

(b) Let

$$y_i = 2\sqrt{n} \sin^{-1} \sqrt{\hat{p}_i}.$$

Then we can fit the regression model  $y_i = \beta_0 + \beta_1 x_i$ . The transformation insures that the variance is stable.

(c) One example is in item response, where subjects are given a series of test questions. The proportion of people that answer correctly is related to the difficulty of the question, but since we are using proportions, this transformation can be used to stabilize the variance.

10.26 Since  $E(R) = \rho$  and

$$\text{Var}(R) \approx (1 - \rho^2)^2 = (1 - \mu^2)^2,$$

then

$$g(\mu) = (1 - \mu^2).$$

The appropriate transformation is

$$h(R) = \int \frac{dR}{(1 - R^2)} = \frac{1}{2} \log_e \left| \frac{R+1}{R-1} \right| = \frac{1}{2} \log_e \left( \frac{1+R}{1-R} \right) = \tanh^{-1}(R).$$

10.27 Since

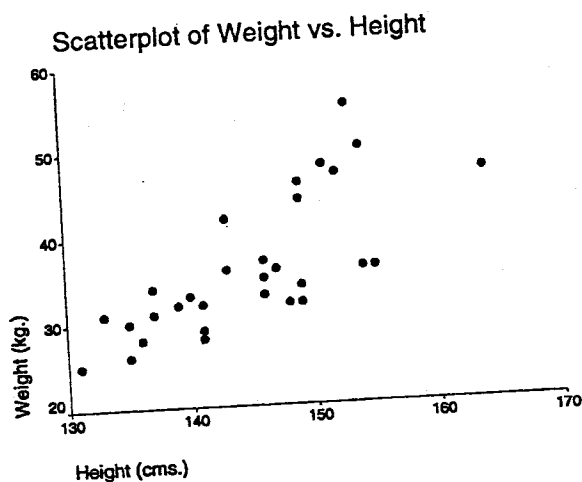
$$g(\mu) = \frac{1}{h'(\mu)} = \frac{1}{-1/\mu^2} = -\mu^2,$$

then

$$V(Y) = g^2(\mu) = \mu^4.$$

Solutions to Section 10.5

10.28 (a)



The strength of the linear relationship is moderate, so the correlation is probably around 0.6 to 0.8.

(b) Using  $\bar{x} = 144.8$  and  $\bar{y} = 36.167$ ,

$$\begin{aligned}
 S_{xy} &= \sum_i x_i y_i - n\bar{x}\bar{y} \\
 &= [135 \times 26 + \dots + 135 \times 30] - 30(144.8)(36.167) \\
 &= 1275. \\
 S_{xx} &= \sum_i x_i^2 - n\bar{x}^2 \\
 &= [135^2 + \dots + 135^2] - 30(144.8)^2 \\
 &= 1716.8. \\
 S_{yy} &= \sum_i y_i^2 - n\bar{y}^2 \\
 &= 26^2 + \dots + 30^2 - 30(36.167)^2 \\
 &= 1718.167.
 \end{aligned}$$

Then the sample correlation is

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \frac{1275}{\sqrt{(1716.8)(1718.167)}} = 0.742.$$

Testing  $H_0: \rho = 0.7$  vs.  $H_1: \rho > 0.7$  is equivalent to testing

$$H_0: \psi = \frac{1}{2} \log_e \left( \frac{1+0.7}{1-0.7} \right) = 0.867 \text{ vs. } H_1: \psi > 0.867.$$

Since

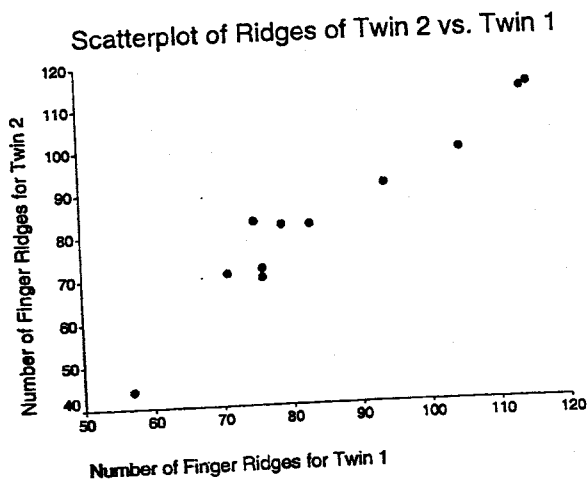
$$\hat{\psi} = \frac{1}{2} \log_e \left( \frac{1+r}{1-r} \right) = \frac{1}{2} \log_e \left( \frac{1+0.742}{1-0.742} \right) = 0.955,$$

the test statistic is

$$z = \sqrt{n-3}(\hat{\psi} - \psi_0) = \sqrt{30-3}(0.955 - 0.867) = 0.457.$$

The  $P$ -value is 0.323, leading to the conclusion that  $\rho$  is not significantly higher than 0.7.

10.29 (a)



This plot shows a very strong linear relationship, so the correlation is probably close to 1.

(b) From Minitab,  $r = 0.971$ . To find a 95% CI for  $\rho$  we must first find a 95% CI for  $\psi$ .

$$\hat{\psi} = \frac{1}{2} \log_e \left( \frac{1+r}{1-r} \right) = \frac{1}{2} \log_e \left( \frac{1+0.971}{1-0.971} \right) = 2.110.$$

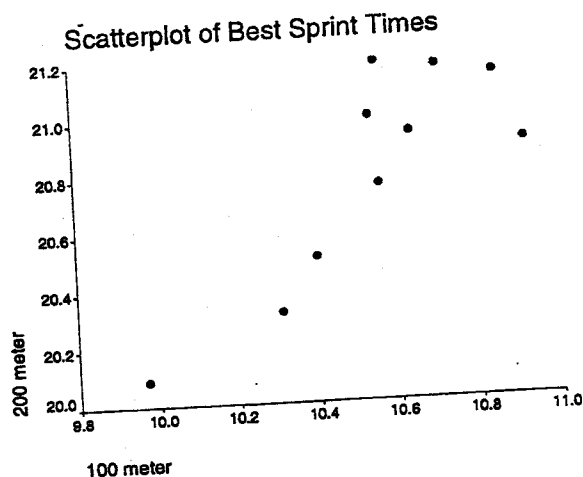
Then a 95% CI for  $\psi$  is

$$\hat{\psi} \pm z_{0.025} \frac{1}{\sqrt{n-3}} = 2.110 \pm 1.96 \frac{1}{\sqrt{12-3}} = [1.456, 2.763].$$

