

R.R. If  $t > t_{\alpha/2; n-(k+1)}$   
 or  $t^2 > t_{\alpha/2; n-(k+1)}^2$

$F_{dup}$  (or  $F_{part}$ )  $\Rightarrow$  R.R. We reject if  $F_{dup}$  (or  $F_{part}$ )  $> F_{\alpha}(1; n-k-1)$

$$t_{n-(k+1)}^2 = F_{part}(1; n-(k+1)) = F_{dup}(1; n-(k+1))$$

### $F_{part}$

- is based on breaking down the SSR into components, i.e. we determine the contribution to SSR made by variable  $X_j$  after all other variables have been included

i.e.  $SS$  (due to  $X_j$  after all other  $X$ 's in the model) -  $SSR(X_j | \text{all other } X\text{'s})$

$$= \overset{\text{"just happened already"}}{SSR(\text{due all } X\text{'s})} - SSR(\text{all } X\text{'s w/out } X_j) =$$

$$= SSR(X_1, X_2, \dots, X_{j-1}, X_{j+1}, \dots, X_k) - SSR(X_1, X_2, \dots, X_{j-1}, X_{j+1}, \dots, X_k)$$

$$= SSR_R - SSR_F$$

✘

$r^2 = \frac{SSR}{TSS}$  = determines how much of the tot. variation in the data can be explained by the reg.  
 ↑ def. of determination line and here, how much is due to error

- if  $r^2$  is 0 (or close to) = when none of the  $X$ 's in the model have strong linear relationship w/  $y$

- if  $r^2 = 1$  (or close to) it does not necessarily imply a good model as there may be only few  $X$ 's that are linearly

related to  $y$  (and hence the rest of the  $x$ 's do not contribute)

- also SSR increases as # of  $x$ 's incs., hence  $r^2 = \frac{SSR}{TSS}$  will incs also (it does not reflect

the goodness of fit)

- better measure is so-called  $r^2$  adjusted which will drop down if unimportant  $x$ 's are entered into the model.

$$r^2_{adj} = 1 - \frac{MSE}{TSS/n-1} = 1 - \frac{SSE/n(k+1)}{TSS/n-1}$$

i.e. - If 96 vs 70 there is an issue

Dummy variables (Indicator variables)

- are used to represent qualitative (categorical) variables such as (sex, marital status, geographic region, mother tongue)

the # of dummy variables needed to represent such a variable is always one less than # of categories the variable in question has.

i.e. sex =  $\begin{cases} m \\ f \end{cases} \Rightarrow 2 \text{ categories} \Rightarrow 1 \text{ dummy } X = \begin{cases} 0 & \text{if } m \\ 1 & \text{if } f \end{cases}$

coffee cup size =  $\begin{cases} S \\ M \\ L \end{cases} \Rightarrow 3 \text{ cat's} \Rightarrow 2 \text{ dum. } X_1 = \begin{cases} 1 & \text{if } S \\ 0 & \text{otherwise} \end{cases}$

$X_2 = \begin{cases} 1 & \text{if } M \\ 0 & \text{otherwise} \end{cases}$   
if  $S$   $X_1 = 1$   $X_2 = 0$  if  $M$   $X_1 = 0$   $X_2 = 1$  if  $L$   $X_1 = 0$   $X_2 = 0$

Ex) want to examine the relationship btwn university salary ( $y$ ), the # of yrs of experience ( $X_1$ ) and the sex of professor ( $X_2$ )

We might suspect a linear relationship btwn  $y$  &  $X_1$  for both males & females i.e. we run SLR model  $y = \beta_0 + \beta_1 X_1 + \epsilon$  2 times (once for males & once for F)

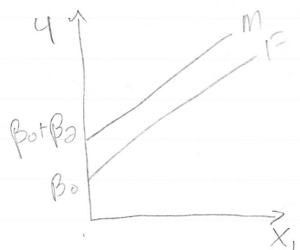
- get 2 lines & can compare them

- But we want  $X_2$  (sex of prof) in the same model w/  $X_1$

i.e.  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$  where  $X_2 = \begin{cases} 0, \text{ if } F \\ 1, \text{ if } M \end{cases}$

if prof M:  $y = \beta_0 + \beta_1 X_1 + \beta_2(1) + \epsilon$   
 or  $y = \beta_0 + \beta_2 + \beta_1 X_1 + \epsilon \dots (M)$

if prof F:  $y = \beta_0 + \beta_1 X_1 + \beta_2(0) + \epsilon$   
 or  $y = \beta_0 + \beta_1 X_1 + \epsilon \dots (F)$



