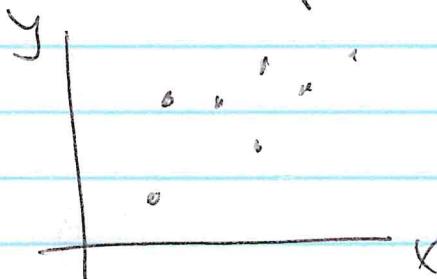


# Ch 11 Regression + Least-Squares

11pt

Inds: best-fit lines (two-ways)

Consider a sample (data)  $S = \{(x_i, y_i)\}$



with some underlying joint  
pd f  $f(x, y)$ .

Assumptions:

- (1)  $Y$  is dependent upon  $X$  (not incl.)
- (2) in general,  $y = mx + b$   
is the relationship btwn  $X$  and  $Y$ .

A) a probabilistic construction.

Given all possible lines  $y = mx + b$  that could describe the data, the most likely line will intersect the point  $(\bar{x}, \bar{y})$ .

(This is the fundamental idea of "regression analysis".)

$$\text{That is } y - \bar{y} = b(x - \bar{x})$$

To find the "best" line, we choose to minimize the vertical distance (error in the y-direction) btwn the points and the line.

1(p2)

For a point  $(x_k, y_k)$  in  $S$ , the vertical distance is

$$\text{dist} = | y_k - (b(x_k - \mu_x) + \mu_y) |$$

actual                    predicted by model.

Minimizing will be a calculus problem.

Old trick from Calc I, min  $\text{dist}^2$  instead,

$$\begin{aligned}\text{dist}^2 &= (y_k - (b(x_k - \mu_x) + \mu_y))^2 \\ &= (y_k - \mu_y - b(x_k - \mu_x))^2\end{aligned}$$

As  $X, Y$  distributed via a PDF, best way to compute this problem is by minimizing the expectation of the parameter  $b$ .

$$\begin{aligned}\text{Define } L(b) &= E(\text{dist}^2(b)) \\ &= E[(y_k - \mu_y - b(x_k - \mu_x))^2]\end{aligned}$$

The problem of finding the "b" that minimizes expectation is called a least-squares problem.

$$\begin{aligned}L(b) &= E[(y_k - \mu_y)^2 - 2b(y_k - \mu_y)(x_k - \mu_x) + b^2(x_k - \mu_x)^2] \\ &= E[(y_k - \mu_y)^2] - 2bE[(x_k - \mu_x)(y_k - \mu_y)] + b^2E[(x_k - \mu_x)^2] \\ &= S_y^2 - 2b \text{Cov}(X, Y) + b^2 S_x^2.\end{aligned}$$

11/3

$$\text{Minimize} \Rightarrow \frac{d}{db}$$

$$L'(b) = -2\text{Cov}(X, Y) + 2b S_x^2$$

$$L'(b) = 0 \quad \text{then} \quad b = \frac{\text{Cov}(X, Y)}{S_x^2}$$

Note that this is a minimum as  $L''(b) = 2S_x^2 > 0$ .

In AP Stats, the slope is moved to a "prettier" form:

$$b = \frac{\text{Cov}(X, Y)}{S_x S_y} \cdot \frac{S_y}{S_x} = \rho \frac{S_y}{S_x}$$

where  $\rho$  is the correlation coef.  $\frac{\text{Cov}(X, Y)}{S_x S_y}$

In last semester,

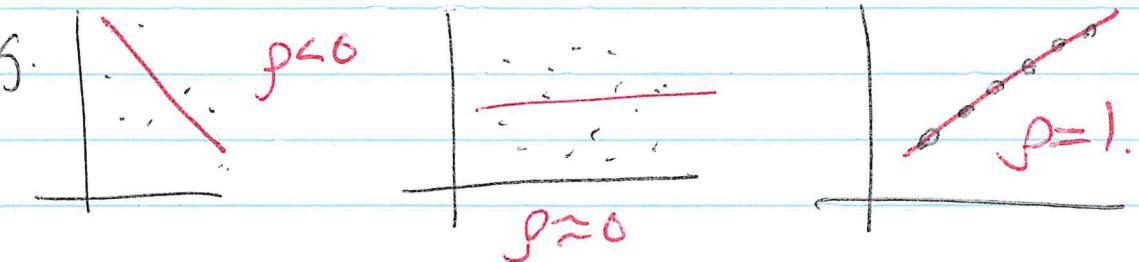
FACT #1: The sign of the slope depends entirely upon  $\rho$ , whose sign is dependent upon  $\text{Cov}(X, Y)$

FACT #2:  $-1 \leq \rho \leq 1$ .

We proved this last semester in HU. Here, much easier to prove simply by evaluating  $L(b)$  at our critical point.

FACT #3: Closer  $|p|$  is to 1, the better the line  $y = mx + b$  "fits" the data.

e.g.



The AP Stats defn of the best fit line

$$\text{is } \hat{y} = \hat{b}_0 + \hat{b}_1 x$$

$$\text{where } \hat{b}_1 = p \frac{\sum x}{\sum x^2} \text{ and } \hat{b}_0 = \bar{y} - \hat{b}_1 \bar{x}$$

B) the linear algebra construction.

If want  $y = mx + b$ .

Using  $S = \{(x_i, y_i)\}$  yields the over-determined  
and always inconsistent system

$$\begin{aligned} y_1 &= mx_1 + b \\ y_2 &= mx_2 + b \\ &\vdots \\ y_n &= mx_n + b \end{aligned} \quad \begin{matrix} \text{as a} \\ \text{vector} \\ \text{eqn} \end{matrix} \quad \Rightarrow$$

$$\vec{y} = M \vec{x} + b \vec{1}$$

where

$$\vec{y} = \langle y_1, \dots, y_n \rangle, \vec{x} = \langle x_1, \dots, x_n \rangle \text{ and } \vec{1} = \langle 1, \dots, 1 \rangle.$$

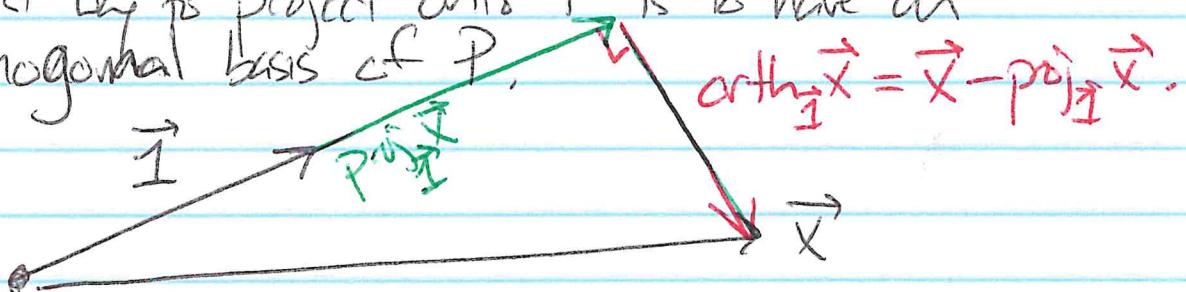
11/25

Consider the plane  $P = \text{span}\{\vec{x}, \vec{1}\}$ .

