

**Ch 10 Inferential tools for multiple regression
Chapter 11 Model Checking & Refinement**

Class 17: 4/13/09 M

HW 12 due Friday 4/17/09

Submit as Myname-HW12.doc (or *.rtf)

- HW 12 10.28: El Niño and Hurricanes
 - Due Friday 4/17/09 Noon
- Read Chapter 11: Model Checking and Refinement
 - Read chapter 11 conceptual problems & solutions
 - Post a question or response about Chapter 11 conceptual problems
- No Class Monday 4/20: Patriot's Day

Case 10.2

Energy of echolocating bats:
Do they require more energy than non-echolocating bats or birds, after accounting for the effects of body mass on energy consumption?

EEOS611

Slide 1 Ch 10 Inferential tools for multiple regression

Chapter 11 Model Checking & Refinement

NOTES:

Slide 2 HW 12 due Friday 4/17/09

NOTES:

Slide 3 Case 10.2

NOTES:

Display 10.3

Mass and in-flight energy expenditure for 4 non-echolocating bats (Type = 1), 12 non-echolocating birds (Type = 2), and 4 echolocating bats (Type = 3)

Species	Mass (g)	Type	Flight Energy Expenditure (W)
<i>Pteropus poliocephalus</i>	779	1	43.7
<i>Pteropus poliocephalus</i>	628	1	34.8
<i>Hypsignathus monstrosus</i>	258	1	23.3
<i>Eidolon helvum</i>	315	1	22.4
<i>Meliphaga virescens</i>	24.3	2	2.46
<i>Melipotis undulatus</i>	35	2	3.93
<i>Corvus vulgaris</i>	72.8	2	9.15
<i>Falco sparverius</i>	120	2	13.8
<i>Falco sparverius</i>	213	2	14.6
<i>Corvus ossifragus</i>	275	2	22.8
<i>Larus atricilla</i>	370	2	26.2
<i>Columba livia</i>	384	2	25.9
<i>Columba livia</i>	442	2	29.5
<i>Columba livia</i>	412	2	43.7
<i>Columba livia</i>	330	2	34.0
<i>Corvus cryptoleucus</i>	480	2	27.8
<i>Phyllostomus hastatus</i>	93	3	8.83
<i>Plecotus auritus</i>	8	3	1.35
<i>Pipistrellus pipistrellus</i>	6.7	3	1.12
<i>Plecotus auritus</i>	7.7	3	1.02

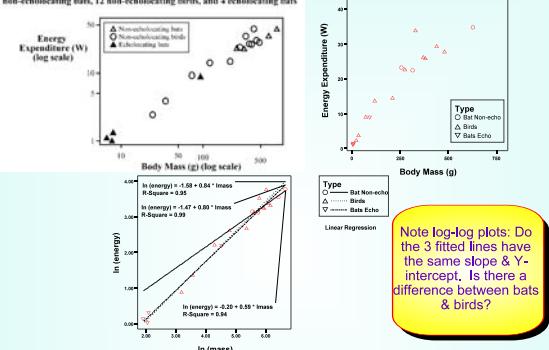
Note that
animal type is
categorical
with 3
categories. It
need 2
(not only 2)
indicator
variables to
code.

Slide 4 Display 10.3

NOTES:

Display 10.4

Log-log scatterplot of in-flight energy expenditure versus body mass for 4 non-echolocating bats, 12 non-echolocating birds, and 4 echolocating bats



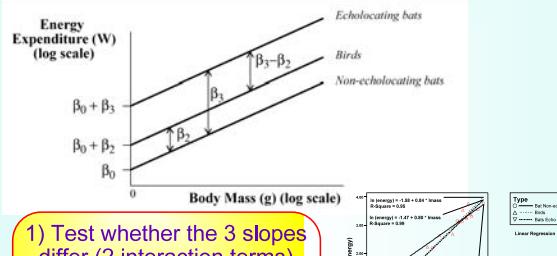
Note log-log plots: Do
the 3 fitted lines have
the same slope &
y-intercept. Is there a
difference between bats
& birds?

Slide 5 Display 10.4

NOTES:

Display 10.5

The parallel regression lines model for the bat echolocation data



- 1) Test whether the 3 slopes differ (2 interaction terms)
- 2) Test whether 3 Y-intercepts differ

Slide 6 Display 10.5

NOTES:

Display 10.12

The extra sum of squares F-test comparing the separate regression lines model to the parallel regression lines model; bat echolocation data

① FIT FULL MODEL: $\mu[\text{energy} \text{Imass, TYPE}] = \beta_0 + \beta_1 \text{Imass} + \beta_2 \text{bird} + \beta_3 \text{ebat} + \beta_4 \text{Imass} \cdot \text{ebat}$					
Source of Variation	Sum of Squares	df	Mean Square	F-Statistic	p-value
Regression	29.46993	5	5.89399	163.4	<.0001
Residual	.50487	14	.03606		
Total	29.97480	19			

② FIT REDUCED MODEL: $\mu[\text{energy} \text{Imass, TYPE}] = \beta_0 + \beta_1 \text{Imass} + \beta_2 \text{bird} + \beta_3 \text{ebat}$					
Source of Variation	Sum of Squares	df	Mean Square	F-Statistic	p-value
Regression	29.42148	3	9.80716	283.6	<.0001
Residual	.55332	16	.03458		
Total	29.97480	19			

...
Residual SS

The models differ in 2 interaction terms

Slide 7 Display 10.12

NOTES:

Display 10.12

The extra sum of squares F-test comparing the separate regression lines model to the parallel regression lines model; bat echolocation data

③ The extra sum of squares is the difference between residual sums of squares	Extra SS = .55332 - .50487 = .04845
⑤ Calculate the F-Statistic	$F\text{-Statistic} = \frac{0.04845}{0.03606} = 1.34$
⑥ Look up $P(F_{2,14} > 1.34)$	$P\text{-value} = 0.53$

Conclusion: There is no evidence that the association between energy expenditure and body size is different for the three types of flying vertebrates ($p\text{-value} = 0.53$).

Slide 8

NOTES:

t tests: Are interaction terms zero?

Parallel slopes model, Can't rely on the t statistic alone to judge whether both interaction terms can be dropped

Model		Unstandardized Coefficients		Standardized Coefficients		t	Sig.	95% Confidence Interval for B		
		B	Std. Error	Beta	t			Lower Bound	Upper Bound	
1	(Constant)	-1.468	.137	.990	-10.705	3.1E-009	-1.76	-1.18		
	In (mass)	.809	.027	.990	30.124	7.4E-017	.75	.86		
	In (mass) * bird	-.170	.077	-.068	-2.207	.030	-.10	-.97		
2	(Constant)	.916	.045	.998	19.297	3.8E-012	.72	.91		
	Birds	.102	.114	.041	.866	.384	-.14	.34		
	Echolocating bats	.079	.203	.028	.388	.703	-.35	.51		
	Inton (Mass, Birds v. Bats)	.246	.213	.536	1.151	.269	-.21	.70		
	Inton (Mass, Echolocating bats)	.215	.224	.204	.961	.353	-.26	.69		

a. Dependent Variable: ln (energy)

Two 1-coefficient-at-a-time t tests can have p values > 0.05, but both together can have p<0.05. Need extra Sum of Squares F test with combined df in numerator.

Slide 9 t tests: Are interaction terms zero?

NOTES:

Extra sum of squares F test (=Partial F test) for both interaction terms

Model Summary ^d										Change Statistics				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square	F Change	df1	df2	Sig. F Change					
1	.999 ^a	.991	.979	1.77	.961	907.838	1	18	.000					
2	.991 ^b	.982	.978	1.659	.001	4.28	2	16	.659					
3	.992 ^c	.983	.977	1.650	.002	.672	2	14	.527					

a. Predictors: (Constant), ln (mass)
b. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo)
c. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo), Intnx (Mass,TwoBats), Intnx (Mass, Birds v. Bats)
d. Dependent Variable: ln (Energy)

Model	Sum of Squares	df	Mean Square	F	Sig.
1	Regression 29.392	1	29.392	907.638	.744E-012*
	Residual .583	18	.032		
Total 29.975	19				
2	Regression 29.421	3	9.807	283.589	.446E-014*
	Residual .553	16	.035		
Total 29.975	19				
3	Regression 29.470	5	5.894	163.440	6.70E-012*
	Residual .505	14	.036		
Total 29.975	19				

a. Predictors: (Constant), ln (mass)
b. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo)
c. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo), Intnx (Mass,TwoBats), Intnx (Mass, Birds v. Bats)
d. Dependent Variable: ln (Energy)

Do echolocating bats differ from non-echolocating bats in energy expenditure?