Then a 95% CI for  $\rho$  is

$$\left[ \frac{e^{2l} - 1}{e^{2l} + 1}, \frac{e^{2u} - 1}{e^{2u} + 1} \right] = \left[ \frac{e^{2(1.456)} - 1}{e^{2(1.456)} + 1}, \frac{e^{2(2.763)} - 1}{e^{2(2.763)} + 1} \right] = [0.897, 0.992].$$

10.30 (a)



This plot shows a moderately strong linear relationship, so the correlation is probably around 0.6 to 0.8.

(b) From Minitab,  $r = 0.836$ . To find a 95% CI for  $\rho$  we must first find a 95% CI for  $\psi$ .

$$\hat{\psi} = \frac{1}{2} \log_e \left( \frac{1+r}{1-r} \right) = \frac{1}{2} \log_e \left( \frac{1+0.836}{1-0.836} \right) = 1.208.$$

Then a 95% CI for  $\psi$  is

$$\hat{\psi} \pm z_{0.025} \frac{1}{\sqrt{n-3}} = 1.208 \pm 1.96 \frac{1}{\sqrt{10-3}} = [0.467, 1.949].$$

Then a 95% CI for  $\rho$  is

$$\left[ \frac{e^{2l} - 1}{e^{2l} + 1}, \frac{e^{2u} - 1}{e^{2u} + 1} \right] = \left[ \frac{e^{2(0.467)} - 1}{e^{2(0.467)} + 1}, \frac{e^{2(1.949)} - 1}{e^{2(1.949)} + 1} \right] = [0.436, 0.960].$$

Testing  $H_0 : \rho = 0.5$  vs.  $H_1 : \rho > 0.5$  is equivalent to testing

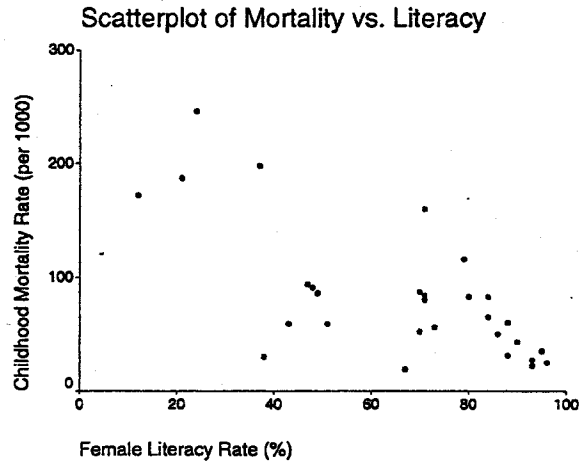
$$H_0 : \psi = \frac{1}{2} \log_e \left( \frac{1+0.5}{1-0.5} \right) = 0.549 \text{ vs. } H_1 : \psi > 0.549.$$

The test statistic is

$$z = \sqrt{n-3}(\hat{\psi} - \psi_0) = \sqrt{10-3}(1.208 - 0.549) = 1.744.$$

The  $P$ -value is  $0.041 < \alpha = 0.05$ , so conclude that  $\rho$  is significantly higher than 0.5.

10.31 (a)



This plot shows a moderately strong negative linear relationship, so the correlation is probably around  $-0.6$  to  $-0.8$ .

(b) From Minitab,  $r = -0.702$ . To find a 95% CI for  $\rho$ , we must first find a 95% CI for  $\psi$ .

$$\hat{\psi} = \frac{1}{2} \log_e \left( \frac{1+r}{1-r} \right) = \frac{1}{2} \log_e \left( \frac{1-0.702}{1+0.702} \right) = -0.871.$$

Then a 95% CI for  $\psi$  is

$$\hat{\psi} \pm z_{0.025} \frac{1}{\sqrt{n-3}} = -0.871 \pm 1.96 \frac{1}{\sqrt{30-3}} = [-1.248, -0.494].$$

Then a 95% CI for  $\rho$  is

$$\left[ \frac{e^{2l} - 1}{e^{2l} + 1}, \frac{e^{2u} - 1}{e^{2u} + 1} \right] = \left[ \frac{e^{2(-1.248)} - 1}{e^{2(-1.248)} + 1}, \frac{e^{2(-0.494)} - 1}{e^{2(-0.494)} + 1} \right] = [-0.848, -0.457].$$

Testing  $H_0 : \rho = -0.7$  vs.  $H_1 : \rho < -0.7$  is equivalent to testing

$$H_0 : \psi = \frac{1}{2} \log_e \left( \frac{1-0.7}{1+0.7} \right) = -0.867 \text{ vs. } H_1 : \psi < -0.867.$$

The test statistic is

$$z = \sqrt{n-3}(\hat{\psi} - \psi_0) = \sqrt{30-3}(-0.871 + 0.867) = -0.020.$$

The  $P$ -value is 0.492, leading to the conclusion that  $\rho$  is not significantly higher than 0.7.