Inconsistent  $\rightarrow \vec{y} \notin P$

But the vector in  $P$  that is closest to  $\vec{y}$  is  $\text{proj}_P \vec{y}$ . (closest under Euclidean distance  
 $\text{dist}^2 = \sum (y_i - r_i)^2 \dots$  i.e. least squares)

Easiest way to project onto  $P$  is to have an orthonormal basis of  $P$ .



$$P = \text{span}\{\vec{1}, \text{orth}_{\vec{1}} \vec{x}\}$$

$$\text{where } \vec{v} = \vec{x} - \frac{\vec{x} \cdot \vec{1}}{\vec{1} \cdot \vec{1}} \vec{1} = \vec{x} - \frac{\vec{x} \cdot \vec{1}}{n} \vec{1}.$$

$$\begin{aligned} \text{Then } \text{proj}_P \vec{y} &= \frac{\vec{y} \cdot \vec{1}}{\vec{1} \cdot \vec{1}} \vec{1} + \frac{\vec{y} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \vec{v} \\ &= \frac{\vec{y} \cdot \vec{1}}{n} \vec{1} + \frac{\vec{y} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \vec{x} - \frac{\vec{y} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \left( \frac{\vec{x} \cdot \vec{1}}{n} \right) \vec{1} \\ &= \underbrace{\frac{\vec{y} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \vec{x}}_{\text{this is slope } m} + \underbrace{\left( \frac{\vec{y} \cdot \vec{1}}{n} - \frac{\vec{y} \cdot \vec{v}}{\vec{v} \cdot \vec{v}} \cdot \frac{\vec{x} \cdot \vec{1}}{n} \right) \vec{1}}_{\text{intercept } b}. \end{aligned}$$

11/26.

$$\begin{aligned}
 \vec{V} \cdot \vec{V} &= \left( \vec{X} - \frac{\vec{X} \cdot \vec{I}}{n} \vec{I} \right) \cdot \left( \vec{X} - \frac{\vec{X} \cdot \vec{I}}{n} \vec{I} \right) \\
 &= \vec{X} \cdot \vec{X} - 2 \frac{(\vec{X} \cdot \vec{I})^2}{n} + \frac{(\vec{X} \cdot \vec{I})^2}{n} \vec{I} \cdot \vec{I} \\
 &= \vec{X} \cdot \vec{X} - \frac{2(\vec{X} \cdot \vec{I})^2}{n} + \frac{(\vec{X} \cdot \vec{I})^2}{n} \\
 &= \frac{n(\vec{X} \cdot \vec{X}) - (\vec{X} \cdot \vec{I})^2}{n}
 \end{aligned}$$

$$\begin{aligned}
 \vec{V} \cdot \vec{Y} &= \vec{X} \cdot \vec{Y} - \frac{(\vec{X} \cdot \vec{I}) (\vec{Y} \cdot \vec{I})}{n} \\
 &= \frac{n(\vec{X} \cdot \vec{Y}) - (\vec{X} \cdot \vec{I})(\vec{Y} \cdot \vec{I})}{n}
 \end{aligned}$$

$$m = \frac{\vec{V} \cdot \vec{Y}}{\vec{V} \cdot \vec{V}} = \frac{n(\vec{X} \cdot \vec{Y}) - (\vec{X} \cdot \vec{I})(\vec{Y} \cdot \vec{I})}{n(\vec{X} \cdot \vec{X}) - (\vec{X} \cdot \vec{I})^2}$$

$$\begin{aligned}
 \text{Also note } b &= \frac{\vec{Y} \cdot \vec{I}}{n} - m \cdot \frac{\vec{X} \cdot \vec{I}}{n} \\
 &= \vec{Y} - m \vec{X}
 \end{aligned}$$

Claim: These are the set of formulas.

11p7

$$m = \frac{\text{Cov}(X, Y)}{S_x^2} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2}$$

Useful formula:  $\sum (x_i - \bar{x})(y_i - \bar{y})$

$$\begin{aligned} &= \sum (x_i y_i - \bar{x}\bar{y} - \bar{x}y_i + \bar{x}\bar{y}) \\ &= \sum x_i y_i - \bar{y} \sum x_i - \bar{x} \sum y_i + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y} - n \bar{x} \bar{y} + n \bar{x} \bar{y} \\ &= \sum x_i y_i - n \bar{x} \bar{y}, \end{aligned}$$

$$\begin{aligned} \text{So } m &= \frac{\text{Cov}(X, Y)}{S_x^2} = \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \\ &= \bar{x} \cdot \bar{y} - n \frac{\bar{x} \cdot \mathbf{I}}{n} \cdot \frac{\bar{y} \cdot \mathbf{I}}{n} \\ &\quad \overline{\bar{x} \cdot \bar{x} - n \left( \frac{\bar{x} \cdot \mathbf{I}}{n} \right)^2} \\ &= \frac{n \bar{x} \cdot \bar{y} - (\bar{x} \cdot \mathbf{I})(\bar{y} \cdot \mathbf{I})}{n \bar{x} \cdot \bar{x} - (\bar{x} \cdot \mathbf{I})^2}. \quad \text{see this!} \end{aligned}$$

11p8

## § 11.4 - the coefficient as estimators.

We have shown the best-fit line to be

where  $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x$   
 $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$ ,  $S_{xy} = \sum (x_i - \bar{x})(y_i - \bar{y})$   
 $S_{xx} = \sum (x_i - \bar{x})^2$ .  
and  $\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$ .

but we recognize  $\hat{\beta}_0$  and  $\hat{\beta}_1$  as stats.  
dependent upon  $\bar{x}$ ,  $S_x$  and  $\bar{y}$ . But what  
exactly are they estimating (Ch 9)?

In theory  $\epsilon$   $x$   $\times$   $y$  distributed via  $f_{xy}$   
and there is some "best" linear over the entire  
probability space.

$$Y = \beta_0 + \beta_1 x$$

That is, when  $\omega x$ ,  $E[Y] = \beta_0 + \beta_1 x$ .

A modern approach to the distribution of  $Y \omega x$   
is to add an error parameter  $\epsilon$

That is  $Y = \beta_0 + \beta_1 x + \epsilon$   
deterministic component the random argument  
of  $Y$

11p9.

Still want  $E(Y) = \beta_0 + \beta_1 x$ , hence  $E(\epsilon) = 0$  is required.

We make the additional assumption that the variance of  $\epsilon$  is independent of  $x$ .

That is  $V(Y) = V(\epsilon) = \sigma^2$

i.e.  $y$  depends on  $x$ , but the spread of  $Y$  does not

Prop:  $\hat{\beta}_0, \hat{\beta}_1$  are unbiased estimators of  $\beta_0, \beta_1$   
where  $Y = \beta_0 + \beta_1 x + \epsilon$ .

$$\begin{aligned}
 \text{Reason: } E(\hat{\beta}_1) &= E\left(\frac{S_{xy}}{S_{xx}}\right) = E\left(\frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{S_{xx}}\right) \\
 &= E\left(\frac{\sum(x_i - \bar{x})y_i - \bar{y}\sum(x_i - \bar{x})}{S_{xx}}\right) \xrightarrow{\text{defn.}} = 0 \\
 &= \frac{\sum(x_i - \bar{x})E(y_i)}{S_{xx}} \\
 &= \frac{\sum(x_i - \bar{x})(\beta_0 + \beta_1 x_i)}{S_{xx}} \\
 &= \frac{\beta_0 \sum(x_i - \bar{x})}{S_{xx}} + \beta_1 \frac{\sum(x_i - \bar{x})x_i}{S_{xx}} \\
 &\quad \star \text{ last day } \sum(x_i - \bar{x})x_i = \sum x_i^2 - n\bar{x}^2 = S_{xx}.
 \end{aligned}$$