Non-echolocating bats are the reference category:
Little evidence (t test, p=0.7) that the difference in Y intercepts ≠ 0, so conclude that there is little evidence that echolocating and non-echolocating bats differ in energy expenditure.

Model	Unstandardized Coefficients			Standardized Coefficients			t	Sig.	95% Confidence Interval for B		
	B	Std. Error	Beta	Beta	t	Sig.	Lower Bound	Upper Bound			
1	(Constant) -1.468	.137			-10.765	.000	-17.766	-1.181			
	In (mass) .809	.227	.390	.30127	3.542	.027	.732	.865			
2	(Constant) -1.576	.287			-5.488	.000	-2.188	.867			
	In (mass) .815	.045	.998	.18297	18.297	.000	.721	.969			
	Birds -.102	.114	.241	.096	.904	.396	.140	.241			
	Echolocating bats .379	.203	.703	.359	1.851	.071	.551	.598			
3	(Constant) -.302	1.261			-1.81	.075	-2.908	2.603			
	In (mass) .590	.208	.722	.2861	.313	.148	.102				
	Birds -.1378	1.295	-.252	-.1044	.205	.416	-.416	.407			
	Echolocating bats -.103	1.238	-.414	-.087	.321	.423	-.423	.449			
	Bird x mass .248	.213	.538	1.151	.389	.212	.703				
	Ebat x mass .215	.224	.304	.061	.353	.365	.694				

a. Dependent Variable: ln (Energy)

Extra sum of squares F test for different Y intercepts

Display 10.10

The extra-sum-of-squares F-test for testing equality of intercepts in the parallel regression lines model; bat echolocation data

① Fit the FULL model: $\ln(\text{energy}) | \ln(\text{mass}, \text{TYPE}) = \beta_0 + \beta_1 \ln(\text{mass}) + \beta_2 \text{bird} + \beta_3 \text{ebat}$
Sum of squared residuals = .55332 d.f. = 16 $\hat{\sigma}^2 = .03458$

② Fit the REDUCED model: $\ln(\text{energy}) | \ln(\text{mass}, \text{TYPE}) = \beta_0 + \beta_1 \ln(\text{mass})$
Sum of squared residuals = .58289 d.f. = 18

③ The extra-sum-of-squares is the difference between the two residual sums of squares → Extra SS = .58289 - .55332 = .02957

④ Numerator degrees of freedom are the number of β 's in the full model minus the number of β 's in the reduced model.

⑤ Calculate the F-Statistic = $\frac{.02957}{2} / \frac{.03458}{16} = .014785 / .00222 = .428$

⑥ Find $P(F_{2,14} > .428)$ from table, computer, or calculator → p-Value = .66

Conclusion: There is no evidence that mean log energy differs for birds, echolocating bats, and non-echolocating bats, after accounting for body mass.

Slide 10 Extra sum of squares F test (=Partial F test) for both interaction terms

NOTES:

Slide 11 Do echolocating bats differ from non-echolocating bats in energy expenditure?

NOTES:

Slide 12 Display 10.10

NOTES:

Display 10.6

Partial summary of the least squares fit to the regression of log energy expenditure on log body mass, an indicator variable for bird, and an indicator variable for echolocating bat

Variable	Coefficient	Standard Error	t-Statistic	2-Sided p-Value
CONSTANT	-1.5764	0.2872	5.4880	<0.0001
<i>lmass</i>	0.8150	0.0445	18.2966	<0.0001
<i>bird</i>	0.1023	0.1142	0.8956	0.3837
<i>ebat</i>	0.0787	0.2027	0.3881	0.7030

Estimate of $\sigma = 0.1860$, df = 16

Coefficients^a

	Unstandardized Coefficients		95% Confidence Interval for B				
	B	Std. Error	t	Sig.	Lower Bound	Upper Bound	
(Constant)	-1.576	.287	5.488	4.96E-005	.219	.87	
Echolocating bats	.079	.203	.388	.703	-.35	.51	
Birds	.102	.114	.896	.384	-.14	.34	
In (mass)	.815	.045	18.297	3.76E-012	.72	.91	

a. Dependent Variable: ln (energy)

Slide 13 Display 10.6

NOTES:

Regression σ estimated from ✓Residual mean square

Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square Change	Change Statistics			
						F Change	df1	df2	Sig. F Change
1	.980 ^b	.981	.979	.17995	.981	907.638	1	18	.000
2	.982 ^b	.978	.976	.16596	.001	4.28	2	16	.659
3	.992 ^c	.983	.977	.16990	.002	.872	2	14	.527

a. Predictors: (Constant), ln (mass)
b. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo)
c. Predictors: (Constant), ln (mass), Birds not Bats, Bats (echo) v. Bats (non-echo), Intnx (Mass,TwoBats), Intnx (Mass, Birds v. Bats)
d. Dependent Variable: ln (Energy)

Note that mean squares are variances, the square root of the residual mean square provides the standard error for the regression.

Slide 14 Regression σ estimated from ✓Residual mean square

NOTES:

Display 10.6

Partial summary of the least squares fit to the regression of log energy expenditure on log body mass, an indicator variable for bird, and an indicator variable for echolocating bat

Variable	Coefficient	Standard Error	t-Statistic	2-Sided p-Value
CONSTANT	-1.5764	0.2872	5.4880	<0.0001
<i>lmass</i>	0.8150	0.0445	18.2966	<0.0001
<i>bird</i>	0.1023	0.1142	0.8956	0.3837
<i>ebat</i>	0.0787	0.2027	0.3881	0.7030

Estimate of $\sigma = 0.1860$, df = 16

How do you estimate sigma, σ , standard error of the estimate, root mean square error for the regression, standard error for the regression?

Slide 15

NOTES:

Model Summaries ^d											Change Statistics				
Model	R	R Square	Adjusted R Square	Std. Error of the Estimate	R Square	F Change	df1	df2	Sig. F Change						
1	.990 ^a	.981	.979	17955	.881	907.638	1	18	.000						
2	.991 ^b	.982	.978	18596	.001	4.28	2	16	.659						
3	.992 ^c	.983	.977	18990	.002	6.72	2	14	.527						

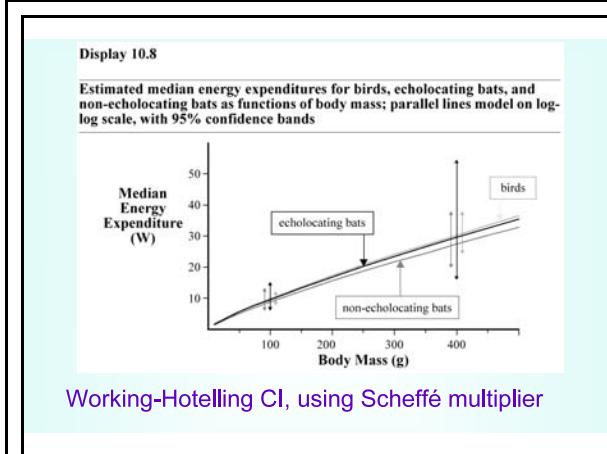
a. Predictors: (Constant), ln(mass)
b. Predictors: (Constant), ln(mass), Birds not Bats, Bats (echo) v. Bats (non-echo)
c. Predictors: (Constant), ln(mass), Birds not Bats, Bats (echo) v. Bats (non-echo), Intxn (Mass, TwoBats), Intxn (Mass, Birds v. Bats)
d. Dependent Variable: ln(energy)

Model		Sum of Squares	df	Mean Square	F	Sig.
1	Regression	29.392	1	29.392	907.638	7.44E-017 ^d
	Residual	.583	18	.032		
	Total	29.975	19			
2	Regression	29.421	3	9.807	283.589	4.46E-014 ^b
	Residual	.553	16	.035		
	Total	29.975	19			
3	Regression	29.470	5	5.894	163.440	6.70E-012 ^c
	Residual	.505	14	.036		
	Total	29.975	19			

