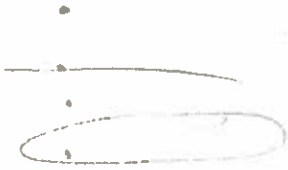


Intro Activity

- Have sharbill pick up as they walk in.
- Ask for theories.
- Show subsequent pics and ask for for revisions

TheoriesAt So, k ft Science

- Make observations in the natural world collect data
- Generate theories/models.
- Evaluate the consistency between the model and the data.
- Gather more data, revise theories
- Iterate on.

Regression: the central statistical model

- Express one variable of interest (response) as a function of some other variables (predictors)

$$Y = f(X_1, X_2, X_3, \dots)$$

- It is a probabilistic model

$$P(Y | X_1, X_2, \dots, X_n)$$

Examples from NYT

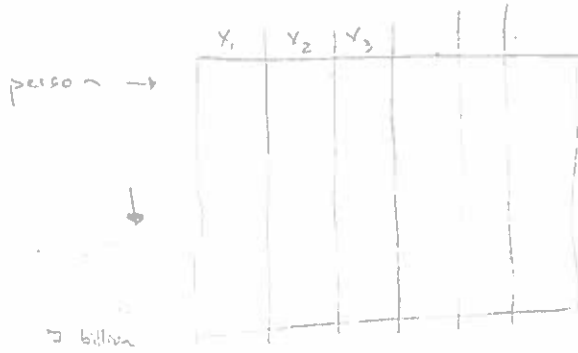
Activity

Let's build a model to predict someone's actual age.

Our only predictor for now is a guess based on a photograph.

Week 2 Day 1: Describing distributions

→ Taxonomy of Data



What variables do we have data on?

What type of data is each?

Write out pseudo-code for the subsetting operator they do:

↑ mean photo.

Taxonomy of Data

numerical

categorical R factor

continuous
numeric

discrete
R: integer

unordered

ordinal

One Variable Numerical Descriptors

center: mean, median, mode

spread: standard deviation, variance, IQR

shape: unimodal/bimodal

skew: right left

Graphical Descriptors

categorical

frequency/
probability

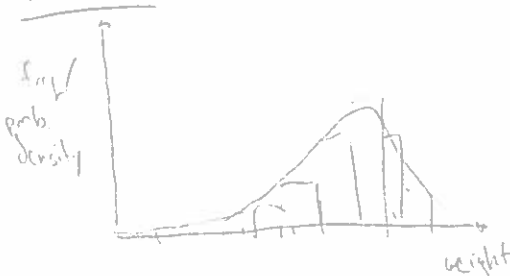
female male

- bar plot

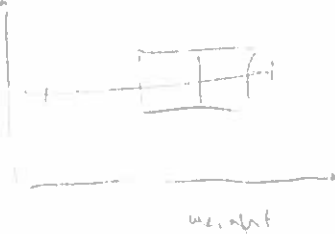
barplot (table (rectangles))



Numerical



Boxplot



- must be aggregated by binning (hist) or smoothing (density)
- beware of binwidth, bandwidth selection!

- Shows median, Q1, Q3, and outliers.
- good for side-by-side comparisons

Two Variables

Numerical summaries (2 num)

- shape: linear, quadratic
- direction: positive/neg, curvature
- strength: how tightly clustered?
correlation coefficient r

Graphical

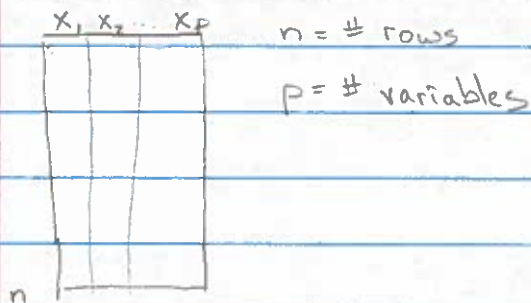


Scatterplot

- time often on x
- beware of overplotting
- response often on y

Week 3C: From Data to Distributions

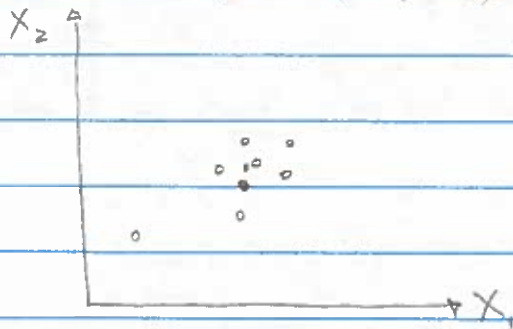
Representations of Data



1D Dotplot $\{x_i\}$



2D Scatterplot $\{x_1, x_2\}$



* In all of these representations, the entire data set is recoverable, zero information loss.