To test that starting salaries are the same  
 $H_0: \beta_2 = 0$   
 $H_a: \beta_2 \neq 0$

$\beta_0$  is y intercept  
 $\beta_1$  = slope of F & M  
 $\beta_2 = (\beta_0 + \beta_2) - \beta_0$   
 difference in y intercepts btwn M & F's

- Different y intercepts, but same (parallel) slopes "DIFF starting salaries but same increase"

- in gen.



- in this case, we say that  $X_1$  &  $X_2$  interact and therefore interaction term (ie. cross-product term) needs to be added to the model

cross product  
i.e. interaction  
term

$$\text{i.e. } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon \leftarrow \text{"interaction model"}$$

Here we have different  $y$  intercepts and slopes

if a prof is a male:  $y = \beta_0 + \beta_1 x_1 + \beta_2(1) + \beta_3 x_1(1) + \epsilon$   
 or  $y = (\beta_0 + \beta_2) + (\beta_1 + \beta_3)x_1 + \epsilon \dots (m)$

if a prof is a female:  $y = \beta_0 + \beta_1 x_1 + \beta_2(0) + \beta_3 x_1(0) + \epsilon$   
 or  $y = \beta_0 + \beta_1 x_1 + \epsilon \dots (f)$

$\beta_0 = y$ -intercept for  $F$

$\beta_2 = (\beta_0 + \beta_2) - \beta_0 = \text{difference in } y\text{-intercepts btwn } m + F$

$\beta_1 = \text{slope of } F$

$\beta_3 = (\beta_1 + \beta_3) - \beta_1 = \text{diff. in slopes for } m + F$

### \* Interaction models

- if the relationship btwn  $y$  +  $x_1$  for example, is different for diff. values of  $x_2$  then there is an interaction effect btwn  $x_1$  +  $x_2$  which should be accounted for in the model by adding the cross product (i.e. interaction term)  $x_1 x_2$

$$\text{i.e. } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_1 x_2 + \epsilon$$

NOTE: when using dummy variable(s) you often need to include terms representing the possible interactions of the dummy variable(s) w/ all the other non-dummy variables.

- Comparing two lines for M's + F's

testing if the two lines have the same y-intercept

(only pair use t-test)

$H_0: \beta_2 = 0; \alpha$

$H_a: \beta_2 \neq 0$

t-test:  $t = \frac{\hat{\beta}_2}{\sqrt{\text{Var}_{\hat{\beta}_2} \cdot \text{MSE}}}$  Rej. if  $t > t_{\alpha/2; n-4}$   
 $t < -t_{\alpha/2; n-4}$

$F_{part} = \frac{[SSR_P - SSR_R] / (df_{SSR_P} - df_{SSR_R})}{\text{MSE}_R}$

Full:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$

R:  $y = \beta_0 + \beta_1 x_1 + \beta_3 x_3 + \epsilon$

$\frac{[SSR_P - SSR_R] / [3 - 2]}{SSE_R / (n - 4)} \Rightarrow$  RR if  $F_{part} > F_{\alpha}(1, n-4)$

test if the 2 lines have the same slopes or test for parallelism

$H_0: \beta_3 = 0; \alpha$  / t-test:  $t = \frac{\hat{\beta}_3}{\sqrt{\text{Var}_{\hat{\beta}_3} \cdot \text{MSE}}}$   $\Rightarrow$  RR if  $t > t_{\alpha/2; n-4}$   
 $H_a: \beta_3 \neq 0$   $t < -t_{\alpha/2; n-4}$

$F_{drop} = \frac{[SSE_R - SSE_P] / (df_{SSE_R} - df_{SSE_P})}{\text{MSE}_P}$

Full: same

R:  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon$

$\frac{[SSE_R - SSE_P] / [n - 3 - (n - 4)]}{SSE_P / (n - 4)}$  RR if  $F_{drop} > F_{\alpha}(1; n-4)$

