

CHAPTER 14

Section 14.1

1.

- a. We reject H_0 if the calculated \mathbf{c}^2 value is greater than or equal to the tabled value of $\mathbf{c}_{a,k-1}^2$ from Table A.7. Since $12.25 \geq \mathbf{c}_{.05,4}^2 = 9.488$, we would reject H_0 .
- b. Since 8.54 is not $\geq \mathbf{c}_{.01,3}^2 = 11.344$, we would fail to reject H_0 .
- c. Since 4.36 is not $\geq \mathbf{c}_{.10,2}^2 = 4.605$, we would fail to reject H_0 .
- d. Since 10.20 is not $\geq \mathbf{c}_{.01,5}^2 = 15.085$, we would fail to reject H_0 .

2.

- a. In the d.f. = 2 row of Table A.7, our \mathbf{c}^2 value of 7.5 falls between $\mathbf{c}_{.025,2}^2 = 7.378$ and $\mathbf{c}_{.01,2}^2 = 9.210$, so the p-value is between .01 and .025, or $.01 < \text{p-value} < .025$.
- b. With d.f. = 6, our \mathbf{c}^2 value of 13.00 falls between $\mathbf{c}_{.05,6}^2 = 12.592$ and $\mathbf{c}_{.025,6}^2 = 14.440$, so $.025 < \text{p-value} < .05$.
- c. With d.f. = 9, our \mathbf{c}^2 value of 18.00 falls between $\mathbf{c}_{.05,9}^2 = 16.919$ and $\mathbf{c}_{.025,9}^2 = 19.022$, so $.025 < \text{p-value} < .05$.
- d. With $k = 5$, d.f. = $k - 1 = 4$, and our \mathbf{c}^2 value of 21.3 exceeds $\mathbf{c}_{.005,4}^2 = 14.860$, so the p-value $< .005$.
- e. The d.f. = $k - 1 = 4 - 1 = 3$; $\mathbf{c}^2 = 5.0$ is less than $\mathbf{c}_{.10,3}^2 = 6.251$, so p-value $> .10$.

Chapter 14: The Analysis of Categorical Data

3. Using the number 1 for business, 2 for engineering, 3 for social science, and 4 for agriculture, let p_i = the true proportion of all clients from discipline i . If the Statistics department's expectations are correct, then the relevant null hypothesis is $H_o : p_1 = .40, p_2 = .30, p_3 = .20, p_4 = .10$, versus H_a : The Statistics department's expectations are not correct. With d.f = $k - 1 = 4 - 1 = 3$, we reject H_o if $\mathbf{c}^2 \geq \mathbf{c}_{.05,3}^2 = 7.815$. Using the proportions in H_o , the expected number of clients are :

Client's Discipline	Expected Number
Business	$(120)(.40) = 48$
Engineering	$(120)(.30) = 36$
Social Science	$(120)(.20) = 24$
Agriculture	$(120)(.10) = 12$

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

test statistic is
$$\mathbf{c}^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i} = \sum_{all\ cells} \frac{(observed - expected)^2}{expected}$$

$$= \left[\frac{(52 - 48)^2}{48} + \frac{(38 - 36)^2}{36} + \frac{(21 - 24)^2}{24} + \frac{(9 - 12)^2}{12} \right] = 1.57$$
, which is not

≥ 7.815 , so we fail to reject H_o . (Alternatively, p-value = $P(\mathbf{c}^2 \geq 1.57)$ which is $> .10$, and since the p-value is not $< .05$, we reject H_o). Thus we have no evidence to suggest that the statistics department's expectations are incorrect.

4. The uniform hypothesis implies that $p_{i0} = \frac{1}{8} = .125$ for $i = 1, \dots, 8$, so

$H_o : p_{10} = p_{20} = \dots = p_{80} = .125$ will be rejected in favor of H_a if

$\mathbf{c}^2 \geq \mathbf{c}_{.10,7}^2 = 12.017$. Each expected count is $np_{i0} = 120(.125) = 15$, so

$$\mathbf{c}^2 = \left[\frac{(12 - 15)^2}{15} + \dots + \frac{(10 - 15)^2}{15} \right] = 4.80$$
. Because 4.80 is not ≥ 12.017 , we fail to

reject H_o . There is not enough evidence to disprove the claim.

Chapter 14: The Analysis of Categorical Data

5. We will reject H_0 if the p-value $< .10$. The observed values, expected values, and corresponding \mathbf{C}^2 terms are :

Obs	4	15	23	25	38	21	32	14	10	8
Exp	6.67	13.33	20	26.67	33.33	33.33	26.67	20	13.33	6.67
\mathbf{C}^2	1.069	.209	.450	.105	.654	.163	1.065	1.800	.832	.265

$\mathbf{C}^2 = 1.069 + \dots + .265 = 6.612$. With d.f. = $10 - 1 = 9$, our \mathbf{C}^2 value of 6.612 is less than $\mathbf{C}_{.10,9}^2 = 14.684$, so the p-value $> .10$, which is not $< .10$, so we cannot reject H_0 . There is no evidence that the data is not consistent with the previously determined proportions.

6. A 9:3:4 ratio implies that $p_{10} = \frac{9}{16} = .5625$, $p_{20} = \frac{3}{16} = .1875$, and $p_{30} = \frac{4}{16} = .2500$. With $n = 195 + 73 + 100 = 368$, the expected counts are 207.000, 69.000, and 92.000, so

$$\mathbf{C}^2 = \left[\frac{(195 - 207)^2}{207} + \frac{(73 - 69)^2}{69} + \frac{(100 - 92)^2}{92} \right] = 1.623$$

. With d.f. = $3 - 1 = 2$, our

\mathbf{C}^2 value of 1.623 is less than $\mathbf{C}_{.10,2}^2 = 4.605$, so the p-value $> .10$, which is not $< .05$, so we cannot reject H_0 . The data does confirm the 9:3:4 theory.

