2,1},\varepsilon_{4,2}] & Cov[\varepsilon_{3,1},\varepsilon_{4,2}] & V[\varepsilon_{4,2}] & \cdot & \cdot & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{5,2}] & Cov[\varepsilon_{2,1},\varepsilon_{5,2}] & Cov[\varepsilon_{3,1},\varepsilon_{5,2}] & Cov[\varepsilon_{4,2},\varepsilon_{5,2}] & V[\varepsilon_{5,2}] & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{5,2}] & Cov[\varepsilon_{2,1},\varepsilon_{6,2}] & Cov[\varepsilon_{3,1},\varepsilon_{6,2}] & Cov[\varepsilon_{4,2},\varepsilon_{5,2}] & V[\varepsilon_{5,2}] & \cdot \\ Cov[\varepsilon_{1,1},\varepsilon_{5,2}] & Cov[\varepsilon_{2,1},\varepsilon_{6,2}] & Cov[\varepsilon_{3,1},\varepsilon_{6,2}] & Cov[\varepsilon_{4,2},\varepsilon_{6,2}] & Cov[\varepsilon_{5,2},\varepsilon_{6,2}] & V[\varepsilon_{6,2}] \end{bmatrix} \end{bmatrix} \\ = \begin{bmatrix} \sigma^2 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & 0 & 0 & 0 \\ \sigma^2 \cdot \rho & \sigma^2 & \sigma^2 \cdot \rho & \sigma^2 & 0 & 0 & 0 \\ \sigma^2 \cdot \rho & \sigma^2 \cdot \rho & \sigma^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma^2 & \sigma^2 \cdot \rho & \sigma^2 \cdot \rho \\ 0 & 0 & 0 & \sigma^2 \cdot \rho & \sigma^2 & \sigma^2 \cdot \rho \end{bmatrix} \end{split}$$

which can be verified as follows:

• 
$$V[\varepsilon_{ij}] = V[v_j + u_{ij}] = V[v_j] + V[u_{ij}] = \rho\sigma^2 + (1 - \rho)\sigma^2 = \sigma^2$$
  
•  $Cov[\varepsilon_{ij}, \varepsilon_{ij}] = E[\varepsilon_{ij}\varepsilon_{ij}] - E[\varepsilon_{ij}]E[\varepsilon_{ij}] = E[\varepsilon_{ij}\varepsilon_{ij}] = E[(v_j + u_{ij})(v_j + u_{ij})]$   
 $= E[v_j^2] + E[v_ju_{ij}] + E[v_ju_{ij}] + E[u_{ij}u_{ij}]$   
 $= E[v_j^2] + E[v_j]E[u_{ij}] + E[v_j]E[u_{ij}] + E[u_{ij}]E[u_{ij}]$   
 $= E[v_j^2] = V[v_j] + (E[v_j])^2 = V[v_j] = \rho\sigma^2$ 

• 
$$Cov[\varepsilon_{ij}, \varepsilon_{lk}] = E[\varepsilon_{ij}\varepsilon_{lk}] - E[\varepsilon_{ij}]E[\varepsilon_{lk}] = E[\varepsilon_{ij}\varepsilon_{lk}] = E[(v_j + u_{ij})(v_k + u_{lk})]$$
  
=  $E[v_jv_k] + E[v_ju_{lk}] + E[v_ku_{ij}] + E[u_{ij}u_{lk}]$   
=  $E[v_j]E[v_k] + E[v_j]E[u_{lk}] + E[v_k]E[u_{ij}] + E[u_{ij}]E[u_{lk}] = 0$ 

#### Error Variance Matrix with Cluster Dependence

The variance-covariance matrix of the error,  $\Sigma$ , is block diagonal:

• By independence, the errors are uncorrelated across clusters:

$$V[\varepsilon] = \Sigma = \begin{bmatrix} \Sigma_1 & 0 & \dots & 0 \\ \hline 0 & \Sigma_2 & \dots & 0 \\ \hline & & \ddots & \\ \hline 0 & 0 & \dots & \Sigma_M \end{bmatrix}$$

• But the errors may be correlated for units within the same cluster:

$$\boldsymbol{\Sigma}_{j} = \begin{bmatrix} \sigma^{2} & \sigma^{2} \cdot \rho & \dots & \sigma^{2} \cdot \rho \\ \sigma^{2} \cdot \rho & \sigma^{2} & \dots & \sigma^{2} \cdot \rho \\ & & \ddots & \\ \sigma^{2} \cdot \rho & \sigma^{2} \cdot \rho & \dots & \sigma^{2} \end{bmatrix}$$

