

Statistics 516- Spring 2002 - Exam 1 (modified for Spring 2003 practice)

Part I: Answer the following three questions. First two are eight points each, the third is seven points.

- 1) Define what is meant by the p-value (or empirical significance level of a test).
- 2) In performing a linear regression to predict y from x, what four assumptions must be satisfied?

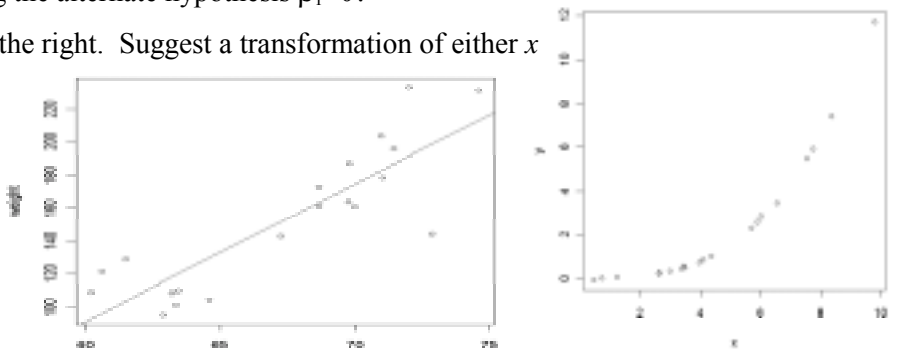
$$3) \sum_{i=1}^n (x_i - \bar{x})^2 = 22.0 \quad \sum_{i=1}^n (y_i - \bar{y})^2 = 21752.0 \quad \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = 236.5 \quad \bar{x} = 22.5 \quad \bar{y} = 332.0 \quad n = 7$$

Complete the Table	Source	SS	DF	MS	F	Prob>F
	Regression	_____	_____	_____	_____	0.4529
	Error	_____	5	3841.9250	_____	_____
	Total	_____	_____	_____	_____	_____

Part II: Answer eleven of the following twelve questions. If you answer more than eleven I will grade only the first eleven. Seven points each.

1) SAS produces a t-statistic of -1.87 and a p-value of 0.075 for testing the null hypothesis $\beta_1=0$ vs. the alternate hypothesis of $\beta_1 \neq 0$. What would the p-value be for testing the alternate hypothesis $\beta_1 < 0$?

2) Consider the scatter plot of x vs. y shown at the right. Suggest a transformation of either x or y that may allow you to perform linear regression on the data.



3) The regression line shown to the right is: $\text{height} = -407.771 + 8.320 \text{ weight}$. Even if all of the assumptions are met, why couldn't we use this regression equation to predict the weight of someone who is 55" tall?

4) In class we gave two situations where it was ok to remove an outlier from a data set. Name one of those two situations.

Questions 5-8 use the attached data set *Minnows* is from a study by Ryan, Hubert, Sprague, and Parrot. It concerned the effect of Zinc and Copper in the water on the amount of protein in minnow larvae. One unit of *Metals* equals a total Copper parts per million + 0.1 Zinc parts per million of one. One unit of *Protein* is a micro-gram. In each case the researchers controlled the level of the metals and randomly assigned the groups of minnow larvae to one of the contaminated tanks. Use the above description and the accompanying SAS output to answer the following questions.

- 5) For each 100 unit increase in the amount of metal, what is the predicted change in the amount of protein?
- 6) Is there a statistically significant relationship ($\alpha=0.05$) between the amount of metal and the amount of protein? How did you check this?
- 7) Give a 95% confidence interval for the range of protein values a new set of minnow larvae would be expected to fall in, if the metal amount was 0.
- 8) Comment on the regression assumptions for this problem. Which are met or not met, and how could you tell?

Questions 9-12 use the attached data set *Galapagos* from an article in the journal *Science* by Johnson and Raven. It concerns the number of native species on the various Galapagos Islands based on the number of non-native species (*NonNative*), the area of the island in km^2 (*Area*), the elevation of the island in m (*Elev*), the distance from the nearest other island in km (*DistNear*), and the distance from Santa Cruz in km (*DistSC*). Use the accompanying SAS output to answer the following questions. You may assume the assumptions for performing regression are met.

