Stat 511, HW 7 Answers

1 point each part, with 3 points for free.

1. Analysis of boy/girl ratios

   (a) Test the hypothesis that $\pi_0 = 0.5$, using test statistic $z = \frac{\pi - \pi_0}{\sqrt{\pi_0(1-\pi_0)/n}} \sim N(0,1)$

   p-value $<0.001$. Very strong evidence that the proportion of girls is not equal to 0.5.

   (b) Estimate $\pi$ and calculate a 95% confidence interval. Approximate (or asymptotic) $100(1-\alpha)$% CI for $\pi$ is $\pi \pm z_{1-\alpha/2}\sqrt{\pi(1-\pi)/n} =$
   $(0.4799, 0.4894)$

   (c) Test whether the proportion of female children is different in 6 child families and 12 child families. Use pooled estimates of $\pi = \frac{y_1 + y_2}{n_1 + n_2}$, then

   $$\frac{\pi_1 - \pi_2}{\sqrt{\pi(1-\pi)(1/n_1 + 1/n_2)}} \sim N(0,1)$$

   p-value $= 0.59$
   Approximate $100(1-\alpha)$% CI for $\pi_1 - \pi_2$ is

   $$\pi_1 - \pi_2 \pm z_{1-\alpha/2}\sqrt{\pi_1(1-\pi_1)/n_1 + \pi_2(1-\pi_2)/n_2} = c(-0.0090, 0.0159)$$

   (d) No. Since variation in $\pi$ between families violates independence and identical distribution assumptions.
   Note: Only need to give one (either independence or no variation in probability between families), since they are two ways to say the same thing.

2. Analysis of boy/girl ratios

   (a) $P(x = 6) = \frac{6!}{6!0!}\pi^6(1-\pi)^0$ ,$\pi = 0.48465$ and E(# 6 girl families)$=NP(x = 6) = 93.3$

   (b) $\sum_i \frac{(O_i - E_i)^2}{E_i} = 12.6775 + \frac{(113-93.3)^2}{93.3} = 16.837 \sim \chi^2_{df=6}$, p-value$=0.0099$. Strong evidence that the distribution is not binomial with probability of a girl $= 0.48465$

3. breast cancer

   (a) Test whether drinking is associated with breast cancer. Two possible approaches:

      i. Define $\pi$ as the probability of light drinking. $\hat\pi_1 = 0.610$, $\hat\pi_2 = 0.510$. Use pooled estimates of $\pi = \frac{y_1 + y_2}{n_1 + n_2} = 250/459 = 0.5446$, then

      $$\frac{\pi_1 - \pi_2}{\sqrt{\pi(1-\pi)(1/n_1 + 1/n_2)}} = 2.047 \sim N(0,1)$$

      . p $= 0.0407$

1

ii. default $\chi^2$ test p-value reported by R is 0.0512

But, if you add "correct='F'" to turn off the continuity correction then the p-value=0.405

Notes: 1) The difference between 0.0405 and 0.0407 is roundoff error in intermediate values.

2) This is a case where the loss of power with the continuity correction matters.

(b) Estimate the odds ratio (as odds of breast cancer in heavy drinkers to odds in light drinkers)

$$\phi = \frac{\pi_{heavy}/(1 - \pi_{heavy})}{\pi_{light}/(1 - \pi_{light})} = 0.666$$

(c) Estimate the standard error of the log odds ratio

$$s.e.(\log(\phi)) = \sqrt{\frac{1}{Y_{11}} + \frac{1}{Y_{12}} + \frac{1}{Y_{21}} + \frac{1}{Y_{22}}} = 0.199$$

(d) Calculate a 95% confidence interval for the odds ratio.

$$\log(\phi) \pm 1.96 \times s.e.(\log(\phi)) = (0.275, 1.056)$$

. So, the 95% ci for the odds ratio is $(\exp 0.275, \exp 1.056) = (1.32, 2.87)$

4. Statistics and Physics

(a) There is no evidence of a difference between colleges in the proportion of Physics students passing the exam. Because of small sample size, Fisher exact test is used and p-value=1.

Notes:

1) You need to look at the expected counts, not the observed counts.

2) The p-value is 1 because R doubles the one-sided probability of a more extreme table.

(b) Small university has the higher percentage of passes if you are not told the department

| college | Pass | Fail |
|---------|------|------|
| Big U   | 20   | 20   |
| Small U | 21   | 19   |

(c) Both colleges have a higher pass rate for statistics classes than physics classes. The small university has the higher proportion of students in Statistics, while the big university has approximately the same proportion of students in both department.

