r 1 Data and Distributions

Section 1.2

(a)

7.

Number		
Nonconforming	Frequency	RelativeFrequency(Freq/60)
0	7	0.117
1	12	0.200
2	13	0.217
3	14	0.233
4	6	0.100
5	3	0.050
6	3	0.050
7	1	0.017
8	1	0.017
doesn't add exactly to 1 bec	ause relative frequencies have	been rounded $\rightarrow 1.001$

- (b) The number of batches with at most 5 nonconforming items is 7+12+13+14+6+3=55, which is a proportion of 55/60 = .917. The proportion of batches with (strictly) fewer than 5 nonconforming items is 52/60 = .867. Notice that these proportions could also have been computed by using the relative frequencies: e.g., proportion of batches with 5 or fewer nonconforming items = 1-(.05+.017+.107) = .916; proportion of batches with fewer than 5 nonconforming items = 1-(.05+.017+.107) = .866.
- 8.
- (a) The following histogram was constructed using MINITAB:





Note: One way to have MINITAB automatically construct a histogram from grouped data such as this is to use MINITAB's ability to enter multiple copies of the same number by typing, for example, 784(1) to enter 784 copies of the number 1. The frequency data in this exercise was entered using the following MINITAB commands:

(b) From the frequency distribution (or from the histogram), the number of authors who published at least 5 papers is 33+28+19+...+5+3+3 = 144, so the proportion who published 5 or more papers is 144/1309 = .11, or 11%. Similarly, by adding frequencies and dividing by n = 1309, the proportion who published 10 or more papers is 39/1309 = .0298, or about 3%. The proportion who published more than 10 papers (i.e., 11 or more) is 32/1309 = .0245, or about 2.5%.

11.

(b) A histogram of this data, using classes of width 1000 separated at 0, 1000, 2000, and 6000 is shown below. The proportion of subdivisions with total length less than 2000 is (12+11)/47 = .489, or 48.9%. Between 2000 and 4000, the proportion is (10 + 7)/47 = .362, or 36.21%. The histogram shows the same general shape as depicted by the stem-and-leaf display in part (a).



12.

(a) A histogram of the y data appears below. From this histogram, the number of subdivisions having no cul-de-sacs (i.e., y = 0) is 17/47 = .362, or 36.2%. The proportion having at least one cul-de-sac ($y \ge 1$) is (47-17)/47 = .30/47 = .638, or 63.8%. Note that subtracting the number of cul-de-sacs with y = 0 from the total, 47, is an easy way to find the number of subdivisions with $y \ge 1$.



(b) A histogram of the z data appears below. From this histogram, the number of subdivisions with at most 5 intersections (i.e., $z \le 5$) is 42/47 = .894, or 89.4%. The proportion having fewer than 5 intersections ($z \le 5$) is 39/47 = .830, or 83.0%.



16. A histogram of the raw data appears below:



After transforming the data by taking logarithms (base 10), a histogram of the log_{10} data is shown below. The shape of this histogram is much less skewed than the histogram of the original data.



Section 1.3

19.

(a) The density curve forms a rectangle over the interval [4, 6]. For this reason, uniform densities are also called **rectangular densities** by some authors. Areas under uniform densities are easy to find (i.e., no calculus is needed) since they are just areas of rectangles. For example, the total area under this density curve is $\frac{1}{2}(6-4) = 1$.



(b) The proportion of values between 4.5 and 5.5 is depicted (shaded) in the diagram below. The area of this rectangle is $\frac{1}{2}(5.5-4.5) = .5$. Similarly, the proportion of *x* values that exceed 4.5 would be $\frac{1}{2}(6-4.5) = .75$.



- (c) The median of this distribution is 5 because exactly half the area under this density sits over the interval [4,5].
- (d) Since 'good' processing times are short ones, we need to find the particular value x_0 for which the proportion of the data less than x_0 equals .10. That is, the area under the density to the left of x_0 must equal .10. Therefore, the area = $.10 = \frac{1}{2}(x_0 4)$, and so $x_0 4 = .20$. Thus, $x_0 = 4.20$.
- (a) The density function is f(x) = 1/[5 (-5)] = 1/10 over the interval [-5, 5] and f(x) = 0 elsewhere. The proportion of x values that are negative is exactly .5 since the value x = 0 sits precisely in the middle of the interval [-5, 5].
- (b) The proportion of values between -2 and 2 is $\frac{1}{10}[2 (-2)] = .4$. The proportion of the *x* values falling between -2 and 3 is $\frac{1}{10}[3 (-2)] = .5$.
- (c) The proportion of the x values that lie between k and k+4 is $\frac{1}{10}[(k+4)-k] = .4$.

20.

(a) The density curve forms an isosceles triangle over the interval [0, 10]. For this reason, such densities are often called **triangular densities**. The total area under this density curve is simply the area of the triangle, which is ¹/₂ (base)(height) = ¹/₂ (10)(.2) = 1. The height of the triangle is the value of f(x) at x = 5; i.e., f(5) = .4 - .04(5) = .2.



- (b) Proportion $(x \le 3) = \frac{1}{2}(3-0)f(3) = \frac{1}{2}(3)(.12) = .18.$ Proportion $(x \ge 7) = \frac{1}{2}(10-7)f(7) = \frac{1}{2}(3)(.12) = .18.$ Proportion $(x \ge 4) = 1$ - Proportion $(x < 4) = 1 - \frac{1}{2}(4-0)f(4) = 1 - \frac{1}{2}(4)(.16) = .68.$
 - Proportion (4 < x < 7) = 1 [Proportion $(x \le 4)$ + Proportion $(x \ge 7)$] = 1 [.32+.18] = .50.

(2)

(a)

$$\lambda = .00004$$

$$\int_{0}^{\infty} .00004e^{-.0004x} dx = \left[\frac{-(.00004)e^{-.0004x}}{(.00004)}\right]_{20,000}^{\infty} - e^{-.0004x}\Big|_{20,000}^{\infty}$$
(b)

$$= 0 - (-e^{-.00004(20,000)}) = e^{-.8} = .449.$$

Note: for any exponential density curve, the area to the right of some fixed constant always equals $e^{-\lambda c}$, as our integration above shows. That is,

Proportion (x >c) = $\int_{c}^{\infty} \lambda e^{-\lambda x} dx$ We will use this fact in the remainder of the chapter instead of repeating the same type of integration as in part (a) Proportion (x ≤ 30,000) = 1 - Proportion (x > 30,000) = 1 - e^{-\lambda c} = 1 - e^{-00004(30,000)} = 1 - e^{-1.2} = .699.

Proportion $(20,000 \le x \le 30,000) =$ Proportion (x > 30,000) - Proportion $(x \le 20,000) = .699 - (1-.449) = .148$.

(c) For the best 1%, the lifetimes must be at least x_0 , where Proportion $(x \ge x_0) = .01$, which becomes $e^{-\lambda x_0} = .01$. Taking natural logarithms of both sides, $-\lambda x_0 = \ln(.01)$, so $x_0 = -\ln(.01)/\lambda = 4.60517/.00004 = 115,129.25$. For the worst 1%, we have Proportion $(x \le x_0) = .01$, which is equivalent to saying that Proportion $(x \ge x_0) = .99$, so $e^{-\lambda x_0} = .99$. Taking logarithms, $-\lambda x_0 = \ln(.99)$, so $x_0 = -\ln(.99)/\lambda = 251.26$.

26.

(a) (ii) qualifies as a distribution because the probabilities add exactly to 1. On the other hand, the probabilities in (i) add to .7, which disqualifies (i). In (iii), notice that the probability associated with x = 4 is negative, which is impossible for a valid distribution.

(b) Using (ii), the proportion of cars having at most 2 (i.e., 1 or less) under-inflated tires is p(0)+p(1)+p(2) = .4 + .1 + .1 = .6. Similarly, Proportion($x \ge 1$) = 1- Proportion(x = 0) = 1-.4 = .6.

- (a) Proportion $(x \le 3) = .10 + .15 + .20 + .25 = .70$. Proportion $(x \le 3) =$ Proportion $(x \le 2) = .10 + .15 + .20 = .45$.
- (b) Proportion $(x \ge 5) = 1$ Proportion(x < 5) = 1 (.10+.15+.20+.25+.20) = .10.
- (c) Proportion $(2 \le x \le 4) = .20 + .25 + .20 = .65$
- (d) At least 4 lines will *not* be in use whenever 2 or fewer lines *are* in use. At most 2 lines are in use .45, or 45%, of the time from part (a) of this exercise.

28.

The sum of all the proportions must equal 1, so $1 = \sum_{y=1}^{5} cy = c[1+2+3+4+5] = 15c$ and c=1/15. Proportion($y \le 3$) = p(1) + p(2) + p(3) = 1/15 + 2/15 + 3/15 = 6/15 = .4.

Proportion $(2 \le y \le 4) = 2/15 + 3/15 + 4/15 = 9/15 = .60$.

Section 1.4

30.

(a) Proportion($z \le 2.15$) = .9842 (Table I). Proportion($z \le 2.15$) will also equal .9842 because the z distribution is continuous.

- (b) Using the symmetry of the z density, Proportion(z > 1.50) = Proportion(z < -1.50) = .0668. Proportion(z > -2.00) = 1 - Proportion($z \le -2.00$) = 1 - .0228 = .9772.
- (c) Proportion($-1.23 \le z \le 2.85$) = Proportion($z \le 2.85$) Proportion($z \le -1.23$) = .9978 .1093 = .8885.
- (d) In Table I, z values range from -3.8 to +3.8. For z < -3.8, the left tail areas (i.e., table values) are .0000 (to 4 decimal places); for z > +3.8, left tail areas equal 1.0000 (to 4 places). Therefore, Proportion(z > 5) = 1 Proportion($z \le 5$) 1 1.0000 = .0000. Similarly, Proportion(z > -5) = 1 Proportion($z \le -5$) = 1 .0000 = 1.0000.
- (e) Proportion(z < |2.50|) = Proportion(-2.50 < z < 2.50)= Proportion(z < 2.50) - Proportion(z < -2.50) = .9938 - .0062 = .9876.

31.

Proportion($z \le 1.78$) = .9625 (Table I)

- (b) Proportion(z > .55) = 1 Proportion($z \le .55$) = 1 .7088 = .2912.
- (c) Proportion(z > -.80) = 1 Proportion $(z \le -.80) = 1$.2119 = .7881.
- (d) Proportion($.21 \le z \le 1.21$) = Proportion($z \le 1.21$) Proportion($z \le .21$) = .8869 - .5832 = .3037.
- (e) Proportion(z ≤ -2.00 or z ≥ 2.00) = Proportion (z≤ -2.00) + [1- Proportion(z < 2.00)]
 = .0228 + [1 .9772] = .0456. Alternatively, using the fact that the z density is symmetric around z = 0, Proportion(z ≤ -2.00) = Proportion(z ≥ 2.00), so the answer is simply

2 Proportion($z \le -2.00$) = 2(.0228) = .0456.

- (f) Proportion($z \le -4.2$) = .0000
- (g) Proportion(z > 4.33) = .0000

32.

- (a) Proportion($z \le z^*$) = .9082 when z^* = 1.33 (Table I).
- (b) Proportion($z \le 1.33$) = .9082 and Proportion($z \le 1.32$) = .9066; the z value for 0.9082 is just under 1.33, and it is sufficient to approximate it with 1.329.
- (c) Proportion($z > z^*$) = .1210 is a right-tail area. Converting to a left-tail area (so that we can use table I), Proportion($z \le z^*$) = 1 .1210 = .8790. From Table I, z^* = 1.17 has a left tail area of .8790 and, therefore, has a right-tail area of .1210.
- (b) Proportion(-z* ≤ z ≤ z*) = .754. Because the z density is symmetric the two tail areas associated with z < z* and z > z* must be equal and they also account for all the remaining area under the z density curve. That is, 2×Proportion(z ≤ z*) = 1 .754 = .246, which means that Proportion(z ≤ z*) = .1230. From Table I, Proportion(z ≤ -1.16) = .1230, so z* = -1.16 and z* = 1.16.
- (c) Proportion(z > z*) = .002 is equivalent to saying that Proportion(z ≤ z*) = .998. From Table I, Proportion(z ≤ 2.88) = .9980, so z* = 2.88. Similarly, you would have to go a distance of -2.88 (i.e., 2.88 units to the left of 0) to capture a left-tail area of .002.

34.

(a) Let z^* denote the 91st percentile, so Proportion($z \le z^*$) = .9100. From Table I, Proportion($z \le 1.34$) = .9099 and Proportion($z \le 1.35$) = .9115 and, so the z^* is just over 1.34, and it is sufficient to approximate it with 1.341.

- (b) Let z^* denote the 9th percentile, Proportion($z \le z^*$) = .0900. From Table I, $z^* \approx -1.34$. Note that the 9th percentile should be the same distance to the *left* of 0 that the 91st percentile is to the right of 0, so we could have simply used the answer to part (a), after attaching a minus sign.
 - (c) The 22nd percentile occurs at $z \approx -.77$.

40.

(d) Proportion(3432-c < x < 3432+c) = .98; standardizing gives Proportion((3432-c-3432)/482 < z < (3432+c-3432)/482) = Proportion(-c/482 < z < c/482) = .98. This means that the right and left tail areas Proportion(z < -c/482) and Proportion(z > c/482) both equal .01. From Table I, z = -2.33 has (approximately) a left tail area of .01, so -c/482 = -2.33, or, c = 1123.

Section 1.6

54.

(a) Let x = number of red lights encountered. Then x has a binomial distribution with n = 10, π = .40. Proportion(x ≤ 2) = .006 + .040 + .121 = .167 (using Table II). Similarly, Proportion(x ≥ 5) = .201 + .111 + .042 + .011 + .002 + .000 = .367. (b) Proportion $(2 \le x \le 5) = .121 + .215 + .251 + .201 = .788$.

55.

(a) Let x = number of bits erroneously transmitted. Then, x is binomial with n = 20, π = .10, so Proportion(x ≤ 2) = .122 + .270 + .285 = .677 (from Table II).

(b) Proportion($x \ge 5$) = .032 + .009 + .002 + .000 ++ .000 = .043.

(c) 'More than half' means 11 or more, so Proportion $(x \ge 11) = .000 + ... + .000 = .000$.

57.

(c) Proportion $(10 \le x \le 20) = .006 + .011 + ... + .089 + .089 = .556.$ Proportion $(10 \le x \le 20) = .011 + ... + .089 = .461.$

59.

Using Table III ($\lambda = 20$), Proportion($x \ge 15$) = 1 - Proportion($x \le 14$) = 1 - (.000 + .000 + ... + .001 + .001 + .003 + .003 + .006 + .011 + .018 + .-27 + .039) = 1 - .106 = .894. Similarly, Proportion($x \le 25$) = 1 - Proportion($x \ge 26$) = 1 - (.034 + .025 + .018 + .013 + .008) = 1 - .098 = .902.

60.

Although x has a binomial distribution with n = 1000 and $\pi = 1/200 = .005$, its distribution can be approximated by a Poisson distribution with $\lambda = n\pi = 1000/200 = 5$. Therefore, using Table III, Proportion($x \ge 8$) = .065 + .036 + .018 + .008 + .003 + .001 = .131. Note: because the Table entries are rounded to 3 places, you would get a slightly different answer (of .135) if you worked the problem by first adding the proportions for $x \le 7$, then subtracting from 1.

Proportion($5 \le x \le 10$) = .175 + .146 + .104 + .065 + .036 + .018 = .544.

Supplementary Exercises

62.

(a) Let x = fracture strength. Then, Proportion(x < 90) = .50 because 90 is the mean of the (assumed) normal distribution of x. Table I must be used for the other proportions: Proportion(x < 95) = Proportion(z < (95-90)/3.75) = Proportion(z < 1.33) = .9082; therefore, Proportion(x ≥ 95) = 1 - Proportion(x < 95) = 1 - .9082 = .0918.