- test if the 2 lines are coincident (i.e. same y intercept and the same slope)

(+ because of parallel)

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

$H_a$ : @ least one of the  $\beta$ 's  $\neq 0$ ;  $\alpha$

$$F_{\text{part}} = \frac{[SSR_R - SSR_{R'}] - [df_{SSR_R} - df_{SSR'}]}{MSE_D}$$

E: same

$$R: y = \beta_0 + \beta_1 x_1 + \epsilon$$

$$R.R. F_{\text{part}} > F_{\alpha}(2; n-4)$$

$$\downarrow \frac{[SSR_R - SSR_{R'}] / (3-1)}{SSE_D / n-4}$$

$$F_{\text{dup}} = \frac{[SSE_R - SSE_{R'}] / \frac{(n-2) - (n-4)}{df_{SSE_R} - df_{SSE_{R'}}}}{MSE_D} \quad \left. \begin{array}{l} \text{or eq. w/} \\ (n-2) - (n-4) \end{array} \right\}$$

R.R. if  $F_{\text{dup}} > F_{\alpha}(2; n-4)$

A3 Q2 An experimenter wished to compare the potencies of 3 different drug products. To do this 12 test tubes were inoculated w/ a culture of the virus under study and incubated for 2 days @ 35°C. Four dosage levels (0.2, 0.4, 0.8 + 1.6 mg/tub) were to be used from each of the 3 different drug products A, B, + C, w/ only one dose-drug combo. for each of the 12 test tube cultures

"how well the drug work"

a) Write a general lin. model relating the response ( $y$ ) to the independent var.  $x$ :

"dose" and "drug product" Make  $X_i$  = Indose  
Identify the parameters in the model

- here drug product - (A, B or C) ← categorical variable. ∴ we need 2 dummy variables:  
 $- X_2 = \begin{cases} 1; & \text{if B} \\ 0; & \text{otherwise} \end{cases}, X_3 = \begin{cases} 1; & \text{if C} \\ 0; & \text{otherwise} \end{cases}$

$$\text{MLR } y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_1 x_2 + \beta_5 x_1 x_3 + \epsilon$$

if we have drug A

$$x_2 = 0$$

$$x_3 = 0$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2(0) + \beta_3(0) + \beta_4 x_1(0) + \beta_5 x_1(0) + \epsilon$$

$$y = \beta_0 + \beta_1 x_1 + \epsilon \dots \dots (A)$$

if drug B

$$x_2 = 1$$

$$x_3 = 0$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2(1) + \beta_3(0) + \beta_4 x_1(1) + \beta_5 x_1(0) + \epsilon$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2 + \beta_4 x_1 + \epsilon$$

$$y = (\beta_0 + \beta_2) + (\beta_1 + \beta_4) x_1 + \epsilon \dots \dots (B)$$

if drug C

$$x_2 = 0$$

$$x_3 = 1$$

$$y = \beta_0 + \beta_1 x_1 + \beta_2(0) + \beta_3(1) + \beta_4 x_1(0) + \beta_5 x_1(1) + \epsilon$$

$$y = (\beta_0 + \beta_3) + (\beta_1 + \beta_5) x_1 + \epsilon \dots \dots (C)$$

A is the base line

$\beta_0 = y$  intercept for drug A

$$\beta_2 = (\beta_0 + \beta_2) - \beta_0 = \text{diff. in } y\text{-interc. for the lines A + B}$$

$$\beta_3 = (\beta_0 + \beta_3) - \beta_0 = \text{diff. in } y\text{-interc. for the lines A + C}$$

$\beta_1 =$  slope of line A

$$\beta_4 = (\beta_1 + \beta_4) - \beta_1 = \text{diff. in slopes b/w A + B}$$

$$\beta_5 = (\beta_1 + \beta_5) - \beta_1 = \text{diff. in slopes b/w A + C}$$

b.) It would be reasonable to assume that the 3 lines have a common y-interc. I.e. test if their y interc. is equal =.

$$H_0: \beta_0 = \beta_3 = 0 \quad n=12$$

$$H_a: \text{at least one of } \beta_i \neq 0 ; <$$

2 pairs  
no F-test

$$\text{test stat. } F_{\text{pair}} = \frac{[SSR_P - SSR_F] / (df_{SSR_P} - df_{SSR_F})}{MSE_F} =$$

$$\text{reduced } = y = \beta_0 + \beta_1 X_1 + \beta_4 X_1 X_2 + \beta_5 X_1 X_3 + \epsilon$$

$$= \frac{[SSR_P - SSR_F] / (5-3)}{SSE_F / n-6}$$

R.R. Rej.  $H_0$  if  $F_{\text{pair}} > F_{\alpha}(2; 6)$

or

$$F_{\text{dup}} = \frac{[SSE_F - SSE_P] / (df_{SSE_F} - df_{SSE_P})}{MSE_P}$$

$$= \frac{[SSE_F - SSE_P] / (n-4) - (n-6)}{SSE_P / n-6}$$

R.R. Rej.  $H_0$  if  $F_{\text{dup}} > F_{\alpha}(2; 6)$

c.) A3, Q2 Test whether the 3 drug lines are parallel (i.e. the same slopes).

in SAS  $\ln x$  does not exist  
use  $\log(x)$

In Real life

$\ln x$   
 $\log_{10} x$   
 $\log_2 x$

In SAS

$\log(x)$   
 $\log_{10}(x)$   
 $\log_2(x)$

Not on exam

## Multicollinearity (not in textbook)

① IF independent (i.e. explanatory) variables are uncorrelated, then  $r^2_{y|x_1, x_2, \dots, x_k} = r^2_{y|x_1} + r^2_{y|x_2} + \dots + r^2_{y|x_k}$

Are  $x$ 's independent?

$$(r^2 = \frac{SSR}{TSS}) = SSR_{(x_1, x_2, \dots, x_k)} = SSR(x_1) + SSR(x_2) + \dots + SSR(x_k)$$

$$= SSR(x_j | \text{all other } x\text{'s}) = SSR_j - SSR_r = SSR(x_1, x_2, \dots, x_k) - SSR(x_1, x_2, \dots, x_k)$$

$$= SSR(x_j)$$

- IF  $x$ 's are related this is not true

② IF the independent (i.e. explanatory) variables are correlated then we have so-called multicollinearity

- IF occurs when some of the explanatory vars are related to each other i.e. when they are correlated among themselves
- For example.

$$x_3 = 2x_1 + 3x_2 - 5x_4 \leftarrow x_3 \text{ is linearly related to } x_1, x_2, \text{ \& } x_4$$

- $x_3$  does not bring anything new to the model, but can't be disregarded

$$\Rightarrow \bullet SSR(x_j | \text{all other } x\text{'s in the model}) = 0 \text{ (or very small \#)}$$

but if  $SSR(x_j)$  - it might be a very large # as  $x_j$  is linearly related to other  $x$ 's

- $X^T X$  is not invertible (since 1 row is a linear combo. of some other row) and hence there are infinitely many sol's to  $\beta = (X^T X)^{-1} (X^T Y)$

not on exam

- ~~By~~ not rejecting  $H_0: \beta_j = 0$  means only that  $X_j$  has no additional contribution (linear relationship) w/  $Y$  when the other  $X$ 's to which it is related are in the model.

### Indications of multicollinearity in a Regression

- ① Large correlation coefficient btwn pairs of  $X$ 's
- ②  $SSR(X_j | \text{all other } X\text{'s present}) \ll SSR(X_j)$   
much, much smaller than
- ③ Adding or deleting an  $X$  variable that is highly related to some of the other  $X$ 's in the model will cause a huge chng in the val. of the  $\beta$  of those  $X$ 's to which  $(X_j)$  is related.
- ④ Some of the regression coefficients ( $\beta_j$ 's) have opposite signs to what you would expect from theory or common sense.
- ⑤ wide confidence intervals for  $\beta_j$ 's, i.e. lrg sampling variances ( $v_j$ ; MSE) for  $\beta_j$ 's
- ⑥ Even though  $F \text{ test} = \frac{MSR}{MSE}$  for testing whether

$\beta_1 = \beta_2 = \beta_k = 0$  (i.e. for testing for a significant linear relationship btwn  $y$  + @ least one of the  $X$ 's) will indicate a high significant linear relationship, if however, we test for individual contributions of  $X$ 's to  $y$  i.e. (testing  $H_0: \beta_j = 0$  ( $j=1 \dots k$ )), the results are not significant

- ⑦ Large (variance inflation factors) VIF's, where  $VIF = \frac{1}{1-R_j^2}$ , where  $R_j^2$  is coefficient of determination for the regression of  $X_j$  on all other  $X$ 's.

Let on  
Exam

$$\text{i.e. } \hat{X}_j = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \dots + \hat{\beta}_{j-1} X_{j-1} + \hat{\beta}_{j+1} X_{j+1} + \dots + \hat{\beta}_K X_K$$

If  $X_j$  is highly related to all other  $X$ 's then  $R_j^2$  is high, then  $1 - R_j^2$  is very small

$\Rightarrow \frac{1}{1 - R_j^2}$  is very large

~~Rule of thumb~~

VIF > 5 - moderate multicollinearity in the data.

VIF > 10 - serious multicollinearity in the data

in SAS proc Reg corr;  
model  $y = X_1 \dots X_K / \text{VIF};$   
run;

(Ex) fuel:  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$  collinearity?

we get  $R = \begin{matrix} & \begin{matrix} X_1 & X_2 & X_3 \end{matrix} \\ \begin{matrix} X_1 \\ X_2 \\ X_3 \end{matrix} & \begin{bmatrix} 1 & 0.9463 & 0.982 \\ & 1 & 0.8761 \\ & & 1 \end{bmatrix} \end{matrix}$  ex. correlation between  $X_1, X_2 = 94.63\%$

high pairwise correlation coefficients

i.e.  $X_1, X_2 = 0.9463$

$X_1, X_3 = .982$

$X_2, X_3 = .8761$

- start w/ testing  $X_1$  (highest correlation)

•  $SSR(X_1 | X_2, X_3) \ll SSR(X_1) \ll ?$  then?

Reduced model -  $y = \beta_0 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$

$SSR(X_2, X_3) = 322,184,328$   
Full:  $y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$  → OVER  
 $SSR(X_1, X_2, X_3) = 325,360,308$

Not an exam

$$y = \beta_0 + \beta_1 x_1 + \epsilon$$

$$SSR(x_1) = 314,544,161$$

$$SSR(x_1 | x_2, x_3) = \overset{\text{"given"}}{SSR_F} - SSR_r = SSR(x_1, x_2, x_3) - SSR(x_2, x_3)$$

$$= 3,325,980 \ll 314,544,161 \quad (SSR(x_1) \checkmark)$$

Confirmation  $x_1$  is creating problems

Full:  $\beta_1 = 0.854$ ,  $\beta_2 = 44.375$ ,  $\beta_3 = 82.285$  } <sup>1 sig fig</sup> log jump  
 reduced:  $\beta_1 = 71.976$ ,  $\beta_2 = 117.454$

- full -  $H_0: \beta_1 = \beta_2 = \beta_3 = 0$   $F_{(3,17)} = 45.093$ ,  $pval = .0001$   
 highly significant linear relationship between  $y$  and @ least one of the  $x$ 's

To find out which  $x$ :  $\leftarrow \rightarrow H_0: \beta_j = 0$

- |         |       |   |  |
|---------|-------|---|--|
| $(x_1)$ | $j=1$ | $t_{(17)} = 1.176 \rightarrow pval = 0.2558$  | } non-significant linear relationship between $y$ + any of the $x$ 's (i.e. don't reject $H_0$ ) |
| $(x_2)$ | $j=2$ | $t_{(17)} = 1.411 \rightarrow pval = 0.1767$  |  |
| $(x_3)$ | $j=3$ | $t_{(17)} = 1.3978 \rightarrow pval = 0.1805$ |  |

- large variance inflation factors  
 $\sqrt{1/R^2}$  12.88, 9.71 + 5.795      710 serious multic.

### Remedies for Multicollinearity

- ① Drop some of the correlated variables from the model if it is reasonable to do so. (Do not drop a variable that is known from theory to belong in the model)
- ② Fit the regression equation using  $(x_j - \bar{x}_j)$  instead of  $x_j$  (especially if we have polynomial model)

not on exam!