7. We test $H_0 : p_1 = p_2 = p_3 = p_4 = .25$ vs. H_a : at least one proportion $\neq .25$, and d.f. = 3. We will reject H_0 if the p-value $< .01$.

Cell	1	2	3	4
Observed	328	334	372	327
Expected	340.25	340.25	340.25	34.025
\mathbf{C}^2 term	.4410	.1148	2.9627	.5160

$\mathbf{C}^2 = 4.0345$, and with 3 d.f., p-value $> .10$, so we fail to reject H_0 . The data fails to indicate a seasonal relationship with incidence of violent crime.

Chapter 14: The Analysis of Categorical Data

8. $H_o : p_1 = \frac{15}{365}, p_2 = \frac{46}{365}, p_3 = \frac{120}{365}, p_4 = \frac{184}{365}$, versus H_a : at least one proportion is not a stated in H_o . The degrees of freedom = 3, and the rejection region is $\mathbf{c}^2 \geq \mathbf{c}_{.01,3} = 11.344$.

Cell	1	2	3	4
Observed	11	24	69	96
Expected	8.22	25.21	65.75	100.82
\mathbf{c}^2 term	.9402	.0581	.1606	.2304

$\mathbf{c}^2 = \sum \frac{(\text{obs} - \text{exp})^2}{\text{exp}} = 1.3893$, which is not ≥ 11.344 , so H_o is not rejected. The data does not indicate a relationship between patients' admission date and birthday.

9.

- a. Denoting the 5 intervals by $[0, c_1), [c_1, c_2), \dots, [c_4, \infty)$, we wish c_1 for which $.2 = P(0 \leq X \leq c_1) = \int_0^{c_1} e^{-x} dx = 1 - e^{-c_1}$, so $c_1 = -\ln(.8) = .2231$. Then $.2 = P(c_1 \leq X \leq c_2) \Rightarrow .4 = P(0 \leq X_1 \leq c_2) = 1 - e^{-c_2}$, so $c_2 = -\ln(.6) = .5108$. Similarly, $c_3 = -\ln(.4) = .9163$ and $c_4 = -\ln(.2) = 1.6094$. the resulting intervals are $[0, .2231), [.2231, .5108), [.5108, .9163), [.9163, 1.6094)$, and $[1.6094, \infty)$.
- b. Each expected cell count is $40(.2) = 8$, and the observed cell counts are 6, 8, 10, 7, and 9, so $\mathbf{c}^2 = \left[\frac{(6-8)^2}{8} + \dots + \frac{(9-8)^2}{8} \right] = 1.25$. Because 1.25 is not $\geq \mathbf{c}_{.10,4}^2 = 7.779$, even at level .10 H_o cannot be rejected; the data is quite consistent with the specified exponential distribution.

10.

- a.
$$\mathbf{c}^2 = \sum_{i=1}^k \frac{(n_i - np_{i0})^2}{np_{i0}} = \sum_{i=1}^k \frac{N_i^2 - 2np_{i0}N_i + n^2 p_{i0}^2}{np_{i0}} = \sum_{i=1}^k \frac{N_i^2}{np_{i0}} - 2 \sum_{i=1}^k N_i + n \sum_{i=1}^k p_{i0}$$

$$= \sum_{i=1}^k \frac{N_i^2}{np_{i0}} - 2n + n(1) = \sum_{i=1}^k \frac{N_i^2}{np_{i0}} - n$$
 as desired. This formula involves only one subtraction, and that at the end of the calculation, so it is analogous to the shortcut formula for s^2 .
- b. $\mathbf{c}^2 = \frac{k}{n} \sum_{i=1}^k N_i^2 - n$. For the pigeon data, $k = 8, n = 120$, and $\sum N_i^2 = 1872$, so
$$\mathbf{c}^2 = \frac{8(1872)}{120} - 120 = 124.8 - 120 = 4.8$$
 as before.

11.

- a. The six intervals must be symmetric about 0, so denote the 4th, 5th and 6th intervals by $[0, a]$, $[a, b]$, $[b, \infty)$. a must be such that $\Phi(a) = .6667\left(\frac{1}{2} + \frac{1}{6}\right)$, which from Table A.3 gives $a \approx .43$. Similarly $\Phi(b) = .8333$ implies $b \approx .97$, so the six intervals are $(-\infty, -.97)$, $[-.97, -.43]$, $[-.43, 0]$, $[0, .43]$, $[.43, .97]$, and $[.97, \infty)$.
- b. The six intervals are symmetric about the mean of .5. From **a**, the fourth interval should extend from the mean to .43 standard deviations above the mean, i.e., from .5 to $.5 + .43(.002)$, which gives $[.5, .50086]$. Thus the third interval is $[.5 - .00086, .5) = [.49914, .5)$. Similarly, the upper endpoint of the fifth interval is $.5 + .97(.002) = .50194$, and the lower endpoint of the second interval is $.5 - .00194 = .49806$. The resulting intervals are $(-\infty, .49806)$, $[.49806, .49914)$, $[.49914, .5)$, $[.5, .50086)$, $[.50086, .50194)$, and $[.50194, \infty)$.
- c. Each expected count is $45\left(\frac{1}{6}\right) = 7.5$, and the observed counts are 13, 6, 6, 8, 7, and 5, so $\mathbf{c}^2 = 5.53$. With 5 d.f., the p-value $> .10$, so we would fail to reject H_0 at any of the usual levels of significance. There is no evidence to suggest that the bolt diameters are not normally distributed.

Section 14.2

12.