- (b) Proportion $(85 \le x \le 95)$ = Proportion $((85-90)/3.75 \le z \le (95-90)/3.75)$ = Proportion $(-1.33 \le z \le 1.33)$ = Proportion $(z \le 1.33)$ - Proportion(z < -1.33) = .9082 - .0918 = .8164. Likewise, Proportion $(80 \le x \le 100)$ = Proportion $(-2.67 \le z \le 2.67)$ = .9962 - .0038 = .9924.
- (c) Let x* denote the value exceeded by 90% of the x data. From Table I, the corresponding value z* for the z distribution is $z^* = -1.28$, which is 1.28 σ 's below the mean of the z distribution. Therefore, x* must be 1.28 σ 's below the mean of the x data: $x^* = 90 1.28(3.75) = 85.20$.
- (b) The corresponding interval for the z distribution is $.99 = Proportion(-z^* \le z \le z^*)$, which means that the left-tail proportion Proportion($z \le z^*$) = .99 + .005 = .9950. From Table I, $z^* \approx 2.58$; i.e., 2.58 σ 's from the mean. Therefore, the two values for the x data must lie 2.59 σ 's on either side of the mean: $90 \pm 2.58(3.75) = 80.325$ and 99.675.

64.

(a) Use the formula for right-tail areas given in the answer to Exercise 23(b): Proportion(x ≤ 100) = 1 - Proportion(x ≥ 100) = 1 - $e^{-\lambda(100)}$ = 1 - $e^{-.01386(100)}$ = 1 - .25007 = .7499, or, .75. Proportion(x ≤ 200) = 1 - $e^{-\lambda(200)}$ = 1 - $e^{-.01386(200)}$ = 1 - .06253 = .9375.

Proportion($100 \le x \le 200$) = Proportion(x > 100) - Proportion(x > 200) = $e^{-.01386(100)} - e^{-.01386(200)} = .25007 - .06253 = .1875.$

- (b) Proportion($x \ge 50$) = $e^{-.01386(50)}$ = .50007, or, *almost* exactly 50%.
- (c) Let ^x denote the median. Then, .50 = Proportion(x ≥ ^x) = e^{-λ(x)}. Taking logarithms of both sides, ln(.50) = -λ^x, so ^x = -.6931472/-.01386 = 50.01. Note that you could have guessed from the answer to (b) that ^x is very close to 50.

66.

(c) Proportion $(5 \le x \le 10) = .196 + .163 + .111 + .062 + .030 + .011 = .573$.

(d) Proportion($5 \le x \le 10$) = .163 + .111 + .062 + .030 = .366

67.

(a) Accommodating everyone who shows up means that $x \le 100$, so Proportion $(x \le 100) = .05 + .10 + ... + .24 + .17 = .82$.

(b) Proportion(x > 100) = 1 - Proportion($x \le 100$) = 1 - .82 = .18.

(c) The first standby passenger will be able to fly as long as the number who show up is $x \le 99$, which leaves one free seat. This proportion is .65. The third person on the standby list will get a seat as long as $x \le 97$, where Proportion($x \le 97$) = .05 + .10 + .12 = .27.

Let X = the number of components that function. X is binomial with n = 5, $\pi = .9$. (Proportion of 3 out of 5 systems that will function.)

$$= \text{Proportion}(x \ge 3) = p(3) + p(4) + p(5) = {5 \choose 3} \pi^{3} (1 - \pi)^{2} + {5 \choose 4} \pi^{4} (1 - \pi)^{4} + {5 \choose 5} \pi^{5} (1 - \pi)^{0} = 10\pi^{3} (1 - \pi)^{2} + 5\pi^{4} (1 - \pi) + \pi^{5} = (10)(.9)^{4} (.1)^{2} + (5)(.9)^{4} (.1) + (.9)^{5} = .99144$$

Alternatively, use the Binomial Table II.

75.

Letting X = "bursting strength", we first find the proportion of all bottles having bursting strength exceeding 300 PSI. Proportion(X > 300) = Proportion(z > ((300-250)/30)) = Proportion(z > 1.67) = .0475 (from Table I). Then, Y = "the number of bottles in a carton of 12 with bursting strength over 300 PSI" is a binomial variable with n = 12 and π = .0475. So the proportion of all cartons with at least one bottle with a bursting strength over 300 PSI is Proportion(Y = 1) = 1 – Proportion(Y = 0) = 1 – (1-p)^{1/2} = 1 – (1 - .0475)^{12} = .4423

Chapter 2 Numerical Summary Measures

Section 2.1

4.

The three quantities of interest are: \overline{x}_n, x_{n+1} , and \overline{x}_{n+1} . Their relationship is as follows:

$$\overline{x}_{n+1} = \left(\frac{1}{n+1}\right) \sum_{i=1}^{n+1} x_i = \left(\frac{1}{n+1}\right) \left[\sum_{i=1}^n x_i + x_{n+1}\right] = \left(\frac{1}{n+1}\right) \left[n\overline{x}_n + x_{n+1}\right] = \left(\frac{n\overline{x}_n + x_{n+1}}{n+1}\right)$$
((10)((10)))

For the strength observations, $\bar{x}_{10} = 640.5$ and $x_{11} = 780$. Therefore, $\bar{x}_{11} = \left(\frac{(10)(640.5)+780}{10+1}\right) = 653.18$

6.

- (a) The *reported* blood pressure values would have been: 120 125 140 130 115 120 110 130 135. When ordered these values are: 110 115 120 120 125 130 130 130 135 140. The median is 125.
- (b) 127.6 would have been reported to be 130 instead of 125. Since this value is the middle value, the median would

change from 125 to 130. This example illustrates that the median is sensitive to rounding in the data.

8.

$$\mu = \int_{-1}^{1} (x)(.75)(1-x^2) dx = \left[\frac{.75x^2}{2} - \frac{.75x^4}{4}\right]_{-1}^{1} = 0$$

The mean value of x equals 0.

9.

(a)
$$\mu = \int_0^2 x(.5x) dx = \int_0^2 \frac{x^3}{3} \Big|_0^2 = 4/3$$
. The mean μ does not equal 1 because the density curve is not symmetric around $x = 1$.

- (b) Half the area under the density curve to the left (or right) of the median $\tilde{\mu}$, so, $.50 = \int_{0}^{\tilde{\mu}} .5x \, dx = .5 \left[\frac{x^2}{2}\right]_{0}^{\tilde{\mu}} = (\tilde{\mu} \cdot \tilde{\mu})^{2}$ $(\tilde{\mu} \cdot \tilde{\mu})^{2} = 4(.5) = 2$ and $\tilde{\mu} = 1.414$. $\mu < \tilde{\mu}$ because the density curve is negatively skewed.
- (c) $\mu \pm \frac{1}{2} = \frac{4}{3} \pm \frac{1}{2} = \frac{5}{6}$ and $\frac{11}{6}$. The area under the curve between these two values is $\int_{5/6}^{11/6} .5x \, dx = \left[\frac{x^2}{4}\right]_{5/6}^{1.914} = .667$. Similarly, the proportion of the times that are within one-half hour of $\widetilde{\mu}$ is: $\int_{.914}^{1.914} .5x \, dx = \left[\frac{x^2}{4}\right]_{.914}^{1.914} = .707$.

12.

$$\mu = \sum_{x=0}^{6} xp(x) = \begin{bmatrix} (0)(.10) + (1)(.15) + (2)(.20) + (3)(.25) + \\ (4)(.20) + (5)(p(5)) + (6)(p(6)) \end{bmatrix} = 2.64$$

 $\Rightarrow 2.1+5(p(5))+6(p(6))=2.64$ $\Rightarrow 5(p(5))+6(p(6))=.54$ Also, we know that $\sum_{x \neq 0} p(x) = 1$ $\Rightarrow p(5) + p(6) = 1 - [(.10) + (.15) + (.20) + (.25) + (.20)]$ $\Rightarrow p(5) + p(6) = 1 - .90 = .10$ Therefore, p(5) = .10 - p(6)Returning to our first equation:

 $\Rightarrow 5(.10 - p(6)) + 6(p(6)) = .54$ $\Rightarrow .5 - 5(p(6) + 6(p(6)) = .54$ $\Rightarrow p(6) = .04$ Therefore, p(5) = .1 - .04 = .06

Finally, p(5) = .06 and p(6) = .04

13.

$$\mu = \sum_{x=0}^{4} x \cdot p(x) = 0(.4) + 1(.1) + 2(.1) + 3(.1) + 4(.3) = 1.8$$

Section 2.2

23.

Let X = the number of drivers who travel between a particular origin and destination during a designated time period. X has a Poisson distribution with $\lambda = 20$.

(a) $\mu_x = \lambda = 20$ Find $P(\mu - 5 \le x \le \mu + 5) = P(15 \le x \le 25)$ Using Table III with $\lambda = 20$, we obtain: $P(15 \le x \le 25) = .052 + .065 + .076 + .084 + .089 + .089 + .085 + .077 + .067 + .056 + + .045 = .785$ (b) $\sigma_x = \sqrt{\lambda} = \sqrt{20} = 4.47$ Find $P(\mu - \sigma \le x \le \mu + \sigma)$ $P(20 - 4.47 \le x \le 20 + 4.47) = P(15.53 \le x \le 24.47)$ But, since X is an integer-valued random variable, only the integers between 16 and 24 satisfy this requirement. So, we find $P(16 \le x \le 24)$ using Table III with $\lambda = 20$, we obtain: $P(16 \le x \le 24) = .065 + .076 + .084 + .089 + .085 + .077 + .067 + .056$ = .688

24.

(a) $\sigma = \sqrt{\lambda} = \sqrt{5} = 2.236$, so x values that lie between 5 - 2.236 = 2.764 and 5 + 2.236 = 7.236 are within one standard deviation from the mean. Because x is integer-valued, only the integers between 3 and 7 satisfy this requirement, i.e., Proportion(2.764 $\leq x \leq 7.236$) = Proportion(3 $\leq x \leq 7$) = p(3) + p(4) + p(5) +p(6) + p(7) = .140 + .175 + .146 + .104 = .740 (using Table III with $\lambda = 5$).

(b) To exceed the mean by *more* than 2 standard deviations, x values must be greater than 5 + 2(2.236) = 9.472. The integer values of x than satisfy this requirement are x = 10 and greater. From Table III, Proportion(x > 9.472) = Proportion(x ≥ 10) = .018 + .008 + .003 + .001 = .030. Note: because table entries are rounded to 3 places, a slightly different answer results if you calculate the proportion as 1 - Proportion(x ≤ 9) = 1 - .966 = .032.

(a)
$$\sigma = \sqrt{\sum (x - \mu)^2 p(x)}$$
$$\sigma = \sqrt{\frac{(0 - 1.8)^2 (4) + (1 - 1.8)^2 (1) + (2 - 1.8)^2 (1) + (3 - 1.8)^2 (.1) + (4 - 1.8)^2 (.3)}{(3 - 1.8)^2 (.1) + (4 - 1.8)^2 (.3)}}$$
$$\sigma = \sqrt{2.96} = 1.72$$
(b)
$$P(\mu - \sigma \le x \le \mu + \sigma) = P(.08 \le x \le 3.52) = P(1 \le x \le 3) = .3$$

(b)
$$P(x \mid \mu + 3\sigma) + P(x \mid \mu - 3\sigma) = P(x \mid 6.96) + P(x \mid -3.36) = 0 + 0 = 0$$

$$\sigma^{2} = \sum (x - \mu)^{2} p(x) = \sum (x^{2} - 2\mu x + \mu^{2}) p(x) = \sum x^{2} p(x) - 2\mu \sum x p(x) + \mu^{2} \sum p(x) = \sum x^{2} p(x) - 2\mu^{2} + \mu^{2} = \sum x^{2} p(x) - \mu^{2}.$$

For the mass function given in Exercise 24, $\sigma^{2} = \sum x^{2} p(x) - \mu^{2} = (0)^{2}(.4) + (1)^{2}(.1) + (2)^{2}(.1) + (4)^{2}(.3) + (1.8)^{2} = 6.2 - 3.24 = 2.96.$

28.

The proportion of x values between $\mu - 1.5\sigma$ and $\mu + 1.5\sigma$ is the same as the proportion of z values between -1.5 and +1.5: Proportion(-1.5 $\leq z \leq 1.5$) = Proportion($z \leq 1.5$) - Proportion($z \leq -1.5$) = .9332-.0668 = .8664. The proportion of x value that exceed μ by more than 2.5 σ 's equals Proportion(z > 2.5) = 1 - Proportion($z \leq 2.5$) = 1 - .9938 = .0062.

29.

X is binomial with $\pi = .2$ and n = 25.

 $\sigma^{2} = n\pi (1 - \pi) = (25)(.20)(.80) = 4 \quad \sigma = 2$ Also, $\mu = n\pi = (25)(.20) = 5$ $P(x) \mu + 2\sigma = P(x) + 2(2) = P(x) = 0$ Using Table II: $P(x) = P(x \ge 10) = .011 + .004 + .002 = .017$

(Notice that $P(x \ge 13) \ge 0$)

Section 2.3

35.

Since 5% of all lengths exceed 3.75 mm, then 3.75 is the 95th percentile of the distribution. Because the circuits are normally distributed, 3.75 is 1.645 standard deviations above the mean; that is, ${}^{3.75 = \mu + 1.645\sigma}$. Furthermore, 3.85 is the 99th percentile of the distribution, and so it is 2.33 standard deviations above the mean: ${}^{3.85 = \mu + 2.33\sigma}$. We then have two equations and two unknowns:

$$\mu + 1.645\sigma = 3.75 \mu + 2.33\sigma = 3.85$$

Subtracting the top equation from the bottom equation yields $.685\sigma = .10$, and so $\sigma = .10/.685 = .146$. Then substituting $\sigma = .146$ into either of the equations gives $\mu = 3.51$.

Section 2.4

The normal quantiles are easy to generate using Minitab or Excel. For example, in Minitab, typing the following commands will generate the normal quantiles:

```
MTB> set c2
MTB> 1:20
MTB> let c3 = (c2-.5)/20
MTB> invcdf c3 c4;
SUBC> norm 0 1.
```

(Note: you don't have to type 'END' to end the input of data into column C2; typing any valid Minitab command will automatically end data input and then will execute the command)

Although it isn't necessary in this exercise, remember to sort (from smallest to largest) the data before plotting it versus the normal quantiles. The quantile plot for this data is shown below. The pattern is obviously nonlinear, so a normal distribution is implausible for this data. The apparent break that appears in the data in the right side of the graph is indicative of data that contains outliers.



Let η_p denote the pth quantile of an exponential distribution with parameter λ . Then, the area to the *right* of η_p is 1-p. Recall from Exercise 23 of Chapter 1, that the right tail area (i.e., the area past x = c) for an

exponential distribution is simply $e^{-\lambda c}$. Therefore, $e^{-\eta_p \lambda} = 1$ -p. Taking logarithms of both sides, $-\eta_p \lambda = \ln(1-p)$, so $\eta_p = \lambda(-\ln(1-p))$. That is, the quantiles η_p are linearly related to the quantities $-\ln(1-p)$, so a plot of the sample quantiles $x_{(i)}$ versus $-\ln(1-p_i)$ is a straight line. In this exercise, n = 16, so the values of p_i equal (i-.5)/16 for i = 1, 2, ... 16. Use Minitab or Excel to compute the plotting values $-\ln(1-p_i)$. The quantile plot for this data appears below. Because the plot exhibits curvature, an exponential distribution would not be appropriate for this data.



53.

(a) A normal quantile plot of x follows.
 Clearly the variable, hourly median power, is not normally distributed, as the normal quantile plot is curvilinear.





Supplementary Exercises

61.

(b)
$$\mu = 100 \Rightarrow \lambda = .01$$

 $\sigma = 100$
 $\mu \pm \sigma \Rightarrow 100 \pm 100 \Rightarrow (0, 200)$; $\int_{0}^{200} 1e^{-.01x} dx = 1 - e^{-.01(200)} = .8647$
 $\mu \pm 2\sigma \Rightarrow 100 \pm 200 \Rightarrow (0, 300)$; $\int_{0}^{300} 1e^{-.01x} dx = 1 - e^{-.01(300)} = .9502$
 $\mu \pm 3\sigma \Rightarrow 100 \pm 300 \Rightarrow (0, 400)$; $\int_{0}^{400} 1e^{-.01x} dx = 1 - e^{-.01(400)} = .9817$

So:

Percentage within	Chebyshev's Rule	Exponential
Ισ	No statement	86.47%
2σ	At least 75%	95.02%
3σ	At least 89%	98.17%

62.