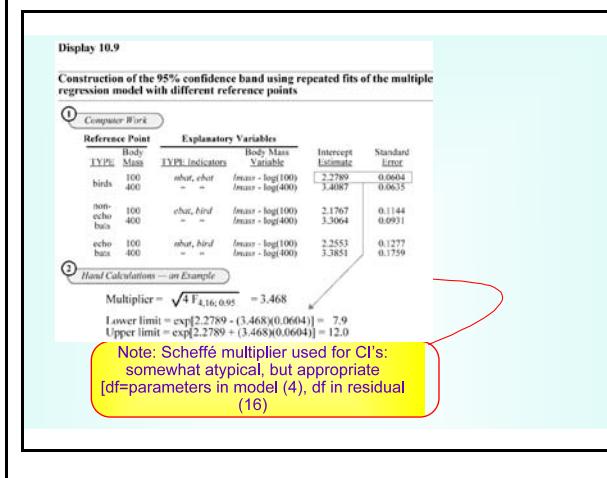
a. Predictors: (Constant), ln(mass)
b. Predictors: (Constant), ln(mass), Birds not Bats, Bats (echo) v. Bats (non-echo)
c. Predictors: (Constant), ln(mass), Birds not Bats, Bats (echo) v. Bats (non-echo), Intxn (Mass, TwoBats), Intxn (Mass, Birds v. Bats)
d. Dependent Variable: ln(energy)

Slide 16

NOTES:

**Slide 17 Display 10.8**

NOTES:

**Slide 18 Display 10.9**

NOTES:

Display 10.8

Estimated median energy expenditures for birds, echolocating bats, and non-echolocating bats as functions of body mass; parallel lines model on log-scale, with 95% confidence bands

Slide 19 Display 10.8

NOTES:

CASE	type	mass	lmass	energy	energylenergy
1	Bat Non-echo	779	6.66	44	3.78
21	Bat Non-echo	100	4.61	-	-
22	Birds	100	4.61	-	-
23	Bats Echo	100	4.61	-	-
24	Bat Non-echo	400	5.99	-	-
25	Birds	400	5.99	-	-
26	Bats Echo	400	5.99	-	-

Slide 20

NOTES:
This Scheffé prediction interval is based on a sample size. There is a further Scheffé adjustment for individual CI's

Syntax: Scheffé multiplier

```
* regpars is the number of parameters in the final model, with 16 df in the residual.
Compute regpars=4.
Compute residdf=16.
exe.
COMPUTE FScheffe = IDF.F(0.95,regpars,residdf) .
EXECUTE .
COMPUTE Scheffemultiplier = sqrt(regpars*FScheffe) .
EXECUTE .
* Scheffé interval is Scheffé multiplier times the standard error for each predicted value, SEP_1 was produced by regression.
COMPUTE SchInt = SEP_1 * Scheffemultiplier .
EXECUTE .
COMPUTE L95S = PRE_1 - SchInt.
COMPUTE U95S = PRE_1 + SchInt.
EXE.
COMPUTE PredE = Exp(PRE_1).
COMPUTE Low95S = Exp(L95S).
COMPUTE Up95S = Exp(U95S) .
EXECUTE .
```

Slide 21 Syntax: Scheffé multiplier

NOTES:

Variance formulae for linear contrasts

Display 10.15

Inference about $\beta_2 - \beta_3$, the coefficient of the indicator variable for birds minus the coefficient of the indicator variable for echolocating bats

① Estimate the linear combination of coefficients as the same linear combination of estimated coefficients.

Estimate of $\beta_2 - \beta_3$, from : .1023 + .0789 = .0826

② Use the estimated variance-covariance matrix of the estimated regression coefficients.

Estimated variance-covariance matrix (from computer):

(Constant)	Intercept	Bird	Echolocating bats
Birds	.03211	.00173	.00096
Non-echolocating bats	.03192	.00173	.00097
Others	.03208	.01304	.01464
Total	.03208	.01304	.01464

Note: (a) the matrix is symmetric; (b) the square roots of the diagonal elements are the standard errors of the estimated coefficients (as reported in).

③ The estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$ is:

$$\hat{\sigma}^2_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3) = (\hat{\beta}_2 - \hat{\beta}_3)^T \hat{\Sigma}_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3)$$

Estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$: .01394 + .04108 - 2 * .01464 = .02484

$\text{SE}(\hat{\beta}_2 - \hat{\beta}_3) = (.02484)^{1/2} = .1576$ (df = 16)

Slide 22 Variance formulae for linear contrasts

NOTES:

Echolocating Bats as Reference

Display 10.15

Inference about $\beta_2 - \beta_3$, the coefficient of the indicator variable for birds minus the coefficient of the indicator variable for echolocating bats

① Estimate the linear combination of coefficients as the same linear combination of estimated coefficients.

Estimate of $\beta_2 - \beta_3$, from : .1023 - .0789 = .02324

② Use the estimated variance-covariance matrix of the estimated regression coefficients.

Estimated variance-covariance matrix (from computer):

(Constant)	Intercept	Bird	Echolocating bats
Birds	.03236	.01378	.0009
Non-echolocating bats	.03207	.02027	.0026
Others	.03208	.01304	.01464
Total	.03208	.01304	.01464

Note: (a) the matrix is symmetric; (b) the square roots of the diagonal elements are the standard errors of the estimated coefficients (as reported in).

③ The estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$ is:

$$\hat{\sigma}^2_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3) = (\hat{\beta}_2 - \hat{\beta}_3)^T \hat{\Sigma}_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3)$$

Estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$: .01394 + .04108 - 2 * .01464 = .02484

$\text{SE}(\hat{\beta}_2 - \hat{\beta}_3) = (.02484)^{1/2} = .1576$ (df = 16)

Slide 23 Echolocating Bats as Reference

NOTES:

Birds vs. Echolocating Bats

Display 10.15

Inference about $\beta_2 - \beta_3$, the coefficient of the indicator variable for birds minus the coefficient of the indicator variable for echolocating bats

① Estimate the linear combination of coefficients as the same linear combination of estimated coefficients.

Estimate of $\beta_2 - \beta_3$, from : .1023 - .0789 = .02324

② Use the estimated variance-covariance matrix of the estimated regression coefficients.

Estimated variance-covariance matrix (from computer):

(Constant)	Intercept	Bird	Echolocating bats
Birds	.03207	.01378	.0009
Non-echolocating bats	.03207	.02027	.0026
Others	.03208	.01304	.01464
Total	.03208	.01304	.01464

Note: (a) the matrix is symmetric; (b) the square roots of the diagonal elements are the standard errors of the estimated coefficients (as reported in).

③ The estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$ is:

$$\hat{\sigma}^2_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3) = (\hat{\beta}_2 - \hat{\beta}_3)^T \hat{\Sigma}_{\text{Var}}(\hat{\beta}_2 - \hat{\beta}_3)$$

Estimated variance of $\hat{\beta}_2 - \hat{\beta}_3$: .01394 + .04108 - 2 * .01464 = .02484

$\text{SE}(\hat{\beta}_2 - \hat{\beta}_3) = (.02484)^{1/2} = .1576$ (df = 16)

Slide 24 Birds vs. Echolocating Bats

NOTES:

Using indicator variables judiciously

- By changing the reference category to echolocating bats, instead of non-echolocating bats, the bird coefficient will test the hypothesis that the bird Y-intercept differs from echolocating bat Y-intercept
- If non-echolocating bats were the reference, the standard error for the difference in Y intercepts could be calculated using the 'propogation of error variance' formula:

$$\text{Var}(a-b) = \text{variance}(a) + \text{variance}(b) - 2 \text{ cov}(a,b)$$
- By coding the flying animal type using non-integer coding, ebats as $\frac{1}{2}$ and nebats as $-\frac{1}{2}$, "Birds vs. Both types of Bats" can be tested. Draper & Smith cover non-integer coding schemes.