$\therefore E(\hat{\beta}_1) = \beta_1$ . Unbiased

11/20

$$\begin{aligned} \text{For } E(\hat{\beta}_0) &= E(\bar{Y} - \hat{\beta}_1 \bar{x}) \\ &= E[\bar{Y}] - \bar{x} E[\hat{\beta}_1] \\ &= E[\bar{Y}] - \beta_1 \bar{x}. \end{aligned}$$

$$\text{But } \bar{Y} = \frac{1}{n} \sum Y_i = \frac{1}{n} \sum (\beta_0 + \beta_1 x_i + \epsilon) = \beta_0 + \beta_1 \bar{x} + \bar{\epsilon}$$

$$\text{and } E[\bar{Y}] = \beta_0 + \beta_1 \bar{x}.$$

$$\text{and } E(\hat{\beta}_0) = \beta_0 + \beta_1 \bar{x} - \beta_1 \bar{x} = \beta_0 \quad \blacksquare.$$

$$\text{Cov: } V(\hat{\beta}_1) = \frac{\sigma^2}{S_{xx}}, \quad V(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n S_{xx}}.$$

$$\begin{aligned} \text{Reason: } V(\hat{\beta}_1) &= V\left[\frac{\sum (x_i - \bar{x}) Y_i}{S_{xx}}\right] = \frac{1}{S_{xx}^2} V\left[\sum (x_i - \bar{x}) Y_i\right] \\ &= \frac{\sum (x_i - \bar{x})^2 V(Y_i)}{S_{xx}^2} = \frac{\sigma^2 \sum (x_i - \bar{x})^2}{S_{xx}^2} = \frac{\sigma^2}{S_{xx}}. \end{aligned}$$

$$\begin{aligned} V(\hat{\beta}_0) &= V(\bar{Y} - \hat{\beta}_1 \bar{x}) \\ &\text{note: in our } \epsilon \text{ set-up, } \bar{Y} \text{ and } \hat{\beta}_1 \text{ are independent.} \\ &= V(\bar{Y}) + V(\hat{\beta}_1 \bar{x}) - 2 \text{Cov}(\bar{Y}, \hat{\beta}_1 \bar{x}) \\ &= V(\bar{Y}) + \bar{x} V(\hat{\beta}_1) - 2 \bar{x} \text{Cov}(\bar{Y}, \hat{\beta}_1) \end{aligned}$$

$$\text{i) } V(\bar{Y}) = V(\bar{\epsilon}) = \frac{1}{n} V(\epsilon) = \frac{\sigma^2}{n}$$

$$\text{ii) } \text{Cov}(\bar{Y}, \hat{\beta}_1) = \text{Cov}\left(\frac{1}{n} \sum Y_i, \frac{\sum (x_i - \bar{x}) Y_i}{S_{xx}}\right)$$

$$\begin{aligned}
 &= \sum_{i=1}^n \frac{(x_i - \bar{x})}{\sqrt{S_{xx}}} \sqrt{V(Y_i)} + \sum_{i \neq j} \frac{(x_i - \bar{x})}{\sqrt{S_{xx}}} \text{Cov}(Y_i, Y_j) \\
 &= \frac{\sigma^2}{\sqrt{S_{xx}}} \sum_{i=1}^n (x_i - \bar{x}) \xrightarrow{\text{red}} = 0 \\
 &= 0.
 \end{aligned}$$

$$\begin{aligned}
 S_{\beta_0} V(\hat{\beta}_0) &= \frac{\sigma^2}{n} + \frac{\bar{x} \sigma^2}{S_{xx}} + 0 \\
 &= \frac{\sigma^2 (S_{xx} + n \bar{x}^2)}{n S_{xx}} \\
 &= \frac{\sigma^2 \sum_{i=1}^n x_i^2}{n S_{xx}} \quad \text{since } S_{xx} = \sum x_i^2 - n \bar{x}^2. \quad \square
 \end{aligned}$$

$$\text{Cov: } \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\frac{\bar{x} \sigma^2}{S_{xx}}$$

Reason: Exercise

Note:  $\hat{\beta}_0, \hat{\beta}_1$  are dependent unless  $\bar{x} = 0$ .

~~Exercise~~

topic  
else: estimating  $\sigma^2$ .

11/12.

We have been working w/ assumption  $V(Y) = V(\epsilon) = \sigma^2$ ,  
but this is also usually unknown.

In the past, we used  $V(Y) = \frac{1}{n-1} \sum_{i=1}^{n-1} (Y_i - \bar{Y})^2$ .

However, in this setting, our estimator of  $E[Y_i]$  is  
 $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$  and not  $\bar{Y}$ .

We want to define  $S^2$  via  $SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$   
 $S^2 := K \sum_{i=1}^n (Y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2$

where  $K$  is chosen to make  $S^2$  an unbiased estimator of  $\sigma^2$ .

(I don't want to do this again pg 580-581)

As we did in Ch8, we can show

$$E[SSE] = (n-d) \sigma^2$$

So we will use

$$\hat{S}^2 := \frac{1}{n-d} SSE = \frac{1}{n-d} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Prop:  $E[\hat{S}^2] = \sigma^2$ .

Cor: (The computation cor.)  $SSE = S_{yy} - \hat{\beta}_1 S_{xy}$   
where  $S_{yy} = \sum_{i=1}^n (Y_i - \bar{Y})^2$ .

1/p13.

ex: (11.16)

potency of antibiotic after storage (Y)	38, 43, 29	32, 26, 33	19, 27, 27	14, 19, 21
temp stored at (X)°F	30°	50°	70°	90°

note n = 12

$$\bar{Y} = 27 \quad , \quad \bar{x} = 60$$

$$S_{xy} = -1900 = \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$S_{xx} = 6000 = \sum (x_i - \bar{x})^2$$

$$S_{yy} = 792 = \sum (y_i - \bar{y})^2.$$

$$\text{then } \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = -\frac{19}{60}$$

$$\text{and } \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{x} = 27 + \frac{19}{60} \cdot 60 = 46$$

So the best-fit line is  $\hat{Y} = 46 - \frac{19}{60} x$ .

$$\begin{aligned} \text{For } S^2 &= \frac{1}{n-2} SSE = \frac{1}{n-2} (S_{yy} - \hat{\beta}_1 S_{xy}) \\ &= \frac{1}{10} [792 - \left(-\frac{19}{60}\right)(-1900)] \\ &= 571/30 \approx 19.03 \end{aligned}$$

11p24,

disc: do we know how  $S^2$  is distributed?

This depends on our assumptions on how  $\epsilon$  is distributed.

We have  $Y = \beta_0 + \beta_1 X + \epsilon$  w/  $E(\epsilon) = 0$  and  $V(\epsilon) = \sigma^2$ , and these are independent of  $X$ .

Note the dist'sns of the point estimators  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are  $\sigma^2$  dependent (which we can estimate w/  $S^2$ ).

However, it is a common assumption that things are normally distributed.

If we also assume  $\epsilon \sim N(0, \sigma^2)$

then we can prove a Fisher's Thm - like result for  $S^2$ .

Namely,

$$\text{Thm: } \frac{(n-2)S^2}{\sigma^2} = \frac{\text{SSE}}{\sigma^2} \sim \chi^2_{n-2}$$

and  $S^2$  is independent of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

Moreover,  $\hat{\beta}_0$  and  $\hat{\beta}_1$  are also normally distributed.

Pf: beyond our scope.

## § 11.5 Inferences concerning the point estimators. 11pt

Under the assumption  $\epsilon \sim N(0, \sigma^2)$ , we have  
 $\hat{\beta}_0, \hat{\beta}_1$  are Normal  $N(\beta_0, V(\beta_0))$

Thus, we can do confidence intervals for the true  $\beta_0, \beta_1$  ( $Y = \beta_0 + \beta_1 X + \epsilon$ ) and hypothesis tests.

Use the same Z / T rules as before.

• If  $n \leq 30$ , use t-dist'n

•  $n > 30$ , use Z.

Ex: Given  $Z_{\alpha/2}$ ,  $\hat{\beta}_1 \pm Z_{\alpha/2} \sqrt{V(\hat{\beta}_1)}$   
 $n > 30$

is an  $\alpha$ -level confidence interval for the true  $\beta_1$ .

Ex: Find a 90% C.I. for the slope  $\hat{\beta}_1$  of the paternity example.

We have  $\hat{\beta}_1 = -\frac{19}{60} \approx -0.31667$

$$V(\hat{\beta}_1) = \frac{s^2}{S_{xx}} \approx \frac{s^2}{S_{xx}} \approx \frac{19.3}{6000} = 0.0317\dots$$

$$\Rightarrow S = 0.0563\dots$$

11/16.

\* Because we are using  $S^2$ , we need to use  $S^2$  df. \*

$$n=12$$

$$S^2 \sim \chi^2(12-2) = \chi^2(10)$$

$$\text{Need } t_{0.05}(10) = 1.812.$$

$$\text{And } \hat{\beta}_1 \pm t_{0.05}(10) S$$

$$\begin{aligned} S &= -0.31667 \pm (0.0563)(1.812) \\ &= -0.31667 \pm 0.102856 \\ \text{or } &(-0.4187, -0.214644) \end{aligned}$$

ex: (potency example again)

Let us assume that it is commonly known that penicillin derivatives are considered "good" if when stored in a deep freezer ( $0^\circ\text{F}$ ) the potency is  $50$  (or better).