The transformation of the x_i 's into the z_i 's is analogous to 'standardizing' a normal distribution. The purpose of standardizing is to reduce a distribution (or, in this exercise, a set of data) to one that has a mean of 0 and a standard deviation of 1. To show that this transformation achieves this goal, note that:

$$\sum_{i=1}^{n-1} \sum_{x=1}^{n-1} \frac{1}{s} (x_i - \overline{x}) = \frac{1}{s} \sum_{x=1}^{n-1} (x_i - \overline{x}) = \frac{1}{s} (0) = 0, \text{ so, dividing this sum by n, } \overline{z} = 0. \text{ Next,}$$
$$\frac{1}{n-1} \sum_{x=1}^{n-1} (z_i - \overline{z})^2 = \frac{1}{n-1} \left(\frac{1}{s^2}\right) \sum_{x=1}^{n-1} (x_i - \overline{x})^2 = \left(\frac{1}{s^2}\right)^2 = 1.$$

- (a) Let y denote the capacitance of a capacitor. Capacitors will conform to specification when y is in the interval from 95 nf to 105 nf. Therefore, Proportion($95 \le y \le 105$) = Proportion($(95-98)/2 \le z \le (105-98)/2$) = Proportion($-1.5 \le z \le 3.5$). Using Table I, this proportion is equivalent to Proportion($z \le 3.5$) Proportion($z \le -1.5$) = .9998 .0068 = .9930, or, about 93.3%.
- (b) The number of capacitors in a batch of 20 that conform to specifications will have a binomial distribution with n = 20 and $\pi = .9330$. Therefore, the proportion of batches containing at least 19 conforming capacitors is Proportion($x \ge 20$) = Proportion(x = 19) + Proportion(x = 20). Using the formula for the binomial mass function:

$$\frac{20!}{19!!!}(.9930)^{19}(1-.9930)^{1} + \frac{20!}{20!0!}(.9930)^{20}(1-.9930)^{0} = .9914$$

Section 3.1

5.

(a) The scatter plot with axes intersecting at (0,0) is shown below.



(b) The scatter plot with axes intersecting at (55,100) appears below. The plot in (b) makes it somewhat easier to see the nature of the relationship between the two variables.



(c) A parabola appears to provide a good fit to both graphs.

Section 3.2

11. (a) $SS_{xy} = 5530.92 - (1950)(47.92)/18 = 339.586667$, $SS_{xx} = 251,970 - (1950)^2/18 = 40,720$, and $SS_{yy} = 130.6074 - (47.92)^2/18 = 3.033711$, so

$$r = \frac{339.586667}{\sqrt{40720}\sqrt{3.033711}} = .9662$$

There is a very strong positive correlation between the two variables.

(b) Because the association between the variables is positive, the specimen with the larger shear force will tend to have a larger percent dry fiber weight.

- (c) Changing the units of measurement on either (or both) variables *will have no effect on the calculated value of r*, because any change in units will affect both the numerator and denominator of r by exactly the same multiplicative constant.
- 14. Using a correlation coefficient to summarize the relationship between the artist (x) and the sales price (y) is not appropriate. To compute and interpret a correlation coefficient both x and y variables must be quantitative variables. While the y variable, sale price, is quantitative, the x variable, artist, is not.
- 15. Let d_0 denote the (fixed) length of the stretch of highway. Then, $d_0 = \text{distance} = (\text{rate})(\text{time}) = xy$. Dividing both sides by x, gives the equation $y = d_0/x$ which means the relationship between x and y is curvilinear (in particular, the curve is a hyperbola). However, for values of x that are fairly close to one another, sections of this hyperbola can be approximated very well by a straight line with a negative slope (to see this, draw a picture of the function d_0/x for a particular value of d_0). This means that r should be closer to -.9 than to any of the other choices.
- 16. The value of the sample correlation coefficient using the squared y values would not necessarily be approximately 1. If the y-values are greater than 1, then the squared y-values would differ from each other by more than the y-values differ from one another. Hence, the relationship between x and y^2 would be less like a straight line, and the resulting value of the correlation coefficient would decrease. [Note: I have yet to find an example where r is less than about .96 for (x, y^2), however.]

17. (a)
$$SS_{xx} = 37695 - (561)^2/9 = 2726$$
, $SS_{yy} = 40223 - (589)^2/9 = 1676.222$, and $SS_{xy} = 38281 - (561)(589)/9 = 1566.666$, so

$$r = \frac{1566.667}{\sqrt{2726}\sqrt{1676.222}} = .733.$$

(b)
$$\overline{x}_1 = (70+72+94)/3 = 78.667, \ \overline{y}_1 = (60+83+85)/3 = 76.$$

 $\overline{x}_2 = (80+60+55)/3 = 65, \ \overline{y}_2 = (72+74+58)/3 = 68.$
 $\overline{x}_3 = (45+50+35)/3 = 43.333, \ \overline{y}_3 = (63+40+54)/3 = 52.333.$

$$S_{xx} = [(78.667)^{2} + (65)^{2} + (43.333)^{2} - (78.667 + 65 + 43.333)^{2}/3] = 634.913,$$

$$S_{yy} = [(76)^{2} + (68)^{2} + (52.333)^{2} - (76 + 68 + 52.333)^{2}/3] = 289.923,$$

$$S_{xy} = [(78.667)(76) + (65)(68) + (43.333)(52.333) - (187)(196.333)/3] = 428.348, \text{ so}$$

$$r = \frac{428.348}{\sqrt{2}} = .9984.$$

$$=\frac{1}{\sqrt{634.913}\sqrt{289.923}}$$

(c) The correlation among the averages is noticeably higher than the correlation among the raw scores, so these points fall much closer to a straight line than do the unaveraged scores. The reason for this is that averaging tends to reduce the variation in data, making it more likely that the averages will fall close to a straight line than the more variable raw data.



Section 3.3

26. (a) Yes, the following plot does suggest the aptness of a simple linear regression model.



(b) We need to calculate
$$S_{xx}$$
 and S_{xy} :

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - (1/n)(\sum x_i)^2 = 5099.2412 - (1/32)(384.26)^2 \approx 485.00$$

$$S_{xy} = \left(\sum (x_i - \bar{x}) \sum (y_i - \bar{y})\right) = \sum x_i y_i - (1/n)(\sum x_i)(\sum y_i)$$

$$= 37,850.7762 - (1/32)(384.26)(3149.04) = 36.71025$$

$$h = S_{xy} - (S_{xy}) = 36.71025 - (485) \approx 0.0756$$

Therefore, $b = S_{xy} / S_{xx} = 36.71025 / 485 \approx .0756$, and $a = \overline{y} - b\overline{x} = (1/32)(3149.04) - .0756[(1/32)(384.26)] = 98.4075 - 13.212(12.008125) \approx 97.5.$

Hence, the least squares line is given by $\hat{y} = 97.5 \pm .0756x$.

The following MINITAB output summarizes the least squares regression. The output gives $\hat{y} = 97.5 \pm .0757x$

The regression equation is Removal = 97.5 + 0.0757 Temp Predictor Coef SE Coef T P

Constant Temp	97.4986 0.075691	0.0889 0.007046	1096.17 10.74	0.000	
S = 0.1552	R-Sq =	79.4% R·	-Sq(adj) = 7	8.7%	
Analysis of V	ariance				
Source	DF	SS	MS	F	Р
Regression	1	2.7786	2.7786	115.40	0.000
Residual Erro	r 30	0.7224	0.0241		
Total	31	3.5010			

Now, the point prediction of removal efficiency at the temperature value of 10.50 is $\hat{y} = 97.5 + .0756(10.5) = 98.2938$. The residual is given by $y - \hat{y} = 98.41 - 98.2938 = 0.1162$.

(c) The size of a typical of deviation of points in the scatter plot from the least squares line is given by the standard deviation about the least squares line: $s_{e} = \sqrt{\frac{\text{SSResid}}{n-2}}$ $\text{SSResid} = \text{SSTo} - bS_{xy}$ $\text{SSTo} = S_{yy} = \sum y_{i}^{2} - (1/n)(\sum y_{i})^{2} = 309,892.6548 - (1/32)(3149.04)^{2} = 3.501$ $\text{SSResid} = \text{SSTo} - bS_{xy} = 3.501 - (.0756)(36.71025) = .7257$ $s_{e} = \sqrt{\frac{\text{SSResid}}{n-2}} = \sqrt{\frac{.7257}{32-2}} = \sqrt{.02419} \approx .1552.$

NOTE: From the MINITAB output above, we could see that the output value of s = 0.1552 gives s_e , and we could have calculated s_e by calculating the square root of the mean square residual error of .0241.

(d) The proportion of observed variation in removal efficiency that can be attributed to the approximate

linear relationship is given by $r^2 = 1 - \frac{\text{SSResid}}{\text{SSTo}} = 1 - \frac{.7257}{3.501} = .7927$, which can also be found from the MINITAB output.

(e) First compute the new summary statistics:

$$\sum x_i = 384.26 + 6.53 = 390.79$$

$$\sum x_i^2 = 5099.241 + 6.53^2 = 5141.8819$$

$$\sum y_i = 3149.04 + 96.55 = 3245.59$$

$$\sum y_i^2 = 309,892.6548 + 96.55^2 = 319,214.5573$$

$$\sum x_i y_i = 37,850.78 + (6.53)(96.55) = 38,481.25$$

Then compute the least squares line:

$$S_{xx} = \sum (x_i - \bar{x})^2 = \sum x_i^2 - (1/n)(\sum x_i)^2 = 5141.8819 - (1/33)(390.79)^2 \approx 514.10$$

$$S_{xy} = \left(\sum (x_i - \bar{x}) \left(\sum (y_i - \bar{y})\right) = \sum x_i y_i - (1/n)(\sum x_i)(\sum y_i)$$

$$= 38,481.25 - (1/33)(390.79)(3245.59) = 46.58$$

Therefore, $b = S_{xy} / S_{xx} = 46.58 / 514.1 \approx .0906$, and $a = \overline{y} - b\overline{x} = (3245.59 / 33) - (.0906)(390.79 / 33) = 97.28$ So our new equation is $\hat{y} = 97.28 + .0906$.

Next, compute
S_e
 and ${r}^2$:
SSTo = $S_{yy} = \sum y_i^2 - (1/n)(\sum y_i)^2 = 319,214.5573 - (1/33)(3245.59)^2 = 6.85$
SSResid = SSTo - $bS_{xy} = 6.85 - (.0906)(46.58) = 2.63$
 $s_e = \sqrt{\frac{\text{SSResid}}{n-2}} = \sqrt{\frac{2.63}{33-2}} = \sqrt{.0848} \approx .2913.$
 $r^2 = 1 - \frac{\text{SSResid}}{\text{SSTo}} = 1 - \frac{2.63}{6.85} = .616$

From above, we see that the intercept of our least squares equation decreased from 97.5 to about 97.28. The change in the slope is noticeable: the new slope increased from .0756 to .0906. The addition of this new results in a much larger s_e of .2913, which is about twice the original s_e value of .1552. Consequently, the r^2 valued decreased; the new r^2 of .616 is down from the original r^2 value of .7257. Moreover, s_e increases.

27. Data set #1





Fitting a straight line seems appropriate here. There is no indication of a problem. There is a fair amount of scatter around the least squares line, however. This fact is quantified by the r^2 value of about 67%.





It is inappropriate to fit a straight line through this data. There is a quadratic relationship between these two variables. So, a model that incorporates this quadratic relationship should be fit instead of a linear fit.



In this data set there is one pair of values that is a clear outlier (13.0, 12.74). The y-value of 12.74 is much larger than one would expect based on the other data. This pair of values exerts a lot of influence on the linear fit. It should be investigated. Perhaps it was a recording error or some other similar problem. With the outlier included, it is inappropriate to fit a straight line. With it excluded, an excellent linear fit is achieved.

Data set #4



It is inappropriate to fit a straight line through this data. There is no linear relationship between the variables. There is one pair of observations that is a clear outlier (19.0, 12.50). It should be investigated.

Section 3.4 Section 3.5 Supplementary Exercises

49 (a) Since stride rate is being predicted, y = stride rate and x = speed. Therefore, $SS_{xx} = \sum_{i=1}^{x_i^2} -(\sum_{i=1}^{x_i})^2/n$ = 3880.08 - (205.4)²/11 = 44.7018, $SS_{yy} = \sum_{i=1}^{y_i^2} -(\sum_{i=1}^{y_i})^2/n$ = 112.681 - (35.16)²/11 = .2969, and $SS_{xy} = \sum_{i=1}^{x_i y_i} -(\sum_{i=1}^{x_i})(\sum_{i=1}^{y_i})/n$ = 660.130 - (205.4)(35.16)/11 = 3.5969. Therefore, b = SS_{xy}/SS_{xx} = 3.5969/44.7018 = .0805

- and a = (35.16/11) (.0805)(205.4/11) = 1.6932. The least squares line is then $\hat{y} = 1.6932 + .0805x$.
 - (b) Predicting speed from stride rate means that y = speed and x = stride rate. Therefore, interchanging the x and y subscripts in the sums of squares computed in part (a), we now have SS_{xx} = .2969 and SS_{xy} = 3.5969 (note that SS_{xy} does not change when the roles of x and y are reversed). The new regression line has a slope of b = SS_{xy}/SS_{xx} = 3.5969/.2969 = 12.1149 and an intercept of a = (205.4/11) (12.1149)(35.16/11) = -20.0514; that is, ŷ = -20.0514 + 12.1149x.
 - (c) For the regression in part (a), $r = 3.5969/[\sqrt{44.7018}\sqrt{.2969}] = .9873$, so $r^2 = (.9873)^2 = .975$. For the regression in part (b), r is also equal to .9873 (since reversing x and y has no effect on the formula for r). So, both regressions have the same coefficient of determination. For the regression in part (a), we conclude that about 97.5% of the observed variation in rate can be attributed to the approximate linear relationship between speed and rate. In part (b), we conclude that about 97.5% of the optimized linear relationship between rate and speed.
- 50. Values for the vertical intercept and the slope will be changed using the same constant that was used to change y. For example, suppose that the least squares equation used to predict speed (in ft/sec) from stride rate was:

 $\hat{y} = -20.0514 + 12.1149x$

Then, since 1 m is about 3.2808 ft, if speed was expressed in m/sec instead of ft/sec, the new least squares equation would be:

$$\hat{y} = \left(\frac{-20.0514}{3.2808}\right) + \left(\frac{12.1149}{3.2808}\right)x$$

 $\Rightarrow \hat{y} = -6.1117 + 3.6927x$

Chapter 3

More generally, if each y value in the sample is multiplied by the same number c, the slope and vertical intercept of the least squares line will also change by multiplying each of them by c.

Chapter 3

53. (a) The curvature that is apparent in the plot of y versus x (see below) indicates that merely fitting a straight-line to the data would not be the best strategy. Instead, one should search for some transformation of the x or y data (or both) that would give a more linear plot.



(b) The plot below shows the graph of ln(y) versus 1/x. Because it appears to be approximately linear, a straight-line fit to such data should provide a reasonable approximation to the relationship between the two variables.



The following Minitab printout shows the results of fitting a regression line to the transformed data. From the printout, the prediction equation is $\ln(y) \approx -7.2557 + 8328.4(1/x)$. The r² value of 95.3% indicates that the fit is quite good. When temperature is 720 (i.e., x = 720), the equation gives a predicted value of $\ln(y) \approx -7.2557 + 8328.4(1/720) = 4.31152$. Exponentiating both sides gives a predicted y value of $y \approx e^{4.31152} = 74.6$.

```
The regression equation is logey = - 7.26 + 8328 recipx
```

Predictor	Coef	StDev	Т	P	
Constant	-7.2557	0.9670	-7.50	0.000	
recipx	8328.4	651.1	12.79	0.000	
S = 0.3882	R–Sq	= 95.3%	R-Sq(adj)	= 94.8%	
Analysis of	Variance				
Source	DF	SS	MS	F	Р
Regression	1	24.663	24.663	163.63	0.000
Error	8	1.206	0.151		
Total	9	25.869			

Chapter 5 Probability and Sampling Distributions

Section 5.1

2. A Venn diagram of these three events is



a) The event 'at least one plant is completed by the contract date' is represented by the shaded area covered by all three circles:



b) The event 'all plants are completed by the contract date' is the shaded area where all three circles overlap:



c) The event 'none of the plants is completed by the contract date' is the complement of the shaded area in (a):



d) The event 'only the plant at site 1 is completed by the contract date' is shown shaded:



e) The event 'exactly one of the three plants is completed by the contract date is:



f) The event 'either the plant at Site 1 or Site 2 or both plants are completed by the contract date' is:



7. The event A and B is the shaded area where A and B overlap in the following Venn diagram. Its complement consists of all events that are either not in A or not in B (or not in both). That is, the complement can be expressed as A' or B'.