Slide 25 Using indicator variables judiciously

NOTES:

Birds vs. Both types of bats

Weighted average of the 2 bat types using $-\frac{1}{2}$ and $\frac{1}{2}$ for the indicator variables for ebats and bats

Model		Unstandardized Coefficients		Standardized Coefficients		95% Confidence Interval for B		
		B	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound
1	(Constant)	-.537	.205		-7.481	.000	-1.973	-1.101
	In (mass)	.815	.045	.998	18.297	.000	.721	.909
	Birds	.063	.093	.025	.676	.509	-.134	.260
	Bats (echo) v.	.079	.203	.020	.388	.703	-.351	.508
	Bats (non-echo)							

a. Dependent Variable: ln (energy)

Statistical summary: After accounting for body mass effects, little evidence that the energy consumption by echolocating bats differs from non-echolocating bats (*t* test, $p=0.7$), nor is there much evidence that birds differ from bats, echolocating or not, in energy expenditure (*t* test, $p=0.51$)

Slide 26 Birds vs. Both types of bats

NOTES:

Chapter 11 Model Checking & Refinement Case Studies

EEOS611

Slide 27 Chapter 11 Model Checking & Refinement

Case Studies

NOTES:

Case Study 11.1 Gender differences in alcohol tolerance

Women have a lower alcohol tolerance & more alcohol related disease (with the same amount of drinking). Why?

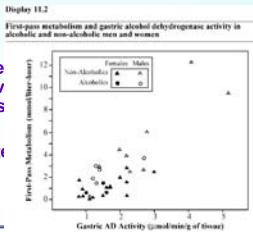
- 18 Men & 14 women volunteers from Trieste
 - 3 women & 5 men were alcoholic
- Ethanol dose 0.3 g/kg orally one day & intravenously another
 - Randomly determined order
- Intravenous- oral= 1st pass metabolism
- Gastric alcohol dehydrogenase activity measured

Display 11.1									
First-pass metabolism of alcohol in the stomach (nmol/litre-hour) and gastric alcohol dehydrogenase activity in the stomach (nmol/min/g of tissue) for 18 control and 14 alcoholics.									
Male (n=8, N=8)					Female (n=7, N=7)				
Metabolism					Metabolism				
Subject					Subject				
1	0.6	1.0	1	3	17	2.5	3.0	1	6
2	1.2	1.5	1	3	18	4.5	1.2	0	1
3	0.4	1.2	1	3	19	2.5	2.0	0	1
4	0.2	1.2	1	0	21	2.0	1.2	0	1
5	0.2	1.2	1	0	22	2.0	1.2	0	1
6	0.2	0.8	1	0	23	3.7	2.5	0	1
7	0.3	0.8	1	0	25	2.5	2.5	0	1
8	0.3	1.9	1	0	25	2.5	2.5	0	0
9	0.0	1.6	1	0	27	3.0	1.4	0	0
10	0.0	1.6	1	0	28	3.0	1.4	0	0
11	1.2	1.6	1	0	29	4.5	2.0	0	0
12	1.2	1.6	1	0	30	4.5	2.0	0	0
13	1.2	1.7	1	0	31	6.5	5.2	0	0
14	1.2	1.7	1	0	32	12.5	4.1	0	0
15	1.8	0.8	1	0					
16	2.0	2.0	1	0					

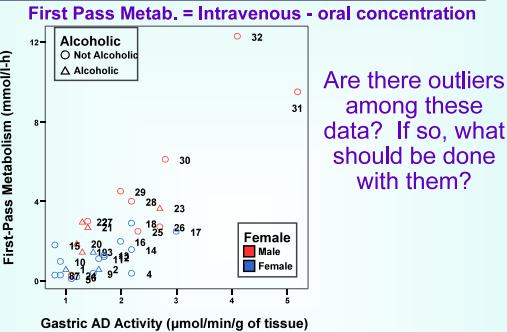
Display 11.2

Research Questions

- Do levels of first-pass metabolism differ between men and women?
- Can differences be explained by postulating that men have more alcohol dehydrogenases activity in their stomachs?
- Are these effects complicated by an alcoholism effect (3 women & 5 men were alcoholics)?



SPSS Analyses



Are there outliers among these data? If so, what should be done with them?

Slide 28 Case Study 11.1 Gender differences in alcohol tolerance

NOTES:

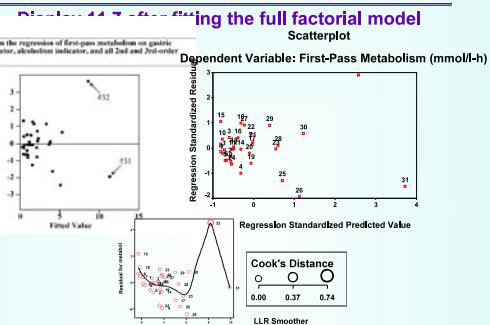
Slide 29 Display 11.2

NOTES:

Slide 30 SPSS Analyses

NOTES:

Identifying outliers that matter

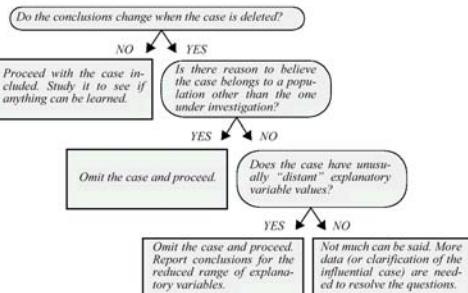


Slide 31 Identifying outliers that matter

NOTES:

Display 11.8

A strategy for dealing with suspected influential cases



Slide 32

NOTES:

Fit gastric model without apparent outliers (Cases 31 & 32)

Restrict analysis to GA < 3.1; Note that Cook's D is not that high. "Some statisticians use a rough guideline that a value of D_i close to or larger than 1 indicates a large influence." Draper & Smith: no test or cutoff for D_i

Display 11.9

Regression parameter estimates, standard errors, and p-values from the regression of first-pass metabolism on gastric activity, an indicator for female, an indicator for alcoholic, and all 2nd and 3rd-order interactions; with (1) all observations and (2) all observations except 31 and 32

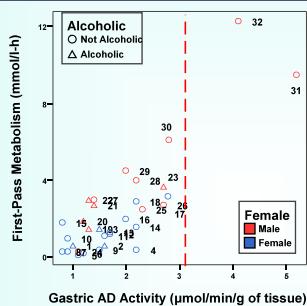
Variable	All 32 Observations			Cases 31 and 32 Removed		
	Estimate	Standard Error	2-sided p-value	Estimate	Standard Error	2-sided p-value
Constant	-1.660	1.000	.11	-0.680	1.309	.61
Gastric activity (G)	1.183	0.123	<.0001	1.901	0.168	.0045
Female (F)	1.466	1.333	.28	0.486	1.462	.74
Alcoholic (A)	2.552	1.946	.20	1.572	1.812	.40
G*F	-1.673	0.870	.011	-1.601	0.721	.15
F*A	-2.552	4.394	.77	-1.272	3.467	.52
G*A	-1.459	1.053	.18	-0.866	0.963	.38
G*F*A	1.199	2.998	.69	0.606	2.316	.80

EEOS611

Slide 33 Fit gastric model without apparent outliers (Cases 31 & 32)

NOTES:

Restrict range of regression to Gastric AD levels < 3.1 $\mu\text{mol}/\text{min/g}$



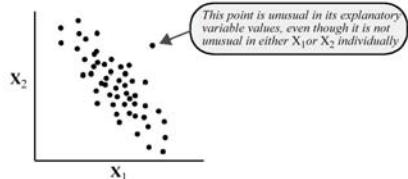
Slide 34 Restrict range of regression to Gastric AD levels < 3.1 $\mu\text{mol}/\text{min/g}$

NOTES:

Leverage

Calculated using only explanatory variables

An illustration of what is meant by "far from the average" of multiple explanatory variables when they are correlated



$1/n < \text{leverage} < p/n$
Potential cutoff leverage > $2p/n$

Slide 35 Leverage

NOTES:

Cook's D

Detecting single observations that matter: there is no critical value for Cook's D, but $D \geq 1$ of concern; D less than 1 could be important too.