* There is value in taking the big picture view of the general, smoother structure at play

3D 3D Scatterplot $\{(x_1, x_2, x_3)\}$



From Data to Distributions

Question of interest: Where is the data dense? Where is it sparse?

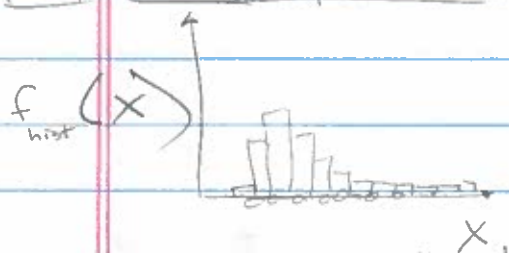
To answer this, we need ^(A) a measure of data density.

proportion of obs.

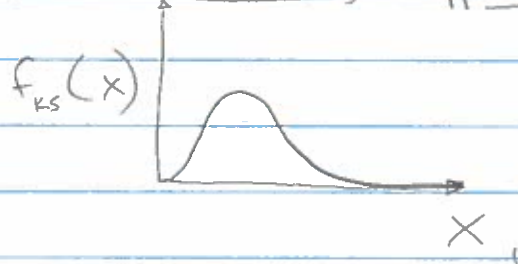
length/area/vol

^(B) a functional map from $X \rightarrow$ density

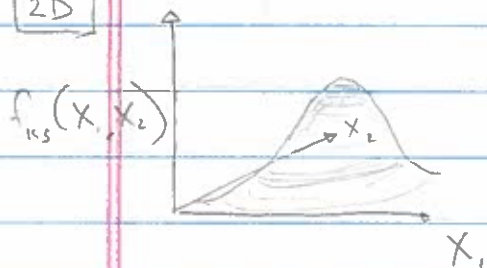
1D Discrete, disjoint bins



Continuous, overlapping bins



2D



3 Types of distributions

Joint Distribution: $f(X_1, X_2)$; what is the dist. of X_1 and X_2 ?

Marginal Distribution: $f(X_1)$, $f(X_2)$ what is the dist of X_1 over all X_2 ?

Conditional Dist: $F(X_1 | X_2 = x_2)$

what is the dist of X_1 given X_2 takes on this particular value?

Week 4A Simple Linear Regression

For Weds: Read 2.1-2.4

Notation

A random variable X takes particular values w/ particular probs.

~~eg1~~ Let X be the value on a rolled fair die.

x	$P(X=x)$
1	$1/6$
2	$1/6$
3	$1/6$
4	$1/6$
5	$1/6$
6	$1/6$

eg2. Let X be the # of heads in 3 fair coin flips

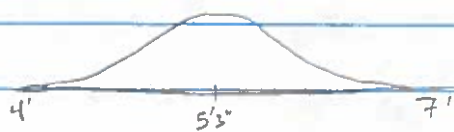
HHH
HTT
HHT
THT
TTH
TTH
TTH

x	$P(X=x)$
0	$1/8$
1	$3/8$
2	$3/8$
3	$1/8$

eg.3

Let X be the height of a randomly selected female college student. X is continuous between 4' and 7', normally distributed w/ mean 5'3" and variance 9".

$$X \sim N(\mu = 5'3", \sigma^2 = 9")$$



Expected Value:

$$E(X) = \sum_{\text{all } x} x P(X=x)$$

eg.2

x	$P(X=x)$
0	$1/8$
1	$3/8$
2	$3/8$
3	$1/8$

$$12/8 = 1.5$$

eg3

$$E(X) = \int_{-\infty}^{\infty} x f(x) dx = \mu$$

Variance

eg 2 $\text{Var}(X) = \sum_{x \in \mathcal{X}} (x - \mu)^2 P(X=x) \cdot \sigma^2$ eg 3. $\text{Var}(X) = \int (x - \mu)^2 f(x) dx$

Describe the relationship

Shape:

Direction:

Strength:

What you're describing is $E(Y|X=x)$ a conditional mean.

Simple Linear Regression

$$E(Y|X=x) = \underset{\substack{\uparrow \\ \text{intercept}}}{\beta_0} + \underset{\substack{\uparrow \\ \text{slope}}}{\beta_1} x$$

— mean function

$$Y_i = E(Y|X=x) + e_i = \beta_0 + \beta_1 x + e_i$$

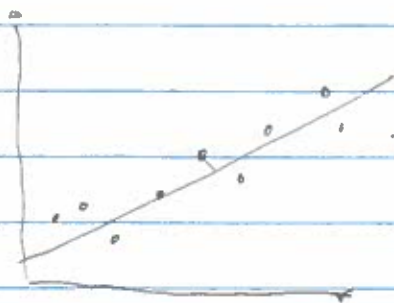
random error

$$E(e_i|X) = 0$$

$$\text{Var}(e_i|X) = \sigma^2$$

— data generating function

Which line?



two options

① minimize distance from (x_i, y_i) to line

* ② minimize distance from (y_i) to line

$$\text{residual: } \hat{e}_i = y_i - \hat{y}_i$$

minimize the sum of squared residuals

$$\sum_{i=1}^n \hat{e}_i^2 \quad (\text{aka RSS})$$

Fitting the Least Squares Line

$$RSS = \sum_{i=1}^n \hat{e}_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2$$

Since we have data for x and y , we can treat them as constant and view the RSS as a function of $\hat{\beta}_0$ and $\hat{\beta}_1$. To find the values that minimize RSS, we can take the derivative and set to zero.

$$\frac{\partial RSS}{\partial \hat{\beta}_0} = -2 \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

and

$$\frac{\partial RSS}{\partial \hat{\beta}_1} = -2 \sum_{i=1}^n x_i (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i) = 0$$

rearrange:

$$\sum_{i=1}^n y_i - \sum_{i=1}^n \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i = 0$$

$$\sum_{i=1}^n y_i = n \hat{\beta}_0 + \hat{\beta}_1 \sum_{i=1}^n x_i$$

normal eqns.

$$\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \hat{\beta}_0 - \sum_{i=1}^n \hat{\beta}_1 x_i^2 = 0$$

$$\sum x_i y_i = \hat{\beta}_0 \sum x_i + \hat{\beta}_1 \sum x_i^2$$

Solve for $\hat{\beta}_0$

$$n \hat{\beta}_0 = \frac{\sum y_i}{n} - \hat{\beta}_1 \frac{\sum x_i}{n} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\begin{aligned} \sum x_i y_i &= (\bar{y} - \hat{\beta}_1 \bar{x}) \sum x_i + \hat{\beta}_1 \sum x_i^2 \\ &= \bar{y} \sum x_i - \hat{\beta}_1 \bar{x} \sum x_i + \hat{\beta}_1 \sum x_i^2 \end{aligned}$$

$$\sum x_i y_i - n \bar{y} \bar{x} = \hat{\beta}_1 (\sum x_i^2 - n \bar{x}^2)$$

rearrange

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Week 6b Outliers and leverage

We want a measure for leverage that is based on

① the distance x_i is away from the \bar{x} .

② the extent to which \hat{y}_i depends on y_i

Recall

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

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum_{j=1}^n \frac{x_j - \bar{x}}{SXX} y_j}{\sum_{j=1}^n \frac{x_j - \bar{x}}{SXX}}$$

centered

$$\hat{y}_i = \bar{y} - \hat{\beta}_1 \bar{x} + \hat{\beta}_1 x_i$$

all the data

$$= \bar{y} + \hat{\beta}_1 (x_i - \bar{x})$$

$$= \frac{1}{n} \sum_{j=1}^n y_j + \sum_{j=1}^n \frac{x_j - \bar{x}}{SXX} y_j (x_i - \bar{x})$$

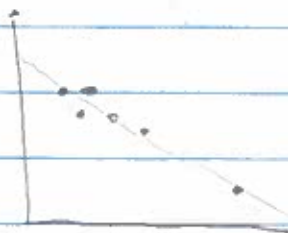
isolate the y

$$= \sum_{j=1}^n \left[\frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{SXX} \right] y_j$$

$$= \sum_{j=1}^n h_{ij} y_j$$

* predicted \hat{y}_i is a weighted sum of all the y 's, with $\sum_{j=1}^n h_{ij} = 1$

$$\hat{y}_i = h_{i1} y_1 + h_{i2} y_2 + \dots + h_{ii} y_i + \dots + h_{in} y_n$$

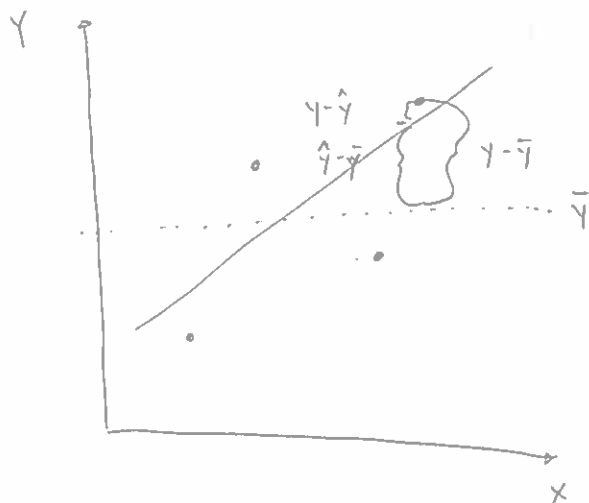


A measure of model fit: R^2 + Adjusted R^2

Recall R ; Pearson's Corr. Coef: measures strength of linear relationship. $R \in [-1, 1]$

R^2 : the proportion of the total variability in the y 's explained by the regression model

$$R^2 = \frac{SS_{reg}}{SST} = 1 - \frac{RSS}{SST}$$



Total Variation in the y

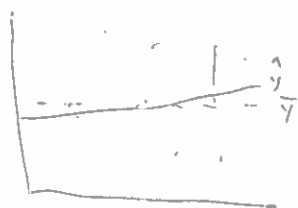
$$SST = \sum_{i=1}^n (y_i - \bar{y})^2$$

SS_{reg} : Total Variation explained by \hat{y}
 $= \sum (\hat{y}_i - \bar{y})^2$

RSS : Total unexplained variation
 $= \sum (y_i - \hat{y}_i)^2$

$$SST = SS_{reg} + RSS$$

No linear trend



RSS : big, SS_{reg} : small $\Rightarrow R^2$ near 0

Strong linear trend



RSS : small, SS_{reg} : big $\Rightarrow R^2$ near 1

Cautions: Adding complexity (i.e. additional terms) to our model will always increase R^2 , even if the term was bogus

Solution: Use adjusted R^2

$$R_{adj}^2 = 1 - \frac{RSS(n-p-1)}{SST(n-1)}$$

n : sample size
 p : # predictors

Now, new terms have to decrease RSS by more than the penalty to increase R_{adj}^2 .

Model Fitting a Metaphor

- △ rice
- quinoa
- salt



← Data
- signal (△, □)
- noise (•)

linear
quadratic



Sieve
(rice stopped)



← Fitted ± 1
Mean, function captures deterministic portion (signal)
 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$



(quinoa + salt
pass through)



← residuals
(noise + some signal)
aka residual plot
w/ structure

Key points

- Your fitted model is a way to separate signal from noise.
- A good model will leave behind residuals w/ no any structure.
- Be wary of building a model that considers everything to be signal.
This is overfitting and will damage prediction

Ockham's Razor aka "Law of Parsimony"

"Plurality should not be posited w/o necessity",



complexity

aka: "Among similar explanations, the simpler is better"
→ useful heuristic in model building.

Wednesday Oct. 22

- Reading: S.3, S.2, Linear Algebra.
- Proposal due Wednesday.
- Survey
- Analysis of Covariance

m1 Questions

- Write out the eqn for the line.
- Is volume a significant predictor?
- How much of the variation in weight is explained by the model containing volume?

ANCOVA Interpretation

- Simple Linear Regression: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$

- Parallel Lines $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$

x_1 : continuous (volume)

x_2 : categorical (cover)
1 if pb
0 if hb

For paperbacks: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$ \nwarrow is 1

$$\hat{y} = (\hat{\beta}_0 + \hat{\beta}_2) + \hat{\beta}_1 x_1$$

new intercept

For Handbooks:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1$$

is 0

same slope
different intercept
 \therefore parallel

- Two intercepts, two slopes

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \hat{\beta}_3 x_1 x_2$$

Pb: $\hat{y} = (\hat{\beta}_0 + \hat{\beta}_2) + (\hat{\beta}_1 + \hat{\beta}_3) x_1$