# Correcting for Clustering

- Including a dummy variable for each cluster (fixed effects)
- (Random effects' models (take above model as true and estimate ρ and σ<sup>2</sup>)
- S Cluster-robust ("clustered") standard errors
- Aggregate data to the cluster-level and use OLS  $\overline{y}_j = \frac{1}{n_i} \sum_i y_{ij}$ 
  - If n<sub>j</sub> varies by cluster, then cluster-level errors will have heteroskedasticity

# Cluster-Robust SEs

• First, let's write the within-cluster regressions like so:

$$\mathbf{y}_j = \mathbf{X}_j \boldsymbol{eta} + \boldsymbol{arepsilon}_j$$

- $\mathbf{y}_j$  is the vector of responses for cluster j, and so on
- We assume that respondents are independent across clusters, but possibly dependent within clusters. Thus, we have

$$\mathsf{Var}[arepsilon_j | \mathbf{X}_j] = \mathbf{\Sigma}_j$$

• Remember our sandwich expression:

$$\mathsf{Var}[\hat{eta}|\mathbf{X}] = \left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'\mathbf{\Sigma}\mathbf{X}\left(\mathbf{X}'\mathbf{X}
ight)^{-1}$$

• Under this clustered dependence, we can write this as:

$$\mathsf{Var}[\hat{eta}|\mathbf{X}] = \left(\mathbf{X}'\mathbf{X}
ight)^{-1} \left(\sum_{j=1}^m \mathbf{X}_j' \mathbf{\Sigma}_j \mathbf{X}_j
ight) \left(\mathbf{X}'\mathbf{X}
ight)^{-1}$$

# Estimating the Variance Components: $\rho$ and $\sigma^2$

The overall error variance  $\sigma^2$  is easily estimated using our usual estimator based on the regression residuals:  $\hat{\sigma}^2 = \frac{\hat{\epsilon}'\hat{\epsilon}}{N-k-1}$ 

The within-cluster correlation can be estimated as follows:

- Subtract from each residual  $\hat{\varepsilon}_{ij}$  the mean residual within its cluster. Call this vector of demeaned residuals  $\tilde{\varepsilon}$ , which estimates the unit error component  $\boldsymbol{u}$
- Some compute the variance of the demeaned residuals as:  $\hat{\tilde{\sigma}}^2 = \frac{\tilde{\epsilon}'\tilde{\epsilon}}{N-M-k-1}$ , which estimates  $(1-\rho)\sigma^2$
- 3 The within cluster correlation is then estimated as:  $\hat{
  ho} = rac{\widehat{\sigma}^2 \widehat{\hat{\sigma}}^2}{\widehat{\sigma}^2}$

# Estimating Cluster Robust Standard Errors

We can now compute the CRSEs using our sandwich formula:

1 Take your estimates of  $\widehat{\sigma^2}$  and  $\widehat{\rho}$  and construct the block diagonal variance-covariance matrix  $\widehat{\Sigma}$ :



<sup>(2)</sup> Plug  $\widehat{\Sigma}$  into the sandwich estimator to obtain the cluster "corrected" estimator of the variance-covariance matrix

$$V[\hat{eta}|\mathbf{X}] = \left(\mathbf{X}'\mathbf{X}
ight)^{-1}\mathbf{X}'\widehat{\Sigma}\mathbf{X}\left(\mathbf{X}'\mathbf{X}
ight)^{-1}$$

 There are multiple implementations in R including multiwayvcov:cluster.vcov and sandwich::vcovCL

#### **Cluster-Robust Standard Errors**

- CRSE do not change our estimates  $\widehat{eta}$ , cannot fix bias
- CRSE is consistent estimator of  $\mathsf{Var}[\widehat{\beta}]$  given clustered dependence
  - ► Relies on independence between clusters, dependence within clusters
  - Doesn't depend on the model we present
  - $\blacktriangleright$  CRSEs usually > conventional SEs—use when you suspect clustering
- Consistency of the CRSE are in the number of groups, not the number of individuals
  - CRSEs can be incorrect with a small (< 50 maybe) number of clusters (often biased downward)
  - Block bootstrap can be a useful alternative (key idea: bootstrap by resampling the clusters)
- There are numerous alternative clustered standard error variants out there.

Concluding Thoughts on Diagnostics

### Residuals are important. Look at them.

# This Week in Review

- We talked about troubleshooting the linear model—few black and white answers but many tools for the toolkit.
- I completely understand than many people won't have all the details of robust regression, generalized additive models or clustered standard errors.
- It is useful to know
  - (a) these tools exist.
  - (b) roughly what problem they help solve.
  - (c) approximately why they work.
- The problem set will give you a chance to practice many of these things.
- There are plenty of other techniques out there (particularly for modeling non-linearity).

Next week: Frameworks for Causal Inference!