

Contingency Tables (Lecture 2)

Heather Weaver & Mark Fulginiti

September 14, 2017

Bayesian Multinomial Example

Adapted from BDA3 by Gelman et al.

Election Polling Results ($n=1,447$)

Obama: $y_1 = 727$

Romney: $y_2 = 583$

Other / NA: $y_3 = 137$

Note: Assuming simple random sampling (i.e. ignoring any biases)

$$(y_1, y_2, y_3 \mid n) \sim \text{Multinomial}(\theta_1, \theta_2, \theta_3)$$

Estimand:

$$\theta_1 - \theta_2$$
$$H_o : \theta_1 = \theta_2$$

Prior: one choice

$$(y_1, y_2, y_3) \sim \text{Dirichlet}(\alpha_1, \alpha_2, \alpha_3)$$
$$\alpha_1 = \alpha_2 = \alpha_3 = 1$$

Posterior:

$$P(\theta_1, \theta_2, \theta_3 \mid \vec{y}) \sim \text{Dirichlet}(y_1 + 1, y_2 + 1, y_3 + 1)$$

$$P(\theta_1, \theta_2, \theta_3 \mid \vec{y}) \sim \text{Dirichlet}(728, 584, 138)$$

- If a closed form solution exists, we could solve for quantiles, CI bounds, etc.

Computationally:

1. Set $i = 1$
2. Draw $\vec{\theta}^{[i]} \sim \text{Dirichlet}()$
3. Calculate $\vec{\theta}_1^{[i]} - \vec{\theta}_2^{[i]}$
4. set $i = i + 1$
5. If $i = 1,000$ skip
6. Compute summary / quantiles of $(\theta_1 - \theta_2)$
 - 80% Credible Interval for $(\theta_1 - \theta_2)$ [10th percentile, 90th percentile]

Bayesian example for small-sample cells

- Exact tests are a frequentist solution for low cell counts

Following example from a Wikipedia article on contingency tables, lingpipe blog (Bob Carpenter)

Gender	Left Handed	Right Handed	Total
Male	43	$y_1 = 9$	$m_1 = 52$
Female	44	$y_2 = 4$	$m_2 = 48$

- Question: Are men and women equally likely to be left-handed?

$$H_o : P(LH \mid M) = P(LH \mid F)$$

$$H_o : P(\pi_1) = P(\pi_2)$$

$$H_o : \pi_1 - \pi_2 = 0$$

We could also look at the risk ratio $\frac{\pi_1}{\pi_2}$ or the odds ratio $\frac{\pi_1/(1-\pi_1)}{\pi_2/(1-\pi_2)}$.

$Y_1 \sim \text{Binomial}(n_1, \pi_1)$ the number of men who are left-handed

$Y_2 \sim \text{Binomial}(n_2, \pi_2)$ the number of women who are left-handed

We are assuming Independence of the above variables

Prior: Assume uniform priors

$$\pi_1 \sim \text{Uniform}(0, 1) = \text{Beta}(1, 1)$$

$$\pi_2 \sim \text{Uniform}(0, 1) = \text{Beta}(1, 1)$$

Posterior for π :

$$\pi_i \mid y_i, n_i \sim \text{Beta}(y_i + 1, n - y_i + 1)$$

$$\begin{pmatrix} \pi_1 \binom{1}{1} \\ \pi_1 \binom{2}{2} \\ \vdots \\ \pi_1 \binom{k}{k} \end{pmatrix} - \begin{pmatrix} \pi_2 \binom{1}{1} \\ \pi_2 \binom{2}{2} \\ \vdots \\ \pi_2 \binom{k}{k} \end{pmatrix} \Rightarrow \begin{pmatrix} (\pi_1 - \pi_2) \binom{1}{1} \\ (\pi_1 - \pi_2) \binom{2}{2} \\ \vdots \\ (\pi_1 - \pi_2) \binom{k}{k} \end{pmatrix}$$

- Advantage: No assumptions about large sample size needed.
- Disadvantage: Results are sensitive to choice to prior.

