

Linear Regression Models

P8111

Lecture 14

Jeff Goldsmith

March 8, 2016



THE DEPARTMENT OF
BIOSTATISTICS



Columbia University
MAILMAN SCHOOL
OF PUBLIC HEALTH

Today's Lecture

- Review session

Exam Thursday

- One hour and twenty minutes (up to two hours)
- Closed book
- Calculators are fine, as long as they don't have internet
- Everything is fair game

What is regression?

$$[y | x]$$

"...to understand as far as possible with the available data how the conditional distribution of the response y varies across subpopulations determined by the possible values of the predictor or predictors." – Cook and Weisberg (1999)

X

Regression model

$$f(x)$$

The process of using data to describe the relationship between outcomes and predictors is called modeling.

- Models are models, not reality. *inc. your model*
- "All models are wrong, but some are useful."
- Introduce structure to $f(x)$ to make the problem of estimation easier (this also introduces elements not found in the data, including judgement calls about important features and assumptions about the world).
- We largely focus on *parametric models* $f(x) = f(x; \underline{\beta})$ and worry about estimating β .

Linear Regression Models

$$f(x; \beta) = \beta_0 + \beta_1 x$$

A linear regression model is a particular type of parametric regression.



- Assume $f(x; \beta) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots$
- Focus is on β_0, β_1, \dots .
- “Linear” refers to the β ’s, not the x ’s:
 - $f(x) = \beta_0 + \beta_1 x + \beta_2 x^2$ is a linear model
 - $f(x) = \beta_0 + x^{\beta_1}$ is not
 - $f^*(x) = \beta_0^* + \beta_1 x^*$

dplyr : Two Other Things

- Grouping (`group_by()`) can make some tasks infinitely easier
- The pipe operator (`%>%`) will change your life

Cheat Sheet

Constructing a ggplot figure

- data: the dataframe you're using to construct your plot
- **aesthetic mappings**: connections between data and visual components (x and y, first; size, color, group, shape, etc)
- **layers**: how the data are actually shown (points, lines, boxplots, densities, smooths)

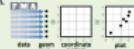
Cheat Sheet

Data Visualization with ggplot2 Cheat Sheet

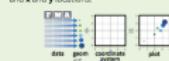


Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same few components: a **data** set, a set of **geoms**—visual marks that represent data points, and a **coordinate system**.



To display data values, map variables in the data set to aesthetic properties of the geom like **size**, **color**, and **x** and **y** locations.



Build a graph with ggplot()

`ggplot(data = cty) + hex(aes(x = cty, y = cty), data = nyc, geom = "point")`
Creates a complete plot with given data, geom, and mappings. Supplies many useful defaults.

`ggplot(data = mtcars, aes(x = cyl, y = hwy))`
Begin a plot that you finish by adding layers to. No defaults, but provides more control than plot().

`data`
`g <- ggplot(mtcars, aes(hwy, cyl)) +
 geom_hex(binwidth = 0.5) +
 coord_cartesian() +
 scale_color_gradient() +
 theme_bw()`
Add a new layer to a plot with a `geom_[]` or `stat_[]` function. Each provides a geom, a set of aesthetic mappings, and a default set and position adjustments.

`last_plot()`
Returns the last plot.
`ggname("plot.png", width = 5, height = 5)`
Saves last plot as 5x5" file named "plot.png" in working directory. Matches file type to file extension.

ggplot2 is a trademark of RStudio, Inc. • [CRAN](#) • [GitHub](#) • [RStudio.com](#) • 844-448-1222 • [support](#)

Geoms – Use a geom to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

One Variable

Continuous

- a -> `geom_area(stat = "bin")`
x,y, alpha, color, fill, linetype, size
- b -> `geom_area(stat = "density")`
x,y, alpha, color, fill, linetype, size, weight
- c -> `geom_density(kern = "gaussian")`
x,y, alpha, color, fill, linetype, size, weight
- d -> `geom_dotplot()`
x,y, alpha, color, fill, shape, size
- e -> `geom_freqpoly()`
x,y, alpha, color, linetype, size
- f -> `geom_histogram(binwidth = 5)`
x,y, alpha, color, fill, size, weight
- g -> `geom_hex()`
x,y, alpha, color, fill, size

Discrete

- h -> `geom_bar(mp, aes(...))`
x, alpha, color, fill, linetype, size, weight
- i -> `geom_bar()`
x, alpha, color, fill, linetype, size, weight

Graphical Primitives

c -> `geom_point(mp, aes(...), lwd)`
x,y, alpha, color, fill, linetype, size

c -> `geom_polygon(pes, group = group)`
x,y, alpha, color, fill, linetype, size

d -> `ggplot(economics, aes(date, unemploy))`

d -> `geom_path(linend = "butt",`

d -> `geom_ribbon(aes(unemployment - 900,`

d -> `unemployment + 900))`

d -> `geom_smooth(span = 0.9))`

d -> `geom_step(direction = "hv")`

d -> `geom_violin(scale = "free")`

e -> `ggplot(seab, aes(x = long, y = lat))`

e -> `geom_segment(aes(x0 = long, y0 = lat, xend = lat, yend = lat))`

e -> `geom_rect(aes(xmin = long, ymin = lat,`

e -> `xmax = long, fill = "#f0f0f0",`

e -> `ymin = lat, ymax = lat))`

e -> `geom_text(aes(x = long, y = lat))`

Coefficient interpretation

$$E(y|x) = \beta_0 + \beta_1 x$$

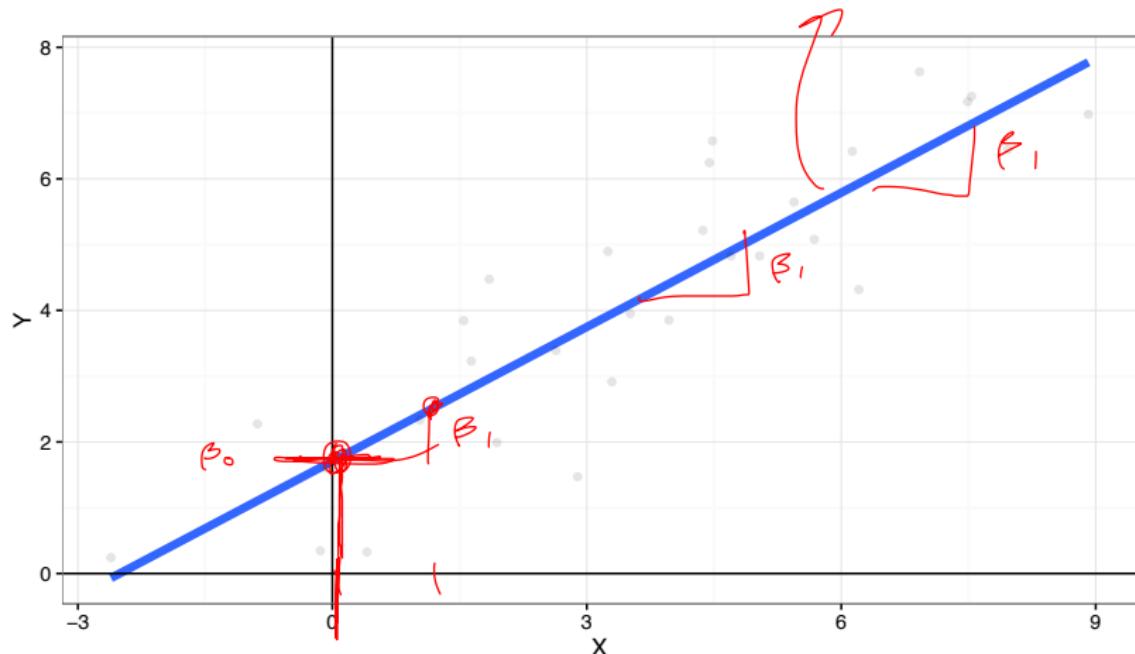
$$\underline{E(y|x=0) = \beta_0} \quad \checkmark$$

$$\beta_1 = (\beta_0 + \beta_1|7) - \beta_0 - \beta_1|6$$

$$= E(y|x=1)_{17} - E(y|x=0)_{16}$$

Coefficient interpretation

$$\mathbb{E}(y|x) = \beta_0 + \beta_1 x$$



Least squares estimation

$$E(y|x) = \beta_0 + \beta_1 x$$

■ Find $\hat{\beta}_0$.

$$RSS(\beta_0, \beta_1) = \sum (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\frac{\partial RSS(\beta_0)}{\partial \beta_0} = \cancel{-2} \sum (y_i - \beta_0 - \beta_1 x_i) = 0$$

$$\sum y_i - n\beta_0 - \beta_1 \sum x_i = 0$$

$$n\beta_0 = \sum y_i - \beta_1 \sum x_i$$

$$\hat{\beta}_0 = \bar{y} - \beta_1 \bar{x}$$

Least squares estimation

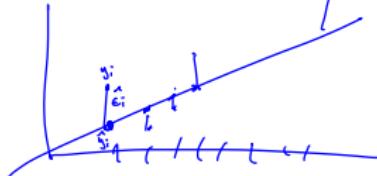
- Now find $\hat{\beta}_1$.

