# 305B PS1

## Isaac Kleisle-Murphy

## 1/17/2022

### 1.14

Recall that the Jeffrey's prior for a binomial distribution follows ~Beta(.5, .5). Hence for Y = 0 ~ Binom(N=25, p), the beta-binomial posterior yields p|Y ~ Beta(0.5, 25.5). The mean of this density is 0.01923077. The 95% equitail interval.

```
c(
  qbeta(.025, .5, 25.5),
  qbeta(.975, .5, 25.5)
)
```

```
## [1] 0.00001944577 0.09468276411
```

Lastly, looking at the posterior density, we see

$$\pi(p|Y) \propto p^{-.5}(1-p)^{24.5}$$

As this density is strictly decreasing on $(0, 1)$, we know that the 95% HPD interval must be pushed all the way to the left, lest a point with higher density (only occuring on the left, due to the decreasingness as $0 \to 1$) be omitted on the LHS. Hence, the 0th and 95th quantiles of the beta are the HPD, i.e.

```
c(
  qbeta(0, .5, 25.5),
  qbeta(.95, .5, 25.5)
)
```

```
## [1] 0.00000000 0.07323939
```

### 2.4

**a.)**

We have, taking sums/ratios over the provided table,

$$P(Fatal|Seatbelt) \approx \frac{703}{703 + 441239} \approx .001591$$

and

$$P(Fatal|NoSeatbelt) \approx \frac{1085}{1085 + 55623} \approx .019133.$$

**b.)**

Similarly,

$$P(Seatbelt|Fatal) \approx \frac{703}{703 + 1085} \approx .39318$$

and

$$P(Seatbelt|Non\ Fatal) \approx \frac{441239}{441239 + 55623} \approx .88805.$$

**c.)**

As seatbelts are intended to prevent them, fatality is the most natural choice of response. Accordingly:

- Difference: $.019133 - .001591 = .017542$ means that the observed proportion of fatalities without a belt was .0175 greater than that with a seatbelt.

- Relative Risk: $.019133/.001591 = 12.02577$ means that the observed proportion of fatalities without a seatbelt was roughly twelve times the observed proportion of fatalities with a seatbelt.

- Odds Ratio: $(441239 * 1085)/(703 * 55623) = 12.24317$ means that the observed odds of a fatality for those without a seatbelt was roughly 12 times the observed odds for those with.

We note that the relative risk and odds ratio are similar here; as outlined in 2.2.7, this stems from the fact that when a proportion is either very close to zero, the $(1 - \pi_2)/(1 - \pi_1)$ is very close to 1 and the odds ratio effectively becomes

$$OR = RR \underbrace{(1 - \pi_2)/(1 - \pi_1)}_{\approx 1} \approx RR \cdot 1.$$

Hence, the OR and RR here are close, around 12.

## 2.13

```
ha_data = matrix(
  c(
    193, 19942 - 193,
    198, 19934 - 198
  ),
  ncol=2,
  byrow=T
)

colnames(ha_data) = c("heart_attack", "no_heart_attack")
rownames(ha_data) = c("placebo", "aspirin")

ha_data
```

```
##          heart_attack no_heart_attack
## placebo           193           19749
## aspirin           198           19736
```

The odds ratio is thus $(19736 \cdot 193)/(198 \cdot 19749) = 0.9741058$, indicating that the odds of a heart attack with aspirin are close, or thereabouts, to the odds of a heart attack without aspirin. This might cast doubt on the hypothesis that aspirin helps to reduce heart attack danger.

## 2.18

This is a classic case of Simpson's paradox. For a hands-on example, consider the following (admittedly outrageous) example: suppose the two age levels are young and old, and we are comparing ME and SC as in the problem.

```r
freqs = matrix(
  c(
    .3, .5,
    .4, .5
  ),
  byrow=T,
  ncol=2
)
colnames(freqs)=c("Maine", "SC")
rownames(freqs)=c("Young", "Old")


counts = matrix(
  c(
    1000, 100,
    4000, 100
  ),
  byrow=T,
  ncol=2
)

colnames(counts)=c("Maine", "SC")
rownames(counts)=c("Young", "Old")

totals = colSums(counts * freqs)
```

As we see, SC has the greater death frequency, as both the RHS columns are greater than the LHS columns.

```r
freqs
```

```
##        Maine  SC
## Young   0.3 0.5
## Old     0.4 0.5
```

However, when we incorporate the overall counts

```r
counts
```

```
##        Maine  SC
## Young   1000 100
## Old     4000 100
```

and then compute overall frequencies, we see that the state-wise death totals are

```
totals
```

```
## Maine    SC
##  1900   100
```

and thus overall frequencies are

```
totals / sum(totals)
```

```
## Maine    SC
##  0.95  0.05
```

with Maine now handily leading the death rate. As stated above, this is classic Simpson's paradox: with subgroups faceted, one group/state appears to be ahead. However, when much of the population is centered in one of the lesser frequencies (i.e. 4,000 in Maine/Old), that cell will dominate the "re-weighting" that is the state vs. state frequency, and hence Maine prevails.