Section 5.2

8. (a) There are 4 possible samples of size 2 that contain item A: {A,B}, (A,C}, {A,D}, and {A,E}. In general, there are exactly 10 possible distinct samples of size 2 that one could draw from a group of 5 items (use a counting procedure similar to that in problem 5.1). Therefore, the probability of a sample containing the defective item A is 4/10 = .40, or, 40%.

(b) There is one chance in 5 that the first inspector will discover item A. Regardless of the outcome, the remaining items will be sent to the second inspector since neither inspector knows in advance how many defectives there may be. If item A was selected during the first inspection, then the second inspector has no chance (probability 0) of finding item A. If item A is not selected during the first inspection, then the second inspector has one chance in 4 of finding it. Using the general definition of probability as the 'percentage of the time that an event occurs', 20% of the time item A will be discovered by the first inspector), the second inspector will catch item A 25% of the time. That is, (.25)(.80) = .20, or 40% of the time item A will be caught during the second inspection. Together, both inspections will catch item A with a probability of 20% + 20% = 40%.

(c) The answers to (a) and (b) are the same, and they will remain the same regardless of the size of the sample taken. Conceptually, the two methods are equivalent since the testing in Method 1 would actually be conducted by examining one of the two items chosen and, after that is done, proceeding to test the second item. Physically, this procedure is no different from that in Method 2 (in which two items are tested, one after the other), except different inspectors are used for the tests in Method 2

9. (a) The total of 724 +751 solder joints overstates the actual number found, since 316 solder joints were found by both inspectors. To avoid such double-counting (i.e., the 316 is part of both the 724 joints found by inspector A and the 751 joints found by Inspector B), 316 should be subtracted from the raw totals, which means that 724 + 751 - 316 = 1,159 distinct joints were identified by the inspectors together. The important point to note in this problem is that the events 'Inspector A finds a defective solder joint' and 'Inspector B finds a defective solder joint' are not necessarily mutually exclusive, so we cannot simply add the numbers of joints (or, equivalently, the probabilities of finding defective joints) for both inspectors.

(b) A and B' contains 724 - 316 = 408 solder joints.

- 10. For i = 1, 2, 3, ..., 10, let Ai denote the event 'component i functions correctly'. The problem indicates that $P(A_i) = .999$ for each component which, by the law of complementary events, means that each $P(A_i') = 1$ -.999 = .001. For a series system built from these 10 components to function correctly, all ten must function, so $P(system functions correctly) = P(A_1 and A_2 and A_3 and ... and A_k)$. According to the problem, this probability exceeds $1 [P(A_1') + P(A_2') + P(A_3') + ... + P(A_k')] = 1 [(.001) + (.001) + ... + (.001)] = 1 10(.001) = 1 .01 = .99$. That is, there is at least a 99% probability that the system will function correctly.
- 11. Letting A_i denote the event that the ith component fails (note that this is different from the definition of A_i used in problem 5.10), the probability that the entire series system fails is denoted by $P(A_1 \text{ or } A_2 \text{ or } A_3 \text{ or } ...$

or A_k). Given that each $P(A_i) = .01$, the problem states that $P(A_1 \text{ or } A_2 \text{ or } A_3 \text{ or } ... \text{ or } A_k) \le P(A_1) + P(A_2) + P(A_3) + ... + P(A_k) = 5(.01) = .05$. That is, there is at most a 5% chance of system failure.

Section 5.3

13. (a) Note that this question is equivalent to asking 'what is the probability A (or B, etc.) is chosen *given* that we know E is not chosen'. That is, the new probabilities are now conditional probabilities, where the conditioning event is that E is not chosen. Therefore,

 $\begin{array}{l} P(A \mid E') = P(A \text{ and } E')/P(E') = P(A)/P(E') = .20/(1-.10) = .20/.90 = 20/.90 \\ P(B \mid E') = P(B \text{ and } E')/P(E') = P(B)/P(E') = .25/(1-.10) = .25/.90 = 25/.90 \\ P(C \mid E') = P(C \text{ and } E')/P(E') = P(C)/P(E') = .15/.(1-.10) = .15/.90 = 15/.90 \\ P(D \mid E') = P(D \text{ and } E')/P(E') = P(D)/P(E') = .30/.(1-.10) = .30/.90 = 30/.90 \\ \end{array}$

Note that the four probabilities above do add exactly to 1, reflecting the fact that A, B, C, and D are the only four choices possible now that company E has been eliminated.

(b) The probability of *not* choosing companies B, D, or E is 1 -(.25+.30+.10). Following the same reasoning as in part (a), the two revised probabilities of choosing A or C are:

P(A | B, D, E not chosen) = P(A)/(1-(.25+.30+.10) = .20/.35 = 20/35. P(C | B, D, E not chosen) = P(C)/(1-(.25+.30+.10) = .15/.35 = 15/.35.

Again, the revised probability sum to 1 since there are now only two choices allowed.

- 17. The probabilities of independent events A and B must satisfy the equation P(A and B) = P(A)P(B). If A and B were also mutually exclusive, then P(A and B) would equal 0, which would mean that P(A)P(B) = P(A and B) = 0. This would require that at least one of A or B have zero probability of occurring. Although this is technically possible, most events of interest have non-zero probabilities, making P(A)P(B) non-zero. It is therefore impossible for independent events with non-zero probabilities to be mutually exclusive.
- 21. Let A denote the event that components 3 and 4 *both* work correctly and let B denote the event that *at least one* of components 1 or 2 works correctly. Then P(systems works) = P(A or B). From the general addition law, P(A or B) = P(A) +P(B) P(A and B). Because all components act independently of one another, P(A) = P(3 and 4 work) = P(3 works)P(4 works) = (.9)(.9) = .81. P(B) = P(1 or 2 works) = P(1 works) + P(2 works) P(1 and 2 work) = .9 + .9 (.9)(.9) = .99. Finally, events A and B are independent since A involves only components 3 and 4, whose actions are independent of components 1 and 2, so P(A and B) = P(A)P(B) = (.81)(.99) = .8019. Therefore, P(A or B) = P(A) +P(B) P(A and B) = .81 + .99 .8019 = .9981.
- 24. (a) Let S_i denote the event that the ith point signals a problem with the manufacturing process. Then, the probability that *none* of the 10 points give such a signal is P(no signals) = P(S_1' and S_2' and ... S_{10}') = P(S_1')P(S_2') ...P(S_{10}') = (1-.01)^{10} = .90438. Therefore, the probability of having at least one point signal a problem is P(at least one signal) = 1 P(no signals) = 1 .90438 = .0956.

(b) P(at least one in 25 signals a problem) = $1 - (1 - .01)^{25} = .2222$.

Section 5.4

Section 5.5

5.45. (a) x has a binomial distribution with n = 5 and $\pi = .05$. Writing $P(.05-.01 \le p \le .05+.01)$ in terms of x, we find $P(.05-.01 \le x/5 \le .05+.01) = P((.04)5 \le x \le (.06)5) = P(.2 \le x \le .3)$. Because x can only have integer values, there is no x between .2 and .3, so the probability of this event is 0.

(b) For n= 25, P(.04 ≤ p ≤ .06) = P((.04)25 ≤ x ≤ (.06)25) = P(1 ≤ x ≤ 1.5) =

$$P(x = 1) = {\binom{25}{1}} (.05)^{1} (.95)^{24} = 0.36498.$$

5.47. (a) x ='disconnect force' has a uniform distribution on the interval [2,4]. M is the maximum of a sample of size n = 2 from the uniform density on [2,4]. The larger of two items randomly selected from the interval [2,4] should, on average, tend to be closer to the upper end of the interval.

- (b) Using the same reasoning as in part (a), the largest value in a sample of n = 100 will, most likely, be even closer to the upper end of the interval [2,4] than is the largest value in a sample of size n =2. So the average of all M's based on n=100 ought to be larger than the average value of all M's based on n = 2.
- (c) For larger samples (e.g., n = 100), the maximum values will usually be fairly close to the upper endpoint of 4, which means that the variability amongst such values will tend to be small. For smaller samples (e.g., n = 2), it is easier for the value of M to wander over the interval [2,4], which means that the variability among these values will be larger than for n = 100.

Section 5.6

48.

x = diameter of a piston ring $\mu = 12 \ cm \ \sigma = .04 \ cm$ (a) $\mu_{\bar{x}} = \mu = 12 \ cm$ $\sigma_{\bar{x}} = \left(\frac{\sigma}{\sqrt{n}}\right) = \left(\frac{.04}{\sqrt{n}}\right)$

(b) When n = 64
$$\mu_{\bar{x}} = 12 \ cm \ \sigma_{\bar{x}} = \left(\frac{.04}{\sqrt{64}}\right) = .005$$

(c) The mean of a random sample of size 64 is more likely to lie within .01 cm of μ , since $\sigma_{\bar{x}}$ is smaller. 5.49.

(a)
$$\mu_p = \pi = .80. \ \sigma_p = \sqrt{\frac{\pi (1 - \pi)}{n}} = \sqrt{\frac{.80(1 - .80)}{25}} = .08$$

(b) Since 20% do not favor the proposed changes (so $\pi = .20$), the mean & standard deviation of the $\pi(1-\pi)$.20(1 - .20)= 1 25

sampling distribution of this proportion are $\mu_p = \pi = .20$ and $\sigma_p = \sqrt{\frac{n}{n}}$ = .08.

(c.) For n = 100 and π = .80, $\mu_p = \pi$ = .80. $\sigma_p = \sqrt{\frac{\pi (1-\pi)}{n}} = \sqrt{\frac{.80(1-.80)}{100}}$ = .04. Notice that it was necessary to <u>quadruple</u> the sample size (from n=25 to n=100) in order to cut σ_p in half (from $\sigma_p = .08$ to $\sigma_p = .04$).

5.52. (a) x = the weight of a bag of fertilizer $\mu = 50 \text{ lbs} \qquad \sigma = 1 \text{ lb} \qquad n = 100$

$$P(49.75 \le \overline{x} \le 50.25) \approx P\left(\frac{49.75 - 50}{1/\sqrt{100}} \le z \le \frac{50.25 - 50}{1/\sqrt{100}}\right)$$

$$= P(-2.5 \le z \le 2.5) = P(z \le 2.5) - P(z \le -2.5)$$

=.9938 - .0062 = .9876

(b) $\mu = 49.8 \text{ lbs} \quad \sigma = 1 \text{ lb} \quad n = 100$

$$P(49.75 \le x \le 50.25) \approx P(-0.5 \le z \le 4.5) = P(z \le 4.5) - P(z \le -0.5)$$

$$=1 - .3085 = .6915$$

5.53

(a)

x = 'lifetime of battery" has a normal density with $\mu = 8$ hours and $\sigma = 1$ hour. Therefore, P(average of 4 9 – 8

exceeds 9 hours) = P($\overline{x} > 9$) = P($z > \frac{1}{\sqrt{4}}$) = P(z > 2) = 1 - P(z < 2) = 1 - .9772 = .0228.

(b) Having the total lifetime of 4 batteries exceeds 36 hours is the same thing as having their average exceed 9, so the probability of this event is .0228, the same as in part (a).

$$\frac{T_0/4-8}{1/\sqrt{1}}$$

(c) $.95 = P(T > T_0) = P(T/4 > T_0/4) = P(\overline{x} > T_0/4) = P(z > \sqrt{4}) =$ $P(z > T_0/2-16)$. For a standard normal distribution, $P(z > -1.645) \approx .95$, so we must have $T_0/2-16 = -1.645$, which gives $T_0 = 28.71$ hours.

5.56. (a)
$$\mu = \sum xp(x) = (0)(8) + (1)(1) + (2)(05) + (3)(05) = .35$$

$$\sigma = \sqrt{\sum (x - \mu)^2 p(x)}$$

= $\sqrt{(0 - .35)^2 (.8) + (1 - .35)^2 (.1) + (2 - .35)^2 (.05) + (3 - .35)^2 (.05)}$
= .7921
(b) $n = 64 \ \mu_{\bar{x}} = \mu = .35$
 $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = (7921/\sqrt{64}) = .099$

(c)
$$P(\bar{x} > 1) \approx P(z > \frac{1 - .35}{.099}) = P(z > 6.57) \approx 0$$

(a) Let p = 'proportion of resistors exceeding 105 Ω '. Then the sampling distribution of p is approximately normal with $\mu_p = \pi = 0.02$ and $\sigma_p = \sqrt{\frac{\pi (1-\pi)}{n}} = \sqrt{\frac{.02(1-.02)}{100}} = 0.014$.

(b)
$$P(p < .03) = P(z < \frac{.03 - .02}{.014}) = P(z < .71) = 0.7611.$$

5.58

x =length of an object $\sigma = 1 mm n = 2$

$$P(-2 \le x - \mu \le 2) = P\left(\frac{-2}{\sigma/\sqrt{n}} \le z \le \frac{2}{\sigma/\sqrt{n}}\right)$$

$$= P(-2.83 \le z \le 2.83) = P(z \le 2.83) - P(z \le -2.83)$$

Supplementary Problems – Chapter 5

5.72. (a) x = 'battery voltage' has a mean value of $\mu = 1.5$ and a standard deviation of $\sigma = .2$ volts. The sampling distribution of \overline{x} (based on n = 4) has a mean value of $\mu_{\overline{x}} = \mu = 1.5$ and a standard error of $\sigma_{\overline{x}} = \frac{\sigma}{\sqrt{n}} = \frac{.2}{\sqrt{4}} = .1$

(a) T is related to \overline{x} by the equation $T = 4\overline{x}$, so the mean of T should be 4 times the mean of \overline{x} . That is, $\mu_T = 4 \mu_{\overline{x}} = 4(1.5) = 6$. Similarly, the standard deviation of T should be 4 times that of \overline{x} , so $\sigma_T = 4$ $\sigma_{\overline{x}} = 4(.1) = .4$

5.70

x = resistance n = 5 μ =100 ohms σ =1.7 ohms

First, we must assume that the population density for x is symmetric. Then, we can use the Central Limit Theorem to proceed.

(a) Note: The question should read: "What is the probability that the *average* resistance in the circuit exceeds 105 ohms?"

$$P(\bar{x} > 105) \approx P\left(\bar{x} > \frac{105 - 100}{1.7/\sqrt{5}}\right) = P(z > 6.58) \approx 0$$

$$P(T > 511) + P(T < 489) = P\left(\overline{x} > \frac{511}{5}\right) + P\left(\overline{x} < \frac{489}{5}\right)$$
(b)

$$P(\overline{x} > 102.2) + P(\overline{x} < 97.8) = P\left(z > \frac{102.2 - 100}{1.7/\sqrt{5}}\right) + P\left(z < \frac{97.8 - 100}{1.7/\sqrt{5}}\right)$$

$$= P(z > 2.89) + P(z < -2.89) = 1 - P(z < 2.89) + P(z < -2.89)$$

$$= 1 - .9981 + .0019 = .0038$$

(c) We know that:
$$P(-1.96 \le z \le 1.96) = .95$$

 $\frac{\overline{x} - 100}{1.7/\sqrt{n}} = 1.96 \Rightarrow \overline{x} = 100 + 1.96 \left(\frac{1.7}{\sqrt{n}}\right)$
So, $T = 510 = n\overline{x} \Rightarrow n = \frac{T}{\overline{x}} = \frac{510}{\overline{x}}$
Also, $n = \frac{510}{100 + 1.96 \left(\frac{1.7}{\sqrt{n}}\right)}$
Thus:

Solving for n produces a value just over 5.

Chapter 7 Estimation and Statistical Intervals

Section 7.1

7.4. (a) Recall:
$$\mu_p = \pi$$
 and $\sigma_p = \sqrt{\frac{\pi(1-\pi)}{n}}$

Even though π is unknown, we can still set an upper bound on σ_p .

When $\pi = .5$, σ_p is maximized. So, In our case, $n = 10 \Rightarrow \sigma_p = .15811$

So,
$$P(-.10 < (p - \pi) < .10) = P(\frac{-.10}{.15811} < z < \frac{.10}{.15811})$$

$$= P(-.63 < z < .63) = .4714$$

Section 7.2

- 7.7. (a) The entry in the 3.0 row and .09 column of the z table is .9990. Similarly, the entry for -3.09 is .0010. Therefore, the area under the z curve between -3.09 and +3.09 is .9990 .0010 = .9980. The confidence level is then 99.8%.
 - (b) Following the example in part (a), the z-table entries corresponding to z = -2.81 and z = +2.81 are .9975 and .0025, respectively. Therefore the area between these two z values is .9975 - .0025 = .9950. The confidence level is then 99.5%.
 - (c) The z-table entries corresponding to z = -1.4 and z = +1.44 are .9251 and .0749, respectively. Therefore, the area between these two z values is .9251 - .0749 = .8502. The confidence level is then 85.02%.
 - (d) The coefficient of s/\sqrt{n} is not written, but is understood to be 1.00. The z-table entries corresponding to z = -1.001 and z = +1.00 are .8413 and .1587, respectively. Therefore, the area between these two z values is .8413 .1587 = .6826. The confidence level is then 68.26%.
 - 7.8. (a) 98% of the standard normal curve area must be captured. This requires that 1% of the area is to be captured in each tail of the distribution. So, P(Z < -z critical value) = .01 and P(Z > z critical value) = .01. Thus, the z critical value = 2.33.

