

$$42.4 \quad 43) \text{ p31 } R^2 = [\text{cor}(Y, \hat{Y})]^2 \quad \text{if } n \geq 5p \quad M484 \quad 8$$

a constant is in the model.

$0 \leq R^2 \leq 1$ . If  $n \geq 5p$ , high  $R^2$

means the least squares "line" is much better at explaining  $Y|X$  than the line  $Y = \bar{Y}$ . (Identity line or hyperplane)

44) p31 Suppose  $X \in \mathbb{R}^p$ , the normal MLR model holds and  $B_2 = \dots = B_p = 0$ .

Then the line  $Y = \bar{Y}$  should be used instead of the OLS "line"

but  $E R^2 = \frac{p-1}{n-1}$ , so

$$\text{if } n=101 \text{ and } p=91 \quad \text{then } E R^2 = \frac{51}{21} = .52$$

This is one reason  $n \geq 5p$  is needed.

(Want  $R^2$  near 0 if  $B_2 = \dots = B_p = 0$ .)

45) know reject  $H_0$  if pvalue <  $\delta$ .  
use  $\delta = 0.05$  if not given.  
 $\delta$  is

46) know for final The Anova F test tests whether the predictors  $x_2, \dots, x_p$  are needed in the model. If needed, use OLS. If not needed, use the line  $\hat{Y} = \bar{Y}$ . Assume a constant is in the model. The 4 step test is

$t_{\text{test}}$   
 $t_{\text{CI}}$   
 $t_{\text{know 483}}$

i)  $H_0: \beta_2 = \dots = \beta_p = 0 \quad H_A: \text{not } H_0$

ii)  $F_0 = \frac{MSR}{MSE}$

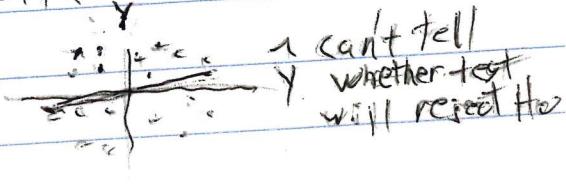
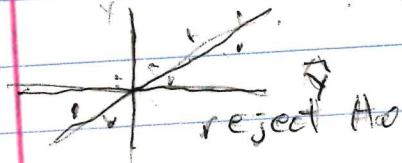
iii)  $P\text{val} = P(F_{p-1, n-p} > F_0)$

iv) Reject  $H_0$  if  $P\text{val} < \alpha$ . Then there is an MLR relationship between  $y$  and the predictors  $x_2, \dots, x_p$ .

Fail to reject  $H_0$  if  $P\text{val} \geq \alpha$ . Then there is not an MLR relationship between  $y$  and  $x_2, \dots, x_p$ .

Replace  $y, x_2, \dots, x_p$  by names from story problem.

47) Response and residual plots should be good with  $n \geq 5p$ . If identity line fits data better than any horiz line, then the Anova F test will reject  $H_0$ . If the response plot looks like the residual plot so there is a horiz line that fits the data about as well as the identity line, then can't tell whether the Anova F test will reject  $H_0$  or fail to reject  $H_0$ , but the MLR relationship is weak.



48) p32 know  
Anova table

Source	SS	df	MS	F	Pval
model or regression	SSR	p-1	MSR = SSR/(p-1)	F = $\frac{MSR}{MSE}$	for $H_0: \beta_2 = \dots = \beta_p = 0$
residual or error	SSE	n-p	MSE = SSE/(n-p)		
total	SS TO	n-1			
often omitted					

See ex 2.5, 2.6

ex) mussel data n=82,  $y = \log(m)$ ,  $m$  = mussel muscle mass

$x_2 = \log(s)$   $x_3 = \log(h)$ ,  $s$  = shell mass

$h$  = shell height

source	df	ss	ms	F	pval
reg	2	33.9208	16.9604	339.05	0.000
residual	79	3.9518	0.05063		

Perform the Anova F test

i)  $H_0: \beta_2 = \beta_3 = 0$   $H_A: \text{not } H_0$

ii)  $F = 339.05$

iii)  $pval = 0.0$

iv) reject  $H_0$ , there is an MLR relationship between  $\log(m)$  and the predictors  $\log(s)$  and  $\log(h)$

49) p30 If there is a constant in the MLR model,  $SS TO = SSE + SSR$

$SSE = \sum_{i=1}^n r_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2$  measures residual variability  
error (hat) sum of squares

$SS TO = \sum_{i=1}^n (y_i - \bar{y})^2$ ,  $SSR = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$   
measures variability of the response  $\uparrow$  regression sum of squares

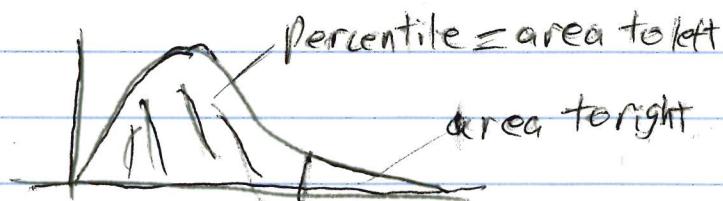
P32, 33, 136 Theorem 2.7 —  $P_{H_0} = P_{F_{n-p}}$ , assuming  $H_0$  is true  
 50) pvalue =  $P_{H_0}$ (obtaining a test statistic  $F_0$ )  
 as extreme as the one observed,  
 where "as extreme" depends on the form of  $H_A$ .  
 $pval = P(F_{p-1, n-p} > F_0)$  is given by  
 the output.