$$\begin{aligned} RSS(\beta_1) &= \sum (y_i - \underbrace{\hat{y}_i}_{\hat{\beta}_0 + \hat{\beta}_1 x_i})^2 \\ &= \sum ((y_i - \bar{y}) - \beta_1(x_i - \bar{x}))^2 \\ \frac{\partial RSS(\beta_1)}{\partial \beta_1} &= -2 \sum ((y_i - \bar{y}) - \beta_1(x_i - \bar{x})) (x_i - \bar{x}) = 0 \\ &= \sum (y_i - \bar{y})(x_i - \bar{x}) - \beta_1 \sum (x_i - \bar{x})^2 = 0 \\ \hat{\beta}_1 &= \frac{\sum (y_i - \bar{y})(x_i - \bar{x})}{\sum (x_i - \bar{x})^2} \end{aligned}$$

Some definitions / SLR products

LSZ $\hat{\beta}_0, \hat{\beta}_1$

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$



- Fitted values: $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$
- Residuals / estimated errors: $\hat{e}_i = y_i - \hat{y}_i$
- Residual sum of squares: $\sum_{i=1}^n \hat{e}_i^2$
- Residual variance: $\hat{\sigma}^2 = \frac{\text{RSS}}{n-2}$
- Degrees of freedom: $n - 2$

{ Notes: residual sample mean is zero; residuals are uncorrelated with fitted values. }

R example

```
> summary(linmod)  
  
Call:  
lm(formula = y ~ x, data = data)
```

Residuals:

Min	1Q	Median	3Q	Max
-1.5202	-0.5050	-0.2297	0.5753	1.8534

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.08743	0.22958	9.092	7.53e-10 ***
x	0.61396	0.05415	11.338	5.61e-12 ***
--				

Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 1

Residual standard error: 0.8084 on 28 degrees of freedom
Multiple R-squared: 0.8211 Adjusted R-squared: 0.8148
F-statistic: 128.6 on 1 and 28 DF, p-value: 5.612e-12

n = 30

Switching to multiple linear regression

- Observe data $(y_i, x_{i1}, \dots, x_{ip})$ for subjects $1, \dots, n$. Want to estimate $\beta_0, \beta_1, \dots, \beta_p$ in the model

$$y_i = \underbrace{\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}}_{E(y|x) = f(x;\beta)} + \epsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- Assumptions (residuals have mean zero, constant variance, are independent) are as in SLR
- Notation is cumbersome. To fix this, let

- $x_i = [1, x_{i1}, \dots, x_{ip}]$
- $\beta^T = [\beta_0, \beta_1, \dots, \beta_p]$
- Then $y_i = x_i \beta + \epsilon_i$

$$\begin{bmatrix} 1 & x_{i1} & \dots & x_{ip} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{bmatrix} = \beta_0 1 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

Matrix notation

$$y_i = x_i \beta + \epsilon$$

- Let

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{np} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}, \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix}$$

Annotations: A blue bracket under y points to the first column of \mathbf{y} . A blue bracket under X points to the first column of \mathbf{X} . A blue bracket under β points to the first column of $\boldsymbol{\beta}$. A blue bracket under ϵ points to the first column of $\boldsymbol{\epsilon}$. A blue arrow labeled x_1 points to the first column of \mathbf{X} .

- Then we can write the model in a more compact form:

$$\underbrace{\mathbf{y}_{n \times 1}}_{\mathbf{y}} = \underbrace{\mathbf{X}_{n \times (p+1)}}_{\mathbf{X}} \underbrace{\boldsymbol{\beta}_{(p+1) \times 1}}_{\boldsymbol{\beta}} + \underbrace{\boldsymbol{\epsilon}_{n \times 1}}_{\boldsymbol{\epsilon}}$$

- \mathbf{X} is called the *design matrix*

$$\mathbf{y} = \mathbf{x}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad \checkmark \checkmark \checkmark$$

Matrix notation

$$E \left(\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} \right) = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}$$

$\epsilon_i \stackrel{\text{iid}}{\sim} (\mu, \sigma^2)$

$$\text{Var} \left(\begin{bmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{bmatrix} \right) = \begin{bmatrix} \sigma^2 & \cdots & \cdots & 0 \\ 0 & \sigma^2 & & \\ 0 & & \ddots & \\ 0 & & & \ddots \\ \vdots & & & \\ 0 & & & \sigma^2 \end{bmatrix}$$

$$y = X\beta + \epsilon$$

- ϵ is a random vector rather than a random variable
- $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma^2 I$
- Note that Var is potentially confusing; in the present context it means the “variance-covariance matrix”

Interpretation of coefficients

$$\boxed{\beta_1} = \underbrace{(\beta_0 + \beta_1)}_{\beta_0} + \underbrace{(\beta_2)}_{\beta_2}$$

$$= (\beta_0 + \beta_1 + \beta_2 + \dots) - (\beta_0 + \beta_2)$$

$\uparrow \quad \uparrow \quad \uparrow \quad \uparrow$

$^{43} \quad ^{43}$

$$E(y|x_1=1, x_2=2, \dots) - E(y|x_1=0, x_2=1)$$

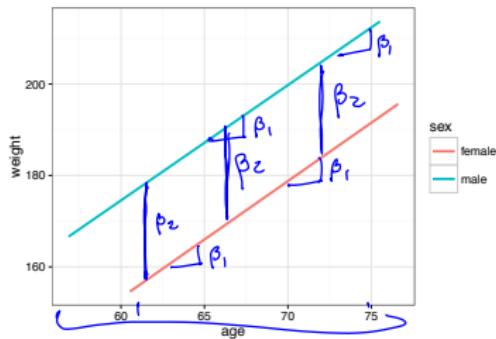
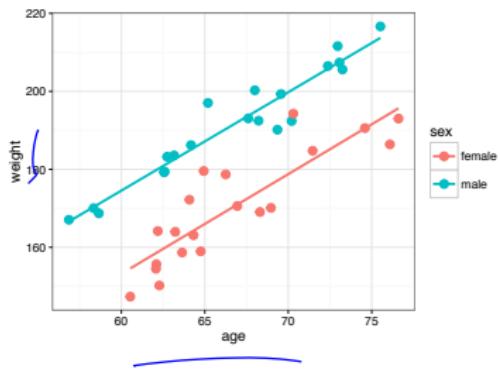
β_1 = the diff in $E(y)$ for a 1-unit Δx_1

Keeping everything else fixed!

Example with two predictors

$$E(y|x) = \beta_0 + \beta_1 x_{age} + \beta_2 x_{sex}$$

$$\beta_0 = E(y | age=0, sex=male)$$



Interpretation of coefficients

$$\underline{x_{i, \text{sex}} = 0}$$

$$\downarrow y_i = \beta_0 + \beta_1 x_{i, \text{age}} + \epsilon_i$$

$$\beta_0 = E[y | \text{age} = 0, \text{men}]$$

$\beta_1 = \Delta E(y) \text{ for a 1 unit } \Delta \text{age for men}$

$$\underline{x_{i, \text{sex}} = 1}$$

$$\hat{y}_i = \underline{\beta_0 + \beta_1 x_{i, \text{age}} + \beta_2 + \beta_3 x_{i, \text{age}}}$$

$$y_i = \underbrace{(\beta_0 + \beta_2)}_{\text{Int for women}} + \underbrace{(\beta_1 + \beta_3)}_{\text{Age slope for women}} x_{i, \text{age}} + \epsilon_i$$

$$\beta_0 + \beta_2 = \text{Int for women}$$

$$\beta_1 + \beta_3 = \text{Age slope for women}$$

$$\left| \begin{array}{l} \beta_2 = (\beta_0 + \beta_2) - (\beta_0) \\ \beta_3 = \overbrace{(\beta_1 + \beta_3)}^{\uparrow} - \overbrace{(\beta_1)}^{\uparrow} \end{array} \right.$$

Indicator variables

- Let x be a categorical variable with k levels (e.g. with $k = 3$, “low”, “med”, “high”).
- Choose one group as the baseline (e.g. “low”)
- Create $(k - 1)$ binary terms to include in the model:

$$\begin{aligned} \underline{x_{\text{med},i}} &= I(x_i = \text{“med”}) \\ \underline{x_{\text{high},i}} &= I(x_i = \text{“high”}) \end{aligned} \quad \left. \begin{array}{l} x_{\text{med}} = 0 \\ x_{\text{high}} = 0 \end{array} \right\} \Rightarrow \text{low!}$$

- For a model with no additional predictors, pose the model

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_{k-1} x_{i,k-1} + \epsilon_i$$

and estimate parameters using least squares

- Note distinction between *predictors* and *terms*

Categorical predictor design matrix

	X_{i1}	X_{i2}	X	has	3 levels
1	0	0			
1	0	0			
1	0	0			
1	0	0			
1	1	0			
1	1	0			
1	1	0			
1	1	0			
1	1	0			
1	0	1			
1	0	1			
1	0	1			
1	0	1			