- Equivalent to performing the regression after deleting each datum one at a time
- The actual cases need not be deleted
- The change in parameters with and without the deleted case are assessed
- A point with high leverage need not affect the regression greatly

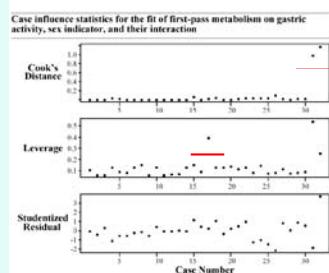
EEOS611

Slide 36 Cook's D

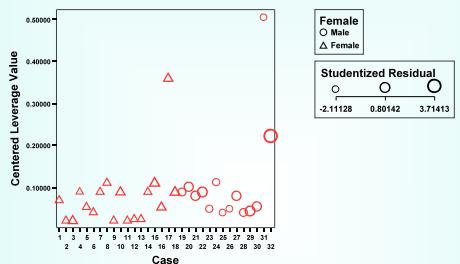
NOTES:

Display 11.11, p 318

Note that this plot is for the REDUCED model (no Display 11.11)

**Slide 37 Display 11.11, p 318**

NOTES:

SPSS analysis: Reduced modelGastric, Gastric x Female (no β_0 , no Female)

EEOS611

Slide 38 SPSS analysis: Reduced model

NOTES:

Model		Coefficients ^{a,b}						
		Unstandardized Coefficients		Standardized Coefficients		95% Confidence Interval for B		
		B	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound
1	(Constant)	-.573	.606		-.946	.352	-1.815	.668
	Gastric:AD Activity (mmol/min/g of tissue)	1.449	.339	.629	4.277	.000	.755	2.143
2	(Constant)	.845	.568		1.488	.148	-.320	2.011
	Gastric:AD Activity (mmol/min/g of tissue)	1.150	.271	.499	4.242	.000	.594	1.706
3	(Constant)	-.528	.345	-.521	-4.434	.000	-2.234	.821
	Female	-1.528	.345	-.521	-4.434	.000	-2.234	.821
4	(Constant)	.070	.802		.087	.932	-1.579	1.718
	Gastric:AD Activity (mmol/min/g of tissue)	1.565	.407	.679	3.843	.001	.728	2.403
	Female	-.267	.993	-.091	-.269	.790	-2.308	1.775
	Gastric:Female	.728	.539	-.449	-1.351	.188	-1.837	.380
	Alcoholic	1.572	1.812	.484	.868	.395	2.185	5.330
	Female:Alcohol	-1.272	3.467	-.266	-.367	.717	-8.462	5.918
	Gastric:Alcohol	-.869	.963	-.428	-.899	.378	-2.963	1.132
	Gastric:Female:Alc	.606	2.316	.177	.262	.796	-4.197	5.403

a. Dependent Variable: First-Pass Metabolism (mmol/L·h)
b. Selecting only cases for which Gastric LE 3 = 1

If there is an interaction, leave the appropriate 1st order effects in the model; Draper & Smith, Harrell (2001)

Slide 39

NOTES:

<p style="text-align: center;">Full model</p> <p>A single female alcoholic plays a big role in fitting models</p> <p>Case</p> <p>Cook's Distance</p> <p>Legend: Female (red circle), Male (blue square), Female (red triangle)</p> <p>Centered Leverage Value: 0.03545, 0.49706, 0.95867</p> <p>EEOS611</p>	<p>Slide 40 Full model</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p>First-Pass Metabolism (mmol/l-h) = $-1.83 + 2.28 \cdot \text{gastric}$ R-Square = 0.68</p> <p>Alcoholic (red square), Not Alcoholic (blue square), Alcoholic (red triangle)</p> <p>Female (red circle), Male (blue square), Female (red triangle)</p> <p>Cook's Distance Means: 0.00001, 1.13987, 2.27973</p> <p>Linear Regression</p> <p>First-Pass Metabolism (mmol/l-h)</p> <p>Gastric AD Activity (nmol/min/g of tissue)</p>	<p>Slide 41</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p>Display 11.8</p> <p>A strategy for dealing with suspected influential cases</p> <p>Do the conclusions change when the case is deleted?</p> <p>NO → YES</p> <p>Proceed with the case included. Study it to see if anything can be learned.</p> <p>Is there reason to believe the case belongs to a population other than the one under investigation?</p> <p>YES → NO</p> <p>Omit the case and proceed.</p> <p>Does the case have unusually "distant" explanatory variable values?</p> <p>YES → NO</p> <p>Omit the case and proceed. Report conclusions for the reduced range of explanatory variables.</p> <p>Not much can be said. More data (or clarification of the influential case) are needed to resolve the questions.</p>	<p>Slide 42</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>

Display 11.12

Least squares estimates for the regression of first-pass metabolism on gastric AD activity, sex, and their interaction (excluding cases 31 and 32)

Variable	Estimate	Standard Error	t-Statistic	2-sided p-value
Constant	0.070	0.802	0.087	.93
gast	1.565	0.407	3.843	.0007
fem	-0.267	0.993	-0.269	.79
gast × fem	-0.728	0.539	0.114	.91

Slide 43 Display 11.12

NOTES:

Sleuth's final gastric model: Force Y intercept to be zero, drop gender main effect & delete cases with GA>3.1

Draper & Smith, Harrell (2001): it is a mistake to leave out the β_0 and female main effect in this type of model! I would recommend leaving them in, even if not significant

In this model, first-pass metabolism is directly proportional to gastric activity, but the estimate of the slope is not zero as we would expect. The P-value for testing this model to the one where it is zero is 0.0007. The t-statistic is 3.843, with 26 degrees of freedom, so the smaller model is adequate.

The results are shown in Display 11.14. In this model, both terms are essential, which is accepted as the final version for reference.

The conclusions stated in the summary of statistical findings are based on Display 11.14. Notice in particular that, for any level of gastric activity, the mean is 0.8 mmol/m² h. At the lowest level of gastric activity, the mean is 1.565 mmol/m² h, which is 2.0 times the mean for males, even after accounting for gender. On the other hand,

Slide 44 Sleuth's final gastric model: Force Y intercept to be zero, drop gender main effect & delete cases with GA>3.1

NOTES:

Sleuth's final Case 11.1 model: an interaction, but without a Y intercept or Female main effect

Model		Coefficients ^{a,b,c}						
		Unstandardized Coefficients B	Standard Error	Standardized Coefficients Beta	t	Sig.	95% Confidence Interval for B Lower Bound	Upper Bound
1	Gastric AD Activity (mmol/m ² of tissue)	1.599	.125	1.218	12.800	.000	1.343	1.855
	GastricFemale	-.873	.174	-.478	-5.019	.000	-1.230	-.517

a. Dependent Variable: First-Pass Metabolism (mmol/l·h)

b. Linear Regression through the Origin

c. Selecting only cases for which Gastric LE 3 = 1

Draper & Smith, Harrell, & I object strongly to fitting an interaction term and leaving the main effect out. It should be left in. Moreover, the intercept should usually be left in the model, even if the 95% CI includes zero

Slide 45 Sleuth's final Case 11.1 model: an interaction, but without a Y intercept or Female main effect

NOTES:

Gallagher's alternate model

Gastric AD & Female main effects, then intnx ns.

Model		Coefficients ^a							
		Unstandardized Coefficients		Standardized Coefficients		t	Sig.	95% Confidence Interval for B	
		B	Std. Error	Beta	t	Sig.	Lower Bound	Upper Bound	
1	(Constant)	-573	.606	-.946	-.352	-1.815	.668		
	Gastric AD Activity	1.449	.339	.629	4.277	.000	.755	2.143	
	(muon/min/g of tissue)								
2	(Constant)	.845	.568	1.488	.148	.320	2.011		
	Gastric AD Activity	1.150	.271	.499	4.242	.000	.594	1.706	
	(muon/min/g of tissue)								
3	(Constant)	.070	.802	-.521	-4.434	.000	-2.234	-.821	
	Gastric AD Activity	1.565	.407	.679	3.843	.001	.728	2.403	
	(muon/min/g of tissue)								
	Female	-.267	.993	-.091	-.269	.790	-2.308	1.775	
	GastricxFemale	-.728	.539	-.449	-1.351	.188	-1.837	.380	
4	(Constant)	-.680	1.309	-.519	.809	.395	-3.395	2.035	
	Gastric AD Activity								
	(muon/min/g of tissue)								
	Female	1.921	.600	.833	3.159	.005	.660	3.183	
	GastricxFemale	1.041	.720	.720	1.431	.146	-2.576	4.15	
	Alcoholic	1.572	1.812	.484	.898	.385	-2.185	5.330	
	FemalexAlcohol	-1.272	3.467	-.266	-.367	.717	-6.462	5.918	
	GastricxAlcohol	-.869	.963	-.428	-.899	.378	-2.863	1.132	
	GastricxFemxAlc	.606	2.316	.177	.262	.796	-4.197	5.408	

a. Dependent Variable: First-Pass Metabolism (mmol/h)

b. Selecting only cases for which Gastric LE 3 = 1

Extra sum of squares F test

For the alternate model, including Y intercept & Females, then gastric x females coefficient not needed