$pval = \text{pvalue}$  if the  $e_i$  are iid  $N(0, \sigma^2)$ ,  
 and if  $n-p$  is large, then  
 for many data sets,  $pval \approx \text{pvalue}$   
 since  $F_0 \approx F_{p-1, n-p}$ .

51)  $\hat{\sigma}^2 = \frac{SSE}{n-p}$  is an estimator  
 of  $\sigma^2 = V(e_i)$ .

52) pvalue=0 means it was impossible  
 for the statistic to have occurred if  
 $H_0$  is true, while pvalue=1 means  
 that it was impossible for the  
 statistic to have occurred if  
 $H_A$  was true

53) F table

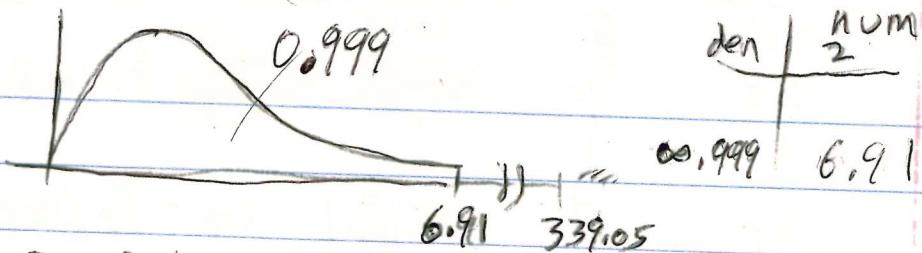


$$pval = 1 - \text{table value}$$

Know If  $n-p > 60$ , use den df =  $\infty$ ,  
 otherwise use closer df.

ex) mussels data  $P(F_{2,79} > 339.05)$

~~SP: P until para X~~

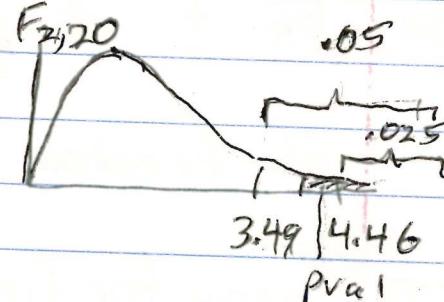


Want  $pval = P(F_{2,18} > 0.4)$

~~den df num df Z~~

20.95	3.49
.975	4.46

} bracket test statistic



so  $.025 < pval < .05$

Want  $pval = P(F_{2,18} > 0.4)$

$$\approx P(F_{2,20} > 0.4) \quad 20.5 \quad | \quad .718$$

so  $pval > .5$



q) For F table a)  $pval > 0.5$  if  $F_0$  is smaller than the 50 value between 1-big table val and 1-small table value  
 b)  $pval \approx 0 < 0.001$  if  $F_0$  value is larger than the .999 value

c)  $pval \approx 0 < 0.001$  if  $F_0$  value is larger than the .999 value

55)	PGI, 63	variable or label	estimate or coef	SE or std err	tvalue	pvalue for H <sub>0</sub>
interceptor constant		$\hat{\beta}_1$	$SE(\hat{\beta}_1)$		$t_{01}$	$\beta_1 = 0$
$x_2$		$\hat{\beta}_2$	$SE(\hat{\beta}_2)$		$t_{02} = \hat{\beta}_2 / SE(\hat{\beta}_2)$	$\beta_2 = 0$
:	:	:	:			:
$x_p$		$\hat{\beta}_p$	$SE(\hat{\beta}_p)$		$t_{0p} = \hat{\beta}_p / SE(\hat{\beta}_p)$	$\beta_p = 0$

56) p50 Know for final A 100 (1- $\delta$ )% confidence interval (CI) for  $\beta_K$  is

$$\hat{\beta}_K \pm t_{n-p, 1-\frac{\delta}{2}} SE(\hat{\beta}_K)$$

use  $\hat{\beta}_K \pm z_{1-\frac{\delta}{2}} SE(\hat{\beta}_K)$  if degrees

of freedom  $d = n - p > 30$ .

57) p50 Know for final A 4 step Wald + test is

i)  $H_0: \beta_K = 0$   $H_A: \beta_K \neq 0$

$$t_{0K} = \frac{\hat{\beta}_K}{SE(\hat{\beta}_K)}$$

$$\text{iii) } p\text{val} = 2 P(t_{n-1} < -|t_{0K}|) \quad \left. \begin{array}{l} \text{usually} \\ \text{from} \\ \text{output} \end{array} \right\}$$

iv) If  $p\text{val} < \delta$ , reject  $H_0$ . So  $x_K$  is needed in the MLR model for  $y$ , given that the other predictors are in the model.

If  $p\text{val} \geq \delta$ , fail to reject  $H_0$ .

So conclude  $x_K$  is not needed in the MLR model for  $y$ , given that the other predictors are in the model. Replace  $x_K$  by variable from story problem.