- a. Let \mathbf{q} denote the probability of a male (as opposed to female) birth under the binomial model. The four cell probabilities (corresponding to $x = 0, 1, 2, 3$) are

$$\mathbf{p}_1(\mathbf{q}) = (1 - \mathbf{q})^3, \mathbf{p}_2(\mathbf{q}) = 3\mathbf{q}(1 - \mathbf{q})^2, \mathbf{p}_3(\mathbf{q}) = 3\mathbf{q}^2(1 - \mathbf{q}), \text{ and } \mathbf{p}_4(\mathbf{q}) = \mathbf{q}^3.$$

The likelihood is $3^{n_2+n_3} \cdot (1 - \mathbf{q})^{3n_1+2n_2+n_3} \cdot \mathbf{q}^{n_2+2n_3+3n_4}$. Forming the log likelihood,

taking the derivative with respect to \mathbf{q} , equating to 0, and solving yields

$$\hat{\mathbf{q}} = \frac{n_2 + 2n_3 + 3n_4}{3n} = \frac{66 + 128 + 48}{480} = .504. \text{ The estimated expected counts are}$$

$$160(1 - .504)^3 = 19.52, 480(.504)(.496)^2 = 59.52, 60.48, \text{ and } 20.48, \text{ so}$$

$$\mathbf{c}^2 = \left[\frac{(14 - 19.52)^2}{19.52} + \dots + \frac{(16 - 20.48)^2}{20.48} \right] = 1.56 + .71 + .20 + .98 = 3.45.$$

The number of degrees of freedom for the test is $4 - 1 - 1 = 2$. H_0 of a binomial

distribution will be rejected using significance level .05 if $\mathbf{c}^2 \geq \mathbf{c}_{.05,2}^2 = 5.992$.

Because $3.45 < 5.992$, H_0 is not rejected, and the binomial model is judged to be quite plausible.

- b. Now $\hat{\mathbf{q}} = \frac{53}{150} = .353$ and the estimated expected counts are 13.54, 22.17, 12.09, and 2.20. The last estimated expected count is much less than 5, so the chi-squared test based on 2 d.f. should not be used.

Chapter 14: The Analysis of Categorical Data

- 13.** According to the stated model, the three cell probabilities are $(1-p)^2$, $2p(1-p)$, and p^2 , so we wish the value of p which maximizes $(1-p)^{2n_1} [2p(1-p)]^{n_2} p^{2n_3}$. Proceeding as in example 14.6 gives $\hat{p} = \frac{n_2 + 2n_3}{2n} = \frac{234}{2776} = .0843$. The estimated expected cell counts are then $n(1-\hat{p})^2 = 1163.85$, $n[2\hat{p}(1-\hat{p})]^2 = 214.29$, $n\hat{p}^2 = 9.86$. This gives
- $$\mathbf{c}^2 = \left[\frac{(1212 - 1163.85)^2}{1163.85} + \frac{(118 - 214.29)^2}{214.29} + \frac{(58 - 9.86)^2}{9.86} \right] = 280.3.$$
- According to (14.15), H_0 will be rejected if $\mathbf{c}^2 \geq \mathbf{c}_{\alpha,2}^2$, and since $\mathbf{c}_{.01,2}^2 = 9.210$, H_0 is soundly rejected; the stated model is strongly contradicted by the data.

14.

- a.** We wish to maximize $p^{\sum x_i - n} (1-p)^n$, or equivalently $(\sum x_i - n) \ln p + n \ln(1-p)$. Equating $\frac{d}{dp}$ to 0 yields $\frac{(\sum x_i - n)}{p} = \frac{n}{(1-p)}$, whence $p = \frac{(\sum x_i - n)}{\sum x_i}$. For the given data, $\sum x_i = (1)(1) + (2)(31) + \dots + (12)(1) = 363$, so
- $$\hat{p} = \frac{(363 - 130)}{363} = .642, \text{ and } \hat{q} = .358.$$
- b.** Each estimated expected cell count is \hat{p} times the previous count, giving
- $$n\hat{q} = 130(.358) = 46.54, \quad n\hat{q}\hat{p} = 46.54(.642) = 29.88, \quad 19.18, \quad 12.31, \quad 7.91, \quad 5.08, \quad 3.26, \dots$$
- Grouping all values ≥ 7 into a single category gives 7 cells with estimated expected counts 46.54, 29.88, 19.18, 12.31, 7.91, 5.08 (sum = 120.9), and $130 - 120.9 = 9.1$. The corresponding observed counts are 48, 31, 20, 9, 6, 5, and 11, giving
- $$\mathbf{c}^2 = 1.87.$$
- With $k = 7$ and $m = 1$ (p was estimated), from (14.15) we need
- $$\mathbf{c}_{.10,5}^2 = 9.236.$$
- Since 1.87 is not ≥ 9.236 , we don't reject H_0 .

15. The part of the likelihood involving \mathbf{q} is $[(1-\mathbf{q})^4]^{n_1} \cdot [\mathbf{q}(1-\mathbf{q})^3]^{n_2} \cdot [\mathbf{q}^2(1-\mathbf{q})^2]^{n_3} \cdot [\mathbf{q}^3(1-\mathbf{q})]^{n_4} \cdot [\mathbf{q}^4]^{n_5} = \mathbf{q}^{n_2+2n_3+3n_4+4n_5} (1-\mathbf{q})^{4n_1+3n_2+2n_3+n_4} = \mathbf{q}^{233} (1-\mathbf{q})^{367}$, so $\ln(\text{likelihood}) = 233 \ln \mathbf{q} + 367 \ln(1-\mathbf{q})$. Differentiating and equating to 0 yields $\hat{\mathbf{q}} = \frac{233}{600} = .3883$, and $(1-\hat{\mathbf{q}}) = .6117$ [note that the exponent on \mathbf{q} is simply the total # of successes (defectives here) in the $n = 4(150) = 600$ trials.] Substituting this \mathbf{q}' into the formula for p_i yields estimated cell probabilities .1400, .3555, .3385, .1433, and .0227. Multiplication by 150 yields the estimated expected cell counts are 21.00, 53.33, 50.78, 21.50, and 3.41. the last estimated expected cell count is less than 5, so we combine the last two categories into a single one (≥ 3 defectives), yielding estimated counts 21.00, 53.33, 50.78, 24.91, observed counts 26, 51, 47, 26, and $\mathbf{c}^2 = 1.62$. With d.f. = $4 - 1 - 1 = 2$, since $1.62 < \mathbf{c}_{.10,2}^2 = 4.605$, the p-value $> .10$, and we do not reject H_0 . The data suggests that the stated binomial distribution is plausible.