Model Summary

Model	Gastric LE 3 = 1 (Selected)	Change Statistics							
		R Square	Adjusted R Square	Std. Error of the Estimate	R Square Change	F Change	df1	df2	Sig. F Change
1	.629^a	.395	.374	1.1557	.395	18.292	1	28	.000
2	.806^b	.650	.624	.8953	.255	19.659	1	27	.000
3	.820^c	.673	.635	.8819	.023	1.824	1	26	.188
4	.828^d	.685	.585	.9411	.012	.209	4	22	.931

a. Predictors: (Constant), Gastric AD Activity (muon/min/g of tissue)

b. Predictors: (Constant), Gastric AD Activity (muon/min/g of tissue), Female

c. Predictors: (Constant), Gastric AD Activity (muon/min/g of tissue), Female, GastricxFemale

d. Predictors: (Constant), Gastric AD Activity (muon/min/g of tissue), Female, GastricxFemale, GastxFemxAlc, GastricxAlcohol, Alcoholic, FemalexAlcohol

EEOS611

Display 11.14 vs. Simpler model

Results for mean metabolism being proportional to gastric activity

Variable	Estimate	Standard Error	t-Statistic	p-value
gast	1.599	0.125	12.800	<.0001
gast*x'fem	-0.873	0.174	-5.019	<.0001
<i>Coefficients^a</i>				
Unstandardized Coefficients				
	B	Std. Error	Beta	t
1	.070	.803	.087	.092
Gastric AD Activity	1.688	.407	.679	4.143
(muon/min/g of tissue)				
Female	-.207	.950	-.051	-.209
GastricxFemale	-.728	.539	-.449	-1.351
(Constant)	-.880	1.309	-.519	.809
Gastric AD Activity	1.921	.608	.833	3.159
(muon/min/g of tissue)				
Female	.496	1.467	.166	.331
GastricxFemale	-.1081	.721	-.096	-1.494
Alcoholic	1.572	1.812	.484	.898
FemalexAlcohol	-1.272	3.467	-.266	-.367
GastricxAlcohol	-.869	.963	-.428	-.899
GastricxFemxAlc	.606	2.316	.177	.262

a. Dependent Variable: First-Pass Metabolism (mmol4h)

b. Selecting only cases for which Gastric LE 3 = 1

Leave in model

No need for interaction terms

Slide 46 Gallagher's alternate model

NOTES:

Slide 47 Extra sum of squares F test

NOTES:

Slide 48 Display 11.14 vs. Simpler model

NOTES:

Case 11.2

Brain barrier disruption

EEOS611

Slide 49 Case 11.2

NOTES:

Case Study 11.2 brain disruption with concentrated sugars

Display 11.3 p. 294

Time line for blood-brain barrier disruption experiment

EEOS611

Slide 50 Case Study 11.2 brain disruption with concentrated sugars

NOTES:

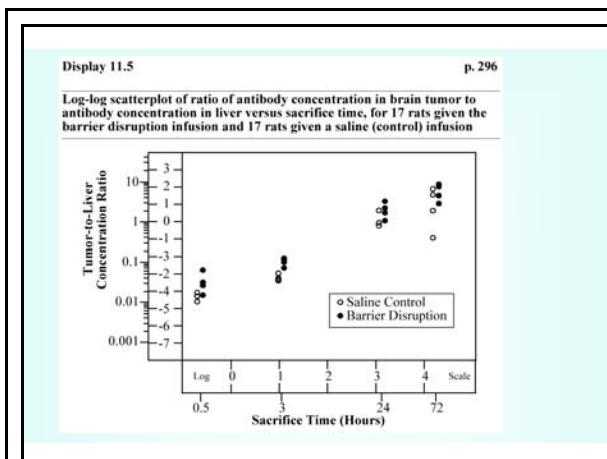
Display 11.4 p. 295

Response variable, design variables, and several covariates for 34 rats in the blood-brain barrier disruption experiment

Case	Response Variable Liver Tumor Count (spicules)	Sacrifice Time (hours)		Covariates	
		Treatment	Days Post Inoculation	Treatment Weight (10 ⁻³ grams)	Weight (g)
1	41097 / 45464a	0.5	10	F	225 4.0 271
2	44246 / 460171	0.5	10	M	200 4.0 246
3	10200 / 10200b	0.5	10	F	225 4.0 271
4	25427 / 2756411	0.5	10	F	184 9.8 168
5	12800 / 12800c	0.5	10	M	200 4.0 246
6	31142 / 296663	0.5	10	F	196 7.7 260
7	23800 / 23800d	0.5	10	M	200 4.0 246
8	16129 / 1432757	0.5	10	F	273 4.0 308
9	22340 / 22340e	0.5	10	M	200 4.0 246
10	77904 / 860057	3	BD	10	F 267 2.6 73
11	79104 / 860057	3	BD	10	F 278 2.6 73
12	76165 / 420145	3	BD	10	F 278 0.0 0.0
13	121009 / 1004812	3	BD	10	F 281 3.4 203
14	21000 / 2000406	3	NS	9	F 281 3.4 203
15	31803 / 3191532	3	NS	10	F 234 0.1 304
16	23200 / 23200f	3	NS	10	F 230 0.1 304
17	30545 / 961097	3	NS	10	F 230 7.0 146
18	30406 / 961097	3	NS	10	F 230 7.0 146
19	84616 / 48815	24	BD	10	F 254 3.9 151
20	55320 / 55320g	24	BD	10	M 200 3.9 150
21	48829 / 22305	24	BD	10	M 247 >2.4 101
22	84616 / 48815	24	BD	10	M 200 3.9 150
23	37928 / 26325	24	NS	10	F 237 2.5 224
24	12000 / 12000h	24	NS	10	M 249 4.4 151
25	23734 / 25895	24	NS	10	M 248 9.7 283
26	13000 / 13000i	24	NS	10	M 249 9.7 283
27	35309 / 41422	72	BD	11	F 251 4.1 39
28	32200 / 32200j	72	BD	10	M 249 4.1 39
29	9625 / 1979	72	BD	10	M 298 12.8 164
30	7490 / 7490k	72	BD	10	M 298 12.8 164
31	4250 / 928	72	NS	10	M 272 11.0 164
32	11310 / 2427	72	NS	10	M 249 4.4 164
33	3104 / 3668	72	NS	10	F 249 4.4 164
34	1334 / 3242	72	NS	10	F 249 4.4 164

Slide 51 Display 11.4

NOTES:



Slide 52 Display 11.5

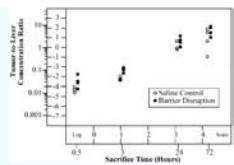
NOTES:

The tentative regression model

Days, sacrifice time, treated as categorical variables

This initial investigation suggests the following tentative regression model (using the shorthand model specification of Section 9.3.5):

$$\mu(\text{antibody} | \text{SAC}, \text{TREAT}, \text{DAYS}, \text{FEM}, \text{weight}, \text{loss}, \text{tumor}) \\ = \text{SAC} + \text{TREAT} + (\text{SAC} \times \text{TREAT}) + \text{DAYS} + \text{FEM} + \text{weight} + \text{loss} + \text{tumor},$$



EEOS611

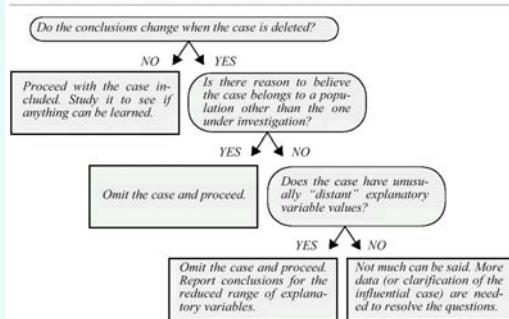
Slide 53 The tentative regression model

NOTES:

Display 11.8

p. 301

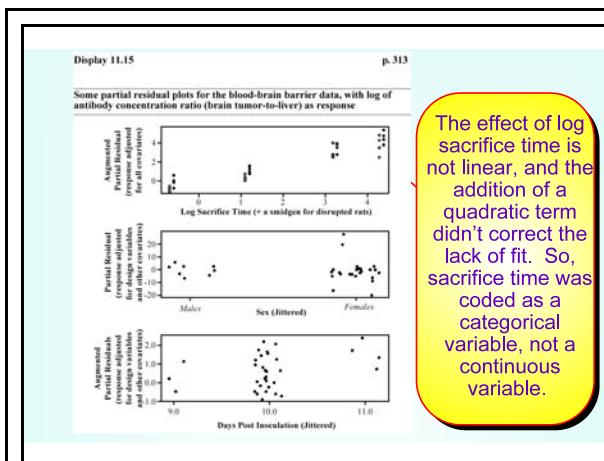
A strategy for dealing with suspected influential cases



Slide 54

NOTES:

<p>Brain Barrier data</p> <p>Studentized residuals: Not influential, so the two data (#31 & #34) left in the analysis</p> <p>Display 11.6</p> <p>Scatterplot of residuals versus fitted values from the fit of the logged response on a rich model for explanatory variables; brain barrier data</p> <p>EOS611</p>	<p>Slide 55 Brain Barrier data</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p>Case 11.02</p> <p>Performing the analysis described on Sleuth p 311.</p> <pre>REGRESSION /MISSING LISTWISE /STATISTICS COEFF CHA CI ANOVA /CRITERIA=PIN(.05) POUT(.10) /NOORIGIN /DEPENDENT lntum21v /METHOD=ENTER fem bd weight loss tumor day10 day11 /METHOD=ENTER fem bd weight loss tumor day10 day11 sac3h sac24h sac72h /METHOD=ENTER fem bd weight loss tumor day10 day11 sac3h sac24h sac72h bdx24 bdx72 /SCATTERPLOT=(*ZRESID ,*ZPRED) /SAVE PRED ZPRED COOK RESID ZRESID .</pre> <p>None of the cases seems to have sufficiently high leverage or Cook's D to be of concern; Some indication of 'trumpet' shaped residual plot</p> <p>Unstandardized Residual</p> <p>Unstandardized Predicted Value</p> <p>Cook's Distance Means 0.00000 0.34750 0.66901</p> <p>LLR Smoother</p>	<p>Slide 56 Case 11.02</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p>Case 11.02</p> <p>No need to delete cases for Case Study 11.02</p> <ul style="list-style-type: none"> Sleuth p. 317: <ul style="list-style-type: none"> p/n is the average leverage, where p is the number of parameters in the model 2*p/n is sometimes used as a lower cutoff for leverage For Full model for Case 11.2 <ul style="list-style-type: none"> 13 parameter full model 2*13/34=0.76 Only 2 data have leverage>0.6 and Max Cook's D is ≈ 0.7 <p>EEOS611</p>	<p>Slide 57 Case 11.02</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>



Slide 58 Display 11.15

NOTES:

Display 11.16 p. 314

Results from the regression of log ratio of antibody concentration (brain tumor-to-liver) on sacrifice time (treated as a factor) and treatment

Variable	Estimate	Standard Error	t-Statistic	2-sided p-value
Constant	-3.505	0.195	-17.94	<.0001
Indicator for time=3	1.341	0.252	5.30	.0001
Indicator for time=24	4.257	0.259	16.43	<.0001
Indicator for time=72	5.154	0.259	19.89	<.0001
Indicator for treatment=BD	0.797	0.183	4.35	.0002

Sacrifice time was coded for as a categorical variable, using 3 (0,1) indicator or 'dummy' variables. Note that changing 'time' to a categorical variable is often an appropriate solution for temporal 'lack of fit,' as shown previously with the Boston Harbor benthic diversity analyses

Slide 59 Display 11.16

NOTES:

Case 11.02 Final model

$e^{(0.422 \ 0.797 \ 1.172)} = [1.53 \ 2.22 \ 3.23]$

The median concentration of antibody was 2.2 times higher with sugar treatment (95% CI: 1.5, 3.2)

Model	Unstandardized Coefficients			Standardized Coefficients			Beta			95% Confidence Interval for B		
	B	Std. Errr	T	Beta	Std. Beta	t	Sig.	Lower Bound	Upper Bound	Lower Bound	Upper Bound	
1	(Constant)	4.002				21.910	4.299E-011	.422	1.172	3.800	4.200	
	Brain Disruption	.797	.183	.430	.422	4.346	1.545E-004	.619	1.650			
	3 hour sacrifice	1.134	.252	.450	.450	4.501	1.038E-004			3.733	4.337	
	24 hour sacrifice	5.341	.259	20.550	20.550	19.692	1.920E-015			5.640	5.984	
	72 hour sacrifice	5.154	.259	19.697	19.697	19.692	1.920E-015			5.654	5.984	
2	(Constant)	-4.930	3.098	-1.593	-1.593	-1.1320	1.459			-1.239	1.239	
	Brain Disruption	1.068	.188	.568	.568	5.125E-005	3.125E-005			1.068	1.239	
	3 hour sacrifice	1.099	.294	.317	.317	3.705	.001	.463	1.696			
	24 hour sacrifice	4.114	.337	.788	.788	12.198	8.895E-012	.3418	4.810			
	72 hour sacrifice	5.341	.341	15.941	15.941	15.940	9.803E-013			5.640	5.984	
	Days	.019	.282	.004	.004	.946	.946	-.053	.001			
	Female	.358	.358	.007	.007	.921	.921	-.703	.774			
	Weight	.002	.002	.017	.017	.754	.754	-.009	.011			
	Loss	-.048	.028	-.091	-.091	.094	.094	-.105	.009			
	Tumor	.001	.001	.083	.083	1.189	1.189	-.001	.004			

a. Dependent Variable: Ln(Tumor to liver)

R Square
Model 1 Change Statistics
Model 1: .301*
Model 2: .006^b
F change: 138.645
df1: 4
df2: 28
Sig. F Change: 1.671E-018

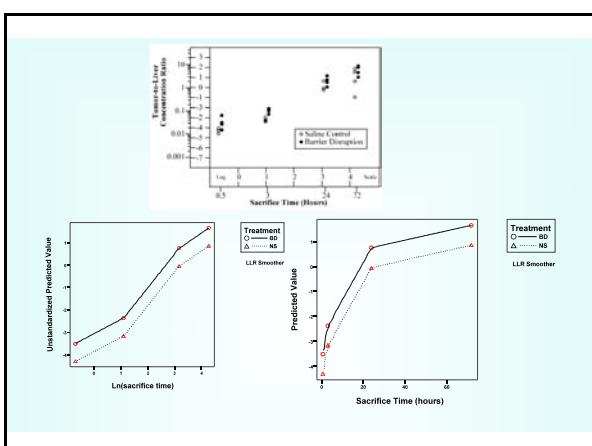
b. Predictors: (Constant), 72 hour sacrifice, Brain Disruption, 24 hour sacrifice, 3 hour sacrifice

c. Predictors: (Constant), 72 hour sacrifice, Brain Disruption, 24 hour sacrifice, 3 hour sacrifice, Loss, Weight, Tumor, Days, Female

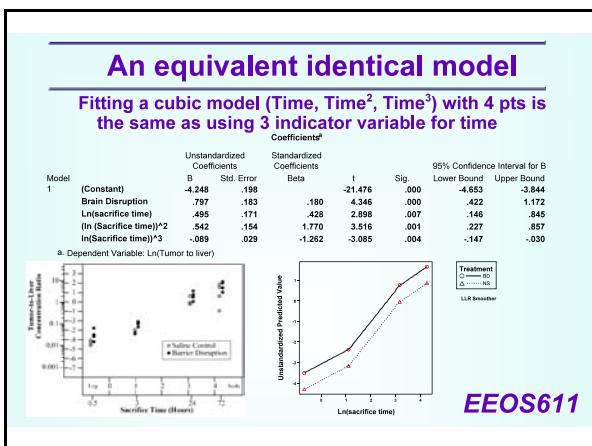
Model Summary^c

Slide 60 Case 11.02 Final model

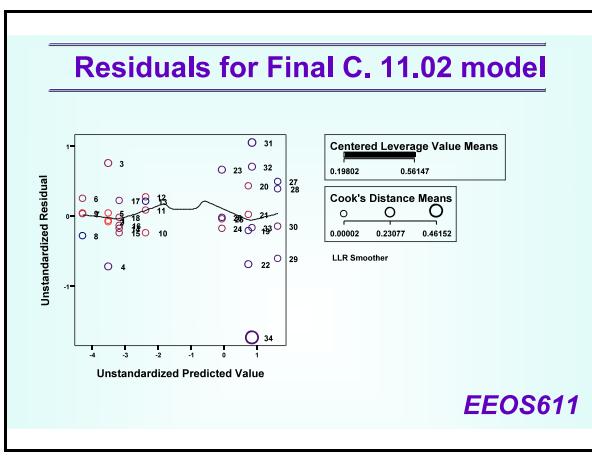
NOTES:

**Slide 61**

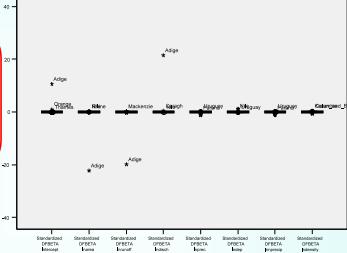
NOTES:

**EEOS611****Slide 62 An equivalent identical model**

NOTES:

**EEOS611****Slide 63 Residuals for Final C. 11.02 model**

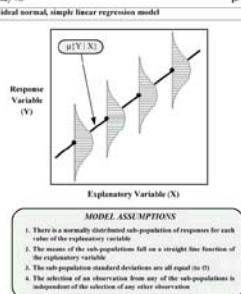
NOTES:

<p style="text-align: center;">SDFBeta</p> <p style="text-align: center;">Change in each parameter estimate caused by the deletion of each data point (in standard errors)</p> <p>Standardized DFBeta</p> $SDBETA_{ij} = \frac{b_j - b_j(i)}{s_{(i)} \sqrt{(\mathbf{X}^t \mathbf{W} \mathbf{X})^{-1}_{jj}}}$ <p>where $b_j - b_j(i)$ is the jth component of $\mathbf{b} - \mathbf{b}(i)$.</p> <p style="text-align: right;">EEOS611</p>	<p>Slide 64 SDFBeta</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p style="text-align: center;">SDFBeta plot, using explore</p> <p>\Data\analyze\explore\ ... These are SDBETA boxplots Using 11.21 River example</p> <p>Also, DFFITS, the change in predicted values with the deletion of each datum</p> 	<p>Slide 65 SDFBeta plot, using explore</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>
<p style="text-align: center;">11.6.1 Related issues: Weighted & Nonlinear regression</p> <ul style="list-style-type: none"> • Weighted Least squares regression available in SPSS <ul style="list-style-type: none"> ‣ /Analyze/regression/weight estimation, or WLS will provide an estimate of the weight function\ • Nonlinear regression is also available in SPSS • Use non-linear regression, available as an advanced regression option in SPSS <p style="text-align: right;">EEOS611</p>	<p>Slide 66 11.6.1 Related issues: Weighted & Nonlinear regression</p> <p>NOTES:</p> <hr/> <hr/> <hr/> <hr/> <hr/>

Review from Chapter 7

Assumptions

- **Linearity**
- **Constant variance**
 - Estimators still unbiased, but p values in err
- **Independence of errors**
 - Cluster & serial tests, Durbin-Watson tests
- **Normality of errors**
 - (not of explanatory variables)
 - Estimators still unbiased if normality assumption violated
 - P values robust to violations of normality
- **[X variable measured without error]**
 - This is an assumption involved in minimlmldn residuals



MODEL ASSUMPTIONS

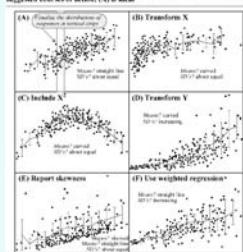
1. There is a normally distributed sub-population of responses for each value of the explanatory variable
2. The means of the sub-populations fall on a straight line function of the explanatory variable
3. The sub-population standard deviations are all equal (σ)
4. The selection of an observation from any of the sub-populations is independent of the selection of any other observation

Slide 67 Review from Chapter 7

NOTES:

When to use Weighted Least Squares or non-linear regression

Display 8.6
Some biometrical assumption of response versus explanatory variable with suggested courses of action: (A) is ideal



EEOS611

Slide 68 When to use Weighted Least Squares or non-linear regression

NOTES:

OLS vs. Generalized Least Squares

Draper & Smith (3rd ed., 1998 p. 224)

"If a generalized least squares analysis were called for [variance not constant] but an ordinary least squares analysis were performed, the estimates obtained would still be unbiased but would not have minimum variance, since the minimum variance estimators are obtained from the corrected generalized least squares analysis"

Slide 69 OLS vs. Generalized Least Squares

NOTES:

The Normal Equations & Matlab

$\mathbf{Y} = \mathbf{X} \boldsymbol{\theta} + \boldsymbol{\epsilon},$
 $E(\boldsymbol{\epsilon}) = \mathbf{0}.$
 $V(\boldsymbol{\epsilon}) = E(\boldsymbol{\epsilon} \boldsymbol{\epsilon}') = \sigma^2 \mathbf{I}.$
 LS method requires minimization of scalar sum of squares, S .
 $S = (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})' (\mathbf{y} - \mathbf{X} \boldsymbol{\theta})$
 To minimize S , set $\frac{\partial S}{\partial \boldsymbol{\theta}} = \mathbf{0}$.
 Differentiating, $2\mathbf{X}'(\mathbf{y} - \mathbf{X} \boldsymbol{\theta}) = \mathbf{0}$.
 $\boldsymbol{\theta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.

Slide 70 The Normal Equations & Matlab

NOTES:

Data with high variance are downweighted, with w_i being the scaling

Weight estimation(WLS) will find the appropriate w

Let us suppose that

$$\mathbf{V}\sigma^2 = \mathbf{V}(\mathbf{Y}) = \begin{bmatrix} 1/w_1 & & & \\ & 1/w_2 & & \\ & & \ddots & \\ & & & 1/w_n \end{bmatrix} \sigma^2,$$

where the w 's are weights to be specified. This means that

$$\mathbf{V}^{-1} = \begin{bmatrix} w_1 & & & \\ & w_2 & & \\ & & \ddots & \\ & & & w_n \end{bmatrix}.$$

Applying the general results above we find, after reduction,

$$\hat{\boldsymbol{\theta}} = \frac{\sum w_i X_i Y_i}{\sum w_i X_i^2},$$

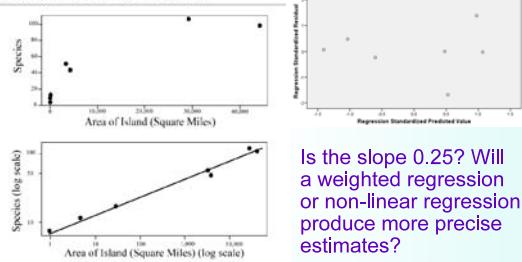
where all summations are from $i = 1, 2, \dots, n$.

EEOS611

Example of WLS & Non-linear regression, Case 8.1

Display 8.2

Scatterplot and log-log-scatterplot of number of reptile and amphibian species versus area for seven islands in the West Indies



Is the slope 0.25? Will a weighted regression or non-linear regression produce more precise estimates?

Slide 71 Data with high variance are downweighted, with w_i being the scaling

NOTES:

Slide 72 Example of WLS & Non-linear regression, Case 8.1

NOTES:

not have minimum variance, since the minimum variance estimates are obtained from the correct generalized least squares analysis. If standard least squares is used, then the estimates are obtained from $\hat{\mathbf{b}}_0 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$ and $E(\hat{\mathbf{b}}_0) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta = \beta$ but $\mathbf{V}(\hat{\mathbf{b}}_0) = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\mathbf{V}(\mathbf{Y})]\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}$ $= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{V}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\sigma^2.$ We recall from Eq. (9.2.13) that if the correct analysis is performed, $\mathbf{V}(\hat{\mathbf{b}}) = (\mathbf{X}'\mathbf{V}^{-1}\mathbf{X})^{-1}\sigma^2$ and, in general, elements of this matrix would provide smaller variances both for individual coefficients and for linear functions of the coefficients. Weighted least squares regression provides higher precision for estimating parameters (and greater power for testing models)	Slide 73 Draper & Smith NOTES: