

# Stat 311: HW 10, not due, solutions to be posted before final

Fritz Scholz

1. Section 13.4, problem 1. Basically you test the hypothesis  $p_1 = \dots = p_8 = 1/8$ , where  $p_i$  denotes the probability that the winner comes from starting position  $i$ . This is based on assigning the horses at random to the starting positions and the hypothesis assumption that starting position has no influence on who the winner is. Thus the winner has equal chance of coming from any of the 8 starting positions. Since no significance level  $\alpha$  is specified, just compute the  $p$ -value and comment on the strength of the evidence against the hypothesis.

```
horses <- function(){
finish <- c(29,19,18,25,17,10,15,11)
k <- length(finish)
n <- sum(finish)
p <- rep(1/k,k)
e <- n*p
X2 <- sum((finish-e)^2/e)
G2 <- 2*sum(finish*log(finish/e))
pX2 <- 1-pchisq(X2,k-1)
pG2 <- 1-pchisq(G2,k-1)
out <- c(X2,pX2,G2,pG2)
names(out) <- c("X2", "p-value (X2)", "G2", "p-value (G2)")
out
}
# with the following result when calling horses()
> horses()
           X2 p-value(X2)           G2 p-value(G2)
16.33333333 0.02223948 16.13810519 0.02388413
```

At level  $\alpha = 0.05$  one should reject the hypothesis that starting position has no influence on which horse is the winner, since the  $p$ -value is  $< 0.05$  for either test statistic.

2. Section 13.4, problem 3. Again, no  $\alpha$  is specified. Thus compute the  $p$ -value and comment on the strength of the evidence against the hypothesis that the cell probabilities in (a) are correct.

(a) The respective cell probabilities (by independence) should be

$$p_1 = P(E_1) = \frac{3}{4} \cdot \frac{3}{4} = \frac{9}{16} \quad p_2 = P(E_2) = \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{16} \quad p_3 = P(E_3) = \frac{1}{4} \cdot \frac{3}{4} = \frac{3}{16} \quad p_4 = P(E_4) = \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{16}$$

(b)

```
recessive <- function(){
obs <- c(926,288,293,104)
n <- sum(obs)
e <- n*c(9/16,3/16,3/16,1/16)
X2 <- sum((obs-e)^2/e)
G2 <- 2*sum(obs*log(obs/e))
pX2 <- 1-pchisq(X2,3)
```

```

pG2 <- 1-pchisq(G2,3)
out <- c(X2,pX2,G2,pG2)
names(out) <- c("X2", "p-value(X2)", "G2", "p-value(G2)")
out
}
# with following result when calling recessive()
> recessive()
      X2 p-value(X2)      G2 p-value(G2)
1.4687220 0.6895079 1.4775868 0.6874529

```

Clearly there is no evidence against the assumed probabilities for the 4 events, since the  $p$ -values are quite large.

**3.** Section 13.4, problem 5. Note that  $\log_{10}(x)$  gives you that number  $y$  such that  $10^y = x$ , e.g.,  $\log_{10}(1000) = 3$  since  $10^3 = 1000$ . Furthermore, we have

$$\log_{10}(x) \geq 0 \quad \text{for } x \geq 1 \quad (1)$$

and

$$y = \log_{10}(x_1 \cdot x_2 \cdot \dots \cdot x_n) = \log_{10}(x_1) + \log_{10}(x_2) + \dots + \log_{10}(x_n) = y_1 + y_2 + \dots + y_n \quad (2)$$

since

$$10^{y_1+y_2+\dots+y_n} = 10^{y_1} \cdot 10^{y_2} \cdot \dots \cdot 10^{y_n} = x_1 \cdot x_2 \cdot \dots \cdot x_n = 10^y$$

Use the properties (1) and (2) to do part (a). In R the command `p <- log10(1+1/(1:9))` would give you the vector of required cell probabilities. You can also check `sum(p)`.

Benford's law (see [http://en.wikipedia.org/wiki/Benford%27s\\_law](http://en.wikipedia.org/wiki/Benford%27s_law)) is quite useful in detecting fraudulent activities when numbers are just made up, such as in accounting when cooking the books, or in faking election results. The latter issue was most recently examined w.r.t. the election results in Iran.

