

# CATEGORICAL DATA ANALYSIS

## ASSIGNMENT III

### Solution Set

3.2

Since the study at hand was a case-control study, it is not possible to estimate the difference in proportions or the relative risk directly from the data.

An appropriate measure of association is the odds ratio.

$$\hat{\theta} = \frac{(688)(59)}{(650)(21)} = 2.9738$$

$$\log \hat{\theta} = 1.0898$$

$$\hat{\sigma}(\log \hat{\theta}) = \left( \frac{1}{688} + \frac{1}{650} + \frac{1}{21} + \frac{1}{59} \right)^{1/2} = 0.2599$$

95% CI for  $\log \theta$ :

$$\begin{aligned} \log \hat{\theta} \pm z_{\alpha/2} (\hat{\sigma}(\log \hat{\theta})) \\ (1.0898) \pm (1.96)(0.2599) \\ (0.5804, 1.5992) \end{aligned}$$

95% CI for  $\theta$ :

$$\begin{aligned} (\exp(0.5804), \exp(1.5992)) \\ (1.7867, 4.9493) \end{aligned}$$

The odds ratio estimate indicates that the odds of developing lung cancer are 2.9738 times higher for a smoker than for a non-smoker.

The confidence interval for the odds ratio suggests that the true odds ratio is at least 1.7867 and at most 4.9493.

3.3

First Shot	Second Shot	
	Made	Missed
Made	251	34
Missed	48	5

We will conduct the Pearson chi-squared test for independence.

$$\begin{aligned} \chi^2 &= \sum_i \sum_j \frac{(m_{ij} - \hat{\mu}_{ij})^2}{\hat{\mu}_{ij}} \\ &= 0.0049 + 0.0378 + 0.0265 + 0.2034 \\ &= 0.2727 \end{aligned}$$

$P = 0.6015$ . There is insufficient evidence to indicate a dependence between the outcomes of successive free throws.

3.7

SAS

	MI		
	Fatal	Non-Fatal	No Attack
Placebo	18	171	10,845
Aspirin	5	99	10,933

$$G^2 = 27.5893 \quad (2 \text{ df}) ; \quad P < .0001$$

	MI	
	Fatal	Non-Fatal
Placebo	18	171
Aspirin	5	99

$$G_1^2 = 2.2173 \quad (1 \text{ df}) ; \quad P = .1365$$

	MI	
	Attack	No Attack
Placebo	189	10,845
Aspirin	104	10,933

$$G_2^2 = 25.3720 \quad (1 \text{ df}) ; \quad P < .0001$$

Note:  $G^2 = G_1^2 + G_2^2$ .

The result of the initial test (on the 2x3 table) indicates that myocardial infarction status is dependent on aspirin intake.

The results of the partitioning support the following conclusions:

- (1) for those experiencing a heart attack, whether the attack is fatal or non-fatal is NOT dependent on aspirin intake,
- (2) whether one experiences a heart attack (either fatal or non-fatal) IS dependent on aspirin intake.

3.11

SAS

$$(a) \chi^2 = 8.8709 \quad (6 \text{ df}); \quad P = 0.1810$$
$$G^2 = 8.9165 \quad (6 \text{ df}); \quad P = 0.1783$$

The Pearson chi-squared and likelihood-ratio chi-squared tests do not provide compelling evidence of a dependence between family income and educational aspirations.

These tests are not optimal for the data at hand because the tests are designed for nominal variables, and yet the variables in the present application are ordinal.

$$(c) M^2 = 0.6316 \quad (1 \text{ df}); \quad P = 0.4268.$$

The CMH linear trend test is more appropriate for the data at hand because it is designed for ordinal variables.

However, in the present application, the CMH test is even less powerful than the Pearson chi-squared or likelihood-ratio chi-squared test, yielding a larger p-value. Thus, the CMH test also fails to provide evidence of a dependence between family income and educational aspirations.

3.12

SAS

$$\hat{\sigma} = 0.3873$$

$$\hat{\alpha}(\hat{\sigma}) = 0.0366 \quad (\text{ASE in SAS})$$

The Wald CI for  $\sigma$  is of the form  $\hat{\sigma} \pm Z_{\alpha/2}(\hat{\alpha}(\hat{\sigma}))$ .

95% CI for  $\sigma$ :

$$\hat{\sigma} \pm Z_{\alpha/2}(\hat{\alpha}(\hat{\sigma}))$$

$$(0.3873) \pm (1.96)(0.0366)$$

$$(0.3156, 0.4590)$$

Based on the point estimate of  $\sigma$  and the CI for  $\sigma$ , it appears that there is a positive association of moderate strength between schooling and attitude on abortion.