(b) 85% of the standard normal curve area must be captured. This requires that 7.5% of the area is to be captured in each tail of the distribution. So, P(Z < -z critical value) = .075 and P(Z > z critical value) = .075. Thus, the z critical value = 1.44.

(c)75% of the standard normal curve area must be captured. This requires that 12.5% of the area is to be captured in each tail of the distribution. So, P(Z < -z critical value) = .125 and P(Z > z critical value) = .125. Thus, the z critical value = 1.15.

(d) 99.9% of the standard normal curve area must be captured. This requires that .05% of the area is to be captured in each tail of the distribution. So, P(Z < -z critical value) = .0005 and P(Z > z critical value) = .0005. Thus, the z critical value is conservatively 3.32.

(Notice that, in this case, there are several possible z critical values listed in the standard normal table.)

7.10. (a) The sample mean is the mid-point of each of the confidence intervals. So, $\left(\frac{115.6+114.4}{2}\right) = 115$ To confirm: $\left(\frac{115.9+114.1}{2}\right) = 115$ The sample mean resonance frequency is 115 Hz.

- (b) The first interval has the 90% confidence level, (114.4, 115.6). We know this because it is the more narrow of the two intervals and a 90% confidence interval will be more narrow than a 99% confidence interval, given the same sample data.
- 7.11. (a) Decreasing the confidence level from 95% to 90% will decrease the associated z value and therefore make the 90% interval narrower than the 95% interval. (*Note: see the answer to Exercise 9 above*)
 - The statement is not correct. Once a particular confidence interval has been created/calculated, then the true mean is either in the interval or not. The 95% refers to the process of creating confidence intervals; i.e., it means that 95% of all the possible confidence intervals you could create (each based on a new random sample of size n) will contain the population mean (and 5% will not).
 - The statement is not correct. A confidence interval states where plausible values of the population mean are, not where the individual data values lie. In statistical inference, there are three types of intervals: **confidence intervals** (which estimate where a population mean is), **prediction intervals** (which estimate where a single value in a population is likely to be), and **tolerance intervals** (which estimate the likely range of values of the items in a population. The statement in this exercise refers to the likely range of all the values in the population, so it is referring to a tolerance interval, not a confidence interval.
 - (d) No, the statement is not exactly correct, but it is close. We *expect* 95% of the intervals we construct to contain μ, but we also expect a little variation. That is, in any group of 100 samples, it is possible to find only, say, 92 that contain μ. In another group of 100 samples, we might find 97 that contain μ, and so forth. So, the 95% refers to the *long run* percentage of intervals that will contain the mean. 100 samples/intervals is <u>not</u> the long run.
- 7.14. (a) A 95% two-sided confidence interval for the true average dye-layer density for all such trees is:

$$\overline{x} \pm (1.96) \left(\frac{s}{\sqrt{n}}\right)$$

$$1.28 \pm (1.96) \left(\frac{.163}{\sqrt{69}}\right)$$

$$1.028 \pm 0.03846$$

$$(.9895, 1.0665)$$

Interpretation 1: We are 95% confident that the (true) average is between 0.9895 and 1.0665. Interpretation 2: There is a 95% probability that a random 95% CI, computed using our formula, will yield an interval that covers the (true) average.

(b)
$$n = \left[\frac{1.96s}{B}\right]^2 = \left[\frac{1.96(.16)}{.025}\right]^2 \approx 158$$

A sample size of 158 trees would be required. (Note: The researchers wanted an interval width of .05. So, the bound on the error of estimation, B, is half of the width. B = .025)

7.18

(a) The entry in the z table corresponding to z = .84 is .7995, so the confidence level is 79.95% or, approximately, 80%.

(b) The entry in the z table corresponding to z = 2.05 is .9798, so the confidence level is 97.98% or, approximately, 98%.

(c) The entry in the z table corresponding to z = .67 is .7486, so the confidence level is 74.86% or, approximately, 75%.

7.19

A 95% upper confidence bound for the true average charge-to-tap time is:

$$\bar{x} + (1.645) \left(\frac{s}{\sqrt{n}} \right)$$

$$382.1 \pm (1.645) \left(\frac{31.5}{\sqrt{36}} \right)$$

$$(382.1 + 8.64) = 390.74 \text{ min}$$

That is, with 95% confidence, the value of μ lies in the interval (0 min, 390.74 min).

7.21

A 90% lower confidence bound for the true average shear strength is:

$$\overline{x} - (1.28) \left(\frac{s}{\sqrt{n}} \right)$$

$$4.25 - (1.28) \left(\frac{1.30}{\sqrt{78}} \right)$$

$$(4.25 - .188) = 4.062 \, kip$$

That is, with 90% confidence, the value of μ lies in the interval (4.062, ∞).

Section 7.3

7.27. (a) Following the same format used for most confidence intervals, i.e., statistic \pm (critical value) (standard error), an interval estimate for π_1 - π_2 is:

$$(p_1 - p_2) \pm z \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}}$$

(b) The response rate for no-incentive sample: $p_1 = 75/110 = .6818$, while the return rate for the incentive sample is $p_2 = 66/98 = .6735$. Using z = 1.96 (for a confidence level of 95%), a two-sided confidence interval for the true (i.e., population) difference in response rates π_1 - π_2 is:

	.6818(16818)	(.6735)(16735)
(.68186735) ±(1.96)	110	98
= .0083 ± .1273 =(119,	.1356).	

The fact that this interval contains 0 implies that we cannot say anything about the equality of the proportions. In other words, including incentives in the questionnaires may or may not have a significant effect on the response rate.

(c) Let p_i denote the sample proportion by adding 1 success and 1 failure to the *i*th sample. We calculate $\tilde{p}_i = (x+1)/(n+2)$, where x is the number of successes (or failures, whichever is desired) in the sample. Then: $\tilde{p}_1 = (75+1)/(110+2) = .67857$ and $\tilde{p}_2 = (66+1)/(98+2) = .67$. Using the format of the equation given in part (a) above, we have the following 95% confidence interval:

$$(\widetilde{p}_{1} - \widetilde{p}_{2}) \pm z \sqrt{\frac{\widetilde{p}_{1}(1 - \widetilde{p}_{1})}{n_{1} + 2}} + \frac{\widetilde{p}_{2}(1 - \widetilde{p}_{2})}{n_{2} + 2}$$

$$(.67857 - .67) \pm 1.96 \sqrt{\frac{.67857(1 - .67857)}{110 + 2}} + \frac{.67(1 - .67)}{98 + 2}$$

$$(.00857) \pm 1.96 \sqrt{.004158} \implies .00857 \pm .1264 \implies (-.1178, .1350)$$

This interval contains 0; so, it is not clear if including incentives in the questionnaire has an effect on the response rate. This is the same conclusion as in part (b).

7.29. (a) As in Exercise 25, the usual confidence interval format *statistic* ± (*critical value*)(*standard error*) gives a confidence interval for:

$$\ln(\pi_1/\pi_2) \text{ of: } \ln(p_1/p_2) \pm z \sqrt{\frac{n_1 - u}{n_1 u} + \frac{n_2 - v}{n_2 v}}$$

(b) Since we want to estimate the ratio of return rates for incentive group to the non-incentive group, we will call group 1 the incentive group (to match the subscripts in the formula above). The number of returns for the non-incentive group is v = 75 out of $n_2 = 110$, so $p_2 = .6818$. For the incentive group, u = 78 out of $n_1 = 100$, so $p_1 = .78$. The 95% confidence interval for $\ln(\pi_1/\pi_2)$ is:

 $\sqrt{\frac{100-78}{100(78)} + \frac{110-75}{110(75)}} = .1346 \pm .1647$ = [-.0301, .2993]. To find a 95% interval for π_1/π_2 , we exponentiate both end points of this interval, [e^{-.0301}, e²⁹⁹³] = [.9702, 1.3489], or, about [.970, 1.349]. Because the 95% confidence interval includes 1, then at 95% confidence level, we cannot say anything about the equality of pi1 and pi2. As in problem 25, there is insufficient evidence from this data to suggest that incentive has an effect on the questionnaire return rate. In other words, we don't know if incentive has an effect or not.

$$n = \pi \left(1 - \pi\right) \left[\frac{2.575}{B}\right]^2$$

Since an estimate of π is not provided, a conservative estimate is .50. (However, it seems improbable that as many as 50% of the coffeepot handles will be cracked.)

$$n = (.50)(.50) \left[\frac{2.575}{.1} \right]^2$$

*n≈*166

A sample of 166 coffeepot handles from the shipment should be inspected.

7.31

For 90% confidence, the associated z value is 1.645. Since nothing is known about the likely values of π we use .25,

the largest possible value of $\pi(1-\pi)$, in the sample size formula: $n = (.25) \left(\frac{1.645}{.05}\right)^2 = 270.6$. To be conservative, we round this value up to the next highest integer and use n = 271.

7.35

Let μ_1 denote the average toughness for the high-purity steel and let μ_2 denote the average toughness for the commercial purity steel. Then, a lower 95% confidence bound for μ_1 - μ_2 is given by: $(\overline{x_1} - \overline{x_2}) - z$

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = (65.6-59.2) - (1.645) \sqrt{\frac{(1.4)^2}{32} + \frac{(1.1)^2}{32}} = 6.4 - .518 = 5.882.$$
 Because this lower interval bound exceeds 5, it gives a reliable indication that the difference between the population toughness levels does exceed 5.

7.36

We need the following equation to be true:

$$B = \left(\text{z critical value}\right) \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}}$$

So,
$$.5 = \left(1.96\right) \sqrt{\frac{2^2}{n} + \frac{2^2}{n}}$$

Solving for n: $n \approx 123$

A sample of 123 batteries of each type should be taken.

Section 7.4

7.42

Recall that this text's definition of upper and lower quartiles may differ a little from the definitions used by some other authors. We define the upper and lower quartiles to be the medians of the lower and upper halves of the data (for n odd, the median is included in both halves). This usually gives results that are very close to the actual (or estimated) quartiles used by other authors, but there are often small differences. For example, in this exercise the medians of the lower and upper halves of the data are: lower quartile = median of lower half (including the median, 437) = 425; upper quartile = median of the upper half (including median) = 448. The software package Minitab, however, estimates the lower quartile as the $.25(n+1)^{th} = .25(17+1)^{th} = 4.5^{th}$ item in the sorted list (here, 422 and 425)

are the 4th and 5th items in the sorted list, and the 4.5th value is defined to be their average, which is 423.5). So, be careful when comparing your answers to those in various software packages.

(a) Using 425 and 448 as the lower and upper quartiles, the IQR = 448 - 425 = 23. To check for outliers, we calculate the values 425 + 1.5(IQR) = 425 - 1.5(23) = 390.5 and 448 + 1.5(IQR) = 448 - 1.5(23) = 482.5. Since the maximum and minimum in the data are 465 and 418, which are inside the limits just calculated, there are no outliers in the data. The median of the data is 437 and a boxplot of the data appears below:



(b) A quantile plot can be used to check for normality. Refer to the answer to Exercise 43 of Chapter 2 for an easy method of creating a quantile plot in Minitab. The quantilile plot below shows a fairly strong linear pattern, supporting the assumption that this data came from a normal population:



(c) Since the sample size is small (n = 17), we use an interval based on the t distribution. The sample mean and standard deviation are 438.29 and 15.144, respectively. For n -1 = 17 -1 = 16 df, the critical t value for a 2-sided 95% confidence interval is 2.120 (from Table IV). Therefore, the desired confidence interval is $438.29 \pm (2.120)^{15.144}/\sqrt{17} = 438.29 \pm 7.787 = [430.5, 446.1]$. This interval suggests that 440 (which is inside the interval) is a plausible value for the mean polymerization. The value of 450, however, is not plausible, since it lies outside the interval.

- 7.48. (a) Confidence level = area between -.687 and 1.725 = .95 - .25 = .70, or, 70%. Note that the value 1.725is in Table IV and is associated with a cumulative area of .95 (for df = 20).
 - (b) Confidence level = area between -.860 and 1.325 = .90 .20 = .70, or, 70%. Note that the value 1.325 is in Table IV and is associated with a cumulative area of .90 (for df = 20).
 - (c) Confidence level = area between -1.064 and 1.064 = (1-2(.15)) = .70, or, 70%. Note that the symmetry of the t distribution allows us to say there is evidence that the area to the left of 1.064 is also .15.

All three intervals have 70% confidence, but they have different widths. The interval in part (c) has the smallest width of 2(1.064) $\sqrt[s]{\sqrt{n}}$, and is therefore the best choice among the three.

Section 7.5

7.51 (a)



The most notable feature of these boxplots is the larger amount of variation present in the midrange data as compared to the high-range data. Otherwise, both boxplots look reasonably symmetric and there are no outliers present.

(b) Minitab output:

> sample n mean

sample standard deviation

Mid-range	17	438.3	15.1
High-range	11	437.45	6.83

A 95% Confidence Interval for (μ mid range - μ high range) is

(Note: df = 23.)

The above analysis was performed by Minitab. The confidence interval was computed as follows:

$$(438.3 - 437.45) \pm (2.069) \sqrt{\frac{(15.1)^2}{17} + \frac{(6.83)^2}{11}}$$

using df = 23, resulting in:

Since plausible values for $(\mu_1 - \mu_2)$ are both positive and negative (i.e., the interval spans zero) we would say that there is not sufficient evidence from data to suggest that μ_1 and μ_2 differ.

7.58

This is paired data with n = 60, $\overline{d} = 4.47$, and $s_d = 8.77$. The critical t value associated with df = n-1 = 60-1 = 59 and 99% confidence is approximately 2.66 (Table IV). Alternatively, since df = 59 is large, we could simply use the 99% z value of 2.58, which we do in the following calculation: $\overline{d} \pm t \sqrt[s]{\sqrt{n}} = 4.47 \pm (2.58)^{\frac{8.77}{\sqrt{60}}} = 4.47 \pm 2.92 = [1.55, 7.39]$. Therefore, we estimate that the average blood pressure in a dental setting exceeds the average blood pressure in a medical setting by between 1.55 and 7.39. The interval does not contain 0, which suggests that the true average pressures are indeed different in the two settings.

Supplementary Exercises

7.77

The center of any confidence interval for μ_1 - μ_2 is always \overline{x}_1 - \overline{x}_2 , so \overline{x}_1 - $\overline{x}_2 = (-473.3 + 1691.9)/2 = 609.3$. Furthermore, the half-width of this interval is [1691.9 - (-473.3)]/2 = 1082.6. Equating this value to

the expression for the half-width of a 95% interval, 1082.6 = (1.96) $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$, we find $\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = 1082.6/1.96 = 552.35$. For a 90% interval, the associated z value is 1.645, so the 90% confidence interval

is then
$$\overline{x}_1 - \overline{x}_2 \pm (1.645) \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = 609.3 \pm (1.645)(552.35) = 609.3 \pm 908.6 = [-299.3, 1517.9].$$

7.91. (a) We shall construct a 95% confidence interval for the true proportion of all American adults who are obese. Here, n = 4115, and $p = 1276/4115 \approx .310085$. With 95% confidence, we use the value $z^* = 1.96$ for a two-sided confidence interval:

$$\frac{p + \frac{(z^*)^2}{2n} \pm z^* \sqrt{\frac{p(1-p)}{n} + \frac{(z^*)^2}{4n^2}}}{1 + \frac{(z^*)^2}{n}} = \frac{.310085 + \frac{1.96^2}{2(4115)} \pm 1.96\sqrt{\frac{.310085(1-.310085)}{4115} + \frac{1.96^2}{4(4115^2)}}}{1 + \frac{1.96^2}{4115}}$$
$$\Rightarrow \frac{.31055 \pm 1.96(.007214)}{1.0009} \Rightarrow (.296, .324)$$

We are therefore 95% confident that between 29.6% and 32.4% of all American adults are obese.