10.32

Let

$$u_i = ax_i + b \text{ and } v_i = cy_i + d.$$

Then

$$\bar{u} = a\bar{x} + b \text{ and } \bar{v} = c\bar{y} + d.$$

Then

$$\begin{aligned}\sum(u_i - \bar{u})(v_i - \bar{v}) &= \sum(ax_i + b - a\bar{x} - b)(cy_i + d - c\bar{y} - d) \\ &= ac \sum(x_i - \bar{x})(y_i - \bar{y}), \\ \sum(u_i - \bar{u})^2 &= \sum(ax_i + b - a\bar{x} - b)^2 \\ &= a^2 \sum(x_i - \bar{x})^2, \\ \sum(v_i - \bar{v})^2 &= \sum(cy_i + d - c\bar{y} - d)^2 \\ &= c^2 \sum(y_i - \bar{y})^2.\end{aligned}$$

Then

$$\begin{aligned}r_{uv} &= \frac{\sum(u_i - \bar{u})(v_i - \bar{v})}{\sqrt{\sum(u_i - \bar{u})^2 \sum(v_i - \bar{v})^2}} \\ &= \frac{ac \sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{a^2 \sum(x_i - \bar{x})^2 c^2 \sum(y_i - \bar{y})^2}} \\ &= r_{xy}.\end{aligned}$$

### Solutions to Chapter 10 Advanced Exercises

10.33 Let

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i = \bar{y} + \hat{\beta}_1(x_i - \bar{x}).$$

Then

$$\begin{aligned}\sum(y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum(y_i - \bar{y} - \hat{\beta}_1(x_i - \bar{x}))(\bar{y} + \hat{\beta}_1(x_i - \bar{x}) - \bar{y}) \\ &= \sum(y_i - \bar{y} - \hat{\beta}_1(x_i - \bar{x}))(\hat{\beta}_1(x_i - \bar{x})) \\ &= \hat{\beta}_1 \sum(y_i - \bar{y})(x_i - \bar{x}) - \hat{\beta}_1^2 \sum(x_i - \bar{x})^2 \\ &= \hat{\beta}_1 S_{xy} - \hat{\beta}_1^2 S_{xx} \\ &= \frac{S_{xy}^2}{S_{xx}} - \frac{S_{xy}^2}{S_{xx}^2} S_{xx} = 0.\end{aligned}$$

10.34 (a) The model is

$$Y_i = \beta_0 + \beta_1 x_i + e_i, \text{ for } i = 1, \dots, n_1 + n_2.$$

If an observation is drawn from population 2, then  $Y_i = \beta_0 + \epsilon_i$ , so that  $\beta_0 = \mu_2$ . Similarly, if an observation is drawn from population 1, then  $Y_i = \beta_0 + \beta_1 + \epsilon_i$ , so that  $\beta_0 + \beta_1 = \mu_1$ , or  $\beta_1 = \mu_1 - \mu_2$ .

(b) Let  $n_1 + n_2 = n$ . Then

$$\begin{aligned}\hat{\beta}_1 &= \frac{\sum x_i y_i - n \bar{x} \bar{y}}{\sum x_i^2 - n \bar{x}^2} \\ &= \frac{n_1 \bar{y}_1 - n \left(\frac{n_1}{n}\right) \left(\frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n}\right)}{n_1 - n(n_1/n)^2} \\ &= \frac{nn_1 \bar{y}_1 - n_1^2 \bar{y}_1 - n_1 n_2 \bar{y}_2}{nn_1 - n_1^2} \\ &= \bar{y}_1 - \bar{y}_2,\end{aligned}$$

and

$$\begin{aligned}\hat{\beta}_0 &= \bar{y} - \hat{\beta}_1 \bar{x} \\ &= \frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n} - (\bar{y}_1 - \bar{y}_2) \frac{n_1}{n} \\ &= \bar{y}_2.\end{aligned}$$

(c)

$$SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^{n_1} (y_i - \bar{y}_1)^2 + \sum_{i=n_1+1}^{n_2} (y_i - \bar{y}_2)^2,$$

and so

$$MSE = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2}{n - 2},$$

with  $n - 2 = n_1 + n_2 - 2$  d.f.

(d)

$$\begin{aligned}t &= \frac{\hat{\beta}_1}{s/\sqrt{S_{xx}}} = \frac{\bar{y}_1 - \bar{y}_2}{s_p/\sqrt{\frac{n_1 n_2}{n}}} \\ &= \frac{\bar{y}_1 - \bar{y}_2}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}.\end{aligned}$$

10.35 (a)

$$E(s^2) = \frac{\sum_{i=1}^k (n_i - 1)E(s_i^2)}{n - k} = \frac{\sigma^2 \sum_{i=1}^k (n_i - 1)}{n - k} = \sigma^2.$$

So  $s^2$  is an unbiased estimate of the common variance  $\sigma^2$ .

(b)

$$\begin{aligned}\sum_i \sum_j (y_{ij} - \hat{y}_{ij})^2 &= \sum_i \sum_j ((y_{ij} - \bar{y}_i) + (\bar{y}_i - \hat{y}_{ij}))^2 \\ &= \sum_i \sum_j (y_{ij} - \bar{y}_i)^2 + \sum_i \sum_j (\bar{y}_i - \hat{y}_{ij})^2 \\ &\quad + 2 \sum_i \sum_j (y_{ij} - \bar{y}_i)(\bar{y}_i - \hat{y}_{ij}) \\ &= \sum_i \sum_j (y_{ij} - \bar{y}_i)^2 + \sum_i n_i (\bar{y}_i - \hat{y}_{ij})^2 \\ &\quad + 2 \sum_i \left( (\bar{y}_i - \hat{y}_{ij}) \sum_j (y_{ij} - \bar{y}_i) \right)\end{aligned}$$

## Chapter 11 Solutions

### Solutions to Section 11.2

11.1

$$Q = \sum (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2)^2.$$

To minimize this, set the partial derivatives equal to 0,

$$\begin{aligned} \frac{\partial Q}{\partial \beta_0} &= -2 \sum (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2) = 0, \\ \frac{\partial Q}{\partial \beta_1} &= -2 \sum x_i (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2) = 0, \\ \frac{\partial Q}{\partial \beta_2} &= -2 \sum x_i^2 (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2) = 0. \end{aligned}$$