3.22

From the class notes from § 1.4, we have the following result for the large-sample distribution of  $\hat{\pi}$ :

$$(\hat{\pi} - \pi) \sim N\left(0, \frac{\pi(1-\pi)}{n}\right).$$

(Here, the variance of  $\hat{\pi}$  is given by the reciprocal of  $I(\pi) = n/(\pi(1-\pi))$ .)

For application of the delta method, we have

$$g(\pi) = \log(\pi/(1-\pi)),$$

$$\begin{aligned} d(\pi) &= \frac{\partial}{\partial \pi} \left[ \log\left(\frac{\pi}{1-\pi}\right) \right] \\ &= \frac{(1-\pi)}{\pi} \left[ \frac{(1-\pi)(1) - (\pi)(-1)}{(1-\pi)^2} \right] \\ &= \frac{(1-\pi)}{\pi} \left[ \frac{1}{(1-\pi)^2} \right] \\ &= \frac{1}{\pi(1-\pi)}, \end{aligned}$$

$$\Sigma(\pi) = \frac{\pi(1-\pi)}{n} \quad (\text{from above}).$$

Thus,

$$\begin{aligned} (d(\pi))^2 \Sigma(\pi) &= \left( \frac{1}{\pi(1-\pi)} \right)^2 \left( \frac{\pi(1-\pi)}{n} \right) \\ &= \frac{1}{n\pi(1-\pi)}. \end{aligned}$$

Therefore, we have the following result for the large-sample distribution of  $\log(\hat{\pi}/(1-\hat{\pi}))$ :

$$\left(\log(\hat{\pi}/(1-\hat{\pi})) - \log(\pi/(1-\pi))\right) \sim N\left(0, \frac{1}{n\pi(1-\pi)}\right).$$

The preceding leads to the Wald CI for  $\log(\pi/(1-\pi))$  provided in the problem statement.

To use the Wald CI to obtain a CI for  $\pi$ , we employ the inverse transformation corresponding to  $g(\pi)$ : i.e.,  $g^{-1}(x) = e^x / (1+e^x)$ . If  $(a, b)$  denotes the Wald CI for  $\log(\pi/(1-\pi))$ , the associated CI for  $\pi$  is given by  $(e^a / (1+e^a), e^b / (1+e^b))$ .

3.34

$$\begin{aligned} (a) \quad \chi^2 &= \sum_i \frac{(m_i - \hat{\mu}_i)^2}{\hat{\mu}_i} \\ &= \sum_i \left( \frac{m_i^2 - 2m_i\hat{\mu}_i + \hat{\mu}_i^2}{\hat{\mu}_i} \right) \\ &= \sum_i \left( \frac{m_i^2}{\hat{\mu}_i} - 2m_i + \hat{\mu}_i \right) \\ &= \sum_i \left( \frac{m_i^2}{\hat{\mu}_i} \right) + \sum_i (\hat{\mu}_i - 2m_i) \\ &= \sum_i \left( \frac{m_i^2}{\hat{\mu}_i} \right) + (-n) \\ &= \sum_i \left( \frac{m_i^2}{\hat{\mu}_i} \right) + (-\sum_i m_i) \\ &= \sum_i m_i \left( \frac{m_i}{\hat{\mu}_i} - 1 \right) \end{aligned}$$

$$= \frac{2}{(1)(1+1)} \sum_i m_i \left( \left( \frac{m_i}{\hat{\mu}_i} \right)^{(1)} - 1 \right)$$

$$\begin{aligned} (b) \quad & \lim_{\lambda \rightarrow 0} \frac{2}{\lambda(\lambda+1)} \sum_i m_i \left( \left( \frac{m_i}{\hat{\mu}_i} \right)^\lambda - 1 \right) \\ &= \left( \lim_{\lambda \rightarrow 0} \frac{2}{(\lambda+1)} \right) \left\{ \sum_i m_i \left[ \lim_{\lambda \rightarrow 0} \left( \frac{\left( \frac{m_i}{\hat{\mu}_i} \right)^\lambda - 1}{\lambda} \right) \right] \right\} \\ &= (2) \left\{ \sum_i m_i \left[ \log \left( \frac{m_i}{\hat{\mu}_i} \right) \right] \right\} \\ &= G^2 \end{aligned}$$