- 9) What percentage of the variation in the number of native species is explained by the five predictor variables?
- 10) Carefully state what null and alternate hypotheses go with the p-value of 0.0008 in the *Type I Tests* box. Do we accept or reject that null hypothesis?
- 11) Assuming the the regression with five variables fit, which simpler models are also acceptable? How could you tell this?
- 12) State which of the islands number of native species is predicted worst by the model and how you found it.

```
DATA Minnows;
INPUT Protein Metals;
CARDS;
201      0
186     37.5
173     75
110    112.5
115    150
202     37.5
161     75
172    112.5
138    150
133    187.5
204     75
165    112.5
148    150
143    187.5
123    225
188    112.5
172    150
157    187.5
115    225
108    262.5
133    150
125    187.5
184    225
135    262.5
114    300
;
```

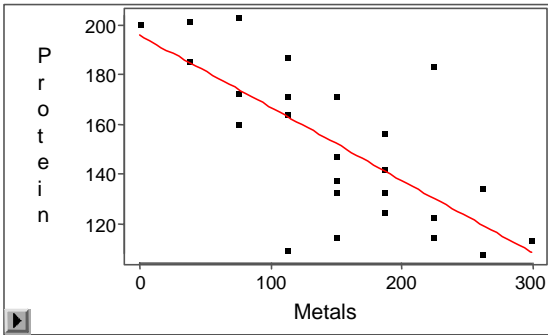
```
PROC INSIGHT;
OPEN Minnows;
FIT Protein = Metals;
RUN;
```

```
PROC GLM DATA=Minnows;
MODEL Protein = Metals / ALPHA=0.05 CLI;
RUN;
```

```
PROC GLM DATA=Minnows;
MODEL Protein = Metals / ALPHA=0.05 CLM;
RUN;
```

Protein = Metals
 Response Distribution: Normal
 Link Function: Identity

Model Equation
 Protein = 195.880 - 0.2912 Metals



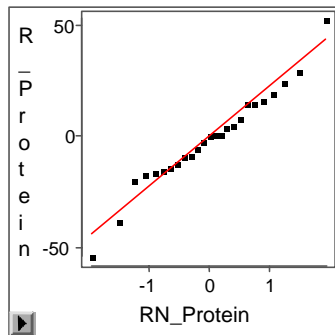
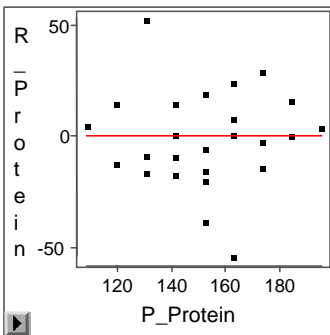
Parametric Regression Fit									
Curve	Degree(Polynomial)	DF	Model Mean Square	DF	Error Mean Square	R-Square	F Stat	Pr > F	
	1	1	11924.6400	23	496.8417	0.5106	24.00	<.0001	

Summary of Fit			
Mean of Response	152.2000	R-Square	0.5106
Root MSE	22.2899	Adj R-Sq	0.4894

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Stat	Pr > F
Model	1	11924.6400	11924.6400	24.00	<.0001
Error	23	11427.3600	496.8417		
C Total	24	23352.0000			

Type III Tests					
Source	DF	Sum of Squares	Mean Square	F Stat	Pr > F
Metals	1	11924.6400	11924.6400	24.00	<.0001

Parameter Estimates							
Variable	DF	Estimate	Std Error	t Stat	Pr > t	Tolerance	Var Inflation
Intercept	1	195.8800	9.9684	19.65	<.0001		0
Metals	1	-0.2912	0.0594	-4.90	<.0001	1.0000	1.0000