## 3.6

The data, in table format, is

```
data = matrix(
  c(
    7, 8,
    0, 15
  ),
  ncol=2,
  byrow=T
)
colnames(data) = c("norm_achieved", "norm_not_achieved")
rownames(data) = c("presnisolone", "control")

data
```

```
##               norm_achieved norm_not_achieved
## presnisolone             7                 8
## control                  0                15
```

Recall the Wald CI, which takes the form

$$\log \hat{\theta} \pm z_{\alpha/2}\sqrt{1/n_{11} + 1/n_{22} + 1/n_{12} + 1/n_{21}}.$$

Unfortunately, since one of the cell counts (norm achieved + control) is zero, this interval cannot be computed (one could interpret an infinity for i.e. $1/n_{2,1} = \infty \implies \sqrt{1/n_{11} + 1/n_{22} + 1/n_{12} + 1/n_{21}} = \infty$.).

We then set up the profile likelihood as follows. Whereas we'll often describe binomial/multinomial 2x2 table model in terms of proportion parameters $\vec{\pi}$ and trials $n$, we can identically express this table/model by a fixed odds ratio $\theta_0$ and marginal proportions $\pi_{i+}$ and $\pi_{+j}$ such that these marginal proportions produce (under assumption of independence, i.e. take products to get cell counts) an odds ratio equal to $\theta_0$. Accordingly, since you can multiply $n$ through an odds ratio (i.e. proportions $\to$ counts), it follows that under this setup

$$\frac{\hat{\mu}_{11}(\theta_0)\hat{\mu}_{22}(\theta_0)}{\hat{\mu}_{12}(\theta_0)\hat{\mu}_{21}(\theta_0)} = \theta_0,$$

in other words, the expected counts (under independence) return an odds ratio of $\theta_0$. In short, we have row and columnwise probabilities $\pi_{+j}, \pi_{i+}$ (presumed independent) such that the expected counts under these row/columnwise probabilities (with independence) give an odds ratio of $\theta_0$.

Now, we proceed to the Profile Likelihood CI by inverting a LRT. Specifically, we choose a null hypothesis/value for $\theta_0$, and compute

$$L\hat{\mu}(\theta_0)$$

(note my notation is a touch different from the book; the key is about evaluating likelihoods at maximized or nearly-maximized values), which is the maximum of the log-likelihood subject to the constraint that the marginals satisfy the

$$\hat{\mu}(\theta_0)$$

property described above; that is, it's the best the constrained marginals could possibly do under the data. We then compute the unrestricted MLE

$$\hat{\theta}$$

, and compute it's likelihood under the data directly as

$$L(\hat{\theta})$$

. Under the standard LRT, the CI is obtained by pivoting

$$-2(L\hat{\mu}(\theta_0) - L(\hat{\theta})) < \chi_1^2,$$

i.e. finding all values $\theta_0$ not rejected under the standard hypothesis test. Note $df = 1$ as this is a 2x2 table, hence $3 - 1 - 1 = 1$.

## 3.7

```
data = matrix(
  c(
    8, 33,
    37, 152
  ),
  ncol=2,
  byrow=T
)
rownames(data) = c("make_1", "miss_1")
colnames(data) = c("make_2", "miss_2")

data
```

```
##        make_2 miss_2
## make_1      8     33
## miss_1     37    152
```

A chi-squared test here returns a p-value of 1:

```
chisq.test(data)
```

```
##
##  Pearson's Chi-squared test with Yates' continuity correction
##
## data:  data
## X-squared = 0.0000000000000000000000000000039336, df = 1, p-value = 1
```

5

in which case we fail to reject the hypothesis that free throw success is independent. This perhaps lends credibility to those who argue the "hot-hand" does not exist, and that successive free throws are in fact independent.

## 3.16

Inits:

```
data = matrix(
  c(
    c(9, 44, 13, 10),
    c(11, 52, 23, 22),
    c(9, 41, 12, 27)
  ),
  byrow=T,
  ncol=4
)
colnames(data) = c("SH", "HS", "SC", "C")
rownames(data) = c("L", "M", "H")
```

**a.)**

A standard Chi-squared test for this data returns the following:

```
chisq.test(data)
```

```
##
##  Pearson's Chi-squared test
##
## data:  data
## X-squared = 8.8709, df = 6, p-value = 0.181
```

i.e. a failure to reject a null hypothesis independence at any reasonable alpha. However, a key potential deficiency is that the test does not account for the ordinal nature of the data, i.e. SH < HS < SC < C or L < M < H. They're just standard (nominal) buckets, and the directed nature of the relationship is not baked into the test. In this way, key ordering information may be overlooked by the regular Chi-squared test.

**b.)**

```
N = sum(data)
# marginal MLEs
pi_row = rowSums(data) / sum(data)
pi_col = colSums(data) / sum(data)
# model expectations
u_hat = pi_row %*% t(pi_col) * N
r_scale = sqrt(
  ((1 - pi_row)%*% t(1 - pi_col)) * u_hat
)
# standardized residuals
(data - u_hat) / r_scale
```

```
##             SH          HS          SC           C
## L  0.4061328  1.5828205 -0.1286367 -2.1078423
## M -0.1898118 -0.5440627  1.3041565 -0.4031584
## H -0.1903291 -0.9459053 -1.2374420  2.4360173
```

As education increases (i.e. columns moving L-> R), we see increasingly large standardized residuals – eventually exceeding the problematic |residual| > 2 or 3 admonished by the book. Per Agresti, these large standardized residuals indicate a lack of fit – specifically, with educational aspiration residuals increasing alongside income – suggesting that the ordinary Chi-Squared is increasingly inadequate here.

**c.)**

In light of the apparent ordinality borking the Chi-Squared test, a Kruskal test may be preferable here.

```
MESS::gkgamma(data)
```

```
##
##  Goodman-Kruskal's gamma for ordinal categorical data
##
## data:  data
## Z = 2.0268, p-value = 0.04268
## 95 percent confidence interval:
##  0.006716385 0.318378906
## sample estimates:
## Goodman-Kruskal's gamma
##               0.1625476
```

With a p-value of .04268, this provides stronger evidence of an (ordinal) association (i.e. non-independence) between family income and educational aspiration.