Notes:

i. This is an example what is now called "Simpson's Paradox". Statistically, it is the same issue as the difference in 2 way ANOVA between "one-bucket" means (aggregating over departments) and "averages of cell means" (keeping departments separate).

ii. The Wikipedia article on Simpson's paradox is very informative. It includes a very good description of the UC Berkeley gender bias case.

2

5. valve failure in a nuclear reactor

   (a) Estimate the effect associated with each level of the five explanatory factors
       Since I didn't specify how to estimate the effect, here are the estimated coefficients in
       the R default contr.treatmetn

       ```
             (Intercept) systemcContainment       systemcNuclear       systemcPower
              -3.64510173         -0.33291634           0.58263966         0.68589050
          systemcSafety     valvecButterfly       valvecDiaphram   valvecDirectional
             0.89017661          0.18532556           0.60673702         1.00891454
             valvecGate         valvecGlobe modecNormally Open     operatorcManual
             2.95894266          1.79318168           0.20934332        -2.47232574
        operatorcMoter    operatorcSolenoid         sizec10-30 in        sizec2-10 in
             -1.19261155          0.70436725           1.61456696        -0.01219413
       ```

       I had not intended to make you copy this information. No points off if you omitted it.
       It is necessary to provide what I specifically asked for, which was the levels with a more
       than 5x increase in failure rate from the most reliable level:
       valve Gate: 19.28 times. Log minimum failure time is for Ball = $\exp(0) = 1$
       valve Globe: 6.01 times
       operator Air: 11.8 times. Log minimum failure time is for Manual = $\exp(-2.47) = 0.0845$
       operator Solenoid: 23.9 times
       size 10-30 inch: 5.023 times. Log minimum failure rate is for 0-2 in = $\exp(0) = 1$
       Notes:

         i. I didn't ask for the minimum failure rate for each factor. That information is included
            in our answers only to help you understand our answer.
        ii. the beta's are the log changes in failure rate when all else is held constant.
       iii. Don't forget that there is a "hidden" level with beta = 0. That is crucial for inter-
            preting the valve coefficients where that 0 is the smallest level = most reliable.
        iv. You can calculate failure rates = $\exp \beta$ for each level or calculate the necessary
            difference of a beta from the minimimum = $\log 5$.

   (b) Test no difference in failure rate between systems when all other factors held constant
       p-value=0.017, evidence of difference in failure rate between systems when all other
       factors held constant.

       Note: Can calculate this from the drop1() approach with helmert contrasts, from an
       explicit model comparison of "with systemc" to "without systemc", or by the usual R
       sequential SS when systemc is the **last** term in the model.

   (c) Evidence of over dispersion, $\hat{\phi} = 195.68/743.41 = 2.64$. If you used Pearson Chi-square
       / df, you got $\hat{\phi} = 4.52$.

   (d) Test the null hypothesis of no difference in failure rate between systems while accounting
       for the over dispersion. $p - value = 0.61$ (using Deviance estimate of OD); $p - value =$
       0.34 (using Chisquare estimate of OD). no evidence of difference in failure rate between

systems.

R code

```
#1
prop.test(t(c(20937,22263)), p=0.5)
prop.test(rbind(c(20937,22263), c(3534,3810)))
#2
pi=0.48465;px6=pi^6;7200*px6
12.6775 + (113-93.3)^2/93.3
1-pchisq(16.8371, 6)  # or pchisq(16.8371, 6, lower=F)
#3
#chisq test
#  prop.test(rbind(c(97,153), c(62,147)))  # WIthout continuity corr., not recomme
prop.test(rbind(c(97,153), c(62,147)),correct=F)
logodd=0.666
se=sqrt(1/62+1/147+1/97+1/153)
c(logodd-1.96*se,logodd+1.96*se)
exp(c(logodd-1.96*se,logodd+1.96*se))
#4
fisher.test(rbind(c(8,14),c(1,3)))
#5
a=read.csv("valves.csv")
o=glm(Failures~systemc+valvec+modec+operatorc+sizec,offset=logtime,data=a,family=p
exp(o$coeff)


o2=glm(Failures~valvec+modec+operatorc+sizec+systemc,offset=logtime,data=a,family=
anova(o2,test="Chisq")
# above is testing systemc with all else held constant by putting it last
quasio=glm(Failures~valvec+modec+operatorc+sizec+systemc,offset=logtime,
    data=a,family=quasipoisson)
anova(quasio,test='F')
# OR, use drop1() on the original model: drop1(o, ~systemc)
drop1(o2,~systemc)
# If you do both, you find that the answers are different
#  anova() or an explicit model comparison estimates overdispersion
#     by Pearson Chisquare/df
#  drop1() estimates overdispersion by deviance/df
# often the two estimates are similar.  For these data, the values are 4.52 and 2.
#   which are quite different
```