The GLM Procedure

Observati on	Observed	Predi cted	Resi dual	95% Confidence Li mi ts for	
				I ndi vi dual	Predi cted Val ue
1	201.0000000	195.8800000	5.1200000	145.36873153	246.39126847
2	186.0000000	184.9600000	1.0400000	135.94411215	233.97588785
3	173.0000000	174.0400000	-1.0400000	126.12080324	221.95919676
4	110.0000000	163.1200000	-53.1200000	115.87103480	210.36896520
5	115.0000000	152.2000000	-37.2000000	105.17656814	199.22343186
6	202.0000000	184.9600000	17.0400000	135.94411215	233.97588785
7	161.0000000	174.0400000	-13.0400000	126.12080324	221.95919676
8	172.0000000	163.1200000	8.8800000	115.87103480	210.36896520
9	138.0000000	152.2000000	-14.2000000	105.17656814	199.22343186
10	133.0000000	141.2800000	-8.2800000	94.03103480	188.52896520
11	204.0000000	174.0400000	29.9600000	126.12080324	221.95919676
12	165.0000000	163.1200000	1.8800000	115.87103480	210.36896520
13	148.0000000	152.2000000	-4.2000000	105.17656814	199.22343186
14	143.0000000	141.2800000	1.7200000	94.03103480	188.52896520
15	123.0000000	130.3600000	-7.3600000	82.44080324	178.27919676
16	188.0000000	163.1200000	24.8800000	115.87103480	210.36896520
17	172.0000000	152.2000000	19.8000000	105.17656814	199.22343186
18	157.0000000	141.2800000	15.7200000	94.03103480	188.52896520
19	115.0000000	130.3600000	-15.3600000	82.44080324	178.27919676
20	108.0000000	119.4400000	-11.4400000	70.42411215	168.45588785
21	133.0000000	152.2000000	-19.2000000	105.17656814	199.22343186
22	125.0000000	141.2800000	-16.2800000	94.03103480	188.52896520
23	184.0000000	130.3600000	53.6400000	82.44080324	178.27919676
24	135.0000000	119.4400000	15.5600000	70.42411215	168.45588785
25	114.0000000	108.5200000	5.4800000	58.00873153	159.03126847