From the first equation,

$$\sum y_i = n\beta_0 + \beta_1 \sum x_i + \beta_2 \sum x_i^2.$$

From the second equation,

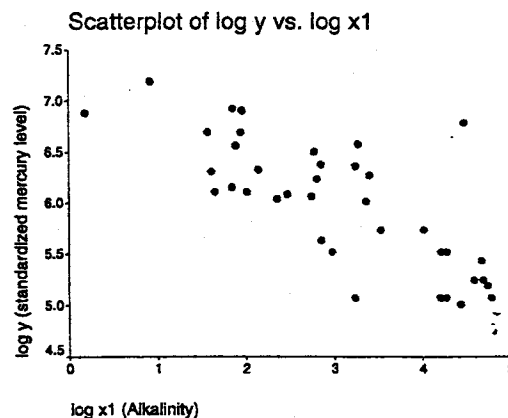
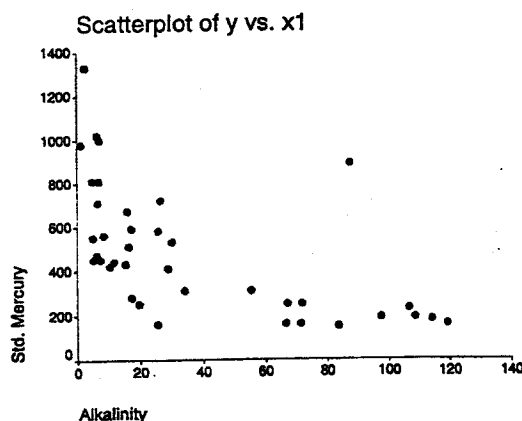
$$\sum x_i y_i = \beta_0 \sum x_i + \beta_1 \sum x_i^2 + \beta_2 \sum x_i^3.$$

From the third equation,

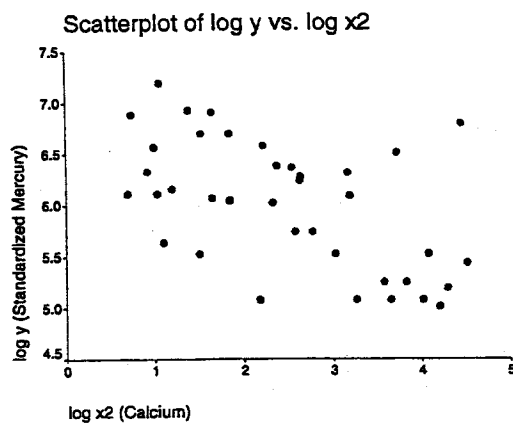
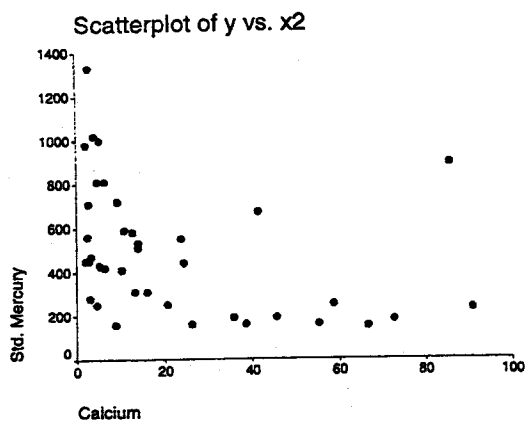
$$\sum x_i^2 y_i = \beta_0 \sum x_i^2 + \beta_1 \sum x_i^3 + \beta_2 \sum x_i^4.$$

These are the normal equations.

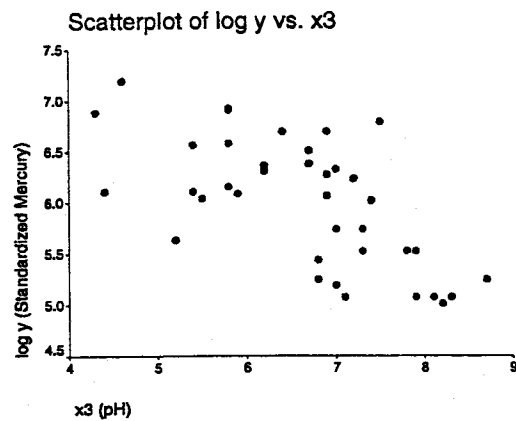
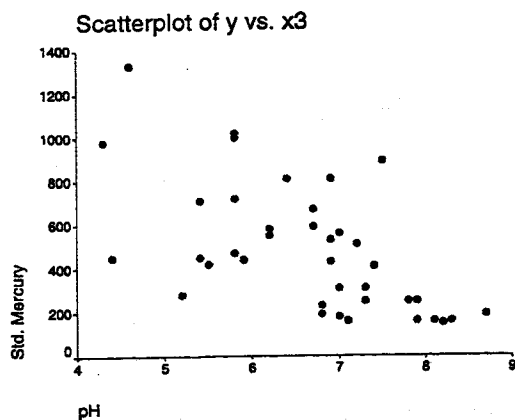
- 11.2 The fitted model is  $\hat{y} = -1.571 + 0.02573 \text{ Verbal} + 0.03361 \text{ Math}$ .  $r^2 = 0.681$ , so 68.1% of the variability in GPA is accounted for by math and verbal scores.
- 11.3 The fitted model is  $\hat{y} = 111.354 + 2.060x_1 - 2.732x_2 + 0.000x_3$ .  $r^2 = 0.295$ , so 29.5% of the variability in PIQ is accounted for by the brain size, height, and weight of a person.
- 11.4 (a)



The log transformation of both  $y$  and  $x_1$  yields an approximately linear relationship.



Similarly, the log transformation of both  $y$  and  $x_2$  yields an approximately linear relationship.



The plot of  $y$  vs.  $x_3$  appears linear to begin with. It remains linear after the transformation.

- (b) The fitted model is  $\log \hat{y} = 7.560 - 0.459 \log x_1 + 0.147 \log x_2 - 0.08x_3$ .  $r^2 = 0.607$ , so 60.7% of the variability in mercury is explained by this model.