hb: $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1$

October 21st

Matrix Multiplication - a review

Ex 1

$$AB = \begin{bmatrix} 2 & 5 & -1 \\ -4 & 3 & -3 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 1 & 1 \\ 0 & 4 \end{bmatrix} = \begin{bmatrix} 2 \cdot 2 + 5 \cdot 1 + (-1) \cdot 0 & 2 \cdot 0 + 5 \cdot 1 + (-1) \cdot 4 \\ (-4) \cdot 2 + 3 \cdot 1 + (-3) \cdot 0 & (-4) \cdot 0 + 3 \cdot 1 + (-3) \cdot 4 \end{bmatrix} = \begin{bmatrix} 9 & 1 \\ -5 & -9 \end{bmatrix}$$

$\begin{matrix} 2 \times 3 & 3 \times 2 \\ \hline & \text{OK} \end{matrix}$

- You can only multiply matrices w/ the same inner dimension.
- The resulting matrix takes its dimensions from the outer dimensions.
- Each entry in the resulting matrix is a dot product of the corresponding row of the 1st matrix and the column of the 2nd.

Ex 2 What is BA^T ?

$$\begin{bmatrix} 2 & 0 \\ 1 & 1 \\ 0 & 4 \end{bmatrix} \begin{bmatrix} 2 & 5 & -1 \\ -4 & 3 & 3 \end{bmatrix} = \begin{bmatrix} 4 & 10 & -2 \\ -2 & 8 & 2 \\ -16 & 12 & 12 \end{bmatrix}$$

$\begin{matrix} 3 \times 2 & 2 \times 3 \\ \hline & \text{OK} \end{matrix}$

Def: Transpose: $A' = \begin{bmatrix} 2 & 4 \\ 5 & 3 \\ -1 & -3 \end{bmatrix}$ The transpose of a matrix is a matrix of flipped dim where the row, col indices have been flipped

Ex 3 What is $B'A'$?

$$\begin{bmatrix} 2 & 1 & 0 \\ 0 & 1 & 4 \end{bmatrix} \begin{bmatrix} 2 & 4 \\ 5 & 3 \\ -1 & -3 \end{bmatrix} = \begin{bmatrix} 9 & -5 \\ 1 & 9 \end{bmatrix}$$

$$\neq B'A' = (AB)'$$

Def Symmetric: a matrix is symmetric if $a_{ij} = a_{ji}$ (must be square)

$$S = \begin{bmatrix} 1 & 3 & 6 \\ 3 & 2 & 1 \\ 6 & 1 & 0 \end{bmatrix} = S'$$

Def Diagonal: a matrix is diagonal if $a_{ij} = 0$ for $i \neq j$.

$$D = \begin{bmatrix} 3 & 0 & 0 \\ 0 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

Def. Identity matrix: a diagonal matrix with $a_{ij} = 1$ for $i=j$

$$I_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Recall from algebra

$$a \cdot 1 = a \quad * 1 \text{ is the identity operator}$$

$$A I = \begin{bmatrix} 2 & 5 & -1 \\ 4 & 3 & -3 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 2 & 5 & -1 \\ 4 & 3 & -3 \end{bmatrix} = A$$

$2 \times 3 \quad 3 \times 3$

$$a \cdot \frac{1}{a} = 1$$

Def. Inverse the inverse of matrix yields I when multiplied.

$$A A^{-1} = I \quad * A \text{ must be square!} *$$

* Not all matrices can be inverted

$$M = \begin{bmatrix} 1 & 3 & -2 \\ 2 & 5 & -3 \\ -3 & 2 & 4 \end{bmatrix} \begin{bmatrix} 14 & -8 & -1 \\ -17 & 10 & 1 \\ 19 & 11 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Matrix Notation for Regression

$$Y = \begin{bmatrix} 1 \\ \vdots \\ y_i \\ \vdots \\ y_n \end{bmatrix} \quad n \times 1$$

$$X = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_p \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_1 & x_2 & x_3 & \dots & x_p \end{bmatrix} \quad n \times (p+1)$$

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \vdots \\ \hat{\beta}_p \end{bmatrix} \quad (p+1) \times 1$$

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad n \times 1$$

To predict the fitted value at x_i

Summation Notation

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i,1} + \hat{\beta}_2 x_{i,2} + \dots + \hat{\beta}_p x_{i,p}$$