### 3.23

As prescribed in Agresti 3.6.2 (p. 97), if we assume row-wise binomial independence, we may put beta priors on the row-wise binomials to obtain a row-wise posterior beta. As further prescribed in 3.6.2, we may further simulate the logs ratio by leveraging the presumed independence, and taking S draws from the row 1 posterior, followed by S draws from the row 2 posterior (independent of row 1), and then go element-by-element through these two sets of draws to compute Monte Carlo odds ratios, for our inferential purposes.

```
set.seed(2022)

### number of simulated draws ###
S = 100000
### interval width ###
Q = .95

### specify priors: uniform(1, 1) ###
ALPHA = c(1, 1)

### make data ###
# first row of data
```
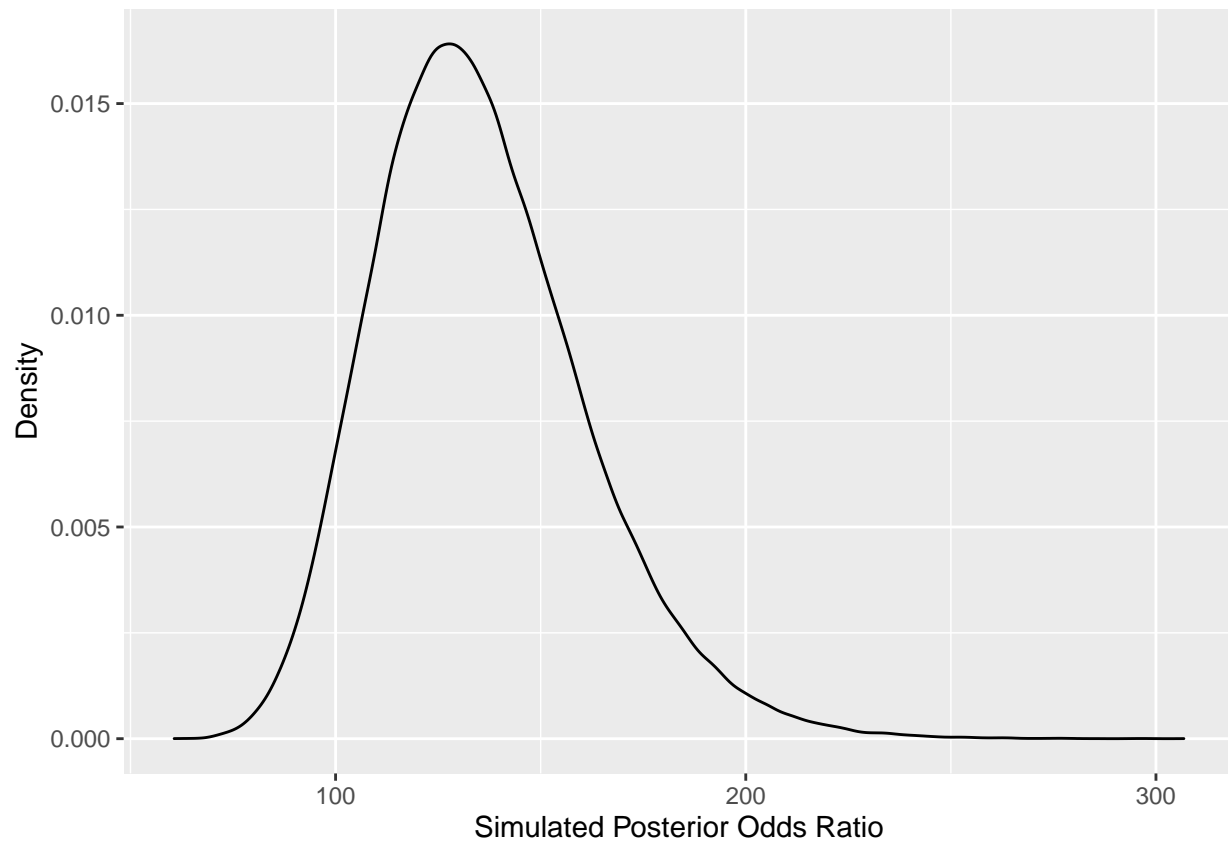
```
y1 = c(763, 65)
# second row of data
y2 = c(59, 680)

### row-wise posteriors ###
y1_post = rbeta(S, y1[1] + ALPHA[1], y1[2] + ALPHA[2])
y2_post = rbeta(S, y2[1] + ALPHA[1], y2[2] + ALPHA[2])

### use simulations for odds ratio, and sort for ECDF ease ###
# strictly due to Bush being on the LHS of the table, define that as the 1 outcome.
odds_ratio_sim = sort(
  (y1_post * (1 - y2_post)) /
  ((1 - y1_post) * y2_post)
)

### plot density ###
ggplot(
  data.frame(odds_ratio=odds_ratio_sim),
  aes(x=odds_ratio)
) +
  geom_density() +
  labs(y="Density", x="Simulated Posterior Odds Ratio")
```



```
### compute CI size and tail size ###
tail_size = as.integer((S * (1 - Q)) / 2) # chose a convenient S
```

```
ci_size = S - 2 * tail_size


### compute equital CI ###
equitail_ci = c(
  odds_ratio_sim[tail_size + 1],
  odds_ratio_sim[S - (tail_size + 1)]
)

hpd_ci = hdi(odds_ratio_sim)
```

The 95% equitail CI (again, under the labeling set forth above) is,

```
equitail_ci
```

```
## [1]   93.01743 193.58020
```

while the HPD CI is

```
hdi(odds_ratio_sim)
```

```
##     lower     upper
##  89.12104 187.50186
## attr(,"credMass")
## [1] 0.95
```

Importantly, the HPD interval here on an odds ratio is dangerous due to its non-invariance under nonlinear parameter transformation (Agresti 3.6.5. For example, suppose we wanted to invert our odds ratio – i.e. in the case we relabeled the data or redefined "success" – we could not just invert the HPD interval, and instead we would have to start from scratch. The example in the second paragraph of 3.6.5 is a perfect cautionary tale of what inversion could do in this problem.

To be thorough, if we had used the alternate labeling (i.e. Kerry is the target), our equitail CI would be

```
c(
  1 / odds_ratio_sim[S - (tail_size + 1)],
  1 / odds_ratio_sim[tail_size + 1]
)
```

```
## [1] 0.005165817 0.010750673
```

while the HPD would be

```
 hdi(1 / odds_ratio_sim)
```

```
##        lower        upper
## 0.004986143 0.010496797
## attr(,"credMass")
## [1] 0.95
```