Polynomial models

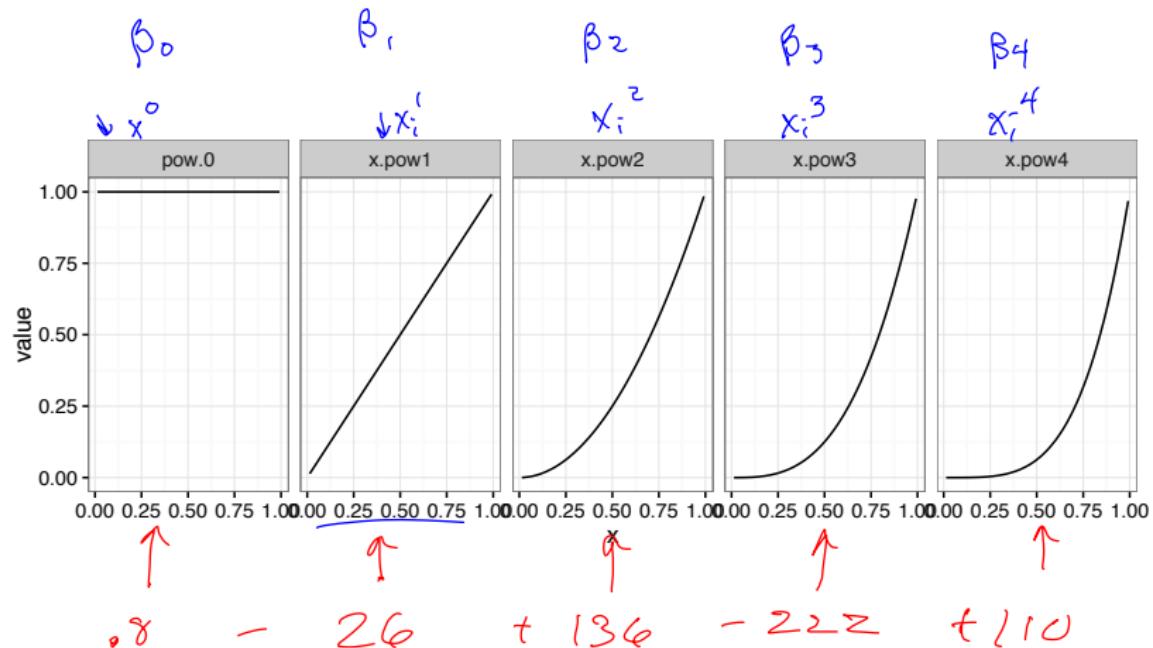
$$y_i = \underbrace{x_i}_{\text{g}} \underbrace{\beta_{\text{z}}}_{\text{g}} + \epsilon_i$$

- Model of the form

$$\underbrace{y_i}_{\text{g}} = \underbrace{\beta_0}_{\text{g}} + \underbrace{\beta_1}_{\text{g}} x_i + \underbrace{\beta_2}_{\text{g}} x_i^2 + \dots + \underbrace{\beta_p}_{\text{g}} x_i^p + \epsilon_i; \epsilon_i \stackrel{iid}{\sim} (0, \sigma^2)$$

- p is the polynomial order
- More polynomial terms can lead to a better approximation of $E(y|x)$, but also higher variability in the fit
- Conversely, smaller p can lead to inability to capture $E(y|x)$, but is often more stable
- Quadratic fits are pretty okay. I don't trust cubic and beyond.

Polynomial models

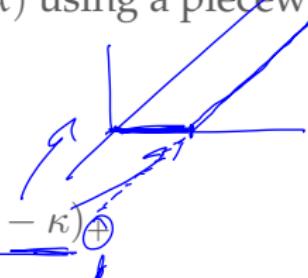


Piecewise linear models

Suppose we want to estimate $E(y|x) = f(x)$ using a piecewise linear model.

- For one knot we can write this as

$$E(y|x) = \underbrace{\beta_0 + \beta_1 x}_{\text{linear part}} + \underbrace{\beta_2(x - \kappa)}_{\text{change point}}$$



where κ is the location of the change point

Interpretation of regression coefficients

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 (x_i - k) + \epsilon_i$$

$x < k$

$$E(y|x+1) - E(y|x)$$

$$= \beta_0 + \beta_1(x+1) - \beta_0 - \beta_1 x$$

$$= \cancel{\beta_0} - \beta_0 + \beta_1 - \cancel{\beta_1 x} + \cancel{\beta_1}$$

$$\beta_1 = E\Delta y \text{ for } \approx 1 \text{ unit } \Delta x, x < k$$

$x > k$

$$E(y|x+1) - E(y|x)$$

$$= \beta_0 + \beta_1(x+1) + \beta_2(x+1-k) - \beta_0 - \beta_1 x - \beta_2(x-k)$$

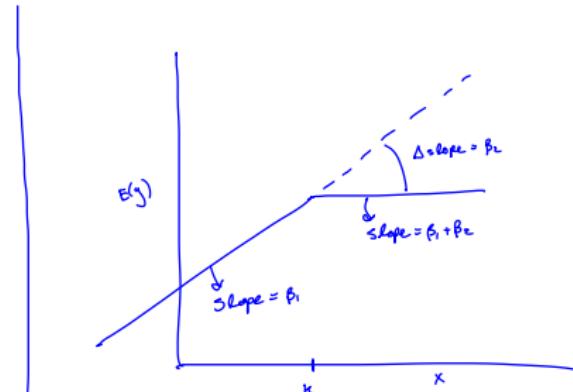
$$= \cancel{\beta_0} - \beta_0 + \beta_1 x - \cancel{\beta_1 x} + \beta_1 + \cancel{\beta_2 x} + \cancel{\beta_2 k} + \cancel{\beta_2 k} + \beta_2$$

$$\beta_1 + \beta_2 = E\Delta y \text{ for } \approx 1 \text{ unit } \Delta x, x > k$$