Matrix Notation

$$= x_{i, \cdot} \hat{\beta} = \begin{bmatrix} \hat{\beta}_0 \cdot 1 + \hat{\beta}_1 \cdot x_{i,1} + \hat{\beta}_2 \cdot x_{i,2} + \dots + \hat{\beta}_p \cdot x_{i,p} \end{bmatrix} = \hat{y}_i$$

$(1 \times (p+1)) \times 1$
 $\quad \quad \quad$
 1×1

True Mean Function

$$E(Y|X=x) = XB$$

Data Gen Fxn

$$Y = XB + e$$

Fitted Mean Fxn

$$\hat{Y} = X \hat{\beta}$$

$$\begin{matrix} \uparrow & \uparrow \\ n \leq \text{fit} & n \leq \text{test} \end{matrix}$$

Wednesday October 29th

$$Y ; X ; \beta ; e ; \hat{e} = Y - X\hat{\beta}$$

$n \times 1 \quad n \times (p+1) \quad (p+1) \times 1 \quad n \times 1 \quad n \times 1$

Estimating $\hat{\beta}$

We want to pick $\hat{\beta}$ to minimize R^2 .

$$\begin{aligned} \text{RSS}(\hat{\beta}) &= \hat{e}'\hat{e} = (Y - X\hat{\beta})'(Y - X\hat{\beta}) \\ &= Y'Y + (X\hat{\beta})'X\hat{\beta} - Y'X\hat{\beta} - (X\hat{\beta})'Y \\ &= Y'Y + \hat{\beta}'(X'X)\hat{\beta} - 2Y'X\hat{\beta} \end{aligned}$$

Take derivative wrt $\hat{\beta}$

$$* (AB') = B'A'$$

$$* B'A = A'B \text{ when result is } 1 \times 1$$

$$= 0 + 2(X'X)\hat{\beta} - 2X'Y$$

$$(X'X)\hat{\beta} = X'Y$$

$$\boxed{\hat{\beta} = (X'X)^{-1}X'Y}$$

SLR

$$X = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ 1 & x_{31} \\ 1 & x_{41} \end{bmatrix}$$

See
page 133

$$\hat{\beta} = (X'X)^{-1}X'Y$$

$$= \begin{bmatrix} \bar{y} - \frac{S_{XY}}{S_{XX}}\bar{x} \\ \frac{S_{XY}}{S_{XX}} \end{bmatrix}$$

Monday 11/3/14

Properties of Least Squares Estimates

$$Y = XB + e, \text{Var}(e) = \sigma^2 I_{n \times n}$$

$$\hat{\beta} = (X'X)^{-1} X'Y$$

$$\text{Var}(e_1) \quad \text{Var}(e_2) \quad \text{Cov}(e_1, e_2)$$

$$\begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

Are our estimates unbiased?

i.e., do they reach the true β in expectation?

$$\begin{aligned} E(cY) &= c E(Y) \\ E(e) &= 0 \end{aligned}$$

$$E(\hat{\beta}|X) = E[(X'X)^{-1} X'Y]$$

$$= (X'X)^{-1} X' E(Y|X) \rightarrow E(Y|X) = E(X\beta) + E(e|X)$$

$$= (X'X)^{-1} X' X\beta = X\beta + 0$$

$$= \beta$$

✓ the least-squared estimates are unbiased

What is their variance?

$$\text{Var}(\hat{\beta}|X) = \text{Var}[(X'X)^{-1} X'Y]$$

$$= (X'X)^{-1} X' \text{Var}(Y|X) X (X'X)^{-1}$$

$$= (X'X)^{-1} X' \sigma^2 I X (X'X)^{-1}$$

$$= \sigma^2 (X'X)^{-1} X' X (X'X)^{-1}$$

$$= \sigma^2 (X'X)^{-1}$$

estimate w/ $S^2 = \frac{RSS}{n-p-1} = \frac{1}{n-p-1} e' e$

$$\text{Var}(aY) = a^2 \text{Var}(Y)$$

as matrix

$$\text{Var}(AX) = A \text{Var}(X) A'$$

$$(AB)' = B' A'$$

Gauss-Markov Theorem

IF $\text{Var}(e) = \sigma^2 I$

the least squares estimates are

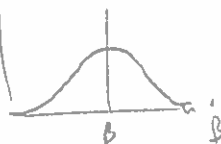
Best

Linear

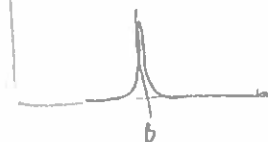
Unbiased

Estimates

method 1



method 2



Wednesday November 5th

- Special seminar + lunch
- Proposals due midnight
- BMI Paper
- Diagnostics I

Leverage

def: the extent to which \hat{y}_i is attracted by y_i .