### Solutions to Section 11.3

#### 11.5



- (e) If there are  $n$  patients in each group, then  $X$  would have the same 4 rows, but repeated  $n$  times. Then

$$X'X = \begin{bmatrix} 4n & 0 & 0 & 0 \\ 0 & 4n & 0 & 0 \\ 0 & 0 & 4n & 0 \\ 0 & 0 & 0 & 4n \end{bmatrix},$$

$$(X'X)^{-1} = \begin{bmatrix} 1/4n & 0 & 0 & 0 \\ 0 & 1/4n & 0 & 0 \\ 0 & 0 & 1/4n & 0 \\ 0 & 0 & 0 & 1/4n \end{bmatrix},$$

and

$$\hat{\beta} = (X'X)^{-1}X'y = \frac{1}{4} \begin{bmatrix} \bar{y}_1 + \bar{y}_2 + \bar{y}_3 + \bar{y}_4 \\ -\bar{y}_1 - \bar{y}_2 + \bar{y}_3 + \bar{y}_4 \\ -\bar{y}_1 + \bar{y}_2 - \bar{y}_3 + \bar{y}_4 \\ \bar{y}_1 - \bar{y}_2 - \bar{y}_3 + \bar{y}_4 \end{bmatrix},$$

where  $\bar{y}_i$  is the sample mean for the  $i$ th group. The error d.f. would now be  $4n - 4 = 4(n - 1)$ .

#### Solutions to Section 11.4

##### 11.11

Analysis of Variance				
Source	SS	d.f.	MS	F
Regression	37.70	3	12.567	2.084
Error	180.90	30	6.03	
Total	218.60	33		

Since  $F < f_{3,30,0.05} = 2.922$ , do not reject  $H_0$  and conclude that the regression is not significant.

- 11.12 A 95% CI for  $\beta_1$  is given by

$$\hat{\beta}_1 \pm t_{22,0.025}SE(\hat{\beta}_1) = 0.11 \pm 2.074 \times 0.055 = [-0.004, 0.224].$$

A 95% CI for  $\beta_2$  is given by

$$\hat{\beta}_2 \pm t_{22,0.025}SE(\hat{\beta}_2) = 1.40 \pm 2.074 \times 0.64 = [0.073, 2.727].$$

$x_2$  should be kept in the model since the CI for  $\beta_2$  is entirely above 0. Since the CI for  $\beta_1$  contains 0,  $\beta_1 = 0$  is plausible, and  $x_1$  could be removed from the model.

##### 11.13

$$t_1 = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)} = \frac{0.06}{0.05} = 1.2,$$

$$t_2 = \frac{\hat{\beta}_2}{SE(\hat{\beta}_2)} = \frac{1.84}{0.89} = 2.067,$$

$$t_3 = \frac{\hat{\beta}_3}{SE(\hat{\beta}_3)} = \frac{0.65}{0.11} = 5.909.$$

Math            0.033615      0.004928      6.82

S = 0.4023      R-Sq = 68.1%      R-Sq(adj) = ~~66.4%~~

Analysis of Variance

Source	DF	SS	MS	F
Regression	2	12.7859	6.3930	<del>39.51</del>
Residual Error	37	5.9876	0.1618	
Total	39	18.7735		

Source	DF	Seq SS
Verbal	1	5.2549
Math	1	7.5311

Then a 95% CI for the Verbal coefficient is given by

$$\hat{\beta}_1 \pm t_{37,0.025}SE(\hat{\beta}_1) = 0.026 \pm 2.021 \times 0.004 = [0.018, 0.034]$$

Similarly, a 95% CI for the Math coefficient is given by

$$\hat{\beta}_2 \pm t_{37,0.025}SE(\hat{\beta}_2) = 0.034 \pm 2.021 \times 0.005 = [0.024, 0.044]$$

11.17 The test statistic is

$$F = \frac{(SSE_{\text{linear}} - SSE_{\text{quad}})/3}{SSE_{\text{quad}}/(40 - 6)} = \frac{(5.9876 - 1.1908)/3}{1.1908/34} = 45.65$$

Since  $F > f_{3,34,0.05} = 2.92$ , we reject  $H_0$  and conclude that the quadratic fit is **significantly** better.

11.18 (a) The regression output is shown below:

Regression Analysis

The regression equation is

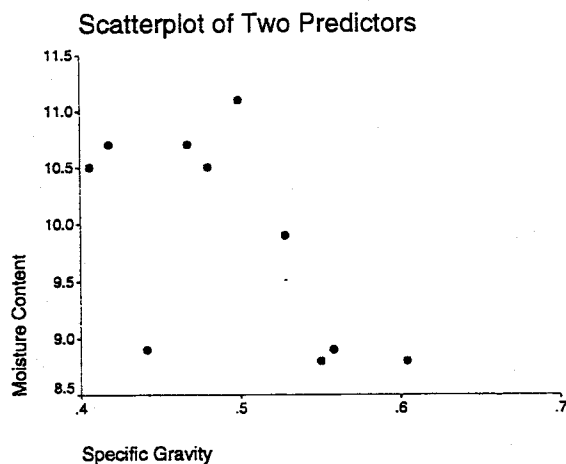
$$PIQ = 111 + 2.06 \text{ MRI} - 2.73 \text{ Height} + 0.001 \text{ Weight}$$

Predictor	Coef	StDev	T	P
Constant	111.35	62.97	1.77	0.086
MRI	2.0604	0.5634	3.66	0.001
Height	-2.732	1.229	-2.22	0.033
Weight	0.0006	0.1971	0.00	0.998

S = 19.79      R-Sq = 29.5%      R-Sq(adj) = 23.3%

Analysis of Variance

(a) — Prob 11,22



From the above plot, beam number 4 is far away from the other data points and appears to be influential.