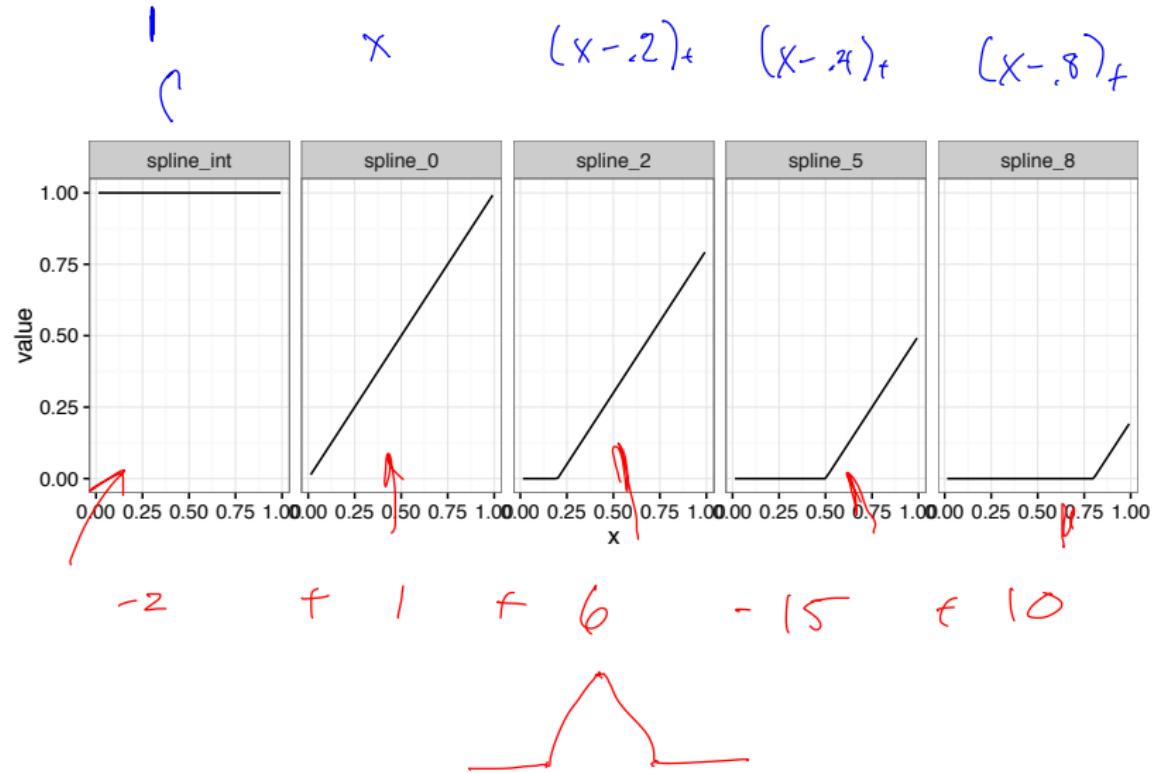
\equiv

$\beta_2 = \text{Change in slope, comparing}$

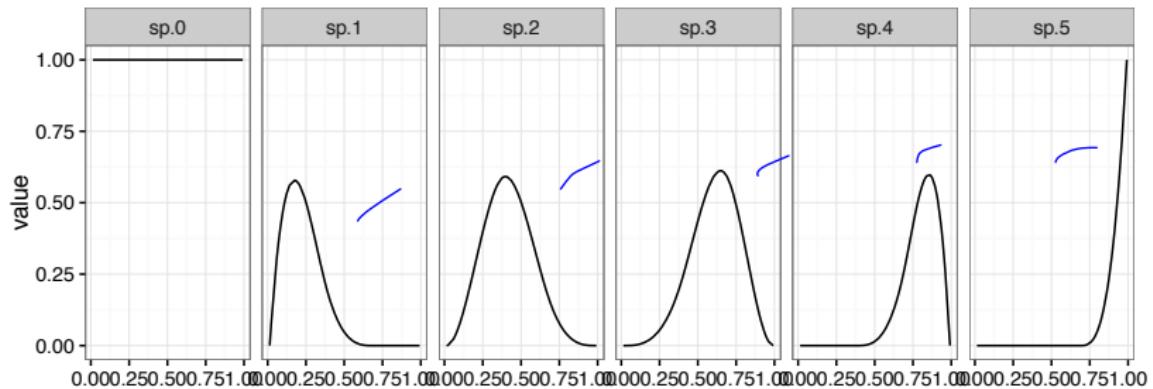
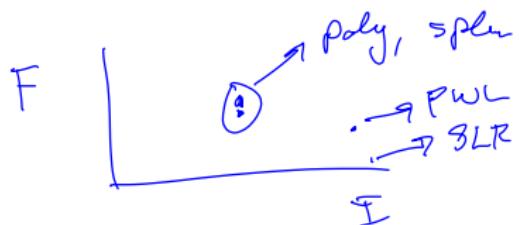
$x > k \text{ to } x < k$



Example



Spline models



$$\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 + \hat{\beta}_3 + \hat{\beta}_4 + \hat{\beta}_5$$



Least squares

As in simple linear regression, we want to find the β that minimizes the residual sum of squares.

$$\text{RSS}(\beta) = \sum_i \epsilon_i^2 = \underbrace{\{ \epsilon_1, \epsilon_2, \epsilon_3, \dots, \epsilon_n \}}_{\mathbf{\epsilon}^\top \mathbf{\epsilon}} \left[\begin{array}{c} \epsilon_1 \\ \vdots \\ \epsilon_n \end{array} \right]$$
$$\underline{(y - x\beta)^\top (y - x\beta)}$$

Least squares

$$\begin{aligned}
 \frac{\partial RSS(\beta)}{\partial \beta} &= \frac{\partial}{\partial \beta} (y - x\beta)^T (y - x\beta) \\
 &= (y^T - \beta^T x^T)(y - x\beta) \\
 &\quad \beta^T x^T x \beta - \beta^T x^T y - \underbrace{y^T x \beta + y^T y}_{RSS(\beta)} \\
 &= (\underbrace{\beta^T x^T x \beta - 2\beta^T x^T y + y^T y}_{RSS(\beta)}) \\
 &\quad \left(x^T x + (x^T x)^T \right) \beta - 2x^T y \\
 &\quad 2x^T x \beta - 2x^T y = 0 \\
 &\boxed{(x^T x)\beta = x^T y} \\
 \hat{\beta} &= (x^T x)^{-1} x^T y
 \end{aligned}$$

Unbiasedness of LSEs

$$E(\hat{\beta}) =$$

$$y = x\beta + \epsilon$$

$$\epsilon \sim (0, \sigma^2 I)$$

$$\checkmark y \sim (\underline{x}\underline{\beta}, \underline{\sigma^2 I}) \quad \checkmark$$

$$E(\underbrace{(x^T x)^{-1} x^T y})$$

$$= (x^T x)^{-1} x^T E(y)$$

~~$$= ((x^T x)^{-1} x^T x) \beta$$~~

$$= \beta$$

Variance of LSEs

$$Var(\hat{\beta}) =$$

$$\begin{aligned} & \text{Var}(\underline{(x^T x)^{-1} x^T y}) \\ &= (x^T x)^{-1} x^T \underbrace{\text{Var}(y)}_{\sigma^2 I} x (x^T x)^{-1} \\ &= \sigma^2 (x^T x)^{-1} \cancel{(x^T x)(x^T x)^{-1}} \\ &= \sigma^2 (x^T x)^{-1} \end{aligned}$$

$$\hat{Var}(\hat{\beta}) = \hat{\sigma}^2 (x^T x)^{-1}$$

$$Var(c\hat{\beta}) =$$

$$\underbrace{c \text{Var}(\hat{\beta}) c^T}_{\sigma^2 c (x^T x)^{-1} c^T}$$

$$\left[\frac{\begin{bmatrix} 1 & x_{i1} & x_{i2} & x_{i3} \end{bmatrix}^T}{c^T} \right]^T = y_i$$

Definitions

"Hat matrix"

- Fitted values: $\hat{y} = X\hat{\beta} = X(X^T X)^{-1}X^T y = Hy$
 - Residuals / estimated errors: $\hat{\epsilon} = y - \hat{y}$
 - Residual sum of squares: $\sum_{i=1}^n \hat{\epsilon}_i^2 = \hat{\epsilon}^T \hat{\epsilon}$
 - Residual variance: $\hat{\sigma}^2 = \frac{\text{RSS}}{n-p-1}$
 - Degrees of freedom: $n - p - 1$
- $y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$
- # pred other than intercept

Hat matrix

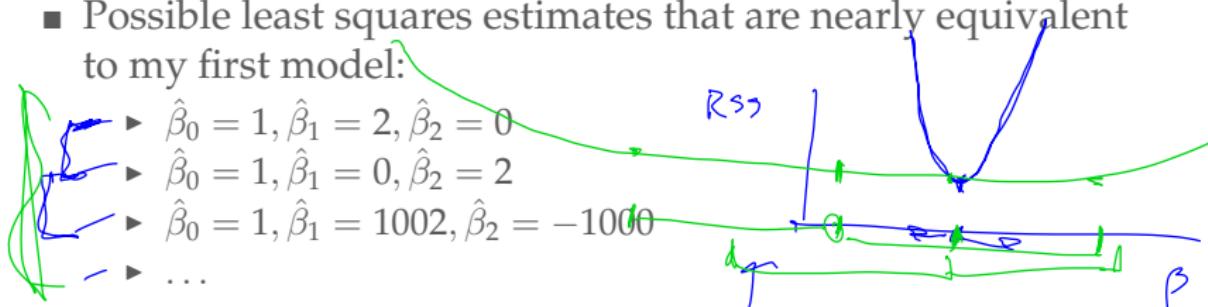
Some properties of the hat matrix:

- It is a projection matrix: $HH = H$
- It is symmetric: $H^T = H$
- The residuals are $\hat{\epsilon} = (\underline{I} - \underline{H})\underline{y}$ $\underline{y} - \underline{H}\underline{y} = \underline{y} - \hat{\underline{y}}$
- The inner product of $(I - H)y$ and Hy is zero (predicted values and residuals are uncorrelated).

Effects of collinearity

Suppose I fit a model $y_i = \beta_0 + \beta_1 x_{i1} + \epsilon_i$.

- I have estimates $\hat{\beta}_0 = 1, \hat{\beta}_1 = 2$
- I put in a new variable $x_2 = x_1 + \text{error}$, where *error* is pretty small
- My new model is $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i$
- Possible least squares estimates that are nearly equivalent to my first model:



- A unique solution exists, but it is hard to find

Sampling distribution of $\hat{\beta}$

$$\hat{\beta} \sim (\beta, \sigma^2(X^T X)^{-1})$$

If our usual assumptions are satisfied and $\epsilon \sim \underline{N}[0, \sigma^2 I]$ then

$$\hat{\beta} \sim N\left[\beta, \sigma^2(X^T X)^{-1}\right]. \quad y \sim N(x\beta, \sigma^2 I)$$

- This will be used later for inference.
- Even without Normal errors, asymptotic Normality of LSEs is possible under reasonable assumptions.

Asymptotic distribution

Assume that

- $E(\epsilon_i | \mathbf{x}_i) = \underline{0} \forall i;$
- $Var(\epsilon_i | \mathbf{x}_i) = \underline{\sigma^2} \forall i;$
- $n \xrightarrow{\text{lim}} \infty \frac{\mathbf{X}^T \mathbf{X}}{n} \rightarrow Q$ where Q is a finite non-singular matrix.

Then

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{} N \left[0, \sigma^2 \underbrace{Q^{-1}}_{\downarrow} \right]$$

(This is essentially an extension of the central limit theorem)

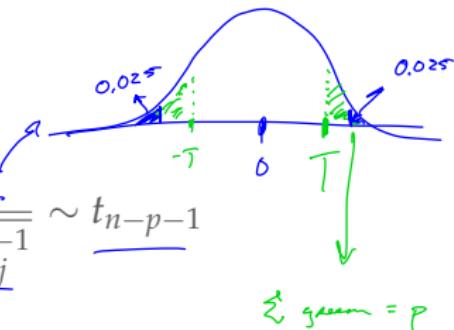
Individual coefficients

est - anal
 \hat{s}_e

For individual coefficients

- We can use the test statistic

$$T = \frac{\hat{\beta}_j - \beta_j}{\hat{s}_e(\hat{\beta}_j)} \stackrel{H_0: \beta_j = 0}{=} \frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\sigma}^2(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}} \sim t_{n-p-1}$$



$$\sum \text{green} = p$$

- For a two-sided test of size α , we reject if

$$|T| > t_{1-\alpha/2, n-p-1}$$

always give $P!!!!!!$

- The p-value gives $P(t_{n-p-1} > T_{obs} | H_0)$

Note that t is a symmetric distribution that converges to a Normal as $n - p - 1$ increases.

Inference for linear combinations

$$H_0: \beta_1 - \beta_2 = 0 \quad ?? \quad | \quad \text{Cat} \quad y_i = \beta_0 + \underbrace{\beta_1}_{\begin{bmatrix} 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}} y_{i1} + \underbrace{\beta_2}_{\begin{bmatrix} 0 & 1 & -1 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}} y_{i2} + \epsilon_i$$

$$\begin{aligned} &= (\beta_1 - \beta_0) - (\beta_2 - \beta_0) \\ &= \beta_1 - \beta_2 \end{aligned}$$

Sometimes we are interested in making claims about $c^T \beta$ for some c .

$$c = \begin{bmatrix} 1 & x_{i1} & x_{i2} \end{bmatrix}$$

- Define $H_0: c^T \beta = c^T \beta_0$ or $H_0: c^T \beta = 0$
- We can use the test statistic

$$\text{Var}(c^T \hat{\beta})$$

$$T = \frac{c^T \hat{\beta} - c^T \beta_0}{\widehat{\text{se}}(c^T \hat{\beta})} = \frac{c^T \hat{\beta} - c^T \beta_0}{\sqrt{\widehat{\sigma}^2 c^T (\mathbf{X}^T \mathbf{X})^{-1} c}}$$

$$\widehat{\sigma}^2 = \widehat{\sigma}^2 c^T (\mathbf{X}^T \mathbf{X})^{-1} c$$

- This test statistic is asymptotically Normally distributed
- For a two-sided test of size α , we reject if

$$|T| > z_{1-\alpha/2}$$

Global tests

$$H_0: \beta_1 = \beta_2 = 0$$

$$y_i = \beta_0 + \underbrace{\beta_1 x_{i1} + \beta_2 x_{i2}} + \beta_3 x_{i3} + \dots + e_i \quad \checkmark$$

$$\downarrow y_i = \beta_0 + 0 + 0 + \beta_3 x_{i3} + \dots + e_i \quad \checkmark$$

Compare a smaller “null” model to a larger “alternative” model

- Smaller model must be nested in the larger model
- That is, the smaller model must be a special case of the larger model
- For both models, the *RSS* gives a general idea about how well the model is fitting
- In particular, something like

$$\frac{RSS_S - RSS_L}{RSS_L}$$

small = Reject
Big = fail to rej.

compares the relative *RSS* of the models

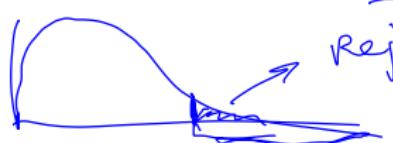
Global F tests

- Compute the test statistic

$$F_{obs} = \frac{(RSS_S - RSS_L) / (df_S - df_L)}{RSS_L / df_L}$$

- If H_0 (the null model) is true, then $F_{obs} \sim F_{\underline{df_S - df_L}, \underline{df_L}}$
- Note $\underline{df_S} = n - \underline{p_S} - 1$ and $\underline{df_L} = n - \underline{p_L} - 1$
- We reject the null hypothesis if the p-value is above α , where

$$\text{p-value} = P(F_{\underline{df_S - df_L}, \underline{df_L}} > F_{obs})$$



Global F tests

Can

$$y_i = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon_i$$

$$H_0: \beta_1 + \beta_2 = 0$$

$$-y_i = \beta_0 + \epsilon_i$$

There are a couple of important special cases for the F test

- The null model contains the intercept only
 - ▶ When people say ANOVA, this is often what they mean (although all F tests are based on an analysis of variance)
- The null model and the alternative model differ only by one term

- ▶ Gives a way of testing for a single coefficient

- ▶ Turns out to be equivalent to a two-sided t-test: $t_{df_L}^2 \sim F_{1, df_L}$

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i \\ y_i &= \beta_0 + \beta_1 x_{i1} + \epsilon_i \end{aligned} \quad \left. \right\} H_0: \beta_2 = 0$$

Alternative global tests: the Wald test

$$\left(\frac{\hat{\beta}_j - \beta_0}{\hat{\sigma}_{\hat{\beta}_j}(\hat{\beta}_j)} \right)^2 = t^2$$

For a vector of coefficients, we can test $H_0 : \underline{\beta} = \beta_0$:

- Use the test statistic

$$W = (\hat{\beta} - \beta_0)^T \underline{\hat{Var}(\hat{\beta})}^{-1} (\hat{\beta} - \beta_0)$$

- Under the null, this test statistic has an asymptotic χ_p^2 distribution
- In practice, we replace $Var(\hat{\beta})$ with $\widehat{Var}(\hat{\beta})$ and use an F distribution

Alternative global tests: the likelihood ratio test

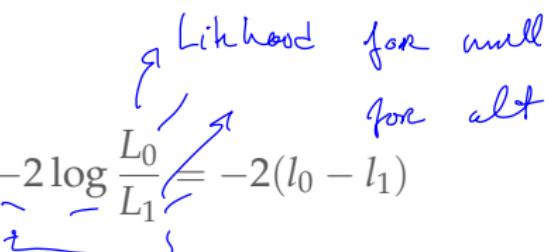
$$H_0: \beta_2 = \beta_3 = 0$$

If we are using maximum likelihood estimation (we'll cover this soon – turns out to be least squares in MLR), we can use a LRT:

- Use the test statistics

$$\Delta = -2 \log \frac{L_0}{L_1} = -2(l_0 - l_1)$$

*Likelihood for null
for alt*



- This test statistic has an asymptotic χ_d^2 distribution where d is the difference in the number of parameters between the two models.
- Must compare nest models

Confidence intervals: individual parameters

$$H_0: \beta_j = \underline{\beta_0} \Rightarrow \frac{\hat{\beta}_j - \beta_0}{\widehat{se}(\hat{\beta}_j)}$$

- A confidence interval with coverage $(1 - \alpha)$ is given by

$$\hat{\beta}_j \pm t_{1-\alpha/2, n-p-1} \widehat{se}(\hat{\beta}_j) \quad \left(\hat{\beta}_j \pm z \widehat{se}(\hat{\beta}) \right)$$

- Assuming all the standard assumptions hold,

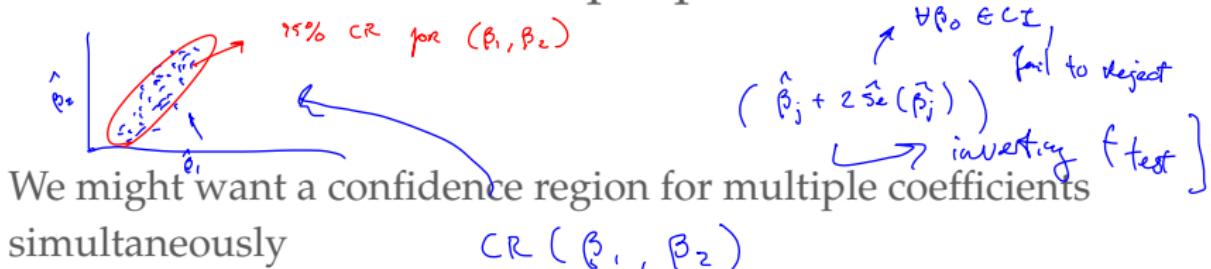
$$(1 - \alpha) \text{ ""=} P(LB < \beta_j < UB)$$

.95 0 / 1

Note there is a one-to-one correspondence between this confidence interval and the hypothesis test.

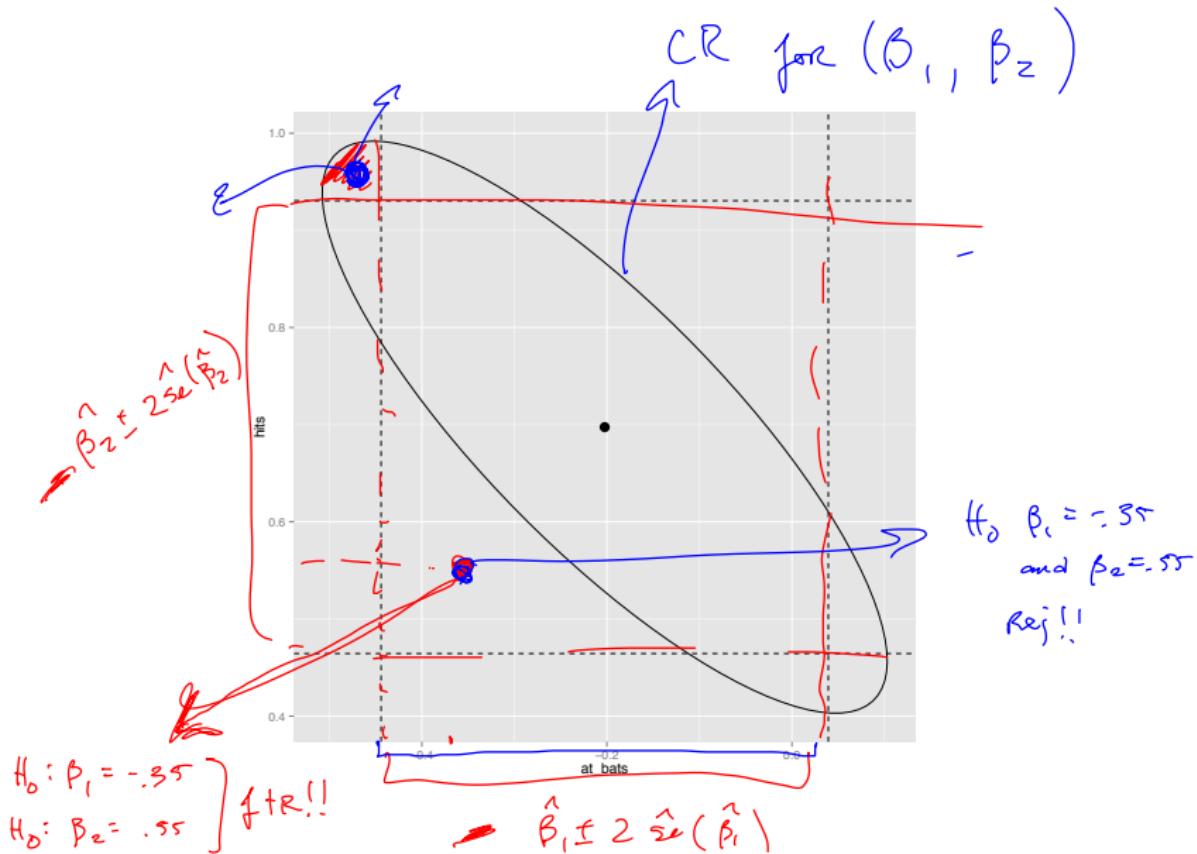
$$\frac{\hat{\beta}_j - \beta_0}{\widehat{se}(\hat{\beta})} < t_{\alpha/2} \Rightarrow \beta_0 \in (\hat{\beta}_j \pm t_{\alpha/2} \widehat{se}(\hat{\beta}))$$

Confidence intervals: multiple parameters



- Invert Wald test for multiple coefficients – find region containing all values β_0 for which p-value from global Wald test is $> \alpha$ $\downarrow H\beta_0, (\hat{\beta} - \beta_0)^T (\underline{\text{Var}(\hat{\beta})}^{-1}) (\hat{\beta} - \beta_0)$
- Then $(1 - \alpha) = P[\beta \in \text{region}] \leq \text{Critical Value}$
- This region is an ellipsoid in higher dimensions; we can visualize in 2D most easily and 3D pretty well.

Confidence intervals: multiple parameters



Predictions and prediction intervals

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2}$$

$$v_{\alpha}(\hat{y}) = x_0^T (\sigma^2 (X^T X)^{-1}) x_0$$

$$\begin{matrix} + \epsilon; \\ \epsilon \sim (0, \sigma^2) \end{matrix}$$

- What is the prediction value y for a given x_0
- What range would you give for the value of a new outcome?
- Two sources of variance to consider: variance in estimates and variance in outcome

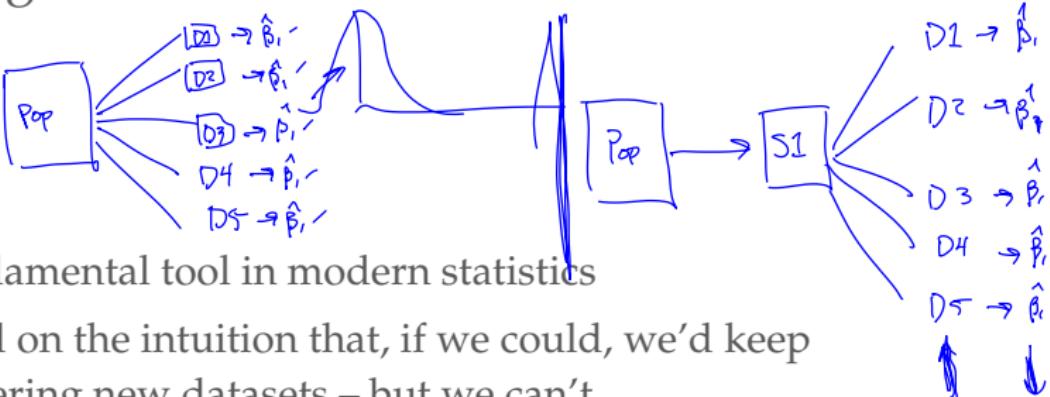
$$Var(y) = \underbrace{Var(E(y|x))}_{\sigma^2 x_0^T (X^T X)^{-1} x_0} + \underbrace{E(Var(y|x))}_{\sigma^2}$$

A prediction interval is given by

$$(\hat{y}|x = x_0) \pm t_{1-\alpha/2, n-p-1} \hat{s}e_{pred}(\hat{y}|x_0)$$

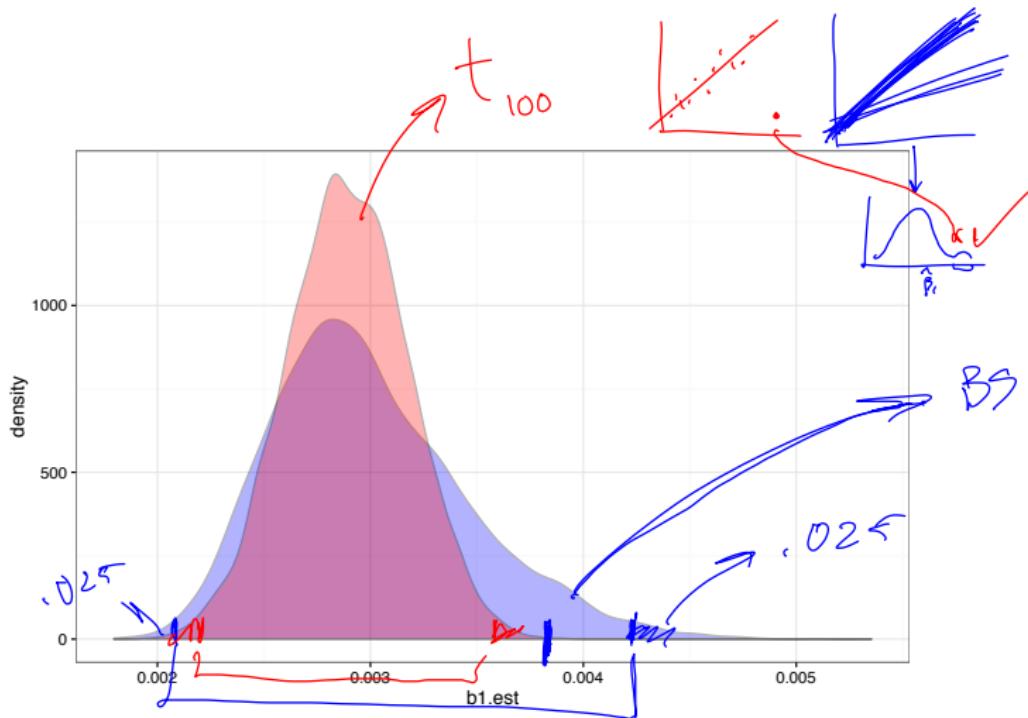
This can be estimated for any x_0 .

Resampling methods



- Fundamental tool in modern statistics
- Build on the intuition that, if we could, we'd keep gathering new datasets – but we can't
- Use repeated samples of a training set to understand variability
- Computationally intensive ... but we have computers

Bootstrap results



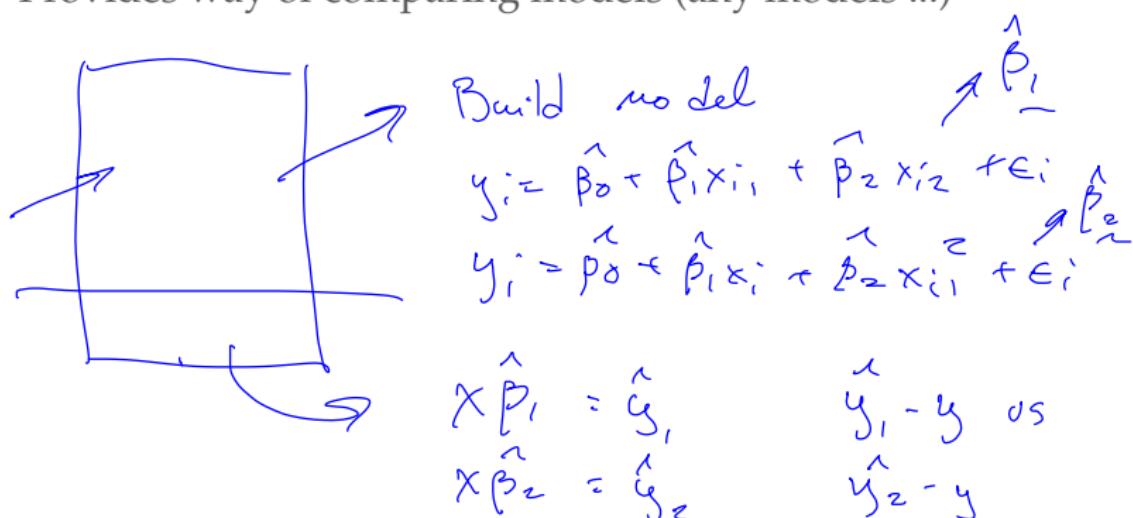
Cross Validation

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \epsilon_i \quad]$$
$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i1}^2 + \epsilon_i \quad]$$

- Focus is on model performance, quantified by prediction error
- We get in-sample performance ...
- But we want generalization to new data
- Most of the time, we don't have an external testing dataset

Cross Validation: one validation set

- Simplest case: create a validation set by randomly splitting the full dataset
- Fit model to training data; compute mean squared prediction error on test set
- Provides way of comparing models (any models ...)



Gauss-Markov theorem

Assume the model

$$\underline{y = X\beta + \epsilon}$$

where $E(\epsilon) = 0$ and $Var(\epsilon) = \sigma^2 I$. Also assume X is a full rank design matrix.

- Among all unbiased linear estimators \underline{Cy} of the regression coefficients β , the LSE has minimum variance and is unique.

We call the LSEs “BLUE”.

Maximum likelihood estimation

Using matrix notation:

$$L(\beta; y) = \underbrace{(2\pi)^{-n/2} (|-\mathbf{c}^T \mathbf{I}|)^{-1/2}}_{\downarrow} \exp \left\{ -\frac{1}{2} \underbrace{\mathbf{y}^T (\mathbf{y} - \mathbf{x}\beta) \mathbf{c}^T \mathbf{I}^{-1} (\mathbf{y} - \mathbf{x}\beta)}_{\text{RSS}(\beta)} \right\}$$

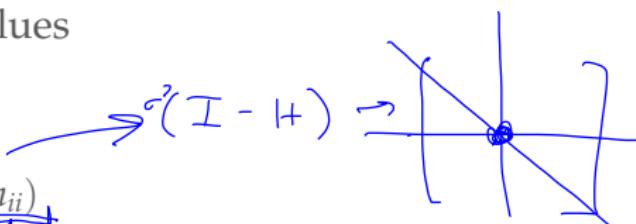
$$\frac{(\mathbf{x}^T \mathbf{x})^{-1} \mathbf{x}^T \mathbf{y}}{\vdots \quad | \quad | \quad | \quad | \quad \vdots}$$

Residuals when model is correct

- Often we plot the residuals against one of the predictors or against the fitted values
- What we look for:

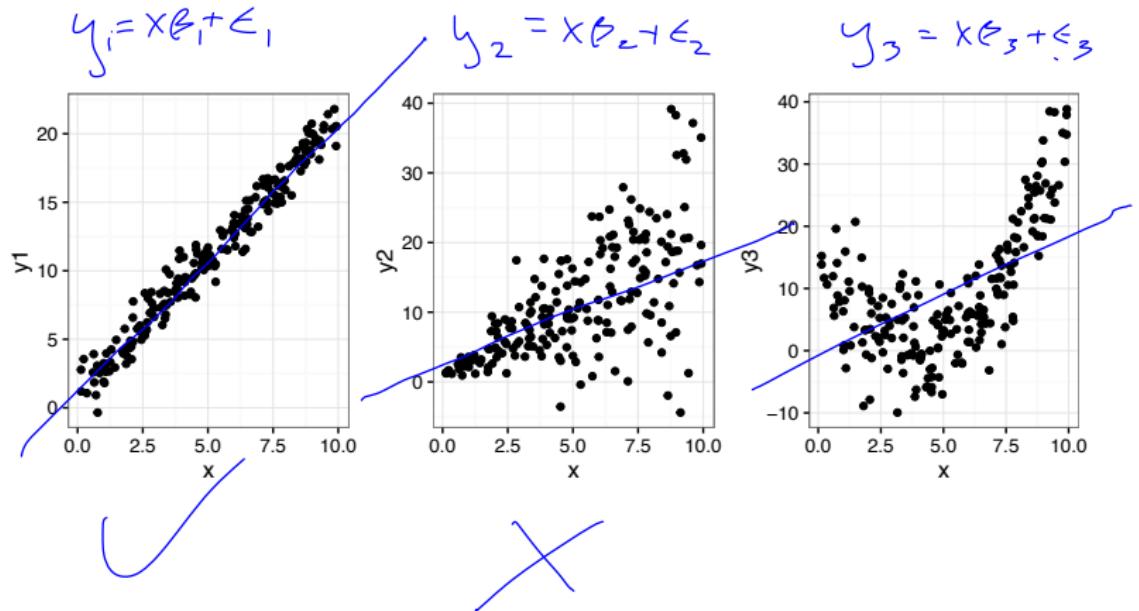
✓ ▶ $E(\hat{\epsilon}|x) = 0$

✓ ▶ $V(\hat{\epsilon}|x) = \sigma^2(1 - h_{ii})$

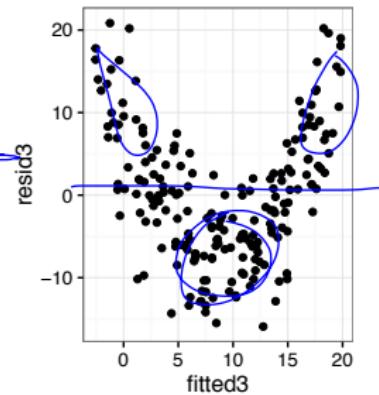
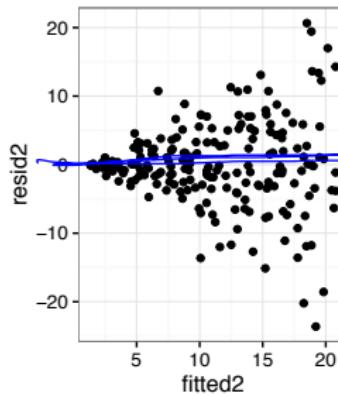
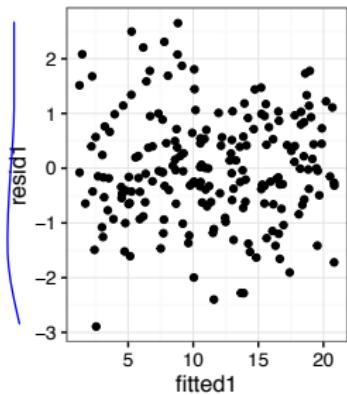


- If the model is incorrect, you may be able to spot:
 - ▶ Patterns in the residuals
 - ▶ Clear non-constant variance

Some data plots



Some residual plots



Model checking

Two major areas of concern:

- Global lack of fit, or general breakdown of model assumptions

- "Linearity" $x\beta$ is "right"
- Unbiased, uncorrelated errors $E(\epsilon|x) = E(\epsilon) = 0$
- Constant variance $\text{Var}(y|x) = \text{Var}(\epsilon|x) = \sigma^2$
- Independent errors
- Normality of errors

$$\epsilon \sim (0, \sigma^2 I)$$

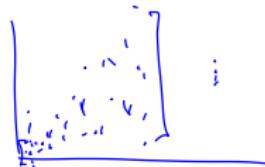
$$y \sim (\beta_0 + \beta_1 x, \sigma^2 I)$$

- Effect of influential points and outliers

Model checking

- Global lack of fit, or general breakdown of model assumptions
 - ▶ Residual analysis – QQ plots, residual plots against fitted values and predictors
 - ▶ Adjusted variable plots
- Effect of influential points and outliers
 - ▶ Measure of leverage, influence, outlying-ness

Variance-stabilizing transformation



Suppose y is strictly positive, $\mu = E(y|x)$, $\text{Var}(y|x) = \sigma^2 g(\mu)$