<http://blog.jgc.org/2009/06/benford-s-law-and-iranian-election.html>

Since no  $\alpha$  is specified, work with the  $p$ -value to assess the evidence.

(a) Since  $1 + 1/x > 1$  for  $x = 1, \dots, 9$  we have  $f(x) = \log_{10}(1 + 1/x) > 0$  and

$$\begin{aligned}
f(1) + \dots + f(9) &= \log_{10}\left(1 + \frac{1}{1}\right) + \log_{10}\left(1 + \frac{1}{2}\right) + \dots + \log_{10}\left(1 + \frac{1}{8}\right) + \log_{10}\left(1 + \frac{1}{9}\right) \\
&= \log_{10}\left(\frac{2}{1}\right) + \log_{10}\left(\frac{3}{2}\right) + \dots + \log_{10}\left(\frac{9}{8}\right) + \log_{10}\left(\frac{10}{9}\right) \\
&= \log_{10}\left(\frac{2}{1} \cdot \frac{3}{2} \cdot \frac{4}{3} \cdot \dots \cdot \frac{8}{7} \cdot \frac{9}{8} \cdot \frac{10}{9}\right) = \log_{10}(10) = 1
\end{aligned}$$

i.e.,  $f(x)$  is a pmf.

(b)

```

Benford <- function(){
obs <- c(107,55,39,22,13,18,13,23,15)
k <- length(obs)
n <- sum(obs)
p <- log10(1+1/(1:9))
e <- n*p

```

```

X2 <- sum((obs-e)^2/e)
G2 <- 2*sum(obs*log(obs/e))
pX2 <- 1-pchisq(X2,k-1)
pG2 <- 1-pchisq(G2,k-1)
out <- c(X2,pX2,G2,pG2)
names(out) <- c("X2", "p-value(X2)", "G2", "p-value(G2)")
out
}
# calling Benford() yields
> Benford()
      X2 p-value(X2)      G2 p-value(G2)
14.75964770 0.06399094 15.55588649 0.04919622

```

The  $p$ -values are borderline significant at  $\alpha = 0.05$ .

**4. Section 14.6, problem 7.** Rather than just calling `cor.test(sister,brother,conf.level=0.9)` for appropriate data vectors `sister` and `brother`, write yourself a function (using the steps on slides 27-28 in Ch. 14)

```

Conf.Int <- function(x,y,conf.level){
  ....
}

```

that computes  $\hat{\rho}$ , then computes  $\hat{\zeta}$ , then computes the interval end points `zeta.L` and `zeta.U` for  $\zeta$ , and from that computes the interval endpoints `rho.L` and `rho.U` for  $\rho$  and returns as output `c(rho.L, rho.U)`. Compare the results from using

```
Conf.Int(sister,brother,conf.level = 0.9)
```

with that when using

```
rho.test(sister,brother,conf.level = 0.9)
```

You may use the fact that `cor(x,y)` returns the sample correlation coefficient for the data vector `x` and `y`.

```

Conf.Int <- function(x,y,conf.level){
rho.hat <- cor(x,y)
n <- length(x)
zeta.hat <- 0.5*log((1+rho.hat)/(1-rho.hat))
alpha <- 1-conf.level
qz <- qnorm(1-alpha/2)
zeta.L <- zeta.hat -qz/sqrt(n-3)
zeta.U <- zeta.hat +qz/sqrt(n-3)
rho.L <- (exp(2*zeta.L)-1)/(exp(2*zeta.L)+1)
rho.U <- (exp(2*zeta.U)-1)/(exp(2*zeta.U)+1)
c(rho.L,rho.U)
}
# with the following results when calling Conf.Int and cor.test
> sister <- c(69,64,65,63,65,62,65,64,66,59,62)
> brother <- c(71,68,66,67,70,71,70,73,72,65,66)
> Conf.Int(sister,brother,.9)
[1] 0.04842215 0.83714300

```

```
> cor.test(sister,brother,conf.level=.9)
```

Pearson's product-moment correlation

data: sister and brother

t = 2.0175, df = 9, p-value = 0.07442

alternative hypothesis: true correlation is not equal to 0

90 percent confidence interval:

0.04842215 0.83714300

sample estimates:

cor

0.5580547

The confidence intervals coincide, i.e., it shows that `cor.test` uses the approximate procedure given in slides 27-28 for Ch. 14.