$$\begin{aligned} (c) \quad & \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \sum_i m_i \left( \left( \frac{m_i}{\hat{\mu}_i} \right)^\lambda - 1 \right) \\ &= \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \sum_i m_i \left( \left( \frac{\hat{\mu}_i}{m_i} \right) \left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} - 1 \right) \\ &= \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \sum_i \left( \hat{\mu}_i \left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} - m_i \right) \\ &= \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \left[ \sum_i \left( \hat{\mu}_i \left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} \right) - \sum_i m_i \right] \\ &= \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \left[ \sum_i \left( \hat{\mu}_i \left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} \right) - \sum_i \hat{\mu}_i \right] \\ &= \lim_{\lambda \rightarrow -1} \frac{2}{\lambda(\lambda+1)} \left[ \sum_i \left( \hat{\mu}_i \left( \left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} - 1 \right) \right) \right] \\ &= \left( \lim_{\lambda \rightarrow -1} \frac{2}{\lambda} \right) \left\{ \sum_i \hat{\mu}_i \left[ \lim_{\lambda \rightarrow -1} \left( \frac{\left( \frac{m_i}{\hat{\mu}_i} \right)^{\lambda+1} - 1}{\lambda+1} \right) \right] \right\} \\ &= (-2) \left\{ \sum_i \hat{\mu}_i \left[ \log \left( \frac{m_i}{\hat{\mu}_i} \right) \right] \right\} \end{aligned}$$



$$= 2 \sum_i \hat{\mu}_i \log \left( \frac{\hat{\mu}_i}{m_i} \right)$$

$$(d) \sum_i \frac{(m_i - \hat{\mu}_i)^2}{m_i}$$

$$= \sum_i \left( \frac{m_i^2 - 2m_i\hat{\mu}_i + \hat{\mu}_i^2}{m_i} \right)$$

$$= \sum_i \left( \frac{\hat{\mu}_i^2}{m_i} - 2\hat{\mu}_i + m_i \right)$$

$$= \sum_i \left( \frac{\hat{\mu}_i^2}{m_i} \right) + \sum_i (m_i - 2\hat{\mu}_i)$$

$$= \sum_i \left( \frac{\hat{\mu}_i^2}{m_i} \right) + (-m)$$

$$= \sum_i \left( \frac{\hat{\mu}_i^2}{m_i} \right) + (-\sum_i m_i)$$

$$= \sum_i m_i \left( \left( \frac{\hat{\mu}_i}{m_i} \right)^2 - 1 \right)$$

$$= \frac{2}{(-2)(-2)+1} \sum_i m_i \left( \left( \frac{m_i}{\hat{\mu}_i} \right)^{-2} - 1 \right)$$

$$(e) 4 \sum_i (\sqrt{m_i} - \sqrt{\hat{\mu}_i})^2$$

$$= 4 \sum_i (m_i + \hat{\mu}_i - 2\sqrt{m_i\hat{\mu}_i})$$

$$= 4 \left[ (-2 \sum_i \sqrt{m_i\hat{\mu}_i}) + (\sum_i (m_i + \hat{\mu}_i)) \right]$$

$$= 4 \left[ (-2 \sum_i \sqrt{m_i\hat{\mu}_i}) + (2m) \right]$$

$$= 4 \left[ (-2 \sum_i \sqrt{m_i\hat{\mu}_i}) + (2 \sum_i m_i) \right]$$

$$= -8 \left( \sum_i (\sqrt{m_i \hat{\mu}_i} - m_i) \right)$$

$$= \frac{2}{\left(-\frac{1}{2}\right)\left(-\frac{1}{2}+1\right)} \sum_i m_i \left( \left(\frac{m_i}{\hat{\mu}_i}\right)^{\left(-\frac{1}{2}\right)} - 1 \right)$$

3.39

$$\hat{u} = - \frac{\sum_i \sum_j p_{ij} \log(p_{ij} / (p_i + p_j))}{\sum_j p_{+j} \log p_{+j}}$$

$$= - \frac{\sum_i \sum_j (m_{ij}/m) \log\left(\frac{m_{ij}/m}{(m_i+m)/m} \cdot \frac{m_j/m}{(m_j+m)/m}\right)}{\sum_j p_{+j} \log p_{+j}}$$

$$= - \frac{2 \sum_i \sum_j m_{ij} \log\left(\frac{m_{ij}}{(m_i+m)_j/m}\right)}{2m \sum_j p_{+j} \log p_{+j}}$$

$$= - \frac{\left(2 \sum_i \sum_j m_{ij} \log\left(\frac{m_{ij}}{\hat{\mu}_{ij}}\right)\right)}{2m \sum_j p_{+j} \log p_{+j}}$$

$$= - \frac{G^2}{2m \sum_j p_{+j} \log p_{+j}}$$