(b) We shall construct a one-sided 95% confidence interval for the true proportion of all American adults who are obese. Here, n = 4115, and $p = 1276/4115 \approx .310085$. With 95% confidence, we use the value $z^* = 1.645$ for a one-sided confidence interval:

$$\frac{p + \frac{(z^*)^2}{2n} - z^* \sqrt{\frac{p(1-p)}{n} + \frac{(z^*)^2}{4n^2}}}{1 + \frac{(z^*)^2}{n}} = \frac{.310085 + \frac{1.645^2}{2(4115)} - 1.645 \sqrt{\frac{.310085(1-.310085)}{4115} + \frac{1.645^2}{4(4115^2)}}}{1 + \frac{1.645^2}{4115}}$$
$$\Rightarrow \frac{.310412 - 1.645(.007213)}{1.0007} \Rightarrow .2983$$

Since the interval value dips below 30%, there is no evidence from data in favor of the claim that the 2002 percentage is more than 1.5 times the 1998 percentage.

Chapter 8 Testing Statistical Hypotheses

Section 8.1

8.1

(a) Yes, $\sigma > 100$ is a statement about a population standard deviation, i.e., a statement about a population parameter.

- (b) No, this is a statement about the statistic \overline{x} , not a statement about a population parameter.
- (c) Yes, this is a statement about the population median μ .
- (d) No, this is a statement about the statistic s (s is the sample, not population, standard deviation.
- (e) Yes, the parameter here is the ratio of two other parameters; i.e, σ_1/σ_2 describes some aspect of the populations being sampled, so it is a parameter, not a statistic.
- (f) No, saying that the difference between two samples means is -5.0 is a statement about sample results, not about population parameters.
- (g) Yes, this is a statement about the parameter λ of an exponential population.
- (h) Yes, this is a statement about the proportion π of successes in a population.

(i) Yes, this is a legitimate hypothesis because we can make a hypothesis about the population distribution [see (4) at the beginning of this section].

(j) Yes, this is a legitimate hypothesis. We can make a hypothesis about the population parameters [see (3) at the beginning of this section].

8.2

The purpose of inspecting pipe welds in nuclear power plants is to determine if the welds are defective (i.e., do not conform to specifications). The benefit of the experimental design where H₀: $\mu = 100$ and H_a: $\mu > 100$ over the design H₀: $\mu = 100$ and H_a: $\mu < 100$ can be understood in terms of Type 1 and Type II error.

In the first design, a type I error corresponds to saying there is evidence from the data that supports the claim that the mean weld population strength is greater than 100 when in fact the mean population strength is less than or equal to 100. Under this design, a type II error corresponds to saying there is not evidence from the data that supports the claim that the mean weld population strength is greater the 100 when in fact such evidence does exist.

In the second design, a type I error corresponds to saying there is evidence from the data that supports the claim that the mean weld population strength is less than 100 when in fact the mean population strength is greater than or equal to 100. Under this design, a type II error corresponds to saying there is not evidence from the data that supports the claim that the mean weld population strength is less the 100 when in fact such evidence does exist.

Thus, under the first design, setting a suitably small significance level will allow you to minimize the chance of claiming that the weld population conforms to specification based on the data, when in fact such

a claim is unsupported. Conversely, a small significance level in the second design will minimize the chance of claiming that the weld population fails to conform to specification based on the data, when in fact such a claim is unsupported. The first instance of Type I error is a danger to public safety; the second only to reputation. Thus, the first design is superior.

8.3 Let μ denote the average amperage in the population of all such fuses. Then the two relevant hypotheses are H₀: $\mu = 40$ (the fuses conform to specifications) and H_a: $\mu \neq 40$ (the average amperage either exceeds 40 is less than 40).

8.4

Let σ denote the population standard deviation of sheath thickness. The relevant hypotheses are:

 $H_0: \sigma = .05 \ versus \ H_a: \sigma < .05$

This is because the company is interested in obtaining conclusive evidence that $\sigma < .05$. A Type I error would be: concluding that the true standard deviation of sheath thickness is less than .05mm when, in fact, it is not.

A Type II error would be: concluding that the true standard deviation of sheath thickness is equal to .05mm when, in fact, it is really less than .05mm.

8.5

Let μ denote the average breaking distance for the new system. The relevant hypotheses are $H_0:\mu = 120$ versus $H_a:\mu < 120$, so implicitly H_0 really says that $\mu \ge 120$. A Type I error would be: *concluding that the new system really does reduce the average breaking distance (i.e., rejecting H₀) when, in fact (i.e., when H₀ is true) it doesn't. A Type II error would be: <i>concluding that the new system does not achieve a reduction in average breaking distance (i.e., not rejecting H₀) when, in fact (i.e., when H₀ is false) it actually does.*

8.6

Let μ denote the true average compressive strength of the mixture. The relevant hypotheses are:

 $H_0:\mu = 1,300 \ versus \ H_a:\mu > 1,300$

A Type I error would be: concluding that the mixture meets the strength specifications when, in fact, it does not.

A Type II error would be: concluding that the mixture does not meet the strength specifications when, in fact, it does.

8.13

(a) This is a test about the population average μ = average silicon content in iron. The null hypothesis value of $\overline{x} - .85$

interest is $\mu = .85$, so the test statistic is of the form $z = \sqrt[s]{\sqrt{n}}$. From the wording of the exercise it seems that a 2-sided test is appropriate (since the silicon content is supposed to average .85 and not be substantially larger or smaller than that number), so the relevant hypotheses are H₀: $\mu = .85$ versus H_a: $\mu \neq .85$. We can verify that a 2-sided test was done by calculating the P-value associated with the z value of -.81 given in the printout: the area to the left of -.81 is .2090 (from Table I), so the 2-sided P-value associated with z = .81 is 2(.2090) = .418 $\approx .42$.

(b) The P-value of .42 is quite large, so we don't expect it to lead to rejecting H₀ for any of the usual values of α used in hypothesis testing. Indeed, P = .42 exceeds both α = .05 and α = .10, so in neither case

would this data lead to rejecting H_0 . It appears to be quite likely that the average silicon content does not differ from .85

- 8.14. (a) If H_a had been $\mu > 750$, the p-value would have been P(z > -2.14), which is clearly not equal to .016. In fact, the p-value would have been .9838.
 - (b) At a significance level of .05, since .016 < .05, we would reject H₀ and say there is evidence from data in favor of the claim that the true average lifetime is smaller than what is advertised. Therefore, the customer should not purchase the light bulbs. Whereas, at a significance level of .01, since .016 > .01, we would fail to reject H₀. That is, there is insufficient evidence to claim that the true average lifetime is smaller than what is advertised. Therefore, the customer should purchase the light bulbs.

8.16

Let μ denote the true average penetration. Since we are concerned about the specifications not being met, the relevant hypotheses are:

H₀: the specifications are met ($\mu = 50$), versus H_a: the specifications are not met ($\mu > 50$).

$$z = \left[\frac{52.7-50}{4.8/\sqrt{35}}\right] = 3.33$$

The test statistic is: $\lfloor \sqrt{35} \rfloor$

The corresponding p-value = P(z > 3.33) = .0004

Since p-value = $.0004 < \alpha$ = .05, we reject H₀ and say that there is evidence from data in favor of the claim that the true average penetration exceeds 50 mils. Thus, the specifications have not been met.

Section 8.2

- 20. (a) Let μ denote the true average writing lifetime. The wording in this exercise indicates that the investigators believe, a priori, that μ can't be less than 10 (i.e., μ≥ 10), so the relevant hypotheses are H₀: μ = 10 versus H_a: μ < 10.
 - (b) The degrees of freedom are d.f. = n-1 = 18-1 = 17, so the P-value = P(t < -2.3) = P(t > 2.3) = .017. Since P-value = $.017 < \alpha = .05$, we should reject H₀ and say that there is evidence from data that the design specification has not been satisfied.
 - (c) For t = -1.8, P-value = P(t<-1.8) = P(t>1.8) = .045, which exceeds $\alpha = .01$. Therefore, in this case H₀ would not be rejected. That is, there is not sufficient evidence from data to claim that the design specification is not satisfied.
 - (d) For t = -3.6, P-value = P(t < -3.6) = P(t > 3.6) = .001. This P-value is smaller than either $\alpha = .01$ or $\alpha = .05$, so in either case H₀ would be rejected. In fact, H₀ would be rejected for any value of α that exceeds .001. For such values of α , we would say that there is evidence from data to support the claim that the design specification has not been satisfied.

Let μ_1 denote the true average gap detection threshold for normal subjects and let μ_2 denote the true 29. average gap detection threshold for CTS subjects. Since we are interested in whether the gap detection threshold for CTS subjects exceeds that for normal subjects, a lower tailed test is appropriate. So, we test:

$$H_0: (\mu_1 - \mu_2) = 0$$
 versus $H_a: (\mu_1 - \mu_2) < 0$

Using the sample statistics provided, the test statistic is:

$$t = \left[\frac{1.71 - 2.53}{\sqrt{\frac{(53)^2}{8} + \frac{(87)^2}{10}}}\right]$$

$$t = -2.46 \approx -2.5$$

Using the equation for df provided in the section, the approximate df = 15.1, which we round down to 15.

The corresponding p-value = P(t < -2.5) = .012.

Since the p-value = $.012 > \alpha = .01$, we fail to reject H₀. We have insufficient evidence to claim that the true average gap detection threshold for CTS subjects exceeds that for normal subjects.

29. Let μ_1 denote the true average proportional stress limit for red oak and let μ_2 denote the average stress limit for Douglas fir. The wording of the exercise suggests that we are interested in detecting any differences between the two averages, which means a 2-sided test is appropriate, so we test H_0 : $\mu_1 - \mu_2 = 0$ versus H_a : μ_1 - $\mu_2 \neq 0$. The test statistic is:

$$t = \frac{\overline{x_1 - \overline{x_2}}}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_1^2}{n_2}}} = \frac{8.48 - 6.65}{\sqrt{\frac{.79^2}{14} + \frac{1.28^2}{10}}} = 1.83/.4565 = 4.01 \approx 4.0.$$

 $\frac{\frac{179^{-1}}{14} + \frac{1.28}{10}}{13} = .04344/.003135 = 13.85, \text{ which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = 13.85, \text{ which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = 13.85, \text{ which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = 13.85, \text{ which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = 13.85, \text{which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = 13.85, \text{which we round } \frac{\text{down}}{10} \text{ to d.f.} = .04344/.003135 = .04344/.00344/.00344/.00344/.00344/.00344/.00344/.00344/.00344/.00344/.00$ The approximate d.f. \approx 13. For a 2-sided test, we use Table VI to find the P-value = 2P(t > 4.0) = 2(.001) = .002. Since P-value = .002 is very small (smaller than the usual significance like .05 or .01), we reject H₀ and say that there is evidence from data that supports the claim that there is a difference between the two average stress limits.

40. Each sample is measured by two different methods so the data is paired. The observed differences (MSI minus SIB) are: 0.03, -0.51, -0.80, -0.57, -0.66, -0.63, -0.18, 0.01. The sample mean and standard deviation of these observations are $\overline{d} = -.4138$ and $s_d = .3210$. The wording in the exercise indicates that we are interested in detecting any difference between the two methods, so a 2-sided test is required. To test H₀: $\mu_d = 0$ versus H_a: $\mu_d \neq 0$, the test statistic is:

$$t = \frac{\frac{\overline{d} - 0}{s_d / \sqrt{n}}}{s_d / \sqrt{n}} = \frac{\frac{-.4138 - 0}{.3210 / \sqrt{8}}}{= -3.64 \approx -3.6.}$$

For d.f. = n-1 = 8-1 = 7, we use Table VI to find the P-value for this 2-sided test: P-value = 2P(t < -3.6) = 2(.004) = .008. At significance levels of either $\alpha = .05$ or $\alpha = .01$, H₀ would be rejected and we would say that there is evidence from data in favor of the claim that there is a difference between the two methods. However, at $\alpha = .001$ we would not reject H₀ since the P-value of .008 exceeds .001.

Section 8.3

45. Using the number 1 (for business), 2 (for engineering), 3 (for social science), and 4 (for agriculture), let π^i = the true proportion of all clients from discipline i. If the Statistics Department's expectations are correct, then the relevant null hypothesis is:

 $H_0: \pi_1 = .40 \ \pi_2 = .30 \ \pi_3 = .20 \ \pi_4 = .10 \ versus$

 H_a : The Statistics Department's expectations are not correct

Using the proportions in H₀, the expected number of clients are:

Client's discipline	Expected number of clients
Business	(120)(.40) = 48
Engineering	(120)(.30) = 36
Social science	(120)(.20) = 24
Agriculture	(120)(.10) = 12

Since all expected counts are at least 5, the chi-squared test can be used. The value of the χ^2 test statistic is:

$$\chi^{2} = \left[\sum \frac{(\text{observed} - \text{expected})^{2}}{\text{expected}} \right]$$
$$= \left[\frac{(52 - 48)^{2}}{48} + \frac{(38 - 36)^{2}}{36} + \frac{(21 - 24)^{2}}{24} + \frac{(9 - 12)^{2}}{12} \right] = 1.57$$

For df = (k - 1) = (4 - 1) = 3, the corresponding p-value = P($\chi^2 > 1.57$). Using Table VII, we find that the p-value > .10. Since the p-value is larger than $\alpha = .05$, we fail to reject H₀. We have no evidence to suggest that the Statistics Department's expectations are incorrect.

46. Using the numbering 1 (for Winter), 2 (Spring), 3(Summer), and 4 (Fall), let π_i = the true proportion of all homicides committed in the ith season. If the homicide rate <u>doesn't</u> depend upon the season, then we would expect the rates to be equal; i.e., H₀: $\pi_1 = \pi_2 = \pi_3 = \pi_4 = .25$. The alternative hypothesis in this case would be H_a: *At least one of the proportions does not equal .25*. Using the proportions in H₀, the expected numbers of homicides (out of n = 1361) are shown below the actual numbers from the problem:

Season	Winter	Spring	Summer	Fall	total #
Observed	328	334	372	327	1361
Expected	340.25	340.25	340.25	340.25	1361
χ^2 contribution	.4410	.1148	2.9627	.5160	

The χ^2 test statistic value is $\chi^2 = .4410 + .1148 + 2.9627 + .5160 = 4.0345$. For d.f = k-1 = 4-1 = 3, the value 4.0345 is smaller than any of the entries in Table VII (column df=3), so the P-value associated with 4.0345 must be larger than .10. Therefore, since P-value > .10 > .05 = α , H₀ is not rejected and we say that this data does not support the belief that there are different homicide rates in the different seasons.

48. Using the equation given for π_i , the relevant hypotheses to test are:

 $H_0: \pi_1 = .033, \pi_2 = .067, \pi_3 = .100, \pi_4 = .133, \pi_5 = .167, \pi_6 = .167$

 $\pi_7 = .133, \pi_8 = .100, \pi_9 = .067, \pi_{10} = .033$

versus

 H_a : The a priori proportions are incorrect

Using the proportions in H₀, the expected number of retrieval requests are:

Storage location	Expected number of retrieval requests
1	(200)(0333) = 6.667
2	(200)(.0666) = 13.333
3.	(200)(.1000) = 20.000
4. (200	(.1333) = 26.667
5. (200	(.1666) = 33.333
6. (200	(.1666) = 33.333
7. (200	(.1333) = 26.667
8. (200	(.1000) = 20.000
9. (200	(.0666) = 13.333
10. (200	(.0333) = 6.667

Since all expected counts are at least 5, the chi-squared test can be used.

location	1	2	3	4	5
observed	4	15	23	25	38
expected	6.667	13.333	20.000	26.667	33.333

location	6	7	8	9	10
observed	31	32	14	10	8
expected	33.333	26.667	20.000	13.333	6.667

The value of the χ^2 test statistic is:

$$\chi^{2} = \left[\sum \frac{(\text{observed} - \text{expected})^{2}}{\text{expected}}\right] = 6.6125$$

For df = (k - 1) = (10 - 1) = 9, the corresponding p-value = P($\chi^2 > 6.6125$). Using Table VIII, we find that the p-value > .10. Since the p-value, which is larger than .10 is also larger than $\alpha = .10$, we fail to reject H₀. We have no evidence to suggest that the a priori proportions are incorrect.

Supplementary Exercises

73. (a) No, it does not appear plausible that the distribution is normal. Notice that the mean value, $\overline{x} = 215$, is not nearly in the middle of the range of values, 5 to 1176. The midrange would be about 585. Since the mean is so much lower than this, one would suspect the distribution is positively skewed.