(b)

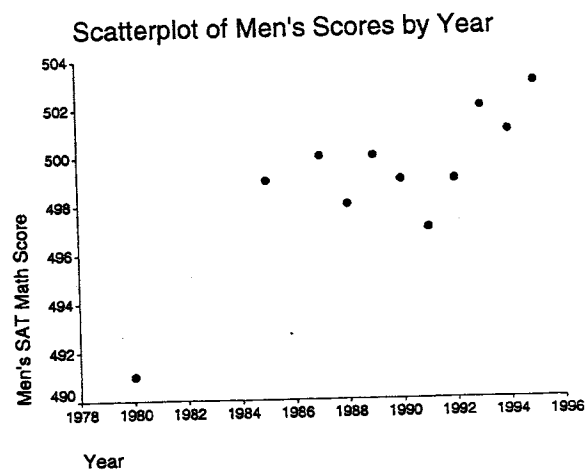
Beam number	$h_{ii}$
1	0.418
2	0.242
3	0.417
4	0.604
5	0.252
6	0.148
7	0.262
8	0.154
9	0.316
10	0.187

0.6

An observation is identified as influential if  $h_{ii} > 2(k + 1)/n = 2(2 + 1)/10 = 0.4$ . Only beam number 4 is influential, which is consistent with our graphical conclusion in (a).

- (c) The LS line using the influential observation is  $\hat{y} = 10.3 + 8.49 \text{ Gravity} - 0.266 \text{ Moisture}$ .  
 The LS line excluding the influential observation is  $\hat{y} = 12.4 + 6.80 \text{ Gravity} - 0.391 \text{ Moisture}$ .  
 Excluding this influential observation dramatically alters the coefficients of the model terms. It would be better to remove the influential observation when fitting the model.

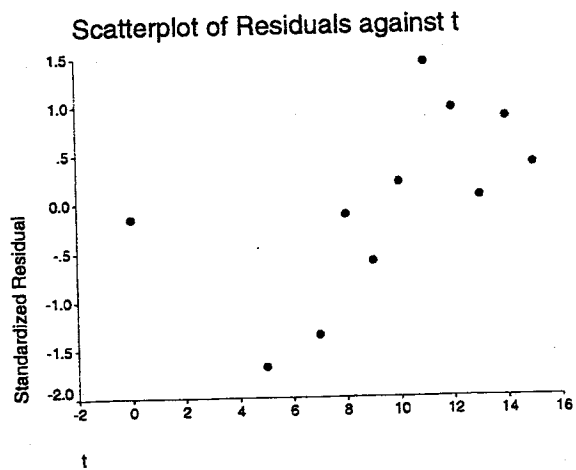
11.23 (a)



This graph indicates an increasing trend in the mens' SAT scores.

(b)  $\hat{y} = 493.399 + 0.592t$ .

(c)



The residuals tend to be negative for lower values of  $t$  and positive for higher values of  $t$ .

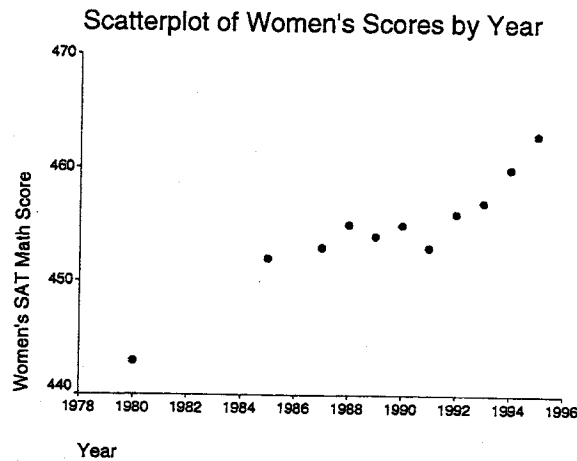
(d)

$t_i$	$e_i^*$
0	-0.157
5	-1.683
7	-1.354
8	-0.123
9	-0.599
10	0.206
11	1.438
12	0.962
13	0.059
14	0.864
15	0.387

Since none of the standardized residuals are  $> 2$ , there do not appear to be any outliers.

- (e) An observation is influential if  $h_{ii} > 2(1 + 1)/11 = 0.364$ . Only  $t = 0$  (1980) satisfies this condition. It is influential because it is much farther to the left than the other observations, according to the plot from (a).

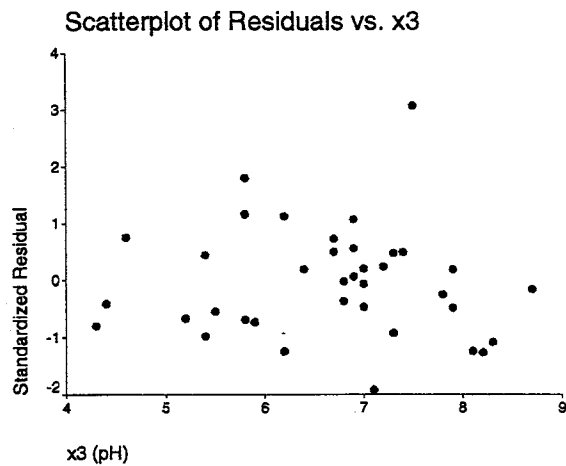
11.24 (a)



This graph indicates an increasing trend in the womens' SAT scores.

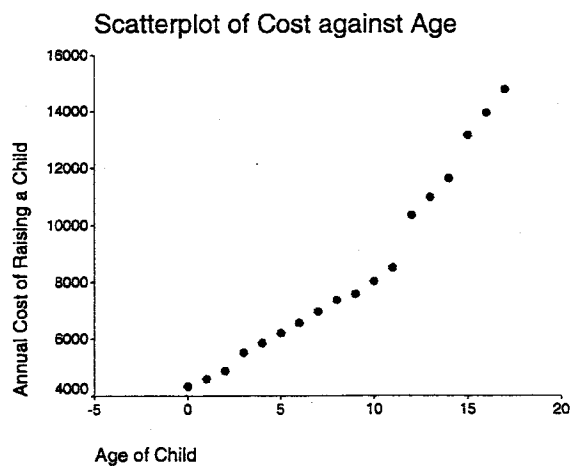
(b)  $\hat{y} = 444.434 + 1.079t$ .