We wish to test if our drug is considered "good".

$$\text{We have } \hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X = 46 - \frac{19}{60}X.$$

$$\text{Then } \hat{Y}(0) = \hat{\beta}_0 = 46$$

We want to test if the true  $\beta_0 = 50$ .  
Given our sample observation  $\hat{\beta}_0 = 46$ .

H<sub>0</sub>

$$H_0: \beta_0 = 50$$

$$H_a: \beta_0 < 50.$$

then  $S^2 \sim \chi^2(10)$  ... need to use t-test

$$\text{p-value } p = P(\hat{\beta}_0 \leq 46 \mid \beta_0 = 50)$$

$$\text{we assume } \hat{\beta}_0 \sim N(\beta_0, V(\hat{\beta}_0))$$

$$P = P\left(\frac{\hat{\beta}_0 - \beta_0}{\sqrt{V(\hat{\beta}_0)}} \leq \frac{46 - 50}{\sqrt{V(\hat{\beta}_0)}}\right)$$

didn't compute  $V(\hat{\beta}_0)$  earlier.

$$\text{Recall } V(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n S_{xx}}$$

Have  $\sigma^2 \approx S^2 \approx 19.03$ ,  $n=10$ ,  $S_{xx}=600$

$$\sum x_i^2 = 30^2 + 28^2 + \dots + 80^2 = \cancel{9740} \quad 16400 \cdot 3 \\ = 49200$$

$$\therefore V(\hat{\beta}_0) = \frac{19.03 \cdot \cancel{9740}}{10 \cdot 600} = \cancel{2.274} \quad 13.003$$

$$\text{Then } \sqrt{V(\hat{\beta}_0)} = 3.606$$

$$p = P(T(10) \leq \frac{46 - 50}{\sqrt{2.274 \cdot 3.606}} = -2.442) = 1.102$$

~~bis~~  $\Rightarrow p\text{-value} \Rightarrow \text{reject } H_0$ . Our drug is not "good".

11/18

§ 11.6 : predictions via  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ .

Our focus so far has been on the coeffs  $\hat{\beta}_i$  and their distribution.

Of course, the original goal was to understand  $y$  as a fun of  $x$

-ex: (potency example again)

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x = 46 - \frac{19}{60} x.$$

If we started the antibiotic at  $X = 20^\circ F$ , what would we expect the potency to be?

$$\hat{y}(20) = 46 - \frac{19}{60}(20) = 119/3 = 39.\bar{6}.$$

But what does this really mean?

This is the expected value of  $y$  when  $x = 20$ . If we stare a bunch of samples at  $20^\circ F$ , we'd expect the sample mean  $\approx 39.\bar{6}$ .

In other words  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  is an estimator for the average value of  $y @ x$ .

$\hat{y}$  is a statistic of its own.

11/19

And where there are point estimators,  
there are interval estimates.

We have that  $\hat{\beta}_i \sim N(\beta_i, V(\beta_i))$ .

$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  is also Normal at any fixed  $x$ .

To avoid formula confusion, we use  $x^*$  for the fixed  $x$ .

$$\text{Note } E(\hat{y}(x^*)) = E(\hat{\beta}_0 + \hat{\beta}_1 x^*) \\ = \beta_0 + \beta_1 x^*.$$

We need  $V(\hat{y})$ .

$$V(\hat{y}) = V(\hat{\beta}_0 + \hat{\beta}_1 x^*) \\ = V(\hat{\beta}_0) + (x^*)^2 V(\hat{\beta}_1) + 2x^* \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) \\ = \frac{\sigma^2 \sum x_i^2}{n S_{xx}} + (x^*)^2 \frac{\sigma^2}{S_{xx}} + 2x^* \left( -\frac{\bar{x} \sigma^2}{S_{xx}} \right) \\ = \frac{\sigma^2}{S_{xx}} \left( \frac{\sum x_i^2 + n(x^*)^2 - 2x^* \bar{x} n}{n} \right)$$

(using  $S_{xx} = \sum x_i^2 - n \bar{x}^2$ )

$$= \frac{\sigma^2}{S_{xx}} \left( \frac{S_{xx} + n \bar{x}^2 + n(x^*)^2 - 2n x^* \bar{x}}{n} \right) \\ = \frac{\sigma^2}{S_{xx}} \left( \frac{S_{xx} + n((x^*)^2 - 2x^* \bar{x} + \bar{x}^2)}{n} \right) \\ = \frac{\sigma^2}{S_{xx}} \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right)$$

11/20

When estimating  $\hat{\sigma}^2$ , make sure to use the new one  $S_{\hat{\sigma}}^2 = \frac{1}{n-d} SSE$ .

As in all things, we standardize  $\frac{\hat{\theta} - \theta}{\sqrt{\hat{\sigma}^2}}$ .

$$\text{Here } \frac{\hat{Y}(x^*) - E(\hat{Y}(x^*))}{\sqrt{V(\hat{Y})}} = \frac{\hat{Y}(x^*) - \hat{\theta}(\hat{\beta}_0 + \hat{\beta}_1 x^*)}{S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$$

$\underbrace{\qquad\qquad\qquad}_{\text{This is a t-stat w/ } n-2 \text{ df!}}$

An  $1-\alpha$  level confidence interval comes from

$$P\left( \left| \frac{\hat{Y}(x^*) - E(\hat{Y}(x^*))}{\sqrt{V(\hat{Y})}} \right| < t_{\alpha/2, (n-2)} \right) = 1-\alpha.$$

$$\Rightarrow \boxed{(\hat{\beta}_0 + \hat{\beta}_1 x^*) \pm t_{\alpha/2, (n-2)} S \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$$

ex: Find a 90% confidence interval for the average potency of our antibiotic when stored at  $20^\circ$ .

$$\text{Here } x^* = 20$$

$$\hat{Y} = 46 - \frac{19}{60} (20) = 39.6$$

11/21

$$\text{need } V(\hat{Y}) = S^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right)$$

data set had  $n=12 \Rightarrow df = 10$

$$\bar{x} = 60$$

$$S^2 = \frac{571}{30}$$

$$, S_{xx} = 6000$$

For  $90\%$ ,  $t_{0.05}(10) = 1.812$ .

$$39.6 \pm 1.812 \sqrt{\frac{571}{30} \left( \frac{1}{12} + \frac{(10-60)^2}{6000} \right)}$$

$$2.5810$$

$$\Rightarrow 39.6 \pm 4.672$$

Rmk: Of course, we could now do hypothesis testing if we'd like.

Something like..  $H_0: \hat{Y}(x^*) = y_0$   
 $H_a: \hat{Y}(x^*) \neq y_0$  etc

Our t-stat  $T = \frac{y_0 - (\hat{\beta}_0 + \hat{\beta}_1 x^*)}{\sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$   
w/  $n-2$  df

11.7

## § 11.7: predictions on $\hat{Y}$ .

The difference btwn this section and the previous is perspective.

In 11.6, we focused on the spread of the average of  $y$ ,  $E[\bar{Y}]$

What if instead we wanted to focus on the distribution of the values that  $\hat{Y}$  can take when  $x=x^*$ .

Recall  $\hat{Y} = \underbrace{\beta_0 + \beta_1 x}_{\text{deterministic part}} + \underbrace{\epsilon}_{\text{the "spread" of } Y \text{ is here}}$

$$V(Y) = V(\epsilon) = \sigma^2.$$

An estimator for  $\hat{Y}$  at  $x^*$  is defined

$$\hat{Y}^* = \hat{\beta}_0 + \hat{\beta}_1 x^* + \epsilon$$

$$\begin{aligned} \text{Since } \epsilon &\sim N(0, \sigma^2), E[\hat{Y}^*] = E[\hat{\beta}_0 + \hat{\beta}_1 x^*] \\ &= \hat{Y}(x^*) \\ &= \beta_0 + \beta_1 x^*. \end{aligned}$$

Computing variance is surprisingly easy.

As before,  $\hat{Y}^*$  is the sum of Normal r.v.  
Hence it is Normal.

11/28

Moreover, the assumption on  $\epsilon$  is that it is independent of  $x$ .

We also have that  $S^2$  is indep of  $\hat{\beta}_0$  and  $\hat{\beta}_1$ .

$$\begin{aligned} \text{So } V(\hat{Y}^*) &= V(\hat{\beta}_0 + \hat{\beta}_1 x^*) + V(\epsilon) \text{ by indep.} \\ &= S^2 \left( \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) + S^2 \\ &= S^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right) \end{aligned}$$

this is the only difference from  $V(\hat{Y})$ .

When  $S^2$  is unknown...

$$V(\hat{Y}^*) = S^2 \left( 1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}} \right)$$

For the  $(1-\alpha)$  level confidence interval

$$P(|\hat{Y}^* - (\hat{\beta}_0 + \hat{\beta}_1 x^*)| \leq t_{\alpha/2}(df)) = 1-\alpha.$$

$$Z = \frac{\hat{Y}^* - E[\hat{Y}^*]}{S_{\hat{Y}^*}} \approx \frac{\hat{Y}^* - (\hat{\beta}_0 + \hat{\beta}_1 x^*)}{S \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}}$$

$t$ -distributed w  $n-d$  df.

and  $(1-\alpha)$  level C.I. is

$$\hat{\beta}_0 + \hat{\beta}_1 x^* \pm t_{\alpha/2} S \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{S_{xx}}}$$

11pd4

disc: confidence + prediction bands.

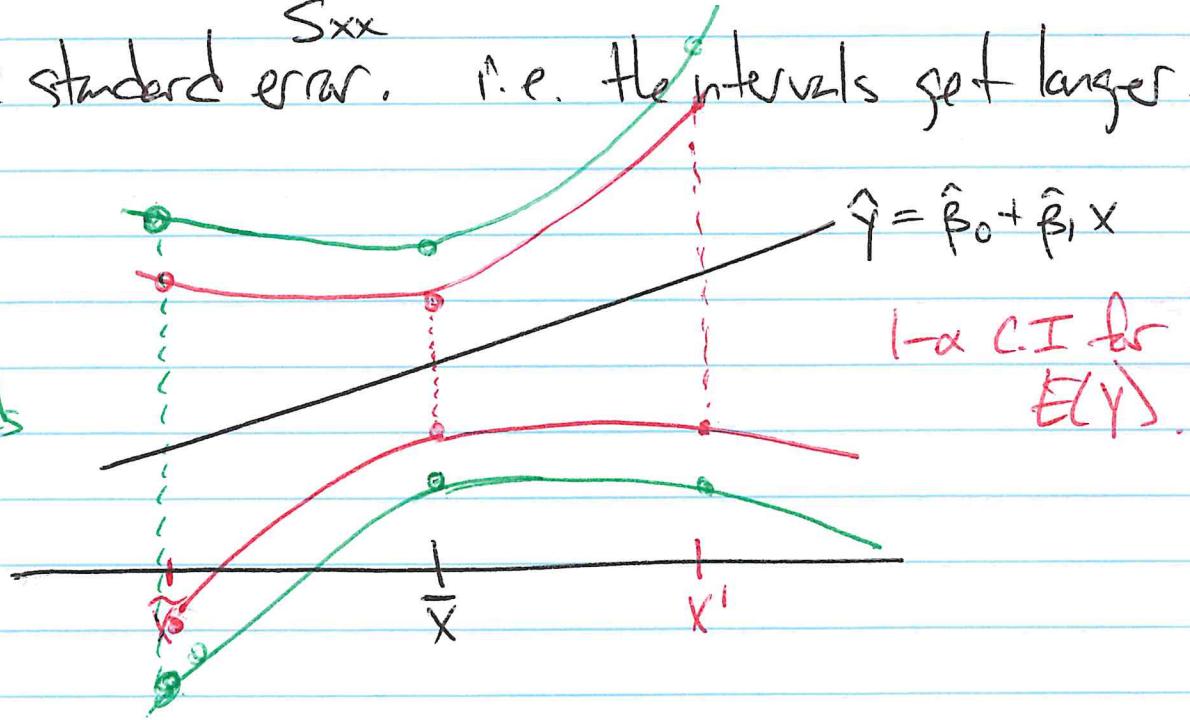
In §11.6, we had C.I. for the mean of  $\hat{Y}(x^*)$

In §11.7, we have C.I. for the spread of  $Y$  at  $x^*$ .

Note that in either case, as  $x^*$  moves away from  $\bar{x}$ , the  $\frac{(x^* - \bar{x})^2}{S_{xx}}$  increases and so does

the standard error. i.e. the intervals get larger.

For level  
prediction bands  
for  $y$ .



Ex: potency of antibiotic when stored at  $20^\circ\text{F}$

Last section we showed  $E[\hat{Y}_{(20)}]$  had the 90% C.I.

$$39.6 \pm 4.677 \quad (\text{red band above})$$

What spread of potency would we expect to see with 90% confidence when stored at  $20^\circ\text{F}$

11/25

(the same calculation but w/ the new standard error)

$$39.6 \pm 1.812 \sqrt{1 + \frac{1}{10} + \frac{1600}{6000}}$$
$$\boxed{39.6 \pm 5.869.}$$
 (green band above)

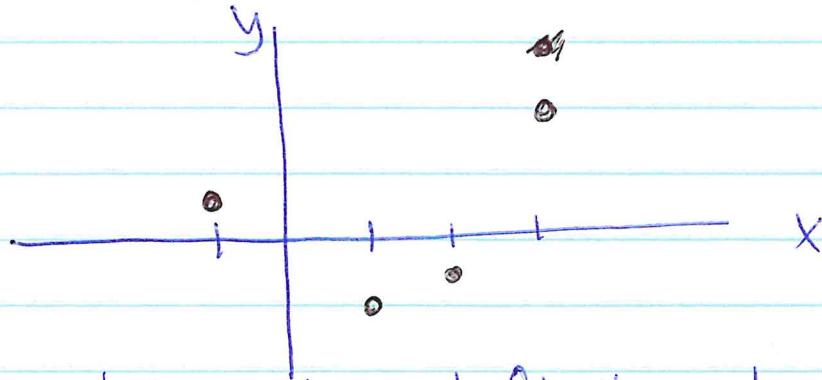
11/26

## § 11.10 Multiple Linear Regression.

Ex: Consider the data

$$(-1, \frac{1}{2}), (1, -1), (2, -\frac{1}{2}), (3, 2)$$

Plot



Seems unlikely a line will fit this well.  
What about a parabola?

Want  $y = b_0 + b_1 x + b_2 x^2$

(of course, in the probabilistic set up  
 $y = \underbrace{f_0 + f_1 x + \beta_2 x^2}_{\text{deterministic}} + \epsilon$ )

Using the data

$$\begin{aligned} \frac{1}{2} &= f_0 + \beta_1(-1) + \beta_2(-1)^2 \\ -1 &= f_0 + \beta_1 + \beta_2 \\ -\frac{1}{2} &= (f_0 + 0)\beta_1 + 4\beta_2 \\ 2 &= f_0 + 3\beta_1 + 9\beta_2 \end{aligned}$$

but this is the same idea as before, just a higher-dimensional problem

11pd7

$$\vec{Y} = \beta_0 \vec{1} + \beta_1 \vec{X} + \beta_2 \vec{X^2}$$
$$\vec{y}_k = \beta_0 \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix} + \beta_1 \begin{pmatrix} -1 \\ 1 \\ 2 \\ 3 \end{pmatrix} + \beta_2 \begin{pmatrix} 1 \\ 1 \\ 4 \\ 9 \end{pmatrix}$$

or, as a matrix-eqn

$$\begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 1/2 \\ -1 \\ -1/4 \\ 2 \end{bmatrix}$$

Remark: This is why it is called "linear" regression.  
The system for the coeffs  $\beta_i$  will be  
a linear system.

But just as before, over-determined  $\Rightarrow$   
no soln. The best we can do is find  
the projection.

We need a way to solve this projection problem  
in general.

11pd8

topic: The Normal Eqns.

$$Y = \beta_0 + \beta_1 f_1(x) + \beta_2 f_2(x) + \dots + \beta_n f_n(x) + \epsilon$$

Let there be  $m$  data points  $(x_i, y_i)$ :  
This will yield a matrix eqn

$$\begin{bmatrix} | & f_1(x_1) & \dots & f_n(x_1) \\ | & \vdots & & \vdots \\ | & f_1(x_m) & \dots & f_n(x_m) \end{bmatrix} \vec{\beta}_0 = \vec{Y}$$

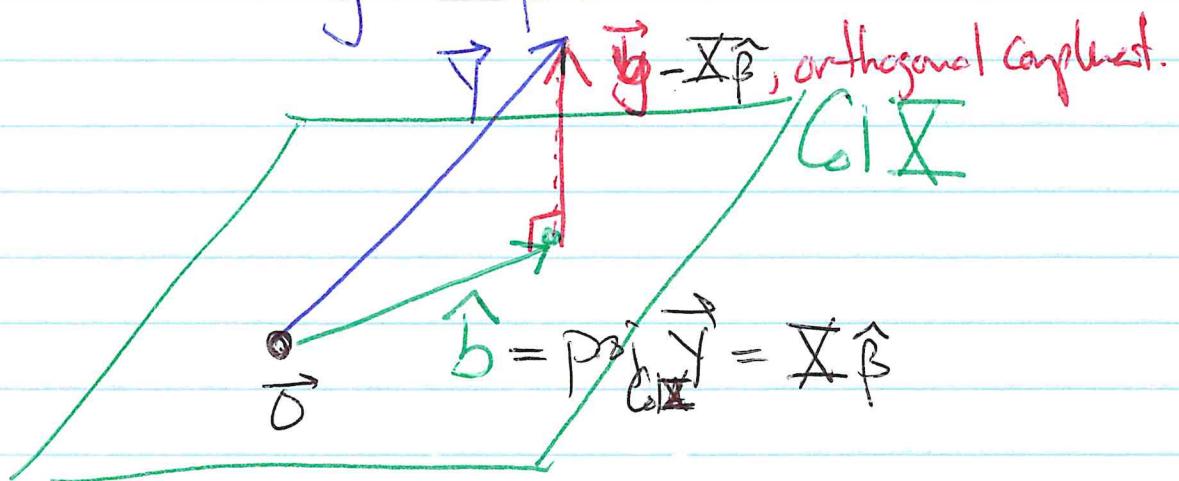
Call this  $m \times n$  matrix  $\vec{X}$

$$\vec{X} \vec{\beta}_0 = \vec{Y}$$

Note, this setup only makes sense when system is  
overdetermined. That is  $M > m$ .

(more rows than columns)

We seek the coefficient vector  $\vec{\beta} = \langle \hat{\beta}_0, \dots, \hat{\beta}_n \rangle$   
such that  $\| \vec{y} - \vec{X} \vec{\beta} \|$  is minimized



11p2A

In the general case, there are now issues to contend with:

- ① no reason the column vectors of  $\vec{X}$  are actually a basis for  $\text{Col } \vec{X}$  (might just be a spanning set).
- ② certainly no reason columns are an orthogonal set.

The implication of ① is that a least squares sol'n  $\hat{\beta}$  need not be unique.

As for ②, ends up we don't need this.

By def'n  $\vec{y} - \vec{X}\hat{\beta}$  is orthogonal to  $\text{Col } \vec{X}$ .  
In other words,  $\vec{y} - \vec{X}\hat{\beta}$  must lie in the null space of  $\vec{X}^T$ .  
(Recall  $(\text{Col } \vec{A})^\perp = \text{Null } \vec{A}^T$ .)

$$\text{So } \vec{X}^T (\vec{y} - \vec{X}\hat{\beta}) = \vec{0}$$

or

$$\vec{X}^T \vec{y} - \vec{X}^T \vec{X} \hat{\beta} = \vec{0}$$

which yields ...

def'n: The Normal Equns  $\vec{X}^T \vec{X} \hat{\beta} = \vec{X}^T \vec{y}$ .

1/p30

## FACTS:

- ①  $\mathbf{X}^T \mathbf{X}$  is a square  $n \times n$  matrix.
- ② The normal equation is always a consistent system,  
That is, there is a sol'n to the least squares problem.
- ③ If the columns of  $\mathbf{X}$  are linearly independent,  
then  $\mathbf{X}^T \mathbf{X}$  is invertible and the unique sol'n is  
$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$
.

Ex: an best-fit parabola.

$$\mathbf{X} = \begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 1 \\ 1 & 2 & 4 \\ 1 & 3 & 9 \end{bmatrix} \quad \mathbf{X}^T = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & 1 & 2 & 3 \\ 1 & 2 & 4 & 9 \end{bmatrix}$$

$$\hat{\beta} = \langle \hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2 \rangle, \vec{y} = \langle 1/2, -1, -1/6, 2 \rangle$$

$$\mathbf{X}^T \mathbf{X} = \begin{bmatrix} 4 & 5 & 15 \\ 5 & 15 & 35 \\ 15 & 35 & 99 \end{bmatrix}$$

$$(\text{and } \det(\mathbf{X}^T \mathbf{X}) = 440 \dots \text{invertible})$$

11/31

$$\text{So } \hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \begin{bmatrix} -41/44 \\ -379/440 \\ 53/88 \end{bmatrix}$$

So the best-fit parabola is

$$y = \frac{-41}{44} - \frac{379}{440}x + \frac{53}{88}x^2$$

$$\approx -0.932 - 0.861x + 0.602x^2$$

(Show picture)

disc: revisit the best-fit line.

Given  $n$   $(x_i, y_i)$ 's,  $y = \beta_0 + \beta_1 x + \epsilon$

$$\Rightarrow [\vec{1} \vec{x}] \hat{\beta} = \vec{y}, \hat{\beta} = \langle \hat{\beta}_0, \hat{\beta}_1 \rangle$$

Here  $\mathbf{X} = [\vec{1} \vec{x}]$  is  $n \times d$ .

$$\mathbf{X}^T = \begin{bmatrix} \vec{1}^T \\ \vec{x}^T \end{bmatrix} \text{ is } 2 \times n.$$

$$\text{Then } \mathbf{X}^T \mathbf{X} = \begin{bmatrix} \vec{1}^T \\ \vec{x}^T \end{bmatrix} [\vec{1} \vec{x}] \text{ is } 2 \times 2$$

$$= \begin{bmatrix} \vec{1} \cdot \vec{1} & \vec{1} \cdot \vec{x} \\ \vec{x} \cdot \vec{1} & \vec{x} \cdot \vec{x} \end{bmatrix} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x^2 \end{bmatrix}$$

$$\text{Note } \det(\mathbf{X}^T \mathbf{X}) = n \sum x^2 - (\sum x)^2 = \Delta$$

11, p3)

$$\begin{aligned} &= n \sum x_i^2 - n \bar{x}^2 \\ &= n (\sum x_i^2 - n \bar{x}) \\ &= n S_{xx} > 0 \dots \text{invertible!} \end{aligned}$$

$$S_0 \hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \vec{y} \text{ unique LSR soln!}$$

$$\mathbf{X}^\top \vec{y} = \begin{bmatrix} \vec{I}^\top \\ \vec{x}^\top \end{bmatrix} \vec{y} = \begin{bmatrix} \vec{I} \cdot \vec{y} \\ \vec{x} \cdot \vec{y} \end{bmatrix} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix}$$

Recall  $\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \frac{1}{ad-bc} \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$

$$S_0 (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix}$$

and

$$\begin{aligned} \hat{\beta} &= \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{bmatrix} \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} \\ &= \frac{1}{n S_{xx}} \begin{bmatrix} (\sum x_i^2)(\sum y_i) - (\sum x_i)(\sum x_i y_i) \\ (-\sum x_i)(\sum y_i) + n \sum x_i y_i \end{bmatrix} \end{aligned}$$

(These are the screencast terms)

But wait! There is more.

$$\text{Look at } (\mathbf{X}^\top \mathbf{X})^{-1} = \frac{1}{n S_{xx}} \begin{bmatrix} \sum x_i^2 & -n \bar{x} \\ -n \bar{x} & n \end{bmatrix}$$

1/p33

$$\sqrt{V(\hat{\beta}_0)} = \left[ \frac{\sum x_i^2}{S_{xx}} - \frac{\bar{x}}{S_{xx}} \right] \cdot \frac{1}{\sqrt{S_{xx}}} \quad \text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{\sum x_i^2}{S_{xx}} \cdot \frac{-\bar{x}}{S_{xx}}$$

$$V(\hat{\beta}_1) = \frac{1}{S_{xx}}$$

We have seen those before.  
They are the cov. of all the  $\text{Cov}(\hat{\beta}_i)$

This actually always happens!

FACT:  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$  is the table of covariances.  
i.e.  $\sigma^2 a_{ij} = \text{Cov}(\hat{\beta}_i, \hat{\beta}_j)$  where  $[a_{ij}] = (\mathbf{X}^T \mathbf{X})^{-1}$ .

$$\text{Lastly, } S^2 = \frac{\text{SSE}}{n-2} = \frac{\sum (y_i - \hat{y}_i)^2}{n-2} \dots$$

$$\text{FACT: } \text{SSE} = \vec{y}^T \vec{y} - \vec{\beta}^T \mathbf{X}^T \vec{y}$$

$$\begin{aligned} \text{Reason: } \text{SSE} &= \sum (y_i - \hat{y}_i)^2 \\ &= \sum (y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i))^2 \\ &= \text{Can be interpreted as the dot product} \\ &= (\vec{y} - \mathbf{X} \vec{\beta})^T (\vec{y} - \mathbf{X} \vec{\beta}) \\ &= \vec{y}^T \vec{y} - 2 \vec{\beta}^T \mathbf{X}^T \vec{y} + \vec{\beta}^T \mathbf{X}^T (\mathbf{X} \vec{\beta}) \end{aligned}$$

Since  $\vec{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$ , we have

$$= \vec{y}^T \vec{y} - 2 \vec{\beta}^T \mathbf{X}^T \vec{y} + \vec{\beta}^T \mathbf{X}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \vec{y}$$

11p34

$$\text{but } (\bar{X}^T \bar{X}) (\bar{X}^T \bar{X})^{-1} = I \text{ so}$$

$$\begin{aligned} SSE &= \bar{Y}^T \bar{Y} - 2 \hat{\beta}^T \bar{X}^T \bar{Y} + \hat{\beta}^T \bar{X}^T \bar{Y} \\ &= \bar{Y}^T \bar{Y} - \hat{\beta}^T \bar{X}^T \bar{Y}. \end{aligned}$$

Ex: (antibiotic, yet again)

temp  $\bar{X} = \langle 30, 30, 30, 50, 50, 50, 70, 70, 70, 90, 90, 90 \rangle$

potency  $\bar{Y} = \langle 38, 43, 41, 32, 26, 33, 19, 27, 23, 14, 19, 21 \rangle$

$$n=12, \bar{X}^T \bar{X} = \begin{bmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix} = \begin{bmatrix} 12 & 720 \\ 720 & 49200 \end{bmatrix}$$

$$\det(\bar{X}^T \bar{X}) = 72000$$

$$(\bar{X}^T \bar{X})^{-1} = \begin{bmatrix} 4/60 & -1/120 \\ -1/120 & 1/6000 \end{bmatrix} \quad \text{Note } V(\hat{\beta}_0) = \frac{4}{60} \sigma^2, \quad V(\hat{\beta}_1) = \frac{1^2}{60000}$$

$$\bar{X}^T \bar{Y} = \begin{bmatrix} \sum y_i \\ \sum x_i y_i \end{bmatrix} = \begin{bmatrix} 324 \\ 17540 \end{bmatrix}$$

$$\hat{\beta} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{Y} = \begin{bmatrix} 46 \\ -19/60 \end{bmatrix} \quad \text{Same as before.}$$

11/35

$$\hat{Y} = 46 - \frac{19}{60} X.$$

$$SSE = \vec{Y}^T \vec{Y} - \hat{\beta}^T \vec{X}^T \vec{Y} = \frac{571}{3}$$

$$S^2 \approx \sigma^2 = \frac{SSE}{n-2} = \frac{573}{30} = 19.1$$

§ 11.11 a big ole thm

Thm Let  $Y_i = \beta_0 + \beta_1 f_1(x_i) + \beta_2 f_2(x_i) + \dots + \beta_k f_k(x_i) + \epsilon_i$

where  $E(\epsilon_i) = 0, V(\epsilon_i) = \sigma^2 \forall i$ .

Then the least-squares estimators are given by

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

provided  $(\vec{X}^T \vec{X})^{-1}$  exists and

$$\textcircled{1} E(\hat{\beta}_i) = \beta_i \quad (\text{unbiased estimators})$$

$$\textcircled{2} V(\hat{\beta}_i) = c_{ii} S^2 \quad \text{where } (\vec{X}^T \vec{X})^{-1} = [c_{ij}]$$

$$\textcircled{3} \text{Cov}(\hat{\beta}_i, \hat{\beta}_j) = c_{ij} S^2$$

$$\textcircled{4} S^2 = \frac{SSE}{n-(k+1)} \quad \text{where } SSE = \vec{Y}^T \vec{Y} - \hat{\beta}^T \vec{X}^T \vec{Y}$$

and  $E(S^2) = \sigma^2$ .

1/p36

If additionally  $\epsilon_i \sim N(0, \sigma^2)$ , then

③  $\hat{\beta}_1$  is normally distributed

⑥  $\frac{(n-(k+1))S^2}{\sigma^2} \sim \chi^2(n-(k+1))$

⑦  $S^2$  and  $\hat{\beta}_1$  are independent for all  $i$ .

## 8.11.2 Hypothesis Testing + C.I. multiple regression!

We have  $Y = \beta_0 + \beta_1 f_1(x) + \beta_2 f_2(x) + \dots + \beta_K f_K(x)$

$$\text{by LS2 soln } \hat{\beta} = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{y}$$

$$= \langle \hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_K \rangle$$

Motivational example: how do we test a specific  $\hat{\beta}_i$ ?

To "pick off" a specific  $\hat{\beta}_i$ , we dot  $\hat{\beta}$  with the  $i^{th}$  standard basis vector  $\vec{e}_{i+} = \langle 0, 0, \dots, 0, 1, 0, \dots, 0 \rangle$   
*i<sup>th</sup> position + 1*

$$\text{then } \hat{\beta}_i = \vec{e}_{i+}^T \hat{\beta} = \vec{e}_{i+}^T \hat{\beta}$$

$$\text{Note } E(\vec{e}_{i+}^T \hat{\beta}) = E(1 \cdot \hat{\beta}_i) = E(\hat{\beta}_i) = \beta_i$$

$$\text{and } V(\vec{e}_{i+}^T \hat{\beta}) = V(\hat{\beta}_i).$$

So a test of the form  $H_0: \beta_i = (\beta_i)_0$   
 $H_{\alpha}: \beta_i > (\beta_i)_0$   
 $<$   
 $=$

yields the statistic

$$Z = \frac{\hat{\beta}_i - (\beta_i)_0}{\sqrt{V(\hat{\beta}_i)}} = \frac{\hat{\beta}_i - (\beta_i)_0}{\sqrt{\lambda_{ii}}}$$

$$[\lambda_{ii}] = (\bar{X}^T \bar{X})^{-1}$$

11/28

$$\text{or } T = \frac{\hat{\beta}_1 - (\beta_1)_0}{S\sqrt{C_{11}}}, \text{ a t-stat of } n-(k+1) \text{ df.}$$

ex: We fit  $(-1, 1/2), (1, -1), (2, -1/2), (3, 2)$

w/  $\hat{Y} = \beta_0 + \beta_1 X + \beta_2 X^2$  because it "looked" non-linear. Is there evidence that the data is in fact non-linear? That is, are we certain  $\beta_2 \neq 0$ ?

$$\begin{aligned} H_0: \beta_2 &= 0 & \text{We had } \hat{\beta} &= (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{y} \\ H_a: \beta_2 &\neq 0 & &= \left\langle \frac{-41}{44}, \frac{-379}{440}, \frac{53}{88} \right\rangle \end{aligned}$$

To pick off  $\hat{\beta}_2$ , we dot w/  $\vec{e}_3 = \langle 0, 0, 1 \rangle$ .

$$\vec{e}_3 \cdot \hat{\beta} = 53/88$$

$$\text{We had } \bar{X}^T \bar{X} = \begin{bmatrix} 4 & 5 & 15 \\ 5 & 15 & 35 \\ 15 & 35 & 99 \end{bmatrix}, \det(\bar{X}^T \bar{X}) = 48$$

$$(\bar{X}^T \bar{X})^{-1} = \frac{1}{480} \begin{bmatrix} 260 & 30 & -50 \\ 30 & 171 & -65 \\ -50 & -65 & 35 \end{bmatrix}$$

$$V(\hat{\beta}_2) = C_{22} \cdot T^2 = \frac{35}{440} \rightarrow$$

1/37

For  $\sigma^2$ , need  $S^2 = \frac{SSE}{n-3}$

$$SSE = \vec{Y}^T \vec{Y} - \vec{\beta}^T \vec{X}^T \vec{Y} = \frac{4791}{440}$$

$$S^2 = \frac{4791}{1320} \approx 3.62955$$

$$\text{and } V(\hat{\beta}_0) \approx \frac{35}{440} \cdot S^2 \approx 0.317585$$

~~$\sqrt{V(\hat{\beta}_0)} = \sqrt{V(\beta_0)} = 1.06872.$~~

$$\sqrt{V(\hat{\beta}_0)} \approx 0.563547$$

$$\begin{aligned} \text{Using } T &= \frac{\vec{e}_3^T \hat{\beta} - (\beta_0)_0}{\sqrt{V(\hat{\beta}_0)}} = \frac{53/88 - 0}{0.563547} \\ &= 1.06872 @ df = 4 - 3 = 1 \end{aligned}$$

$$P(|T| > 1.06872 | df=1) = 0.4788602.$$

Surprising! Supports  $H_0$ . What went wrong?

Probably the small # of data points  $\Rightarrow df=1$ .

This is why statisticians use multiple tests

for example, the SSE of the quadratic fit will be much smaller than the SSE of the linear fit

11/24/10

ex: What about the polynoy example?  
Would quadratic have been a better fit?

ie have line of best fit  $y = 46 - \frac{19}{100}x$

$$w/ SSE \text{ line} = \frac{571}{3} = 190.\overline{3}$$

Instead by  $y = f_0 + f_1 x + f_2 x^2$

$$\Rightarrow \mathbf{X} = [1 \ \vec{x} \ \vec{x^2}]$$

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{1.92 \times 10^6} \begin{bmatrix} 1.09120 \times 10^6 & -391200 & 3100 \\ -391200 & 147200 & -120 \\ 3100 & -120 & 1 \end{bmatrix}$$

$$\hat{\beta} = \left\langle \frac{582}{12}, -\frac{5}{12}, \frac{1}{1200} \right\rangle$$

$$\text{any } y = \frac{582}{12} - \frac{5}{12}x + \frac{1}{1200}x^2$$

$$\approx 48.58 - 0.4166x + 0.00083x^2$$

note  $\hookrightarrow \sim 46 \quad \hookrightarrow \sim -\frac{19}{60} = -0.3166$ )  
pretty small!

11/04/1

To test hypo that model is better as a linear model consider  
H<sub>0</sub>:  $\beta_2 = 0$   
H<sub>1</sub>:  $\beta_2 \neq 0$ .

Need  $V(\hat{\beta}_2) = \frac{1}{1920000} S^2 \approx \frac{S^2}{1920000}$

For  $S^2 = \frac{SSE}{n-3} = \frac{SSE}{9}$ ,

$$SSE = \vec{Y}^T \vec{Y} + \vec{\beta}^T \vec{X}^T \vec{Y} = 189$$

and  $S^2 = \frac{189}{9} = 21$ .

So  $V(\hat{\beta}_2) \approx \frac{21}{1920000} = 1.09375 \times 10^{-5}$

and  $\sqrt{V(\hat{\beta}_2)} \approx 0.00330719$ .

for T-stat,  $T = \frac{\hat{\beta}_2 - (\beta_2)_0}{\sqrt{V(\hat{\beta}_2)}}$  @ df=9  
 $= \frac{1/200}{0.00330719}$   
 $= 0.251976.$

Well... this is no more conclusive than the last

11/24/20

$$P(|T| > 0.251976 \mid df=9)$$

$$= 0.8067194.$$

Again does not support th.

disc: generalization.

Of course, this is math and we have to generalize. Don't have to use an SBV  $\vec{\beta}_i$ .

Let  $\vec{a} = (a_0, \dots, a_k)$  be constant.

Then  $\vec{a}^T \vec{\beta} = a_0 \hat{\beta}_0 + a_1 \hat{\beta}_1 + \dots + a_k \hat{\beta}_k$   
a weighted linear combo of the  $\hat{\beta}_i$ 's.

$$\begin{aligned} \text{Note } E(\vec{a}^T \vec{\beta}) &= \sum_i^k a_i E(\hat{\beta}_i) \text{ by linearity} \\ &= \sum_i^k a_i f_i \text{ by unbiased estimators} \\ &= \vec{a}^T \circ \hat{\beta} \end{aligned}$$

Variance is more complicated.

$$\begin{aligned} V(\vec{a}^T \vec{\beta}) &= V(a_0 \hat{\beta}_0 + a_1 \hat{\beta}_1 + \dots + a_k \hat{\beta}_k) \\ &= \sum_{i=1}^k \sum_{j=1}^k a_i a_j \text{Cov}(a_i \hat{\beta}_i, a_j \hat{\beta}_j) \quad \text{by Chapter 5.} \end{aligned}$$

11/pcB

$$= \sum_{i=1}^k \sum_{j=1}^k a_i a_j \text{cov}(\hat{\beta}_i, \hat{\beta}_j)$$

these are just  $c_{ij}\sigma^2$  like  $[c_{ij}] = (\mathbf{X}^T \mathbf{X})^{-1}$   
In fact, using matrix algebra, this is the same as

$$\sqrt{(\hat{\beta}^T \hat{\beta})} = \hat{\beta}^T (\mathbf{X}^T \mathbf{X})^{-1} \hat{\beta} \cdot \sigma^2$$

yielding  $\hat{\beta} - E(\hat{\beta}) = \frac{\hat{\beta} - \hat{\beta}}{\sqrt{(\hat{\beta}^T \hat{\beta})}} = \frac{\hat{\beta} - \hat{\beta}}{\sigma \sqrt{\hat{\beta}^T (\mathbf{X}^T \mathbf{X})^{-1} \hat{\beta}}}$

(need to know)

For a test,  $H_0: \hat{\beta}^T \hat{\beta} = (\hat{\beta}^T \hat{\beta})_0$

H<sub>a</sub>:  $\hat{\beta}^T \hat{\beta} (\neq, >, <) (\hat{\beta}^T \hat{\beta})_0$

we use

$$T = \frac{\hat{\beta}^T \hat{\beta} - (\hat{\beta}^T \hat{\beta})_0}{S \sqrt{\hat{\beta}^T (\mathbf{X}^T \mathbf{X})^{-1} \hat{\beta}}} \quad \text{on } n-(k+1) \text{ df}$$

where  $S^2 = \frac{SSE}{n-(k+1)}$

Aside: Dr. Ahanda says a common usage of this is  
to test collections of  $\beta_i$  to be zero.  
In other words, to test the necessity of  
variables in your model (like we did in previous chapters)