Observati on	Observed	Predi cted	Resi dual	95% Confidence Li mi ts for	
				Mean Predi cted Val ue	
1	201.0000000	195.8800000	5.1200000	175.25886100	216.50113900
2	186.0000000	184.9600000	1.0400000	168.33470623	201.58529377
3	173.0000000	174.0400000	-1.0400000	160.99804656	187.08195344
4	110.0000000	163.1200000	-53.1200000	152.80943050	173.43056950
5	115.0000000	152.2000000	-37.2000000	142.97794628	161.42205372
6	202.0000000	184.9600000	17.0400000	168.33470623	201.58529377
7	161.0000000	174.0400000	-13.0400000	160.99804656	187.08195344
8	172.0000000	163.1200000	8.8800000	152.80943050	173.43056950
9	138.0000000	152.2000000	-14.2000000	142.97794628	161.42205372
10	133.0000000	141.2800000	-8.2800000	130.96943050	151.59056950
11	204.0000000	174.0400000	29.9600000	160.99804656	187.08195344
12	165.0000000	163.1200000	1.8800000	152.80943050	173.43056950
13	148.0000000	152.2000000	-4.2000000	142.97794628	161.42205372
14	143.0000000	141.2800000	1.7200000	130.96943050	151.59056950
15	123.0000000	130.3600000	-7.3600000	117.31804656	143.40195344
16	188.0000000	163.1200000	24.8800000	152.80943050	173.43056950
17	172.0000000	152.2000000	19.8000000	142.97794628	161.42205372
18	157.0000000	141.2800000	15.7200000	130.96943050	151.59056950
19	115.0000000	130.3600000	-15.3600000	117.31804656	143.40195344
20	108.0000000	119.4400000	-11.4400000	102.81470623	136.06529377
21	133.0000000	152.2000000	-19.2000000	142.97794628	161.42205372
22	125.0000000	141.2800000	-16.2800000	130.96943050	151.59056950
23	184.0000000	130.3600000	53.6400000	117.31804656	143.40195344
24	135.0000000	119.4400000	15.5600000	102.81470623	136.06529377
25	114.0000000	108.5200000	5.4800000	87.89886100	129.14113900

```

DATA Galapagos;
INPUT Native NonNative Are Elev DistNear DistSC Island $;
CARDS;
23 35 25.09 332 0.60 0.60 Baltra
21 10 1.24 109 0.60 26.30 Bartolome
3 0 0.21 114 2.80 58.70 Caldwell
9 16 0.10 46 1.90 47.40 Champion
1 1 1.05 130 1.90 1.90 Coamano
11 7 0.34 119 8.00 8.00 Daphne_Major
12 12 0.08 93 6.00 12.00 Daphne_Minor
7 3 2.33 168 34.10 290.20 Darwin
4 4 0.03 46 0.40 0.40 Eden
2 0 0.18 112 2.60 50.20 Enderby
26 71 58.27 198 1.10 88.30 Wspanola
35 58 634.49 1494 4.30 95.30 Fernandina
17 41 0.57 49 1.10 93.10 Gardner_(Esp)
4 1 0.78 227 4.60 62.20 Gardner_(SM)
19 21 17.35 76 47.40 92.20 Genovesa
89 258 4669.32 1707 0.70 28.10 Isabela
23 28 129.49 343 29.10 85.90 Marchena
2 0 0.01 25 3.30 45.90 Onslow
37 67 59.56 777 29.10 119.60 Pinta
33 75 17.95 458 10.70 10.70 Pinzon
9 3 0.23 84 0.50 0.60 Las_Plazas
30 40 4.89 367 4.40 24.40 Rabida
65 215 551.62 716 45.20 66.60 San_Cristobal
81 156 572.33 906 0.20 19.80 San_Salvador
95 349 903.82 864 0.60 0.00 Santa_Cruz
28 34 24.08 259 16.50 16.50 Santa_Fe
73 212 170.92 640 2.60 49.20 Santa_Maria
16 28 1.84 154 0.60 9.60 Seymour
8 8 1.24 186 6.80 50.90 Tortuga
12 9 2.85 253 34.10 254.70 Wolf
;

```

```

PROC INSIGHT;
OPEN Galapagos;
RUN;

```

```

PROC REG DATA=Galapagos;
MODEL Native = Area DistNear DistSC Elev NonNative /
SELECTION = RSQUARE ADJRSQ CP;
RUN;

```


▶ Native = Area DistNear DistSC Elev NonNative
 Response Distribution: Normal
 Link Function: Identity

▶ Model Equation
 Native = 6.7437 - 0.0025 Area + 0.0982 DistNear - 0.0265 DistSC
 + 0.0184 Elev + 0.2409 NonNative

▶ Summary of Fit
 Mean of Response 26.5000 R-Square 0.9428
 Root MSE 7.1070 Adj R-Sq 0.9308

▶ Analysis of Variance

Source	DF	Sum of Squares	Mean Square	F Stat	Pr > F
Model	5	19963.2724	3992.6545	79.05	<.0001
Error	24	1212.2276	50.5095		
C Total	29	21175.5000			

▶ Type I Tests

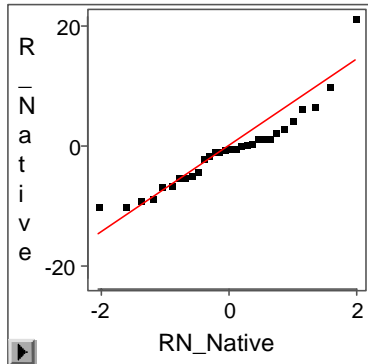
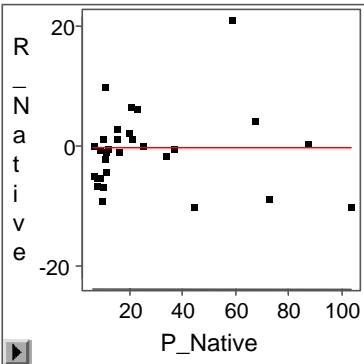
Source	DF	Sum of Squares	Mean Square	F Stat	Pr > F
Area	1	8123.3224	8123.3224	160.83	<.0001
DistNear	1	106.6735	106.6735	2.11	0.1591
DistSC	1	736.3862	736.3862	14.58	0.0008
Elev	1	5165.2902	5165.2902	102.26	<.0001
NonNative	1	5831.6001	5831.6001	115.46	<.0001

▶ Type III Tests

Source	DF	Sum of Squares	Mean Square	F Stat	Pr > F
Area	1	56.0507	56.0507	1.11	0.3026
DistNear	1	34.0055	34.0055	0.67	0.4200
DistSC	1	54.0551	54.0551	1.07	0.3112
Elev	1	554.4144	554.4144	10.98	0.0029
NonNative	1	5831.6001	5831.6001	115.46	<.0001

▶ Parameter Estimates

Variable	DF	Estimate	Std Error	t Stat	Pr > t	Tolerance	Var Inflation
Intercept	1	6.7437	2.1945	3.07	0.0052	.	0
Area	1	-0.0025	0.0024	-1.05	0.3026	0.4104	2.4365
DistNear	1	0.0982	0.1197	0.82	0.4200	0.5965	1.6765
DistSC	1	-0.0265	0.0256	-1.03	0.3112	0.5735	1.7437
Elev	1	0.0184	0.0056	3.31	0.0029	0.3174	3.1510
NonNative	1	0.2409	0.0224	10.75	<.0001	0.4422	2.2615



The REG Procedure
 Model: MODEL1
 Dependent Variable: Native

R-Square Selection Method

Number in Model	R-Square	Adjusted R-Square	C(p)	Variables in Model
1	0.9133	0.9102	10.3452	NonNative
1	0.6253	0.6119	131.0884	Elev
1	0.3836	0.3616	232.4104	Area
1	0.0276	-.0071	381.6643	DistSC
1	0.0000	-.0357	393.2369	DistNear

2	0.9372	0.9326	2.3191	Elev NonNative
2	0.9153	0.9090	11.5007	Area NonNative
2	0.9137	0.9073	12.1833	DistNear NonNative
2	0.9133	0.9069	12.3417	DistSC NonNative
2	0.6489	0.6229	123.2110	DistSC Elev
2	0.6265	0.5989	132.5713	Area Elev
2	0.6254	0.5977	133.0412	DistNear Elev
2	0.3945	0.3496	229.8550	Area DistSC
2	0.3887	0.3434	232.2985	Area DistNear
2	0.0450	-.0257	376.3713	DistNear DistSC

3	0.9401	0.9332	3.1207	Area Elev NonNative
3	0.9379	0.9308	4.0141	DistSC Elev NonNative
3	0.9376	0.9304	4.1650	DistNear Elev NonNative
3	0.9160	0.9063	13.2280	Area DistNear NonNative
3	0.9153	0.9056	13.4967	Area DistSC NonNative
3	0.9141	0.9041	14.0277	DistNear DistSC NonNative
3	0.6666	0.6282	117.7638	DistNear DistSC Elev
3	0.6491	0.6086	125.1227	Area DistSC Elev
3	0.6268	0.5837	134.4600	Area DistNear Elev
3	0.4234	0.3569	219.7193	Area DistNear DistSC

4	0.9411	0.9317	4.6733	Area DistSC Elev NonNative
4	0.9402	0.9306	5.0702	Area DistNear Elev NonNative
4	0.9401	0.9305	5.1097	DistNear DistSC Elev NonNative
4	0.9166	0.9032	14.9764	Area DistNear DistSC NonNative
4	0.6674	0.6141	119.4555	Area DistNear DistSC Elev

5	0.9428	0.9308	6.0000	Area DistNear DistSC Elev NonNative