Polynomial models

in general,  $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_k x^k$   
 - only one variable (high correlation)  
 or

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$$

where  $x_1 = x$   
 $x_2 = x^2$   
 $\vdots$   
 $x_k = x^k$

instead of doing  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_k x_k$  fit  
 $\hat{y} = \hat{\beta}_0^* + \hat{\beta}_1^* (x - \bar{x}) + \hat{\beta}_2^* (x - \bar{x})^2 + \dots + \hat{\beta}_k^* (x - \bar{x})^k$   
 \* means small difference in  $\beta$ 's

- Which x's should be in model? (not in book)  
 - model (or variable) selection

- how do we decide which variable should be used in the model?  
 - how do we select (choose) 'best-fitting' model?  
 (select highest,  $R^2_{adj}$  if different results)  
 - to tot. procedures may give different results

$x_1, x_2, x_3, x_4$

<p>Full is best select 2nd highest <math>r^2</math></p>	$x_1 - r^2 - mse$ - choose and smallest full is 1st	$x_1 x_2 x_3$
	$x_2 - r^2 - mse$	$x_1 x_2 x_4$
	$x_3 - r^2$ "	<u><math>x_2 x_3 x_4</math></u>
	<u><math>x_4 - r^2</math> "</u>	$x_1 x_2 x_3 x_4 \rightarrow r^2 \rightarrow$ must be highest
	$x_1 x_2 - r^2$ "	$r^2 = \frac{SSR}{TSS}$
	$x_1 x_3 - r^2$ "	$\rightarrow mse \rightarrow$ smallest always
	$x_1 x_4 - r^2$ "	
	$x_2 x_3 - r^2$ "	

use if close to  $r^2$  or full

Not on exam

• max  $R^2$  criterion

- we calculate  $r^2 = \frac{SSR}{TSS}$  for each possible model

and then select "best-fitting" model, the one w/ the "highest"  $r^2$  (other than the full model) if possible)

NOTE Choose as best-fitting model the one w/ higher  $r^2$  only if deciding between two models.

• min. MSE criterion

- is equivalent to max  $R^2$

- we calculate MSE for each possible model and then select as the "best-fitting" model the one w/ the smallest MSE (other than the full model if possible)

Ex)  $n=20$  independent pharmacies. Want to predict the sales based on which variables among the following:  $X_1 = \text{tot flr. space}$

$X_2 = \text{\# of parking spots}$

$X_3 = \text{location of the pharmacy}$

$X_4 = \text{space allocated to prescription dept}$

\# of variables in the model

one-var.  
 $Y = \beta_0 + \beta_1 X_1 + \epsilon$       $R^2$      variable  
 .439      $X_1$

two-var.  
 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \epsilon$      .666      $X_1, X_2$

3-var.  
 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon$      .691      $X_1, X_2, X_3$

Full  
 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4 + \epsilon$      .694      $X_1, X_2, X_3, X_4$

since these are close we could calculate  $R^2$  and use model w/ highest value.

← 2nd highest "best-fitting model" w/  $X_1, X_2, X_3, X_4$

Not an exam

• Mallow's Cp Criterion

- It measures the difference of a fitted regression line from the true model, along w/ random errors
- Theory suggests that the "best-fitting" model is one w/ Cp closest to p (where p is the tot. # of parameters in the model including  $\beta_0$ )

We need to calculate Cp statistic for each one of the possible models where Cp stat. is

$$C_p = \frac{SSE_p}{MSE_p} - (n-2p), \text{ where } SSE_p = SSE \text{ for the model w/ } p \text{ parameters}$$

Note: When  $p=1$  i.e. no X's ( $y = \beta_0 + \epsilon$ )  
 $C_p = \frac{ISS}{MSE_p} - (n-2p) \rightarrow$  if this is closest p, all X's are no good.

no. var.	p	Cp	variable.
one var.	1	(closest to 2) 0.17	$X_4$
two var.	2	(closest to 3) 2.47	$X_3, X_4$
3 var.	3	(closest to 4) 2.96	$X_1, X_2, X_4$
4 var.	4	(closest to 5) 4.04	$X_1, X_2, X_3, X_4$

- must draw graph to find shortest distance.