$$\hat{Y} = X\hat{\beta} = \underbrace{X(X'X)^{-1}X'}_H Y = HY$$

$X(X'X)^{-1}X'$ is known as the "hat matrix" since it puts a hat on Y .

$$H = \begin{bmatrix} h_{11} & & \\ & \ddots & \\ h_{12} & & \ddots \end{bmatrix} \quad \text{--- } n \times n$$

$$\hat{y}_i = H_{i \cdot} \cdot \hat{Y} = h_{i1}y_1 + h_{i2}y_2 + \dots + \underbrace{h_{ii}}_{\text{row } i \text{ column } i} y_i + \dots + h_{in}y_n$$

the effect of y_i on \hat{y}_i

The leverages are the diagonal elements of the hat matrix.

* Rule of thumb: an obs has high leverage if its greater than twice the avg. lev.

$$h_{ii} > 2 \frac{(p+1)}{n}$$

- average leverage

Properties of Residuals

$$\hat{e} = Y - \hat{Y} = Y - HY = (I - H)Y$$

$$\begin{aligned} & \bullet \text{Var}(CY) = C \text{Var}(Y) C' \\ & \bullet H \text{ is idempotent: } HH = X(X'X)^{-1}X'X(X'X)^{-1}X' = X(X'X)^{-1}X' = H \end{aligned}$$

Expected value

$$\begin{aligned} E(\hat{e} | X) &= E((I - H)Y) \\ &= (I - H)E(Y) \\ &= (I - H)XB \\ &= XB - X(X'X)^{-1}X'XB \\ &= XB - XB \\ &= 0 \end{aligned}$$

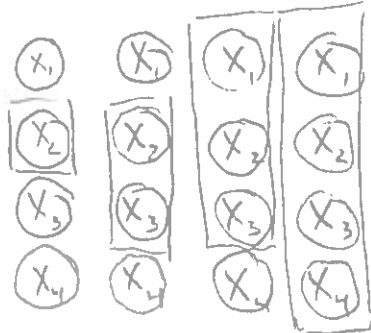
Variance

$$\begin{aligned} \text{Var}(\hat{e} | X) &= \text{Var}((I - H)Y) \\ &= (I - H) \text{Var}(Y) (I - H)' \\ &= (I - H) \sigma^2 I (I - H)' \\ &= \sigma^2 (I - H)(I - H)' \\ &= \sigma^2 (II' - IH' - HI' - HH') \\ &= \sigma^2 (I - H - H + H) \\ &= \sigma^2 (I - H) \end{aligned}$$

Wednesday Nov 19th

- Hack Ebola
- Presentation Schedule
- Exam parameters
- Searching the Model Space

Best Subsets

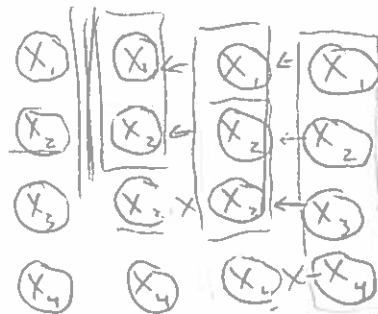


1 2 3 4

Predictors

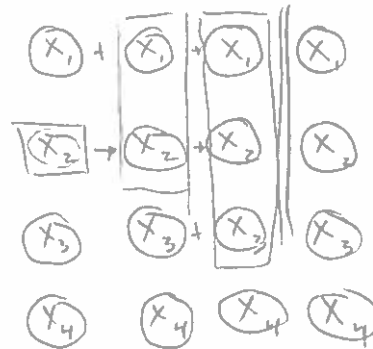
- 1) Find the model of each size that maxes R^2_{adj} .
- 2) Compare those best models by all criteria

Backwards Elimination



- 1) Fit the full model and calculate criteria of choice.
- 2) Remove predictor w/ largest p-val and calculate new criterion.
- 3) Stop iteration when the criterion is no longer improving.

Forward Selection



- 1) Fit the simple model w/ the predictor w/ the smallest p-val. Compute criterion.
- 2) Fit new model w/ that predictor + new predictor, compute criterion.
- 3) Stop iteration w/ criterion is no longer improving.