- Replace y with $y^* = T(y)$ such that $\text{Var}(y^*|x)$ is approximately constant
- Delta method says $\text{Var}(T(y)) = (T'(\mu))^2 \sigma^2 g(\mu)$

$$\mu_0 \in e^{\beta} = 0$$

Variance-stabilizing transformation

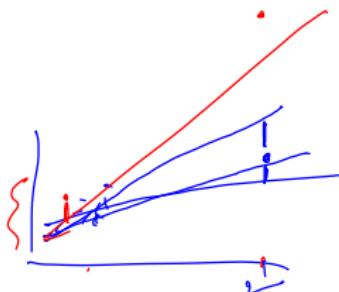
To get constant variance, we want

$$\text{Var}(T(y)) = \frac{(T'(\mu))^2 g(\mu)}{\sqrt{g(\mu)}} = k^2 \text{ (constant)}$$
$$\Rightarrow T'(\mu) = \frac{k}{\sqrt{g(\mu)}} \quad \checkmark$$
$$\Rightarrow T(\mu) = \int \frac{k}{\sqrt{g(\mu)}} d\mu$$
$$\Rightarrow T(y) = \int \frac{k}{\sqrt{g(y)}} dy$$

So the transformation necessary to stabilize the variance really depends on the variance function itself, e.g. $g(\cdot)$

Isolated points

Points can be isolated in three ways



- Leverage point – outlier in x
- Outlier – outlier in $y|x$
- Influential point – a point that largely affects β
 - ▶ Deletion influence; $|\hat{\beta} - \hat{\beta}_{(-i)}|$
 - ▶ Basically, a high-leverage outlier

Leverage is measured by the hat matrix, outlying-ness by the residual

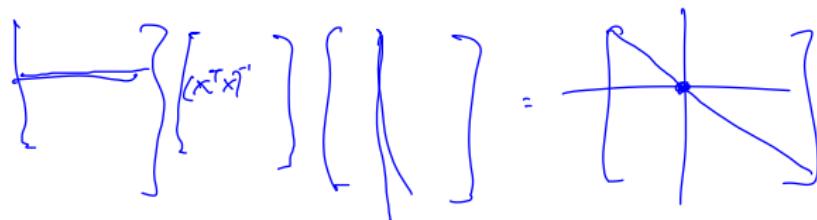
Quantifying leverage

We measure leverage (the “distance” of x_i from the distribution of x) using

$$h_{ii} = \underline{x_i^T} \underline{(\mathbf{X}^T \mathbf{X})^{-1}} \underline{x_i}$$

where h_{ii} is the $(i, i)^{th}$ entry of the hat matrix.

$$\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$



Leverage

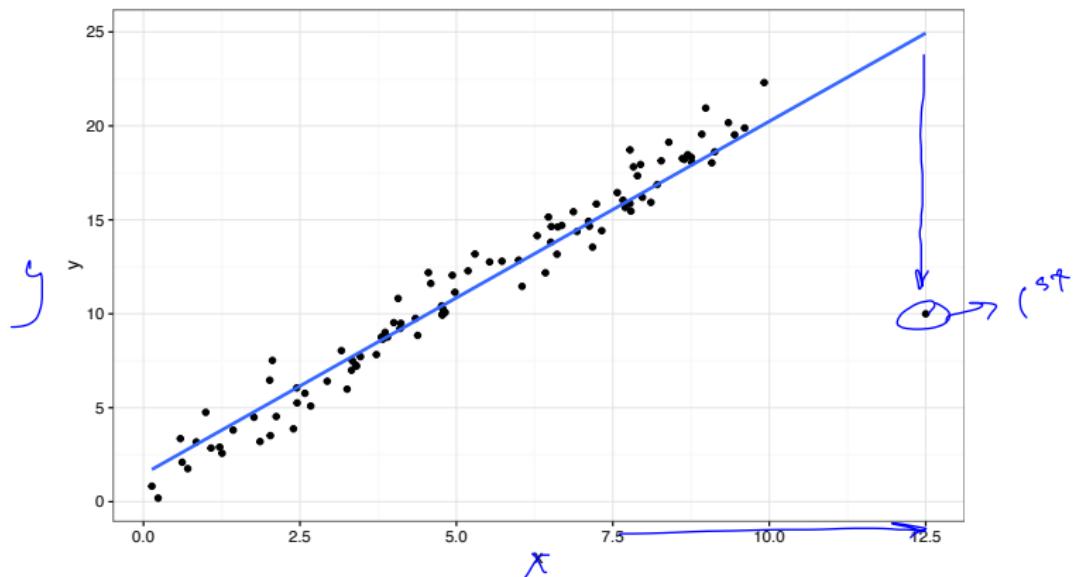
Some notes about the hat matrix

- $\sum_i h_{ii} \stackrel{\text{def}}{=} \text{tr}(\mathbf{H}) = (p + 1)$
$$\text{tr}(\underbrace{\mathbf{x}(\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top}_{\mathbf{I}}) = \text{tr}(\underbrace{\mathbf{x}^\top (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}}_{\mathbf{I}})$$

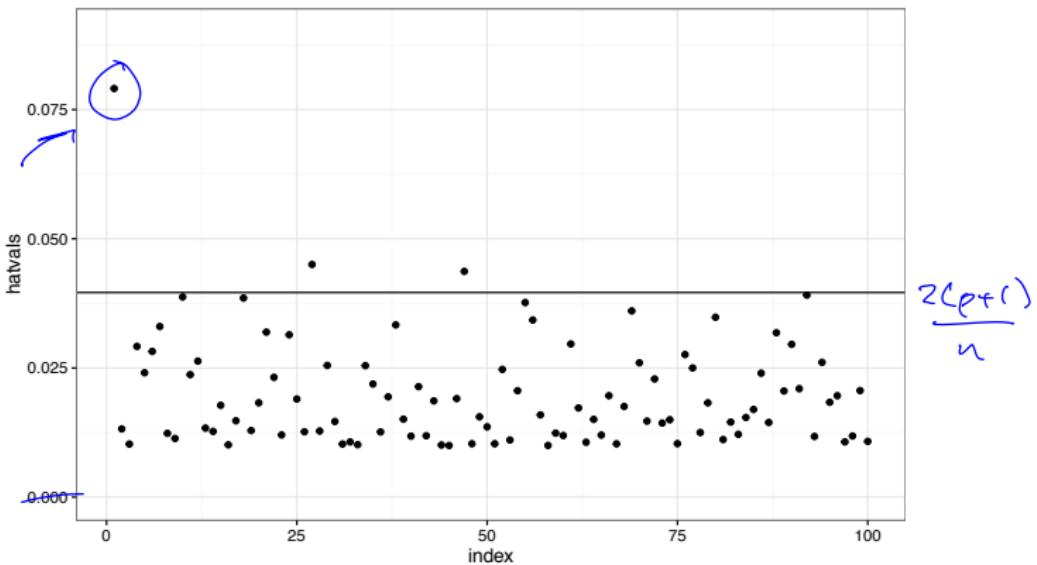
(Note – the trace of the hat matrix generalizes to non-parametric methods, where you don't have a specific number of parameters to count. This is a useful measure of "model size" or "effective degrees of freedom" in these cases.)

$$S_y = \hat{y}$$

Leverage plot



Leverage plot



Outliers

- When we refer to “outliers” we typically mean “points that don’t have the same mean structure as the rest of the data”
- Residuals give an idea of “outlying-ness”, but we need to standardize somehow
- Remember (from last lecture) $Var(\hat{\epsilon}_i) = \sigma^2(1 - h_{ii}) \dots$

$$\frac{1}{\sqrt{1 - h_{ii}}}$$

Outliers

The standardized residual is given by

$$\hat{\epsilon}_i^* = \frac{\hat{\epsilon}_i}{\sqrt{Var(\hat{\epsilon}_i)}} = \frac{\hat{\epsilon}_i}{\hat{\sigma} \sqrt{(1 - h_{ii})}}$$

The *Studentized* residual is given by

$$t_i = \frac{\hat{\epsilon}_i}{\hat{\sigma}_{(-i)} \sqrt{(1 - h_{ii})}} = \hat{\epsilon}_i^* \left(\frac{n - (p + 1)}{n - (p + 1) - \hat{\epsilon}_i^{*2}} \right)^{1/2}$$

Studentized residuals follow a $t_{\underline{n-(p+1)-1}}$ distribution.

Influence

Specifically, deletion influence

$$\|\hat{\beta} - \hat{\beta}_{(-i)}\| = \|(\hat{\beta} - \hat{\beta}_{(-i)})^T (\hat{\beta} - \hat{\beta}_{(-i)})\|^{1/2}$$

Cook's distance is

$$\begin{aligned} D_i &= \frac{(\hat{\beta} - \hat{\beta}_{(-i)})^T (\underline{X^T X}) (\hat{\beta} - \hat{\beta}_{(-i)})}{(p+1)\hat{\sigma}^2} \\ &= \frac{(\hat{y} - \hat{y}_{(-i)})^T (\hat{y} - \hat{y}_{(-i)})}{(p+1)\hat{\sigma}^2} \\ &= \underbrace{\frac{1}{p+1} \hat{\epsilon}_i^2}_{\text{Cook's Distance}} \underbrace{\frac{h_{ii}}{1-h_{ii}}}_{\text{Influence}} \end{aligned}$$