However, it is not necessary to assume normality if the sample size is "large enough", due to the central limit theorem. Since this problem has a sample size which is "large enough" (i.e., 47 > 30), we can proceed with a test of hypothesis about the true mean consumption.

(b) Let μ denote the true mean consumption. Since we are interested in determining if there is evidence to contradict the prior belief that μ was at most 200 mg, the following hypotheses should be tested.

H₀: $\mu = 200$ versus H_a: $\mu > 200$.

The value of the test statistic is:

$$z = \left(\frac{\overline{x} - 200}{s / \sqrt{n}}\right) = \left(\frac{215 - 200}{235 / \sqrt{47}}\right) = .44$$

The corresponding p-value = P(z > .44) = .33.

Since p-value = .33 > most any choice of α , we fail to reject H₀. There is insufficient evidence to suggest that the true mean caffeine consumption of adult women exceeds 200 mg per day.

- (a) The uniformity specification is that σ not exceed .5, so the relevant hypotheses are H₀:σ = .5 versus H_a:σ > .5 (i.e, we want to see if the data shows that the specified uniformity has been exceeded). The test statistic is χ² = (n-1)s²/σ² = (10-1)(.58)²/(.5)² = 12.1104. From the χ² table (Table VII) with d.f. = n-1 = 10-1 = 9, we note that 12.1104 is smaller than the smallest entry (14.68) in the d.f. = 9 column, so the P-value > .10. Therefore, H₀ should not be rejected at any of the usual significance levels (e.g., .05, .01) and we say that there is no evidence from data that contradicts the uniformity specification.
 - (b) No. The calculated test statistic is $\chi^2 = (n-1)s^2/\sigma^2 = (10-1)(.58)^2/(.7)^2 = 6.1788$ under the new null hypothesis. Since 6.1788 is smaller than 12.1104, the smallest entry in the d.f. = 9 column, so the right-tail area > .100. All we know is that the left-tail area < .9. Since we are having the alternative hypothesis H_a: σ < .7, we have to do a left-tail test and with the given information we cannot decide whether to reject the null at the usual levels (e.g., .05, .01).
- 77. Let π denote the true proportion of front-seat occupants involved in head-on collisions, in a certain region, who sustain no injuries. Given the wording of the exercise, the relevant hypotheses are:

H₀:

$$\pi = \left(\frac{1}{3}\right)^{\text{versus Ha:}} \pi < \left(\frac{1}{3}\right)$$
[Note:

$$\left(319\right)\left(\frac{1}{3}\right) = 106.3^{\text{and}} \left(319\right)\left(\frac{2}{3}\right) = 212.7^{\text{are each}} \ge 5^{\text{J}}$$
So the test statistic is:

So, the test statistic is:

$$z = \left[\frac{p - \pi_0}{\sqrt{\frac{(\pi_0)(1 - \pi_0)}{n}}}\right] = \left[\frac{\left(\frac{95}{319}\right) - \left(\frac{1}{3}\right)}{\sqrt{\frac{(\frac{1}{3})(\frac{2}{3})}{319}}}\right] = -1.35$$

The corresponding p-value = P(z < -1.35) = .0885

Since p-value = $.0885 > \alpha = .05$, we fail to reject H₀. We have insufficient evidence to claim that less than one-third of all such accidents result in no injuries.

81. (a) Let μ_1 denote the true mean strength for males and μ_2 denote the true mean strength for females. The hypotheses tested here were:

$$H_0: (\mu_1 - \mu_2) = 0$$
 versus $H_a: (\mu_1 - \mu_2) \neq 0$

If one assumes equal population variances, and uses the pooled sample variance you will obtain: $s_p =$ 31.77, t = 2.47, and df = 24. The corresponding p-value for this test is: 2P(t > 2.5) = 2(.010) = .02. These values are quite close to those reported in the exercise.

Notice, however, that the assumption of equal population variances and the t-test statistic that accompanies this assumption is not described in the body of the chapter.

If one uses the method described in the body of the chapter, then t = 2.84 and df = 18. This results in a p-value = 2(P(t > 2.8) = 2(.006) = .012.

(b) Revise the hypotheses:

$$H_0: (\mu_1 - \mu_2) = 25$$
 versus $H_a: (\mu_1 - \mu_2) > 25$

The test statistic (without assuming equal population variances) is:

$$t = \left[\frac{(129.2 - 98.1) - 25}{\sqrt{\frac{(39.1)^2}{15} + \frac{(14.2)^2}{11}}}\right] = .556$$

With df = 18, the p-value = P(t > .6) = .278

Since the p-value is greater than any sensible choice of α , we fail to reject H₀. There is insufficient evidence that the true average strength for males exceeds that for females by more than 25N.

86. The relevant hypotheses are H₀:*the leader's winning* % *is homogeneous across all 4 sports* versus H_a:*the leader's winning* % *differs among the 4 sports*. The appropriate test is a χ^2 test of homogeneity of several proportions. The following table shows the actual observations along with the expected values (underneath each observation) for the χ^2 test:

	Leader wins	Leader loses	totals
Basketball	150 155.28	39 33.72	189
Baseball	86 75.59	6 16.41	92
Hockey	65 65.73	15 14.27	80
Football	72 76.41	21 16.59	93
Totals	373	81	454

The test statistic value is:

 $\chi^2 = 0.180 + 0.827 +$ 1.435 + 6.607 +0.008 + 0.037 +0.254 + 1.171 = 10.518

From Table VII, with d.f. = (2-1)(4-1) = 3, the P-value associated with $\chi^2 = 10.518$ is P-value $\approx .015$. Since the P-value of .015 is smaller than the specified significance level of $\alpha = .05$, H₀ is rejected and we say that there is evidence from data in favor of the claim that the leader's winning % is not the same across all 4 sports. In particular, the win percentages are 79.4% (basketball), 93.5% (baseball), 81.3% (hockey), and 77.4% (football), so it appears that the leader's winning percentage is much higher for baseball than for the other three sports.

88. Let μ_d denote the true mean difference in retrieval time. We shall test $H_0: \mu_d = 10$ versus $H_a: \mu_d > 10$ at the $\alpha = .05$ significance level. We will use a paired *t* test, which assumes that the paired differences are normally distributed. From the data, we have $\overline{d} = 20.538$ and $s_d = 11.9625$

By using the Ryan-Joiner test of normality, it appears plausible that this normality condition is satisfied. The following probability plot also appears fairly linear.



For the paired-*t* test, the appropriate test statistic is
$$t = \frac{d - \Delta_0}{s_d / \sqrt{n}} = \frac{20.538 - 10}{11.9625 / \sqrt{13}} \approx 3.176.$$
 With $n = 13$, we

have df = n - 1 = 13 - 1 = 12. So the appropriate *t* critical value is 1.782. So we reject H_0 and say there is evidence from data in favor of the claim that the true mean difference in retrieval time does exceed 10 seconds.

Chapter 9 The Analysis of Variance

Section 9.2

17. (a) Changing units of measurement amounts to simply multiplying each observation by an appropriate conversion constant, c. In this exercise, c = 2.54. Next, note that replacing each x_i by cx_i causes any $x_i = 1$ sample mean to change from x to c^x while the grand mean also changes from x to c^x . Therefore, in the formulas for SSTr and SSE, replacing each x_i by cx_i will introduce a factor of c^2 . That is,

SSTR(for the cx_i data) = $n_1 (\frac{z}{z_1 - cx})^2 + ... + n_k (\frac{z}{z_k - cx})^2 = n_1 c^2 (\overline{x_1 - x})^2 + ... + n_k c^2 (\frac{z}{z_k - x})^2 = 0$

 (c^2) SSTr(for the original x_i data). The same thing happened for SSE; i.e., SSE(for the cx_i data) = (c^2) SSE(for the original x_i data). Using these facts, we also see that SST(for the cx_i data) = SSTR(for the cx_i data) = (c^2) SSTr(for the original x_i data) + (c^2) SSE(for the original x_i data) = (c^2) [SSTR(for the original x_i data) + SSE(for the original x_i data) + (c^2) SST(for the original x_i data) = (c^2) [SSTR(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the original x_i data)] = (c^2) SST(for the original x_i data) + SSE(for the number of treatments nor the number of observations is altered, the entries in the degrees of freedom column of the ANOVA table is not changed. Notice also that the F-ratio remains unchanged: F(for the cx_i data) = MSTR(for the cx_i data)/MSE(for the cx_i data)/[($c)^2$ MSE(for the original x_i data)] = MSTR(for the original x_i data)/[($c)^2$ MSE(for the original x_i data)] = MSTR(for the original x_i data) = F-ratio(for the original x_i data). This makes sense, for otherwise we could change the significance of an ANOVA test by merely changing the units of measurement.

- (b) The argument in (a) holds for *any* conversion factor c, not just for c = 2.54. We can say then, that *any* change in the units of measurement will change the 'Sum of Squares' column in the ANOVA table, but the degrees of freedom and F ratio will remained unchanged.
- 23. (a) Let x_i denote the true value of an observation. Then $x_i + c$ is the measured value reported by an instrument which is consistently off (i.e., out of calibration) by c units. Therefore, if \overline{x} denotes the mean of the true measurements, then $\overline{x} + c$ is the mean of the measured values. Similarly, the grand = = = = mean of the measured values equals x + c, where x is the grand mean of the true values. Putting these results in the formula for SSTR, we find SSTr(measured values) = $n_1(\overline{x_1 + c - (x + c)})^2 + ... + n_k(x + c)$

 $= \frac{1}{\overline{x}_k + c - (\overline{x} + c)})^2 = n_1(\overline{x}_1 - \overline{x})^2 + \dots + n_k(\overline{x}_k - \overline{x})^2 = \text{SSTR} \text{ (true values)}.$ Furthermore, note that any

sample variance is unchanged by the calibration problem since the deviations from the mean for the measured data are identical to the deviations from the mean for the true values; i.e., $(x_i + 2.5) - (\overline{x} + 2.5) = (x_i - \overline{x})$. Therefore, SSE (measured values) = $(n_1-1)s_1^2 + ... + (n_k-1)s_k^2 = SSE$ (true values). Finally, because SSTR and SSE are unaffected, so too will SST be unaffected by the calibration error since SST = SSTR + SSE. Thus, *none* of the sums of squares are changed by the calibration error. Obviously, the degrees of freedom are unchanged too, so the net result is that there will be <u>no change</u> in the entire ANOVA table.

(b) Calibration error will not change any of the ANOVA table entries and therefore will not affect the results of an ANOVA test. That is, if all data points are shifted (up or down) by the same amount c, the ANOVA entries will not be affected. However, the mean of each sample *will* shift by an amount equal to c.

Supplementary Exercises

45. (a) $n_1 = 9$ and $n_2 = 4$, so k = 2 and $n = n_1 + n_2 = 9 + 4 = 13$. The grand mean is the weighted average of the sample means, $\overline{x} = [9(-.83) + 4(-.70)]/[9+4] = -.79$. SSE = $(n_1-1)s_1^2 + (n_2-1)s_2^2 = (9-1)(.172)^2 + (4-1)s_2^2 = (9-1)(.172)^2 + (1-1)s_2^2 = (9-1)(.172)^2 + (1-1)$

1)(.184)² = .33824 and SSTr = $n_1(\overline{x_1} - x)^2 + n_2(\overline{x_2} - x)^2 = 9(-.83-(-.79))^2 + 4(-.70-(-.79))^2 = .0468$. Therefore, MSTR = SSTR/(k-1) = .0468/(2-1) = .0468, MSE = SSE/(n-k) = .33824/(13-2) = .03075, and F = MSTR/MSE = .0468/.03075 = 1.522. For df₁ = 1 and df₂ = 11, the P-value associated with F = 1.522 exceeds .10 (Table VIII). Therefore, since the P-value is larger than α = .01, we can not reject H₀: *no difference between average diopter measurements for the symptom 'present' and 'absent'*. The data does not show a significant difference between the two groups of pilots.

Chapter 8 is the independent samples t-test. The test statistic is $t = \frac{\overline{x_1 - \overline{x}_2}}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$

(b) The equivalent test from Chapter 8 is the independent samples t-test. The test statistic is $t = \sqrt{n_1}$

$$\frac{-.83 - (-.70)}{\sqrt{\frac{0.172^2}{9} + \frac{0.184^2}{4}}} = -1.199. \text{ For df} = \frac{\frac{(se_1^2 + se_2^2)^2}{n_1 - 1} + \frac{se_2^4}{n_2 - 1}}{\frac{1}{9} - 1} = \frac{(\frac{0.172^2}{9} + \frac{0.184^2}{4})^2}{\frac{0.172^4}{81}} = 5$$
(where set²)

= s_1^2/n_1 and $se_2^2 = s_2^2/n_2$), the p-value associated with t = -1.199 (for a 2-sided test) is approximately 2(.142) = .284 (from Table VI). Since the p-value \approx .284 exceeds α = .01, we do not reject H₀: $\mu_1 - \mu_2 = 0$. That is, the test does not show a significant difference between the two groups of pilots. This is the same conclusion as in part (a).

Chapter 11 Inferential Methods in Regression and Correlation

Section 11.1

- 1. (a) The slope of the estimated regression line ($\beta = .095$) is the expected <u>change</u> of in the response variable y for each one-unit <u>increase</u> in the x variable. This, of course, is just the usual interpretation of the slope of a straight line. Since x is measured in inches, a one-unit increase in x corresponds to a one-inch increase in pressure drop. Therefore, the expected change in flow rate is .095 m³/min.
 - (b) When the pressure drop, x, changes from 10 inches to 15 inches, then a 5 unit increase in x has occurred. Therefore, using the definition of the slope from (a), we expect about a 5(.095) = .475 m³/min. increase in flow rate (it is an *increase* since the sign of $\beta = .095$ is *positive*).
 - (c) For x = 10, $\mu_{y,10} = -.12 + .095(10) = .830$. For x = 15, $\mu_{y,15} = -.12 + .095(15) = 1.305$.
 - (d) When x = 10, the flow rate y is normally distributed with a mean value of $\mu_{y.10} = .830$ and a standard deviation of $\sigma_{y.10} = \sigma = .025$. Therefore, we standardize and use the z table to find: P(y > .835) = P(z > .835 .830) = P(z > .20) = 1 P(z \le .20) = 1 .5793 = .4207 (using Table I).
- 2. (a) The slope of the estimated regression line ($\beta = -.01$) is the expected change in reaction time for a one degree Fahrenheit increase in the temperature of the chamber.

So, with a one degree Fahrenheit increase in temperature, the true average reaction time will decrease by .01 hours.

With a 10 degree increase in temperature, the true average reaction time will decrease by (10)(.01) = 1 hour.

(b) When x = 200, $\mu_{y^{\bullet}200} = 5 - .01(200) = 3$

When x = 250, $\mu_{v^{\bullet}250} = 5 - .01(250) = 2.5$

(c) P(2.4 < y < 2.6 when x = 250) = $P\left(\frac{2.4 - 2.5}{.075} < z < \frac{2.6 - 2.5}{.075}\right) =$ P(-1.33 < z < 1.33) = P(z < 1.33) - P(z < -1.33) =.9082 - .0918 = .8164

Next, the probability that all five observed reaction times are between 2.4 and 2.6 is $(.8164)^5 = .3627$

4. (a)



The scatterplot appears linear, so a simple linear regression model seems reasonable.

(b) The following quantities are needed:

$$\overline{x} = 53.200$$

$$\overline{y} = 42.867$$

$$S_{xy} = \left[51232 - \left(\frac{(798)(643)}{15}\right)\right] = 17024.4$$

$$S_{xx} = \left[63040 - \left(\frac{(798)^{2}}{15}\right)\right] = 20586.4$$

$$S_{yy} = \left[41999 - \left(\frac{(643)^{2}}{15}\right)\right] = 14435.7$$

$$b = \left(\frac{S_{xy}}{S_{xx}}\right) = \left(\frac{17024.4}{20586.4}\right) = .82697$$

$$a = (\overline{y} - b\overline{x}) = (42.867 - (.82697)(53.2)) = -1.1278$$

(c)
$$\mu_{y \cdot 50} = -1.1278 + (.82697)(50) = 40.2207$$

(d) $SSResid = S_{yy} - bS_{xy} = 14435.7 - (.82697)(17024.4) = 357.07$

$$s_e = \sqrt{\frac{\text{SSResid}}{n-2}} = \sqrt{\frac{357.07}{15.2}} = 5.24$$

(e)

$$r^{2} = 1 - \left(\frac{\text{SSResid}}{\text{SSTo}}\right) = 1 - \left(\frac{\text{SSResid}}{\text{S}_{yy}}\right) = 1 - \left(\frac{357.07}{14435.7}\right) = .9753$$

So, 97.53% of the observed variation in runoff volume can be attributed to the simple linear regression relationship between runoff and rainfall.