(c)



There is no unusual pattern in this residual plot against the omitted predictor, pH. There is no reason that pH should be included in the model, as it contributes little to the fit of the model, nor is it strongly associated with the residuals of the fit. This conclusion is consistent with the previous conclusion from Exercise 11.19.

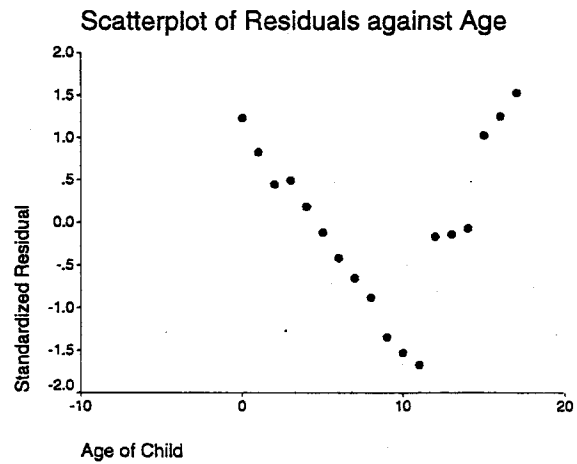
11.28 (a)



A curve would better describe the relationship.

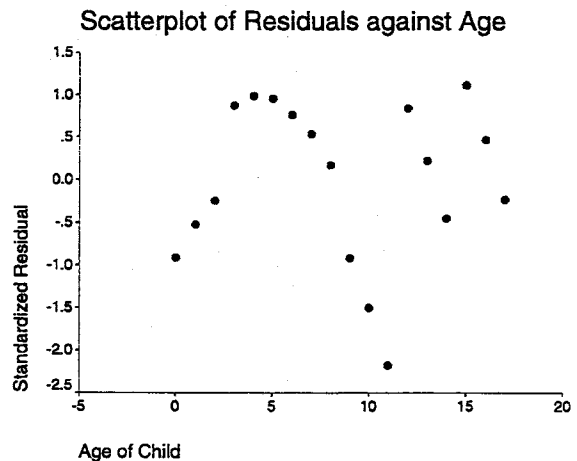
(b) The LS fitted line for the straight line model is  $\hat{y} = 3296 + 600t$ .  $r^2 = 0.940$ .

(c)



The residuals of the straight line fit have a V shape, indicating that a linear fit is not sufficient.

- (d) The LS line for the quadratic model is  $\hat{y} = 4656 + 90t + 30t^2$ .  $r^2 = 0.990$ , so that an additional 5% of the variation in  $y$  is accounted for by the quadratic term.
- (e)



This residual plot still exhibits some patterns, indicating that the quadratic fit is not sufficient. There is an outlier around  $t = 11$ , with a standardized residual around  $-2.2$ .

### Solutions to Section 11.6

11.29 The rate of change is  $\beta_1 + 2\beta_2 t$ .  $\beta_2$  represents how the rate of change in cigarette consumption is changing over time. If  $\beta_2$  is negative, then cigarette consumption is flattening out.

11.30

(a)

Coefficient	Interpretation	Sign			
		(a)	(b)	(c)	(d)
$\beta_0$	Intercept(A)	+	+	+	+
$\beta_1$	Slope(A)	+	+	+	-
$\beta_2$	Intercept(B)-Intercept(A)	+	+	0	-
$\beta_3$	Slope(B)-Slope(A)	0	+	+	+

(b)

Scatterplot	Group	Model
(a)	A	$y = \beta_0 + \beta_1 x$
	B	$y = (\beta_0 + \beta_2) + \beta_1 x$
(b)	A	$y = \beta_0 + \beta_1 x$
	B	$y = (\beta_0 + \beta_2) + (\beta_1 + \beta_3)x$
(c)	A	$y = \beta_0 + \beta_1 x$
	B	$y = \beta_0 + (\beta_1 + \beta_3)x$
(d)	A	$y = \beta_0 + \beta_1 x, \beta_1 < 0$
	B	$y = (\beta_0 + \beta_2) + (\beta_1 + \beta_3)x, \beta_1, \beta_2 < 0$

11.31 Using method (i),

$$\hat{\beta}_1^* = \hat{\beta}_1 \left( \frac{s_{x_1}}{s_y} \right) = 0.0257 \left( \frac{16.10}{0.694} \right) = 0.596, \text{ and}$$

$$\hat{\beta}_2^* = \hat{\beta}_2 \left( \frac{s_{x_2}}{s_y} \right) = 0.0336 \left( \frac{13.15}{0.694} \right) = 0.637.$$

Using method (ii),

$$R = \begin{bmatrix} 1 & -0.107 \\ -0.107 & 1 \end{bmatrix}, R^{-1} = \begin{bmatrix} 1.0116 & 0.1082 \\ 0.1082 & 1.0116 \end{bmatrix}, r = \begin{bmatrix} 0.529 \\ 0.573 \end{bmatrix}.$$

Then

$$\hat{\beta} = R^{-1}r = \begin{bmatrix} 1.0116 & 0.1082 \\ 0.1082 & 1.0116 \end{bmatrix} \begin{bmatrix} 0.529 \\ 0.573 \end{bmatrix} = \begin{bmatrix} 0.596 \\ 0.637 \end{bmatrix}.$$

SAT-M has a slightly larger effect on GPA than does SAT-V.

11.32

$$\hat{\beta}_1^* = \hat{\beta}_1 \left( \frac{s_{x_1}}{s_y} \right) = 2.060 \left( \frac{7.25}{22.60} \right) = 0.661,$$

$$\hat{\beta}_2^* = \hat{\beta}_2 \left( \frac{s_{x_2}}{s_y} \right) = -2.732 \left( \frac{3.99}{22.60} \right) = -0.482,$$

$$\hat{\beta}_3^* = \hat{\beta}_3 \left( \frac{s_{x_3}}{s_y} \right) = 0.000 \left( \frac{23.48}{22.60} \right) = 0.000.$$

MRI brain size has the largest effect on performance IQ, followed by height.