Confidence intervals from likelihood ratios

Define a $(1 - \alpha)\%$ credible interval/region. The set of Θ for which the Likelihood Ratio Test Statistic (LRST)

$$LRTS(\Theta) < \chi_{df}^2(1 - \alpha)$$

$$LRTS(\Theta) = -2[L(\Theta) - L(\hat{\Theta})]$$

where $L(\Theta)$ is the log-likelihood from Θ , across a subset of components of Θ , and $L(\hat{\Theta})$ is the fixed log-likelihood for Θ calculated at the MLE

- $\chi_{df}^2(1 - \alpha)$ is the $(1 - \alpha)^{th}$ quantile of χ_{df}^2
- The df is the number of parameters varying in Θ

e.g. 1-Parameter CI

$$\Theta = \mu$$

$$LRTS = -2(L(\mu) - L(\hat{\mu}))$$

- 95% CI for $\mu = \{ \mu : LRTS < \chi_1^2(0.95) \}$

Example

$$y_i \sim Poisson(\mu), \quad i = 1, 2, 3$$

$$P(y \mid \mu) = \frac{(e^{-\mu} \mu^y)}{y!}$$

$$\begin{aligned} L(\mu \mid y) &= \log \left(\prod_{i=1}^3 \frac{e^{-\mu} \mu^{y_i}}{y_i!} \right) = \sum_{i=1}^3 \log \left(\frac{e^{-\mu} \mu^{y_i}}{y_i!} \right) = \\ &= \sum (-\mu + y_i \log(\mu) - \log(y_i!)) \propto -n\mu + \log(\mu) * \sum y_i + C \\ L(\mu \mid \vec{y}) &\propto -3\mu + 7\log(\mu) \end{aligned}$$

$$\begin{aligned} LRTS &= -2(L(\mu) - L(\hat{\mu})) < \chi_1^2(0.95) \\ &\rightarrow (L(\mu) > L(\hat{\mu})) - \frac{\chi_1^2(0.95)}{2} \\ \frac{\chi_1^2(0.95)}{2} &= \frac{3.84}{2} \end{aligned}$$

Example with Multiple parameters

$$\begin{aligned} \Theta &= (\theta_1, \theta_2, \dots, \theta_r, \theta_{r+1}, \dots, \theta_p) \\ \text{of particular interests are} &(\theta_{r+1}, \dots, \theta_p) \end{aligned}$$

$$LRST = -2 \left\{ L(\theta_{r+1}, \dots, \theta_p \mid \hat{\theta}_1, \dots, \hat{\theta}_r) - L(\hat{\theta}_{r+1}, \dots, \hat{\theta}_p \mid \hat{\theta}_1, \dots, \hat{\theta}_r) \right\}$$

Calculate the likelihood on a grid of $(\theta_{r+1}, \dots, \theta_p) = \theta^*$

There are p - r parameters, so the df is p - r

$$95 \% \text{ CI for } \theta^* = \{ \theta^* : LRTS < \chi_{p-r}^2(0.95) \}$$

- Note: This whole assumption is based on large sample approximation

Diagnostic Tests

Disease Presence	Test Positive	Test Negative	Total
Yes	π_{11}	π_{12}	1
No	π_{21}	π_{22}	1

- Sensitivity: $\Pr(+ | D) = \pi_{1|1} =$ probability the test is positive given that you have the disease
- Specificity: $\Pr(- | D^C) = \pi_{2|2} =$ probability that you test negative given that you do not have the disease

Breast Cancer Example

Cancer Presence	Test Positive	Test Negative	Total
Yes	86	14	$100 = n_1$
No	12	88	$100 = n_2$

Gold standard would be to have $\alpha = 0$, $\beta = 0$.

- Sensitivity: $\Pr(+ | D) = \pi_{1|1} = 86\%$
- Specificity: $\Pr(- | D^C) = \pi_{2|2} = 88\%$

From a clinical perspective we are much more interested in positive/negative predicted value.

$P(D | +)$ sometimes = $\frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$

e.g. $P(\text{Breast Cancer} | +) = \frac{86}{86+12}$

Note: This is only true if $(\frac{n_1}{n_1+n_2}) \approx P(D) =$ prevalence of the disease

$$\begin{aligned}
 P(D | +) &= \frac{P(D \cap +)}{P(+)} = \frac{P(+ | D) * P(D)}{P(+ \cap D)} * P(+ \cap D^C) \\
 &= \frac{\text{Sensitivity} * \text{Population Prevalance}}{P(+ | D)P(D) + P(+ | D^C)P(D^C)} \\
 &= \frac{\text{Sensitivity} * \text{Population Prevalance}}{(\text{Sensitivity} * \text{Population Prevalance}) + ((1 - \text{Specicifity}) * (1 - \text{Population Prevalance}))}
 \end{aligned}$$