16. $\hat{\mathbf{I}} = \bar{x} = \frac{(0)(6) + (1)(24) + (2)(42) + \dots + (8)(6) + (9)(2)}{300} = \frac{1163}{300} = 3.88$, so the estimated cell probabilities are computed from $\hat{p} = e^{-3.88} \frac{(3.88)^x}{x!}$.

x	0	1	2	3	4	5	6	7	≥ 8
np(x)	6.2	24.0	46.6	60.3	58.5	45.4	29.4	16.3	13.3
obs	6	24	42	59	62	44	41	14	8

This gives $\mathbf{c}^2 = 7.789$. To see whether the Poisson model provides a good fit, we need $\mathbf{c}_{.10,9-1-1}^2 = \mathbf{c}_{.10,7}^2 = 12.017$. Since $7.789 < 12.017$, the Poisson model does provide a good fit.

17. $\hat{\mathbf{I}} = \frac{380}{120} = 3.167$, so $\hat{p} = e^{-3.167} \frac{(3.167)^x}{x!}$.

x	0	1	2	3	4	5	6	≥ 7
\hat{p}	.0421	.1334	.2113	.2230	.1766	.1119	.0590	.0427
$n\hat{p}$	5.05	16.00	25.36	26.76	21.19	13.43	7.08	5.12
obs	24	16	16	18	15	9	6	16

The resulting value of $\mathbf{c}^2 = 103.98$, and when compared to $\mathbf{c}_{.01,7}^2 = 18.474$, it is obvious that the Poisson model fits very poorly.

18. $\hat{p}_1 = P(X < .100) = P\left(Z < \frac{.100 - .173}{.066}\right) = \Phi(-1.11) = .1335$,
 $\hat{p}_2 = P(.100 \leq X \leq .150) = P(-1.11 \leq Z \leq -.35) = .2297$,
 $\hat{p}_3 = P(-.35 \leq Z \leq .41) = .2959$, $\hat{p}_4 = P(.41 \leq Z \leq 1.17) = .2199$, and
 $\hat{p}_5 = .1210$. The estimated expected counts are then (multiply \hat{p}_i by $n = 83$) 11.08, 19.07,
 24.56, 18.25, and 10.04, from which $\mathbf{c}^2 = 1.67$. Comparing this with
 $\mathbf{c}_{.05,5-1-2}^2 = \mathbf{c}_{.05,2}^2 = 5.992$, the hypothesis of normality cannot be rejected.

19. With $A = 2n_1 + n_4 + n_5$, $B = 2n_2 + n_4 + n_6$, and $C = 2n_3 + n_5 + n_6$, the likelihood is proportional to $\mathbf{q}_1^A \mathbf{q}_2^B (1 - \mathbf{q}_1 - \mathbf{q}_2)^C$, where $A + B + C = 2n$. Taking the natural log and equating both $\frac{\partial}{\partial \mathbf{q}_1}$ and $\frac{\partial}{\partial \mathbf{q}_2}$ to zero gives $\frac{A}{\mathbf{q}_1} = \frac{C}{1 - \mathbf{q}_1 - \mathbf{q}_2}$ and $\frac{B}{\mathbf{q}_2} = \frac{C}{1 - \mathbf{q}_1 - \mathbf{q}_2}$, whence $\mathbf{q}_2 = \frac{B\mathbf{q}_1}{A}$. Substituting this into the first equation gives $\mathbf{q}_1 = \frac{A}{A + B + C}$, and then $\mathbf{q}_2 = \frac{B}{A + B + C}$. Thus $\hat{\mathbf{q}}_1 = \frac{2n_1 + n_4 + n_5}{2n}$, $\hat{\mathbf{q}}_2 = \frac{2n_2 + n_4 + n_6}{2n}$, and $(1 - \hat{\mathbf{q}}_1 - \hat{\mathbf{q}}_2) = \frac{2n_3 + n_5 + n_6}{2n}$. Substituting the observed n_i 's yields $\hat{\mathbf{q}}_1 = \frac{2(49) + 20 + 53}{400} = .4275$, $\hat{\mathbf{q}}_2 = \frac{110}{400} = .2750$, and $(1 - \hat{\mathbf{q}}_1 - \hat{\mathbf{q}}_2) = .2975$, from which $\hat{p}_1 = (.4275)^2 = .183$, $\hat{p}_2 = .076$, $\hat{p}_3 = .089$, $\hat{p}_4 = 2(.4275)(.275) = .235$, $\hat{p}_5 = .254$, $\hat{p}_6 = .164$.

Category	1	2	3	4	5	6
np	36.6	15.2	17.8	47.0	50.8	32.8
observed	49	26	14	20	53	38

This gives $\mathbf{c}^2 = 29.1$. With $\mathbf{c}_{.01,6-1-2}^2 = \mathbf{c}_{.01,3}^2 = 11.344$, and $\mathbf{c}_{.01,6-1}^2 = \mathbf{c}_{.01,5}^2 = 15.085$, according to (14.15) H_0 must be rejected since $29.1 \geq 15.085$.

20. The pattern of points in the plot appear to deviate from a straight line, a conclusion that is also supported by the small p-value ($< .01000$) of the Ryan-Joiner test. Therefore, it is implausible that this data came from a normal population. In particular, the observation 116.7 is a clear outlier. It would be dangerous to use the one-sample t interval as a basis for inference.

Chapter 14: The Analysis of Categorical Data

21. The Ryan-Joiner test p-value is larger than .10, so we conclude that the null hypothesis of normality cannot be rejected. This data could reasonably have come from a normal population. This means that it would be legitimate to use a one-sample t test to test hypotheses about the true average ratio.

22.

x_i	y_i	x_i	y_i	x_i	y_i
69.5	-1.967	75.5	-.301	79.6	.634
71.9	-1.520	75.7	-.199	79.7	.761
72.6	-1.259	75.8	-.099	79.9	.901
73.1	-1.063	76.1	.000	80.1	1.063
73.3	-.901	76.2	.099	82.2	1.259
73.5	-.761	76.9	.199	83.7	1.520
74.1	-.634	77.0	.301	93.7	1.967
74.2	-.517	77.9	.407		
75.3	-.407	78.1	.517		

n.b.: Minitab was used to calculate the y_I 's. $\sum x_{(i)} = 1925.6$, $\sum x_{(i)}^2 = 148,871$, $\sum y_i = 0$,

$\sum y_i^2 = 22.523$, $\sum x_{(i)}y_i = 103.03$, so

$$r = \frac{25(103.03)}{\sqrt{25(148,871) - (1925.6)^2} \sqrt{25(22.523)}} = .923. \text{ Since } c_{.01} = .9408, \text{ and } .923 < .9408,$$

even at the very smallest significance level of .01, the null hypothesis of population normality must be rejected (the largest observation appears to be the primary culprit).

23. Minitab gives $r = .967$, though the hand calculated value may be slightly different because when there are ties among the $x_{(i)}$'s, Minitab uses the same y_I for each $x_{(i)}$ in a group of tied values. $C_{10} = .9707$, and $c_{.05} = .9639$, so $.05 < p\text{-value} < .10$. At the 5% significance level, one would have to consider population normality plausible.

Section 14.3

24. H_0 : TV watching and physical fitness are independent of each other

H_a : the two variables are not independent

$$Df = (4 - 1)(2 - 1) = 3$$

With $\alpha = .05$, RR: $c^2 \geq 7.815$

Computed $c^2 = 6.161$

Fail to reject H_0 . The data fail to indicate an association between daily TV viewing habits and physical fitness.

Chapter 14: The Analysis of Categorical Data

25. Let P_{ij} = the proportion of white clover in area of type i which has a type j mark ($i = 1, 2$; $j = 1, 2, 3, 4, 5$). The hypothesis $H_0: p_{1j} = p_{2j}$ for $j = 1, \dots, 5$ will be rejected at level .01 if $\mathbf{c}^2 \geq \mathbf{c}_{.01, (2-1)(5-1)}^2 = \mathbf{c}_{.01, 4}^2 = 13.277$.

\hat{E}_{ij}	1	2	3	4	5		
1	449.66	7.32	17.58	8.79	242.65	726	$\mathbf{c}^2 = 23.18$
2	471.34	7.68	18.42	9.21	254.35	761	
	921	15	36	18	497	1487	

Since $23.18 \geq 13.277$, H_0 is rejected.

26. Let p_{i1} = the probability that a fruit given treatment i matures and p_{i2} = the probability that a fruit given treatment i aborts. Then $H_0: p_{i1} = p_{i2}$ for $i = 1, 2, 3, 4, 5$ will be rejected if $\mathbf{c}^2 \geq \mathbf{c}_{.01, 4}^2 = 13.277$.

Observed		Estimated Expected		
Matured	Aborted	Matured	Aborted	
141	206	110.7	236.3	347
28	69	30.9	66.1	97
25	73	31.3	66.7	98
24	78	32.5	69.5	102
20	82	32.5	69.5	102
		238	508	746

Thus $\mathbf{c}^2 = \frac{(141 - 110.7)^2}{110.7} + \dots + \frac{(82 - 69.5)^2}{69.5} = 24.82$, which is ≥ 13.277 , so H_0 is rejected at level .01.

27. With $i = 1$ identified with men and $i = 2$ identified with women, and $j = 1, 2, 3$ denoting the 3 categories $L > R$, $L = R$, $L < R$, we wish to test $H_0: p_{1j} = p_{2j}$ for $j = 1, 2, 3$ vs. $H_a: p_{1j}$ not equal to p_{2j} for at least one j . The estimated cell counts for men are 17.95, 8.82, and 13.23 and for women are 39.05, 19.18, 28.77, resulting in $\mathbf{c}^2 = 44.98$. With $(2 - 1)(3 - 1) = 2$ degrees of freedom, since $44.98 > \mathbf{c}_{.005, 2}^2 = 10.597$, p -value $< .005$, which strongly suggests that H_0 should be rejected.

Chapter 14: The Analysis of Categorical Data

28. With p_{ij} denoting the probability of a type j response when treatment i is applied, $H_0: p_{1j} = p_{2j} = p_{3j} = p_{4j}$ for $j = 1, 2, 3, 4$ will be rejected at level .005 if $\mathbf{c}^2 \geq \mathbf{c}_{.005,9}^2 = 23.587$.

\hat{E}_{ij}	1	2	3	4
1	24.1	10.0	21.6	40.4
2	25.8	10.7	23.1	43.3
3	26.1	10.8	23.4	43.8
4	30.1	12.5	27.0	50.5

$\mathbf{c}^2 = 27.66 \geq 23.587$, so reject H_0 at level .005

29. $H_0: p_{1j} = \dots = p_{6j}$ for $j = 1, 2, 3$ is the hypothesis of interest, where p_{ij} is the proportion of the j^{th} sex combination resulting from the i^{th} genotype. H_0 will be rejected at level .10 if $\mathbf{c}^2 \geq \mathbf{c}_{.10,10}^2 = 15.987$.

\hat{E}_{ij}	1	2	3		\mathbf{c}^2	1	2	3	
1	35.8	83.1	35.1	154		.02	.12	.44	
2	39.5	91.8	38.7	170		.06	.66	1.01	
3	35.1	81.5	34.4	151		.13	.37	.34	
4	9.8	22.7	9.6	42		.32	.49	.26	
5	5.1	11.9	5.0	22		.00	.06	.19	
6	26.7	62.1	26.2	115		.40	.14	1.47	
	152	353	149	654					6.46

(carrying 2 decimal places in \hat{E}_{ij} yields $\mathbf{c}^2 = 6.49$). Since $6.46 < 15.987$, H_0 cannot be rejected at level .10.

Chapter 14: The Analysis of Categorical Data

30. H_0 : the design configurations are homogeneous with respect to type of failure vs. H_a : the design configurations are not homogeneous with respect to type of failure.

\hat{E}_{ij}	1	2	3	4	
1	16.11	43.58	18.00	12.32	90
2	7.16	19.37	8.00	5.47	40
3	10.74	29.05	12.00	8.21	60
	34	92	38	26	190

$$c^2 = \frac{(20 - 16.11)^2}{16.11} + \dots + \frac{(5 - 8.21)^2}{8.21} = 13.253. \text{ With 6 df,}$$

$c_{.05,6}^2 = 12.592 < 13.253 < c_{.025,6}^2 = 14.440$, so $.025 < \text{p-value} < .05$. Since the p-value is $< .05$, we reject H_0 . (If a smaller significance level were chosen, a different conclusion would be reached.) Configuration appears to have an effect on type of failure.

31. With I denoting the I^{th} type of car ($I = 1, 2, 3, 4$) and j the j^{th} category of commuting distance, $H_0: p_{ij} = p_i \cdot p_j$ (type of car and commuting distance are independent) will be rejected at level .05 if $c^2 \geq c_{.05,6}^2 = 12.592$.

\hat{E}_{ij}	1	2	3	
1	10.19	26.21	15.60	52
2	11.96	30.74	18.30	61
3	19.40	49.90	29.70	99
4	7.45	19.15	11.40	38
	49	126	75	250

$c^2 = 14.15 \geq 12.592$, so the independence hypothesis H_0 is rejected at level .05 (but not at level .025!)

32.
$$c^2 = \frac{(479 - 494.4)^2}{494.4} + \frac{(173 - 151.5)^2}{151.5} + \frac{(119 - 125.2)^2}{125.2} + \frac{(214 - 177.0)^2}{177.0} + \frac{(47 - 54.2)^2}{54.2}$$

$$= \frac{(15 - 44.8)^2}{44.8} + \frac{(172 - 193.6)^2}{193.6} + \frac{(45 - 59.3)^2}{59.3} + \frac{(85 - 49.0)^2}{49.0} = 64.65 \geq c_{.01,4}^2 = 13.277$$

so the independence hypothesis is rejected in favor of the conclusion that political views and level of marijuana usage are related.

33.
$$c^2 = \sum \sum \frac{(N_{ij} - \hat{E}_{ij})^2}{\hat{E}_{ij}} = \sum \sum \frac{N_{ij}^2 - 2\hat{E}_{ij}N_{ij} + \hat{E}_{ij}^2}{\hat{E}_{ij}} = \frac{\sum \sum N_{ij}^2}{\hat{E}_{ij}} - 2\sum \sum N_{ij} + \sum \sum \hat{E}_{ij},$$
 but

$\sum \sum \hat{E}_{ij} = \sum \sum N_{ij} = n$, so $c^2 = \sum \sum \frac{N_{ij}^2}{\hat{E}_{ij}} - n$. This formula is computationally efficient

because there is only one subtraction to be performed, which can be done as the last step in the calculation.

34. This is a $3 \times 3 \times 3$ situation, so there are 27 cells. Only the total sample size n is fixed in advance of the experiment, so there are 26 freely determined cell counts. We must estimate $p_{..1}, p_{..2}, p_{..3}, p_{.1.}, p_{.2.}, p_{.3.}, p_{1.}, p_{2.},$ and $p_{3.}$, but $\sum p_{i..} = \sum p_{.j.} = \sum p_{..k} = 1$ so only 6 independent parameters are estimated. The rule for d.f. now gives c^2 df = $26 - 6 = 20$.

35. With p_{ij} denoting the common value of $p_{ij1}, p_{ij2}, p_{ij3}, p_{ij4}$ (under H_0), $\hat{p}_{ij} = \frac{N_{ij}}{n}$ and $\hat{E}_{ijk} = \frac{n_k N_{ij}}{n}$. With four different tables (one for each region), there are $8 + 8 + 8 + 8 = 32$ freely determined cell counts. Under H_0 , p_{11}, \dots, p_{33} must be estimated but $\sum \sum p_{ij} = 1$ so only 8 independent parameters are estimated, giving c^2 df = $32 - 8 = 24$.

36.

a.

Observed				Estimated Expected		
13	19	28	60	12	18	30
7	11	22	40	8	12	20
20	30	50	100			

$$c^2 = \frac{(13 - 12)^2}{12} + \dots + \frac{(22 - 20)^2}{20} = .6806$$
. Because $.6806 < c^2_{.10,2} = 4.605$, H_0 is not rejected.

b. Each observation count here is 10 times what it was in a, and the same is true of the estimated expected counts so now $c^2 = 6.806 \geq 4.605$, and H_0 is rejected. With the much larger sample size, the departure from what is expected under H_0 , the independence hypothesis, is statistically significant – it cannot be explained just by random variation.

c. The observed counts are $.13n, .19n, .28n, .07n, .11n, .22n$, whereas the estimated expected $\frac{(.60n)(.20n)}{n} = .12n, .18n, .30n, .08n, .12n, .20n$, yielding $c^2 = .006806n$.

H_0 will be rejected at level .10 iff $.006806n \geq 4.605$, i.e., iff $n \geq 676.6$, so the minimum $n = 677$.

Supplementary Exercises

37. There are 3 categories here – firstborn, middleborn, (2nd or 3rd born), and lastborn. With p_1 , p_2 , and p_3 denoting the category probabilities, we wish to test $H_0: p_1 = .25, p_2 = .50$ ($p_2 = P(2^{\text{nd}}$ or 3^{rd} born) = $.25 + .25 = .50$), $p_3 = .25$. H_0 will be rejected at significance level .05 if

$$c^2 \geq c^2_{.05,2} = 5.992. \text{ The expected counts are } (31)(.25) = 7.75, (31)(.50) = 15.5, \text{ and } 7.75,$$

$$\text{so } c^2 = \frac{(12 - 7.75)^2}{7.75} + \frac{(11 - 15.5)^2}{15.5} + \frac{(8 - 7.75)^2}{7.75} = 3.65. \text{ Because } 3.65 < 5.992, H_0 \text{ is not}$$

rejected. The hypothesis of equiprobable birth order appears quite plausible.

38. Let p_{i1} = the proportion of fish receiving treatment i ($i = 1, 2, 3$) who are parasitized. We wish to test $H_0: p_{1j} = p_{2j} = p_{3j}$ for $j = 1, 2$. With $df = (2 - 1)(3 - 1) = 2$, H_0 will be rejected at level .01 if $c^2 \geq c^2_{.01,2} = 9.210$.

Observed			Estimated Expected	
30	3	33	22.99	10.01
16	8	24	16.72	7.28
16	16	32	22.29	9.71
62	27	89		

This gives $c^2 = 13.1$. Because $13.1 \geq 9.210$, H_0 should be rejected. The proportion of fish that are parasitized does appear to depend on which treatment is used.

39. H_0 : gender and years of experience are independent; H_a : gender and years of experience are not independent. $Df = 4$, and we reject H_0 if $c^2 \geq c^2_{.01,4} = 13.277$.

Gender	Years of Experience				
	1 – 3	4 – 6	7 – 9	10 – 12	13 +
Male Observed	202	369	482	361	811
Expected	285.56	409.83	475.94	347.04	706.63
$\frac{(O-E)^2}{E}$	24.451	4.068	.077	.562	15.415
Female Observed	230	251	238	164	258
Expected	146.44	210.17	244.06	177.96	362.37
$\frac{(O-E)^2}{E}$	47.680	7.932	.151	1.095	30.061

$$c^2 = \sum \frac{(O-E)^2}{E} = 131.492. \text{ Reject } H_0. \text{ The two variables do not appear to be independent.}$$

In particular, women have higher than expected counts in the beginning category (1 – 3 years) and lower than expected counts in the more experienced category (13+ years).

Chapter 14: The Analysis of Categorical Data

40.

- a. H_0 : The probability of a late-game leader winning is independent of the sport played; H_a : The two variables are not independent. With 3 df, the computed $\chi^2 = 10.518$, and the p-value $< .015$ is also $< .05$, so we would reject H_0 . There appears to be a relationship between the late-game leader winning and the sport played.
- b. Quite possibly: Baseball had many fewer than expected late-game leader losses.

41. The null hypothesis $H_0: p_{ij} = p_i \cdot p_j$ states that level of parental use and level of student use are independent in the population of interest. The test is based on $(3 - 1)(3 - 1) = 4$ df.

	Estimated Expected		
119.3	57.6	58.1	235
82.8	33.9	40.3	163
23.9	11.5	11.6	47
226	109	110	445

The calculated value of $\chi^2 = 22.4$. Since $22.4 > \chi^2_{.005,4} = 14.860$, p-value $< .005$, so H_0 should be rejected at any significance level greater than .005. Parental and student use level do not appear to be independent.

42. The estimated expected counts are displayed below, from which $\chi^2 = 197.70$. A glance at the 6 df row of Table A.7 shows that this test statistic value is highly significant – the hypothesis of independence is clearly implausible.

	Estimated Expected			
	Home	Acute	Chronic	
15 – 54	90.2	372.5	72.3	535
55 – 64	113.6	469.3	91.1	674
65 – 74	142.7	589.0	114.3	846
> 74	157.5	650.3	126.2	934
	504	2081	404	2989

Chapter 14: The Analysis of Categorical Data

43. This is a test of homogeneity: $H_0: p_{1j} = p_{2j} = p_{3j}$ for $j = 1, 2, 3, 4, 5$. The given SPSS output reports the calculated $\mathbf{c}^2 = 70.64156$ and accompanying p-value (significance) of .0000. We reject H_0 at any significance level. The data strongly supports that there are differences in perception of odors among the three areas.
44. The accompanying table contains both observed and estimated expected counts, the latter in parentheses.

	Age					
Want	127 (131.1)	118 (123.3)	77 (71.7)	61 (55.1)	41 (42.8)	424
Don't	23 (18.9)	23 (17.7)	5 (10.3)	2 (7.9)	8 (6.2)	61
	150	141	82	63	49	485

- This gives $\mathbf{c}^2 = 11.60 \geq \mathbf{c}_{.05,4}^2 = 9.488$. At level .05, the null hypothesis of independence is rejected, though it would not be rejected at the level .01 ($.01 < \text{p-value} < .025$).
45. $(n_1 - np_{10})^2 = (np_{10} - n_1)^2 = (n - n_1 - n(1 - p_{10}))^2 = (n_2 - np_{20})^2$. Therefore
- $$\mathbf{c}^2 = \frac{(n_1 - np_{10})^2}{np_{10}} + \frac{(n_2 - np_{20})^2}{np_{20}} = \frac{(n_1 - np_{10})^2}{n_2} \left(\frac{n}{p_{10}} + \frac{n}{p_{20}} \right)$$
- $$= \left(\frac{n_1}{n} - p_{10} \right)^2 \cdot \left(\frac{n}{p_{10}p_{20}} \right) = \frac{(\hat{p}_1 - p_{10})^2}{\frac{p_{10}p_{20}}{n}} = z^2.$$

46. a.

obsv	22	10	5	11
exp	13.189	10	7.406	17.405

H_0 : probabilities are as specified.