11.33



$$\hat{\beta}_{\log x_1}^* = \hat{\beta}_{\log x_1} \left( \frac{s_{\log x_1}}{s_{\log y}} \right) = -0.459 \left( \frac{1.199}{0.627} \right) = -0.878,$$

$$\hat{\beta}_{\log x_2}^* = \hat{\beta}_{\log x_2} \left( \frac{s_{\log x_2}}{s_{\log y}} \right) = 0.147 \left( \frac{1.173}{0.627} \right) = 0.275,$$

$$\hat{\beta}_{x_3}^* = \hat{\beta}_{x_3} \left( \frac{s_{x_3}}{s_{\log y}} \right) = -0.080 \left( \frac{1.095}{0.627} \right) = -0.140.$$

Log(Alkalinity) has the largest effect on predicting standardized mercury, followed by log(Calcium).

11.34 The correlation matrix is

$$R = \begin{bmatrix} 1.000 & 0.588 & 0.513 \\ 0.588 & 1.000 & 0.700 \\ 0.513 & 0.700 & 1.000 \end{bmatrix}.$$

Height and weight have the highest correlation. From regressing  $x_1$  on  $x_2$  and  $x_3$ ,  $r_1^2 = 0.366$ . From regressing  $x_2$  on  $x_1$  and  $x_3$ ,  $r_2^2 = 0.561$ . From regressing  $x_3$  on  $x_1$  and  $x_2$ ,  $r_3^2 = 0.505$ . Then

$$\text{VIF}_1 = \frac{1}{1 - r_1^2} = \frac{1}{1 - 0.366} = 1.577,$$

$$\text{VIF}_2 = \frac{1}{1 - r_2^2} = \frac{1}{1 - 0.561} = 2.278,$$

$$\text{VIF}_3 = \frac{1}{1 - r_3^2} = \frac{1}{1 - 0.505} = 2.020.$$

These are all less than 10 and are therefore acceptable. There does not appear to be high multicollinearity.

11.35 The correlation matrix is

$$R = \begin{bmatrix} 1.000 & 0.827 & 0.795 \\ 0.827 & 1.000 & 0.705 \\ 0.795 & 0.705 & 1.000 \end{bmatrix}.$$

Log(Alkalinity) and log(Calcium) have the highest correlation. All pairs of variables have moderately high correlations. To find the variance inflation factors, invert  $R$  to get

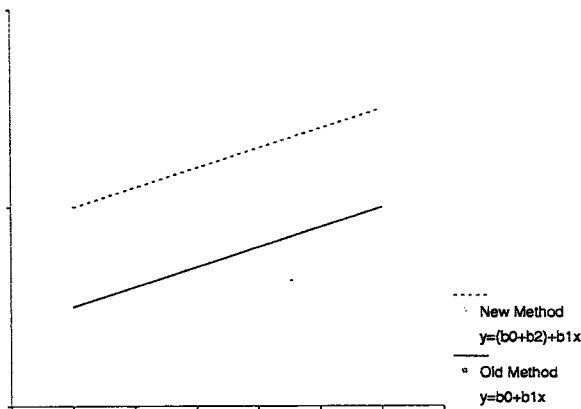
$$R^{-1} = \begin{bmatrix} 4.410 & -2.337 & -1.859 \\ -2.337 & 3.227 & -0.417 \\ -1.859 & -0.417 & 2.771 \end{bmatrix}.$$

Then the variance inflation factors are the diagonal elements of  $R^{-1}$ , namely  $\text{VIF}_{\log x_1} = 4.410$ ,  $\text{VIF}_{\log x_2} = 3.227$ , and  $\text{VIF}_{x_3} = 2.771$ . These are all less than 10 and are therefore acceptable. There does not appear to be high multicollinearity.

11.36 (a)  $\beta_0$  is the yield at a temperature of 0 for the standard method,  $\beta_1$  tells you the slope that both methods change with temperature, and  $\beta_2$  is the difference in yields at a temperature of 0 between the new method and the standard method.

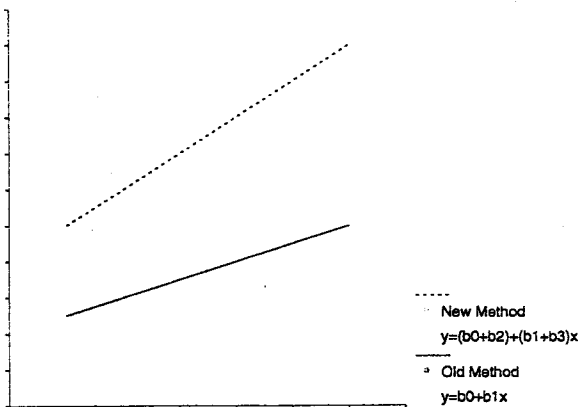
(b)

Plot of Yield Model with Two Methods



(c)

Plot of Yield Model with Two Methods



**11.37** (a) The LS fitted model is  $\hat{y} = 4271.026 + 379.510t - 5269.407z + 549.918tz$ .  $r^2 = 0.998$  compared to 0.940 before, so this model improves the amount of variation explained by almost 6%.

(b)  $SSE_{Full} = 315660$  and  $SSE_{Red.} = 11240400$ . Then

$$F = \frac{(SSE_{Red.} - SSE_{Full})/m}{SSE_{Full}/(n - k + 1)} = \frac{(11240400 - 315660)/2}{315660/14} = 242.264.$$

Since  $F > f_{2,16,0.05} = 3.634$ , we reject  $H_0$  and conclude that the  $z$  terms improve the fit of the model.

(c) The final model is  $\hat{y} = 4271.026 + 379.510t - 5269.407z + 549.918tz$ . For younger children, the model is  $\hat{y} = 4271.026 + 379.510t$ . For older children the model is  $\hat{y} = -998.381 + 929.428t$ .

**11.38**