6. (a) Using the formulas for the various sums of squares, we find:

$$SS_{xy} = \sum_{i} x_{i} y_{i} - \frac{1}{n} \left(\sum_{i} x_{i} \sum_{j} y_{i} \right)_{=40.968 - (12.6)(27.68)/9} = 2.216$$

$$SS_{xx} = \sum_{i} x_{i}^{2} - \frac{1}{n} \left(\sum_{i} x_{i} \right)_{=18.24 - (12.6)^{2}/9} = .600.$$

Therefore, the estimated slope is: $b = SS_{xy}/SS_{xx} = 2.216/.600 = 3.6933$. The estimated intercept is a = \overline{y} - b $\overline{x} = (27.68)/9 - (3.6933)(12.6)/9 = -2.0951$. The estimated regression line is then: $\hat{y} = a + bx =$ -2.0951 + 3.6933x.

(b) For x = 1.5, the point estimate of the *average* y value is: $\mu_{y,1.5} \approx -2.0951 + 3.6933(1.5) = 3.445$. For another measurement made when x = 1.5, the point estimate of the y value (for this x value) would be the same: i.e., $\hat{y} = 3.445$.

(c)
$$SSTo = SS_{yy} = \sum y_i^2 - \frac{1}{n} \left(\sum y_i \right)^2 = 93.3448 - (27.68)^2 / 9 = 8.2134$$
. Therefore, $SSResid = SSTo - b \cdot SS_{xy} = 8.2134 - (3.6933)(2.216) = .0290$. The estimate of σ is then $s_e \approx \sqrt{\frac{SSResid}{n-2}} = \sqrt{\frac{.0290}{9-2}} = .0644$.

of the observed variation in diffusivity can be attributed to the simple linear regression model between diffusivity and temperature.

Section 11.2

10. Let β denote the true average change in runoff (for each 1m³ increase in rainfall). To test the hypotheses H_0 : $\beta = 0$ versus H_a : $\beta \neq 0$, the calculated t statistic is $t = b/s_b = .82697/.03652 = 22.64$ which (from the printout) has an associated P-value of P = 0.000. Therefore, since the P-value is so small, H₀ is rejected and we say that there is evidence from data in favor of the claim that there is a useful linear relationship between runoff and rainfall (no surprise here!).

A confidence interval for β is based on n-2 = 15-2 = 13 degrees of freedom. The t critical value for, say, 95% confidence is 2.160 (from Table IV), so the interval estimate is: $b \pm (t \text{ critical}) \cdot s_b = .82697 \pm .8267 \pm .8267 \pm .8267 \pm .8277 \pm .827700 \pm .8277 \pm .8277000000000000000000000000000000$ $(2.160)(.03652) = .82697 \pm .07888 = [.748, .906]$ Therefore, we can be confident that the true average change in runoff, for each 1m³ increase in rainfall, is somewhere between .748m³ and .906m³.

11. H₀: $\alpha = 0$ versus H_a: $\alpha \neq 0$.

The calculated t statistic is:

$$t = \left(\frac{a}{s_a}\right) = \left(\frac{-1.128}{2.368}\right) = -.48$$

The corresponding p-value = .642. Therefore, since the p-value is large, we would not reject the null hypothesis. We cannot say that there is evidence from data in favor of the claim that the vertical intercept of the population line is nonzero.

13. (a) Method 1: Hypothesis Test, H0: $\beta = 0$ versus Ha: $\beta \neq 0$, t = 54.56, P-value < .0001, reject H0 and

conclude that there is a useful linear relationship between these two variables. Method 2: A confidence interval for $\beta = b \pm (t \operatorname{critical value}) s_b$. A 95% confidence interval for β is: .87825 ± (2.179)(.01610) = (0.8432, 0.9133), using t critical value for df = (n - 2) = (14 - 2) = 12. The plausible values are all positive so we conclude there is a useful linear relationship between the two variables.

(b) The t ratio for testing model utility would be the same value regardless of which of the two variables was defined to be the independent variable. This can be easily seen by looking at the t test statistic for testing if the population correlation coefficient is equal to zero. In that equation the only values required are the sample size (n) and the sample correlation coefficient (r). Both r and n are not dependent on which variable was the independent variable.

(a) SSRegr = 0.071422; SSResid = 0.020178; SSTo = .0916

(b) $R^2 = \frac{SSRegr}{SSTo} = \frac{0.071422}{0.0916} = 0.7797$. Thus, 77.97% of the observed variation in dielectric constant can be accounted for by the simple linear regression with air void as the explanatory variable. $r = -\sqrt{R^2} = -\sqrt{0.7797} = -0.883$; correlation coefficient is negative because the slope is negative.

(c) A 95% confidence interval for is: $-0.0747 \pm (2.12)(.0099) = (-0.09571, -0.05364)$, using t critical

value for df = 16. We are 95% confident the decrease in mean dielectric constant associated with a 1 percentage point increase in air void is between 0.05364 and 0.09571.

(d) $H_0: \beta \ge 0.05$ versus $H_a: \beta < 0.05$ Yes, we have evidence to contradict the null hypothesis, because the 95% CI for the slope is (-0.09571, -0.05364)

17.

(d) By deleting the observation (x, y) = (6.0, 2.50), our new summary statistics become:

$$\begin{split} n &= 17 - 1 = 16; \quad \sum x_i = 221.1 - 6.0 = 215.1; \quad \sum y_i = 193 - 2.50 = 190.50; \\ \sum x_i^2 &= 3056.69 - 6.0^2 = 3020.69; \quad \sum x_i y_i = 2759.6 - (6.0)(2.50) = 2744.6^{-1}; \\ \sum y_i^2 &= 2975 - 2.50^2 = 2968.75 \end{split}$$

$$\begin{split} S_{xy} &= \sum x_i y_i - (1/n) (\sum x_i) (\sum y_i) = 2744.6 - (1/16)(215.1)(190.50) = 183.5656; \\ S_{xx} &= \sum x_i^2 - (1/n) (\sum x_i)^2 = 3020.69 - (1/16)(215.1)^2 = 128.9394 \end{split}$$
So
$$b &= S_{xy} / S_{xx} = 183.5656 / 128.9894 \approx 1.42366 \end{split}$$

By deleting (x, y) = (6.0, 2.50), our new estimate for β is b = 1.42366, which is still fairly close to the value of b = 1.37758 that was computed in part (a) above. Moreover, this new estimate of b = 1.42366 falls well within the 95% CI that was obtained in part (a) above. Therefore, the point (6.0, 2.50) does not appear to exert any undue influence on our regression analysis.

Section 11.3

24. n=15, so the error d.f. = n-2 = 15-2 = 13 and, therefore, the t-critical value for a 95% prediction interval is 2.160 (from Table IV). The prediction interval for x = 40 is centered at $\hat{v} = -1.128 + .82697(40) =$

31.9508. The prediction interval is then: $\hat{y} \pm (\text{t-critical}) \sqrt{s_e^2 + s_{\hat{y}}^2} = 31.9508 \pm (2.160) \sqrt{(5.24)^2 + (1.44)^2} = 31.95 \pm 11.74$

= [20.21, 43.69]. Even though the r^2 value is large ($r^2 = 73.8\%$), the prediction interval is rather wide, so precise information about future runoff levels can not be obtained from this model.

14.

27. (a) Let β denote the true average change in milk protein for each 1 kg/day increase in milk production. The relevant hypotheses to test are H₀: $\beta = 0$ versus H_a: $\beta \neq 0$. The test statistic is t = b/s_b based on n-2 = 14-2 = 12 degrees of freedom. In order to find s_b, we first find s_e:

$$s_e^2 = \frac{SResid}{n-2} = \frac{.02120}{.02120} = .001767$$
, so se = .0420.

Then,
$$s_b = \frac{s_e}{\sqrt{SS_{v_v}}} = \frac{.0420}{\sqrt{762.012}} = .00152$$
, which then gives a calculated t value of $t = b/s_b = \frac{.0420}{\sqrt{762.012}}$

 $.024576/.00152 \approx 16.2$. The 2-sided P-value associated with t = 16.2 is approximately 2(.000) = .000 (Table VI), so H₀ is rejected in favor of the conclusion that there is a useful linear relationship between protein and production. We should not be surprised by this result since the r² value for this data is .956.

(b) For a 99% confidence interval based on d.f. = 12, the t-critical value is 3.055 (from Table IV). The estimated regression line gives a value of $\hat{v} = .175576 + .024576(30) = .913$ when x = 30. Therefore,

$$s_{\hat{y}} = (.0420) \sqrt{\frac{1}{14} + \frac{(30-29.56)^2}{762.012}} = .01124 \text{ and the 95\% confidence interval is then:} .913$$

$$\pm (3.055)(.01124) = .913 \pm .034 = [.879, .947].$$

(c) The 99% prediction interval for protein from a single cow is: $\hat{y} \pm (\text{t-critical}) \sqrt{s_e^2 + s_{\hat{y}}^2} = .913 \pm (3.055) \sqrt{(.0420)^2 + (.01124)^2}$ $= .913 \pm .133 = [.780, 1.046].$

```
28.
```

$$s_a = s_{a+b(0)} = s_e \sqrt{\frac{1}{n} + \frac{(0 - \overline{x})^2}{S_{xx}}}$$

$$=.0420\sqrt{\frac{1}{14} + \frac{(29.564)^2}{762.012}} = .0464$$

H₀: $\alpha = 0$ versus H_a: $\alpha \neq 0$

The value of the test statistic is:

$$t = \left(\frac{.175576 - 0}{.0464}\right) = 3.78$$

With 12 degrees of freedom, the corresponding p-value obtained from Minitab is .0026.

Since the p-value is so small, we reject the null hypothesis and say that there is evidence from data in favor of the claim that the vertical intercept is a nonzero value.

Section 11.4

29. (a) The mean value of y, when $x_1 = 50$ and $x_2 = 3$ is -.800 + .060(50) + .900(3) = 4.9 hours.

- (b) When the number of deliveries (x₂) is held fixed, then average change in travel time associated with a one-mile (i.e., one unit) increase in distance traveled (x₁) is .060 hours. Similarly, when the distance traveled (x₁) is held fixed, then the average change in travel time associated with one extra delivery (i.e., a one-unit increase in x₂) is .900 hours.
- (c) Under the assumption that y follows a normal distribution, the mean and standard deviation of this distribution are 4.9 (because $x_1 = 50$ and $x_2 = 3$) and $\sigma = .5$ (since σ is assumed to be constant regardless of the values of x_1 and x_2). Therefore, $P(y \le 6) = P(z \le (6-4.9)/.5) = P(z \le 2.20) = .9861$ (from Table I). That is, in the long run, about 98.6% of all days will result in a travel time of at most 6 hours.

Section 11.5

37. (a) The appropriate hypotheses are $H_0:\beta_1=\beta_2=\beta_3=\beta_4=0$ versus $H_a:$ at least one of the β_i 's is not zero. The test statistic is $F = \frac{R^2/k}{(1-R^2)/(n-(k+1))} = \frac{.946/4}{(1-.946)/(25-(4+1))} = 87.6$. The test is based on $df_1 = 4$, df_2

= 20. From Table XII, the P-value associated with F = 6.59 is .001, so the P-value associated with 87.6 is obviously .000. Therefore, H₀ can be rejected at any reasonable level of significance. We say that there is evidence from data in favor of the claim that at least one of the four predictor variables appears to provide useful information about tenacity.

(c) The estimated average tenacity when $x_1 = 16.5$, $x_2 = 50$, $x_3 = 3$, and $x_4 = 5$ is: $\hat{y}_2 = 6.121 - .082x_1$

+.113 x_2 +.256 x_3 -.219 x_4 = 6.121 -.082(16.5) +.113(50) +.256(3) -.219(5) = 10.091. For a 99% confidence interval based on 20 d.f., the t-critical value is 2.845. The desired interval is: 10.091 ± (2.845)(.350) = 10.091 ± .996, or, about [9.095, 11.087]. Therefore, when the four predictors are as specified in this problem, the true average tenacity is estimated to be between 9.095 and 11.087.

38. (a) Yes, there does appear to be a useful linear relationship between repair time and the two model predictors. We determine this by conducting a model utility test.

H₀: $\beta_1 = \beta_2 = 0$ versus H_a: At least one of these two β 's are not zero.

The test statistic requires the following quantities.

SSRegr = (SSTo - SSResid) = (12.72 - 2.09) = 10.63

MSRegr = (SSRegr/k) = (10.63/2) = 5.315

MSResid = (SSResid/n - k - 1) = (2.09/9) = .2322

So, F = (MSRegr/MSResid) = (5.315/.232) = 22.91

With 2 numerator degrees of freedom and 9 denominator degrees of freedom, the F critical value at α = .05 is 4.26. Since 22.91 > 4.26, we reject the null hypothesis and say that there is evidence from data in favor of the claim that at least one of the two predictor variables is useful.

(b) $H_0: \beta_2 = 0$ versus $H_a: \beta_2 \neq 0$

The t test statistic = $\left(\frac{1.250 - 0}{.312}\right) = 4.01$

The corresponding p-value = 2P(t > 4.01). With df = 9 and using Minitab, the exact p-value is .003. Since the p-value < α =.01, we reject the null hypothesis and say that there is evidence from data in favor of the claim that the "type of repair" variable does provide useful information about repair time, given that the "elapsed time since the last service" variable remains in the model.

(c) A 95% confidence interval for β_2 is : 1.250 ± (2.262)(.312) 1.250 ± .7057 (.5443, 1.9557) (Note: df = n - (k + 1) = 12 - (2 + 1) = 9) We estimate, with a high degree of confidence, that when an electrical repair is required the repair time will be between .54 and 1.96 hours longer than when a mechanical repair is required, while the "elapsed time" predictor remains fixed.

(d) To compute the prediction interval we need the following quantities. $\hat{y} = .950 + .400(6) + 1.250(1) = 4.6$

 $s_{a}^{2} = MSResid = .23222$

For a 99% prediction interval, the t critical value with df = 9 equals 3.250. So, our 99% prediction interval is:

 $4.6 \pm (3.250)\sqrt{.23222 + (.192)^2}$ 4.6 ± 1.69 (2.91, 6.29)

The prediction interval is quite wide, suggesting a variable estimate for repair time under these conditions.

- (a) The negative value of b_2 , which is the coefficient of x^2 in the model, indicates that the parabola $b_0 + b_1x + b_2x^2$ opens downward.
 - (b) $R^2 = 1$ SSResid/Ssto = 1 .29/202.87 = .9986, so about 99.86% of the variation in output power can be attributed to the relationship between power and frequency.
 - (c) With an R² this high, it is very likely that the test statistic will be significant. The relevant hypotheses are $H_0:\beta_1=\beta_2=0$ versus H_a : at least one of the β_i 's is not zero. The test statistic is: F =

 $\frac{SS \operatorname{Re} gr/k}{SS \operatorname{Re} sid/(n - (k + 1))} = \frac{(202.87 - 2.9)/2}{.29/5} = 1746.$ Clearly, the P-value associated with F = 1746 is

0, so H_0 is rejected and we say that there is evidence from data in favor of the claim that the model is useful for predicting power.

(d) The relevant hypotheses are $H_0:\beta_2=0$ versus $H_a:\beta_2 \neq 0$. The test statistic is: $t = b_2/s_b = -\frac{s_b}{s_b}$

.00163141/.00003391 = 48. The P-value for this statistic is 0 and H₀ is rejected in favor of the conclusion that the quadratic predictor provides useful information.

(e) The estimated average power when x = 150 is $\hat{y} = -1.5127 + .391902x - .00163141x^2 =$

 $1.5127 + .391902(150) - .00163141(150)^2 = 20.57$. The t-critical value based on 5d.f. is 4.032, so the 99% confidence interval is: $20.57 \pm (4.032)(.1410) = 20.57 \pm .57$ or, about [20.00, 21.14] To find the prediction interval, we must first find se. $s_e^2 = \text{SSResid}/(n-3) = .29/5 = .058$, so $s_e = .241$. Therefore,

the 99% prediction interval is:

 $20.57 \pm (4.032) \sqrt{(.241)^2 + (.141)^2} = 20.57 \pm 1.13$, or, about [19.44, 21.70].

39.