H_a : probabilities are not as specified.

Test Statistic: $\mathbf{c}^2 = \frac{(22 - 13.189)^2}{13.189} + \frac{(10 - 10)^2}{10} + \frac{(5 - 7.406)^2}{7.406} + \frac{(11 - 17.405)^2}{17.405}$

$= 5.886 + 0 + 0.782 + 2.357 = 9.025$. Rejection Region: $\mathbf{c}^2 > \mathbf{c}_{.05,2}^2 = 5.99$

Since $9.025 > 5.99$, we reject H_0 . The model postulated in the exercise is not a good fit.

Chapter 14: The Analysis of Categorical Data

b.

p_i	0.45883	0.18813	0.11032	0.24272
exp	22.024	9.03	5.295	11.651

$$c^2 = \frac{(22 - 22.024)^2}{22.024} + \frac{(10 - 9.03)^2}{9.03} + \frac{(5 - 5.295)^2}{5.295} + \frac{(11 - 11.651)^2}{11.651}$$

$$= .0000262 + .1041971 + .0164353 + .0363746 = .1570332$$

With the same rejection region as in a, we do not reject the null hypothesis. This model does provide a good fit.

47.

a. Our hypotheses are H_0 : no difference in proportion of concussions among the three groups. Vs H_a : there is a difference ...

Observed	Concussion	No Concussion	Total
Soccer	45	46	91
Non Soccer	28	68	96
Control	8	45	53
Total	81	159	240

Expected	Concussion	No Concussion	Total
Soccer	30.7125	60.2875	91
Non Soccer	32.4	63.6	96
Control	17.8875	37.1125	53
Total	81	159	240

$$c^2 = \frac{(45 - 30.7125)^2}{30.7125} + \frac{(46 - 60.2875)^2}{60.2875} + \frac{(28 - 32.4)^2}{32.4} + \frac{(68 - 63.6)^2}{63.6}$$

$$+ \frac{(8 - 17.8875)^2}{17.8875} + \frac{(45 - 37.1125)^2}{37.1125} = 19.1842$$

The df for this test is $(I - 1)(J - 1) = 2$, so we reject H_0 if $c^2 > c^2_{.05,2} = 5.99$. $19.1842 > 5.99$, so we reject H_0 . There is a difference in the proportion of concussions based on whether a person plays soccer.

b. We are testing the hypothesis $H_0: \rho = 0$ vs $H_a: \rho \neq 0$. The test statistic is

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{-.22\sqrt{89}}{\sqrt{1-.22^2}} = -2.13$$

At significance level $\alpha = .01$, we would fail to reject and conclude that there is no evidence of non-zero correlation in the population. If we were willing to accept a higher significance level, our decision could change. At best, there is evidence of only weak correlation.

Chapter 14: The Analysis of Categorical Data

- c. We will test to see if the average score on a controlled word association test is the same for soccer and non-soccer athletes. $H_0: \mu_1 = \mu_2$ vs $H_a: \mu_1 \neq \mu_2$. We'll use test statistic

$$t = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{\frac{s_1^2}{m} + \frac{s_2^2}{n}}}. \text{ With } \frac{s_1^2}{m} = 3.206 \text{ and } \frac{s_2^2}{n} = 1.854,$$

$$t = \frac{(37.50 - 39.63)}{\sqrt{3.206 + 1.854}} = -.95. \text{ The df} = \frac{(3.206 + 1.854)^2}{\frac{3.206^2}{25} + \frac{1.854^2}{55}} \approx 56. \text{ The p-value will}$$

be $> .10$, so we do not reject H_0 and conclude that there is no difference in the average score on the test for the two groups of athletes.

- d. Our hypotheses for ANOVA are H_0 : all means are equal vs H_a : not all means are equal.

The test statistic is $f = \frac{MSTr}{MSE}$.

$$SSTr = 91(.30 - .35)^2 + 96(.49 - .35)^2 + 53(.19 - .35)^2 = 3.4659$$

$$MSTr = \frac{3.4659}{2} = 1.73295$$

$$SSE = 90(.67)^2 + 95(.87)^2 + 52(.48)^2 = 124.2873 \text{ and}$$

$$MSE = \frac{124.2873}{237} = .5244. \text{ Now, } f = \frac{1.73295}{.5244} = 3.30. \text{ Using df 2,200 from}$$

table A.9, the p value is between .01 and .05. At significance level .05, we reject the null hypothesis. There is sufficient evidence to conclude that there is a difference in the average number of prior non-soccer concussions between the three groups.

48.

- a. $H_0: p_0 = p_1 = \dots = p_9 = .10$ vs H_a : at least one $p_i \neq .10$, with $df = 9$.
- b. $H_0: p_{ij} = .01$ for i and $j = 1, 2, \dots, 9$ vs H_a : at least one $p_{ij} \neq .01$, with $df = 99$.
- c. For this test, the number of p 's in the Hypothesis would be $10^5 = 100,000$ (the number of possible combinations of 5 digits). Using only the first 100,000 digits in the expansion, the number of non-overlapping groups of 5 is only 20,000. We need a much larger sample size!
- d. Based on these p-values, we could conclude that the digits of p behave as though they were randomly generated.