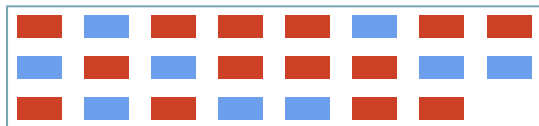


Analyzing categorical data

Chapter 13 of Rice

07/20/2021

In the previous lecture,



If $F=G$, every assignment of ranks to the pooled observations is equally likely.

- Comparing two independent samples:
 - Under Normal assumption, we can use t test.
 - Without any distributional assumption, use Mann-Whitney test:

$$U_1 = mn + \frac{n(n+1)}{2} - R_1$$

$$U_2 = mn + \frac{m(m+1)}{2} - R_2$$

U statistic is $U = \min\{U_1, U_2\}$, and rejection region is $R = \{U \leq c\}$.

- Comparing paired samples:

- Under Normal assumption, we can use t test.

$$\text{Under } H_0, T(\mathbf{X}_n) = \frac{\bar{D}_n}{\sqrt{\frac{1}{n} S_D^2}} \sim t_{n-1}.$$

- Without any distributional assumption, use Wilcoxon ranked sum test:

W_+ = positive rank sum,

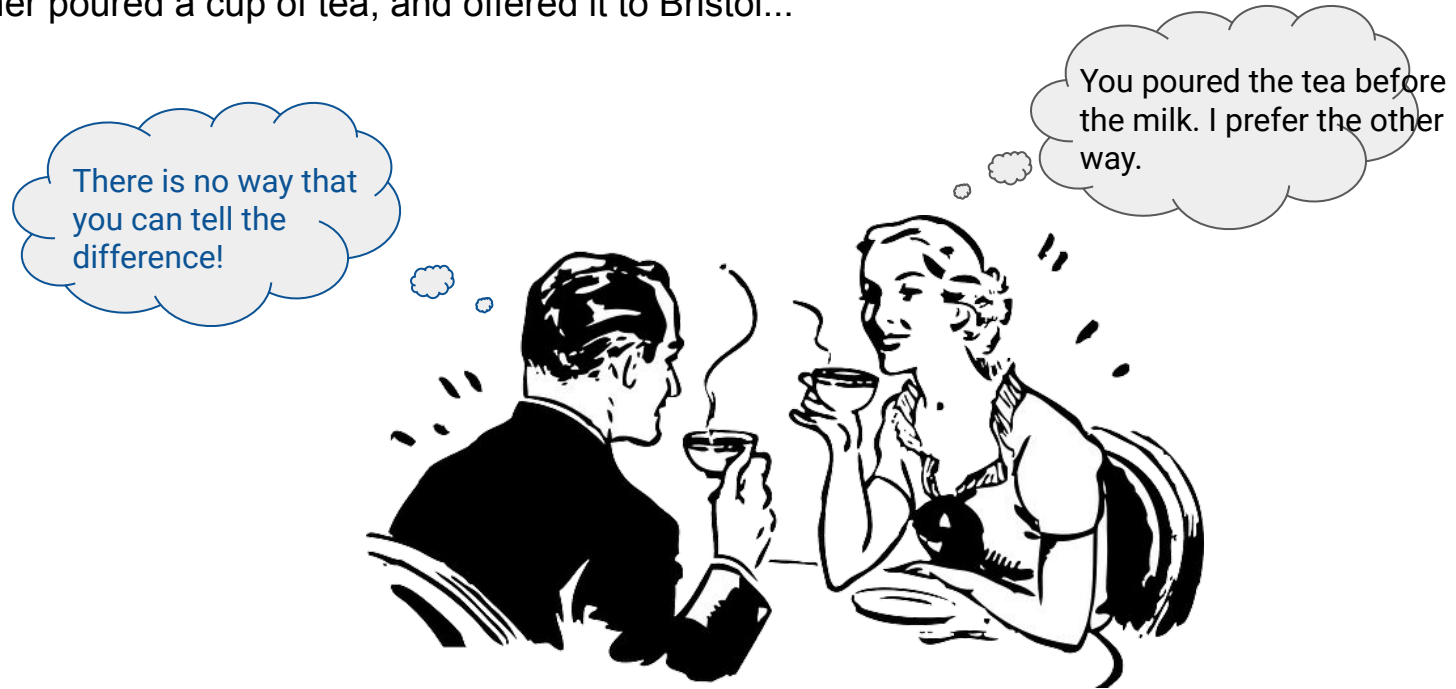
W_- = negative rank sum,

$R = \{W_+ \leq c \text{ or } W_- \leq c\}$ for two-sided hypothesis.

- `wilcox.test(x,y, paired=TRUE)`

Tea-tasting experiment

In 1919, Fisher poured a cup of tea, and offered it to Bristol...



Ronald Fisher

Dr. Muriel Bristol

Tea-tasting experiment



Randomly sort the cups and serve to Bristol, and let her **choose 4** that were prepared by the second method.

Ho = Bristol does not know how to choose the cups prepared by the 2nd method.

4 prepared by first pouring the tea, then adding milk

4 prepared by first pouring the milk, then adding tea

$$\binom{8}{4} = 70 \text{ combinations} \quad \swarrow$$

Tea-tasting experiment

H_0 : Bristol has no skills in determining the order.

→ she has correct skill

→ she has a skill resulting in totally opposite

Success count	Combinations of selection	Number of Combinations
0	0000	$1 \times 1 = 1$
1	000X, 00X0, 0X00, X000	$4 \times 4 = 16$
2	00XX, 0X0X, 0XX0, X0X0, XX00, X00X	$6 \times 6 = 36$
3	0XXX, X0XX, XX0X, XXX0	$4 \times 4 = 16$
4	XXXX	$1 \times 1 = 1$
Total		70

Handwritten notes: A bracket on the left groups rows 0-4. An arrow points from the 'Total' row to the text 'Hypergeometric distribution:'. To the right of the table, there are handwritten binomial coefficients: $\binom{4}{4}$, $\binom{4}{3}$, $\binom{4}{2}$, $\binom{4}{1}$, $\binom{4}{0}$.

Bristol correctly selected out all 4 cups of the second method. classification

$$p\text{-value} = P(\text{count} \geq 4 | H_0) + P(\text{count} \leq 0 | H_0) = \frac{2}{70} \approx 0.028.$$

< 0.05

Hypergeometric distribution:

$$P(X = k) = \frac{\binom{n}{k} \binom{m}{r-k}}{\binom{n+m}{r}}, k = 0, \dots, r.$$

Binomial distribution



n coin flips

p captures the fairness of the coin

$$P(Y = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

↓

number of heads

Multinomial distribution



n tosses

p_i is the probability that i comes up

$$p_1 + \dots + p_6 = 1$$

$$P(Y_1 = k_1, \dots, Y_6 = k_6) = \frac{n!}{k_1! \dots k_6!} p_1^{k_1} \dots p_6^{k_6}, \text{ with } k_1 + \dots + k_6 = n.$$

$\downarrow \qquad \qquad \downarrow$

number of
 i in the n
tosses

$$\underline{E Y_i = n p_i}, \quad i=1, 2, \dots, 6.$$

LRT for multinomial distribution

9.5 of Rice

07/20/2021

LRT for multinomial distribution

Example 1. Let $\theta = (p_1, \dots, p_6)$, and X_1, \dots, X_n are the results of n tosses. Consider testing $H_0 : p_1 = p_2, p_3 = p_4 = p_5 = p_6$ versus $H_1 : H_0$ is not true.

Solution. Denote $Y_j = \#$ of observations equal to j .

We know that $(Y_1, \dots, Y_6) \sim \text{multinomial}(p_1, \dots, p_6)$

$$L(p_1, \dots, p_6 | Y_1, \dots, Y_6) = \frac{n!}{Y_1! \dots Y_6!} p_1^{Y_1} \dots p_6^{Y_6}$$

$$\mathcal{H}_0 = \left\{ \underbrace{p_1 = p_2}_{\text{arrow}}, \underbrace{p_3 = p_4 = p_5 = p_6}_{\text{arrow}}, \sum_{i=1}^6 p_i = 1 \right\}, \quad \mathcal{H} = \mathcal{H}_0 \cup \mathcal{H}_1 = \left\{ \sum_{i=1}^6 p_i = 1 \right\}$$

$$2p_1 + 4p_3 = 1, \quad p_3 = \frac{1 - 2p_1}{4}$$

$$\dim \mathcal{H}_0 = 1$$

$$\dim \mathcal{H} = 5$$

$$\nu = \dim \mathcal{H} - \dim \mathcal{H}_0 = 4$$

$$\sup_{\mathcal{H}} L(p_1, \dots, p_6 | Y_1, \dots, Y_6) = \sup_{p_1, \dots, p_5} \frac{n!}{Y_1! \dots Y_6!} p_1^{Y_1} \dots p_5^{Y_5} \underbrace{(1 - p_1 - \dots - p_5)^{Y_6}}_{\text{Denote } L_0(p_1, \dots, p_5)}$$

$$\text{Then } \ell_0(p_1, \dots, p_5) = \log L_0(p_1, \dots, p_5) = Y_1 \log p_1 + \dots + Y_5 \log p_5 + Y_6 \log (1 - p_1 - \dots - p_5)$$

LRT for multinomial distribution

Solution cont'd. Take the derivatives with respect to p_1, \dots, p_5 :

$$\left\{ \begin{array}{l} \frac{\partial \ell_0}{\partial p_1} = \frac{\gamma_1}{p_1} - \frac{\gamma_6}{1-p_1-\dots-p_5} = 0 \rightarrow \cancel{\gamma_1 (1-p_1-\dots-p_5)} - \gamma_6 \cdot \underbrace{(p_1)}_{\hat{p}_1} = 0 \\ \vdots \\ \frac{\partial \ell_0}{\partial p_5} = \frac{\gamma_5}{p_5} - \frac{\gamma_6}{1-p_1-\dots-p_5} = 0 \rightarrow \cancel{\gamma_5 (1-p_1-\dots-p_5)} - \gamma_6 \cdot \underbrace{(p_5)}_{\hat{p}_5} = 0 \end{array} \right.$$

$\hat{p}_1 = \frac{\gamma_1 (1-p_1-\dots-p_5)}{\gamma_6} = \frac{\gamma_1 p_6}{\gamma_6}$
 $= \frac{\gamma_1}{\gamma_6} \cdot \frac{\gamma_6}{\gamma_1+\dots+\gamma_6} = \frac{\gamma_1}{\gamma_1+\dots+\gamma_6}$
 $\hat{p}_2 = \frac{\gamma_2}{\gamma_1+\dots+\gamma_6}$
 $\hat{p}_5 = \frac{\gamma_5}{\gamma_1+\dots+\gamma_6}$

Sum over these 5 equations:

$$\begin{aligned} & (\gamma_1 + \dots + \gamma_5) (1-p_1-\dots-p_5) - \gamma_6 (p_1 + \dots + p_5) = 0 \\ \Leftrightarrow & \gamma_1 + \dots + \gamma_5 - \underbrace{(p_1 + \dots + p_5) (\gamma_1 + \dots + \gamma_5)}_{= (\gamma_1 + \dots + \gamma_5) (p_1 + \dots + p_5)} - (p_1 + \dots + p_5) \gamma_6 = 0 \\ \Leftrightarrow & \underbrace{(\gamma_1 + \dots + \gamma_5) - (p_1 + \dots + p_5) (\gamma_1 + \dots + \gamma_5)}_{= (\gamma_1 + \dots + \gamma_5) (1 - p_1 - \dots - p_5)} - (p_1 + \dots + p_5) \gamma_6 = 0 \\ \Leftrightarrow & \hat{p}_6 = 1 - \frac{\gamma_1 + \dots + \gamma_5}{\gamma_1 + \dots + \gamma_6} = \frac{\gamma_6}{\gamma_1 + \dots + \gamma_6} \end{aligned}$$

Then $\sup_{(p_1, \dots, p_6) \in \Delta} L(p_1, \dots, p_6 | \gamma_1, \dots, \gamma_6)$
 $= \frac{n!}{\gamma_1! \dots \gamma_6!} \hat{p}_1^{\gamma_1} \dots \hat{p}_6^{\gamma_6}$

$$\sup_{(H)_0} L(p_1, \dots, p_b | \gamma_1 \dots \gamma_b) = \sup_{\substack{p_1 = p_2 \\ p_3 = \dots = p_b \\ 2p_1 + 4p_3 = 1}} L(p_1, \dots, p_b | \gamma_1 \dots \gamma_b) = \sup_{p_1} \frac{n!}{\gamma_1! \dots \gamma_b!} \underbrace{p_1^{\gamma_1} p_1^{\gamma_2}}_{p_1^{\gamma_1 + \gamma_2}} \underbrace{p_2^{\gamma_3} p_3^{\gamma_4} p_3^{\gamma_5} p_3^{\gamma_6}}_{\left(\frac{1-2p_1}{4}\right)^{\gamma_3 + \dots + \gamma_6}}$$

$$= \sup_{p_1} \frac{n!}{\gamma_1! \dots \gamma_b!} \underbrace{p_1^{\gamma_1 + \gamma_2} \left(\frac{1-2p_1}{4}\right)^{\gamma_3 + \dots + \gamma_6}}_{L_0(p_1)}$$

$$\ell_0(p_1) = \log L_0(p_1) = (\gamma_1 + \gamma_2) \log p_1 + (\gamma_3 + \dots + \gamma_6) \log \left(\frac{1-2p_1}{4}\right)$$

$$\frac{\partial \ell_0(p_1)}{\partial p_1} = \frac{\gamma_1 + \gamma_2}{p_1} - 2 \frac{\gamma_3 + \dots + \gamma_6}{\frac{1-2p_1}{4}} = 0 \Rightarrow \tilde{p}_1 = \frac{\gamma_1 + \gamma_2}{\gamma_1 + \dots + \gamma_6}$$



$$\sup_{(H)_0} L(p_1, \dots, p_b | \gamma_1 \dots \gamma_b) = \frac{n!}{\gamma_1! \dots \gamma_b!} \tilde{p}_1^{\gamma_1 + \gamma_2} \tilde{p}_3^{\gamma_3 + \dots + \gamma_6}$$

$$\tilde{p}_1 = \frac{(\gamma_1 + \gamma_2) / 2}{\gamma_1 + \dots + \gamma_6}$$

$$\tilde{p}_3 = \frac{(\gamma_3 + \dots + \gamma_6) / 4}{\gamma_1 + \dots + \gamma_6}$$

$$\lambda(\mathcal{E}_n) = \frac{\sup_{(H)_0} L}{\sup_{(H)} L} = \left(\frac{\tilde{p}_1}{\hat{p}_1}\right)^{\gamma_1} \left(\frac{\tilde{p}_1}{\hat{p}_2}\right)^{\gamma_2} \dots \left(\frac{\tilde{p}_3}{\hat{p}_6}\right)^{\gamma_6}$$

$$-2 \log \lambda(\mathcal{E}_n) = 2 \sum_{i=1}^6 -\log \left(\frac{\tilde{p}_i}{\hat{p}_i}\right)^{\gamma_i} = 2 \sum_{i=1}^6 \gamma_i \log \frac{\hat{p}_i}{\tilde{p}_i}$$

$$\mathcal{P} = \{ \lambda(\mathcal{E}_n) \leq c \}$$

$$= \{ -2 \log \lambda(\mathcal{E}_n) \geq \chi_{4(\alpha)}^2 \}$$

$$\Rightarrow \chi_{\nu}^2 = \chi_4^2$$

$$\underline{-2 \log \lambda(\Sigma_n)} = 2 \sum_{i=1}^b \gamma_i \log \frac{\hat{p}_i}{\widetilde{p}_i} = 2 \sum_{i=1}^b \gamma_i \log \frac{n \hat{p}_i}{n \widetilde{p}_i} = \underline{2 \sum_{i=1}^b O_i \log \frac{O_i}{E_i}} \quad \leftarrow$$

$$\gamma_1 + \dots + \gamma_b = n$$

$$\hat{p}_i = \frac{\gamma_i}{\gamma_1 + \dots + \gamma_b} = \frac{\gamma_i}{n}$$

$$\Rightarrow n \hat{p}_i = \gamma_i = O_i \quad \leftarrow \text{observed count}$$

$$\widetilde{p}_1 = \frac{(\gamma_1 + \gamma_2) / 2}{\gamma_1 + \dots + \gamma_b}$$

$$\Rightarrow n \widetilde{p}_1 = \frac{\gamma_1 + \gamma_2}{2} = E_1 \quad \leftarrow \begin{array}{l} \text{expected count} \\ \text{under } H_0 \end{array}$$

$$\widetilde{p}_3 = \frac{(\gamma_3 + \dots + \gamma_b) / 4}{\gamma_1 + \dots + \gamma_b}$$

$$\Rightarrow n \widetilde{p}_3 = \frac{\gamma_3 + \dots + \gamma_b}{4} = E_3$$

$$E_1 = E_2, \quad E_3 = \dots = E_b$$

LRT for multinomial distribution

Example 2. Let $\theta = (p_1, \dots, p_4)$, and X_1, \dots, X_n are the outcomes of n experiments. Consider testing $H_0 : p_1 = 2\theta, p_2 = \theta, p_3 = 3\theta, p_4 = 1 - 6\theta$ versus $H_1 : H_0$ is not true. in which $\theta \in (0, 1/6)$.

Solution. Denote $Y_j = \#$ of observations equal to j .

$$L(p_1, \dots, p_4 | Y_1, \dots, Y_4) = \frac{n!}{Y_1! \dots Y_4!} p_1^{Y_1} \dots p_4^{Y_4}$$

$$\Theta_0 = \{p_1 = 2\theta, p_2 = \theta, p_3 = 3\theta, p_4 = 1 - 6\theta\}$$

$$\dim \Theta_0 = 1 \rightarrow \nu = \dim \Theta - \dim \Theta_0 = 2 \leftarrow \dim \Theta = 4 - 1 = 3.$$

$$\Theta = \Theta_0 \cup \Theta_1 = \left\{ \sum_{i=1}^4 p_i = 1 \right\}$$

$$\sup_{\Theta} L(p_1, \dots, p_4 | Y_1, \dots, Y_4) = \frac{n!}{Y_1! \dots Y_4!} \hat{p}_1^{Y_1} \dots \hat{p}_4^{Y_4} \quad \text{where } \hat{p}_i = \frac{Y_i}{n}, i=1, 2, 3, 4.$$

$$\begin{aligned} \sup_{\Theta_0} L(p_1, \dots, p_4 | Y_1, \dots, Y_4) &= \sup_{\theta \in (0, 1/6)} \frac{n!}{Y_1! \dots Y_4!} (2\theta)^{Y_1} \theta^{Y_2} (3\theta)^{Y_3} (1-6\theta)^{Y_4} \\ &= \sup_{\theta \in (0, 1/6)} \frac{n!}{Y_1! \dots Y_4!} 2^{Y_1} 3^{Y_3} \theta^{Y_1 + Y_2 + Y_3} (1-6\theta)^{Y_4} \\ &\quad \underbrace{\theta^{Y_1 + Y_2 + Y_3}}_{L_0(\theta)} \end{aligned}$$

$$\ell_0(\theta) = \log L_0(\theta) = (Y_1 + Y_2 + Y_3) \log(\theta) + Y_4 \log(1 - 6\theta)$$

$$\frac{\partial \ell_0}{\partial \theta} = \frac{Y_1 + Y_2 + Y_3}{\theta} - 6 \frac{Y_4}{1 - 6\theta} = 0 \Rightarrow \hat{\theta} = \frac{(Y_1 + Y_2 + Y_3) / 6}{1}$$

$$\sup_{\textcircled{H_0}} L(P_1, \dots, P_4 | Y_1, \dots, Y_4) = \frac{n!}{Y_1! \dots Y_4!} \underbrace{(2\hat{\theta})^{Y_1}}_{P_1(\hat{\theta})} \underbrace{\hat{\theta}^{Y_2}}_{P_2(\hat{\theta})} \underbrace{(3\hat{\theta})^{Y_3}}_{P_3(\hat{\theta})} \underbrace{(1 - 6\hat{\theta})^{Y_4}}_{P_4(\hat{\theta})}$$

$$\lambda(\mathcal{E}_n) = \frac{\sup_{\textcircled{H_0}} L}{\sup_{\textcircled{H}} L} = \left(\frac{P_1(\hat{\theta})}{\hat{P}_1} \right)^{Y_1} \dots \left(\frac{P_4(\hat{\theta})}{\hat{P}_4} \right)^{Y_4}$$

$$-2 \log \lambda(\mathcal{E}_n) = 2 \sum_{i=1}^4 Y_i \log \frac{\hat{P}_i}{P_i(\hat{\theta})} = 2 \sum_{i=1}^4 O_i \log \frac{O_i}{E_i}$$

$$\xrightarrow{d} \chi_{\nu}^2 = \chi_2^2$$

$$O_i = Y_i$$

$$E_i = n P_i(\hat{\theta})$$

LRT for multinomial distribution

Theorem A. Let X_1, \dots, X_n be the results of n experiments. Each experiment has m possible outcomes. Consider testing

$$H_0 : (p_1, \dots, p_m) = (p_1(\theta), \dots, p_m(\theta)) \text{ versus } H_1 : H_0 \text{ is not true.}$$

Then the likelihood ratio under the null hypothesis satisfies

$$\underbrace{-2 \log \lambda(\mathbf{X}_n)}_{\sim} = 2 \underbrace{\sum_{i=1}^m O_i \log \frac{O_i}{E_i}}_{\sim} \xrightarrow{d} \chi_{m-1}^2, \text{ as } n \rightarrow \infty.$$

$O_j = \#$ of observations equal to j .

$\hat{\theta}$ = the MLE in H_0 .

$$\underbrace{E_i = np_i(\hat{\theta})}$$

$$\Theta = \left\{ \sum_{i=1}^m p_i = 1 \right\} \quad \dim \Theta = m-1$$

$$\Theta_0 = \left\{ p_i = p_i(\Theta) \right\} \quad \dim \Theta_0 = 1$$

LRT for multinomial distribution

Corollary A. Let X_1, \dots, X_n be the results of n experiments. Each experiment has m possible outcomes. Consider testing

$H_0 : (p_1, \dots, p_m) = (p_1(\theta), \dots, p_m(\theta))$ versus $H_1 : H_0$ is not true.

Then the likelihood ratio under the null hypothesis satisfies $\rightarrow \chi^2$ statistic

$$2 \sum_{i=1}^m O_i \log \frac{O_i}{E_i} \approx \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \xrightarrow{d} \chi_{m-2}^2 \text{ as } n \rightarrow \infty.$$

$O_j = \#$ of observations equal to j .

$\hat{\theta}$ = the MLE in H_0 .

$$E_i = np_i(\hat{\theta})$$

$$\sum_{j=1}^m O_j = n$$

$$\begin{aligned} \sum_{i=1}^m E_i &= \sum_{i=1}^m n p_i(\hat{\theta}) \\ &= n \sum_{i=1}^m p_i(\hat{\theta}) = 1 = n \end{aligned}$$

$$\begin{aligned} f(x) &= x \log \frac{x}{x_0}, \quad f'(x) = 1 + \log \frac{x}{x_0}, \quad f''(x) = \frac{1}{x} \\ f(x) - f(x_0) &\approx \underbrace{f'(x_0)}_1 (x - x_0) + \underbrace{\frac{f''(x_0)}{2}}_{\frac{1}{2x_0}} (x - x_0)^2 \\ \Rightarrow \underbrace{f(O_i)}_1 - \underbrace{f(E_i)}_0 &\approx (O_i - E_i) + \frac{1}{2E_i} (O_i - E_i)^2 \\ \Rightarrow 2 \sum_{i=1}^m O_i \log \frac{O_i}{E_i} &\approx 2 \sum_{i=1}^m (O_i - E_i) + 2 \sum_{i=1}^m \frac{(O_i - E_i)^2}{2E_i} \Rightarrow 2 \sum_{i=1}^m O_i \log \frac{O_i}{E_i} \approx \sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} \end{aligned}$$

Application: Goodness of fit test

0, 1, ..., n, ∞

all integers
pos

Example 3. Researcher observed the number of emissions of α particles in many 10-sec intervals. If we fit the data using a $\text{Poisson}(\lambda)$ model, what is the best parameter estimate for λ ?

How good is the model assumption? Are the observations really Poisson distributed?

n	Observed
0-2	18
3	28
4	56
5	105
6	126
7	146
8	164
9	161
10	123
11	101
12	74
13	53
14	23
15	15
16	9
17+	5

$$e^{-\lambda} \left(\frac{\lambda^0}{0!} + \frac{\lambda^1}{1!} + \frac{\lambda^2}{2!} \right) \leftarrow p_1(\lambda)$$

$$e^{-\lambda} \frac{\lambda^3}{3!} \leftarrow p_2(\lambda)$$

$$e^{-\lambda} \frac{\lambda^4}{4!} \leftarrow p_3(\lambda)$$

$$\vdots$$

$$e^{-\lambda} \frac{\lambda^{16}}{16!} \leftarrow$$

$$1 - e^{-\lambda} \left(\frac{\lambda^0}{0!} + \dots + \frac{\lambda^{16}}{16!} \right) \leftarrow p_{16}(\lambda)$$

Average count $\bar{X}_n = 8.392$.

$$\hat{\lambda}_{MLE} = \bar{X}_n$$

H_0 = Obs subject to $\text{Poisson}(\lambda)$

vs. H_1 = H_0 is not true

$$(p_1, \dots, p_{16}) = (p_1(\lambda), \dots, p_{16}(\lambda))$$

$$2 \sum_{i=1}^{16} O_i \log \frac{O_i}{E_i} \rightarrow \chi^2_{m-2} = \chi^2_{14}$$

O_i = in the table

$$E_i = n p_i(\hat{\lambda}) = n p_i(\bar{X}_n)$$

Application: Goodness of fit test

Example 3. Researcher observed the number of emissions of α particles in many 10-sec intervals. If we fit the data using a $\text{Poisson}(\lambda)$ model, what is the best parameter estimate for λ ?

How good is the model assumption? Are the observations really Poisson distributed?

$$p\text{-value} = P(-2\log\lambda(\mathbf{X}_n) \geq 8.70854 \mid H_0) \\ = pchisq(8.70854, df=14, \text{lower.tail} = \text{FALSE}) \\ = \underline{0.84926}$$

n	Observed		Expected
0-2	18	$n e^{-\hat{\lambda}} \left(\frac{\hat{\lambda}^0}{0!} + \frac{\hat{\lambda}^1}{1!} + \frac{\hat{\lambda}^2}{2!} \right)$	12.2
3	28	$n e^{-\hat{\lambda}} \frac{\hat{\lambda}^3}{3!}$	27.0
4	56	$n e^{-\hat{\lambda}} \frac{\hat{\lambda}^4}{4!}$	56.5
5	105		94.9
6	126		132.7
7	146	\vdots	159.1
8	164	\vdots	166.9
9	161		155.6
10	123		130.6
11	101		99.7
12	74		69.7
13	53		45.0
14	23		27.0
15	15		15.1
16	9	$n e^{-\hat{\lambda}} \frac{\hat{\lambda}^{16}}{16!}$	7.9
17+	5	$n \left[1 - e^{-\hat{\lambda}} \left(\frac{\hat{\lambda}^0}{0!} + \dots + \frac{\hat{\lambda}^{16}}{16!} \right) \right]$	7.1

Average count $\bar{X}_n = 8.392$.

H_0 : Observations are Poisson distributed.

$$-2 \log \lambda(\mathbf{X}_n) = 2 \sum_{i=1}^m O_i \log \frac{O_i}{E_i} = \underline{8.70854.}$$

$$P = \left\{ -2 \log \lambda(\mathbf{X}_n) \geq \chi_{14}^2(\alpha) \right\}$$

$$\downarrow \\ qchisq(\alpha, df=14, \text{lower.tail} = \text{FALSE})$$

$= 23.685$
we fail to reject $H_0 \Rightarrow \text{Poisson}(\lambda)$ is a good fit.

Application: Goodness of fit test

$$p = \{ -2 \log \lambda(\hat{\theta}_n) \geq \chi^2_1(\alpha) \} = 3.841$$

We fail to reject H_0
 HW model is a good fit.

Example 4. If the gene frequencies are in equilibrium, the genotypes AA, Aa, aa occur in the population with probability θ^2 , $2\theta(1-\theta)$, $(1-\theta)^2$, according to the Hardy-Weinberg equilibrium model.

H_0 : HW model is true. versus H_1 : HW model is not true.

p-value = pchisq

(0.03250, df=1,
 lower.tail = FALSE)
 = 0.8569.

We observed the phenotypes:

AA	Aa	aa
342	500	187

$$H_0 : (p_1, p_2, p_3) = (\theta^2, 2\theta(1-\theta), (1-\theta)^2)$$

vs. H_1 : H_0 is not true.

$$\nu = 3 - 2 = 1$$

O_i

E_i

$$\begin{array}{ccc} n \hat{\theta}^2 & n \cdot 2 \hat{\theta} (1 - \hat{\theta}) & n (1 - \hat{\theta})^2 \\ \text{340.587} & 502.826 & 185.587 \end{array}$$

$$\sup_{H_0} L(p_1, p_2, p_3 | H_0) = \sup_{\theta} \frac{n!}{\gamma_1! \gamma_2! \gamma_3!} (\theta^2)^{\gamma_1} (2\theta(1-\theta))^{\gamma_2} (1-\theta)^{\gamma_3}$$

$$= \sup_{\theta} \frac{n! 2^{\gamma_2}}{\gamma_1! \gamma_2! \gamma_3!} \theta^{2\gamma_1 + \gamma_2} (1-\theta)^{\gamma_2 + \gamma_3}$$

Denote $L_0(\theta)$

$$\log L_0(\theta) = \log L(\theta) = (\gamma_1 + \gamma_2) \log \theta + (\gamma_2 + \gamma_3) \log (1-\theta)$$

$$\sum_{i=1}^m \frac{(O_i - E_i)^2}{E_i} = 0.03250$$

$$2 \sum_{i=1}^m O_i \log \frac{O_i}{E_i} = 0.03249$$

$$\Rightarrow \hat{\theta} = \frac{2\gamma_1 + \gamma_2}{2n} = \frac{2 \times 342 + 500}{2 \times 1029} = 0.57532$$

Application: Goodness of fit test

1. The population is a categorical variable so that the grouped cells subject to a multinomial distribution;
2. Tickets in each cell are independent; ✓ i.i.d assumption
3. Large sample size n so that no more than 20% of expected counts less than 5

Fisher's exact test

13.2 of Rice

07/20/2021

Dependencies between row and column classifications

Example 5. During phase 3 trial, some vaccine recipients were asked to complete diaries of their symptoms during the 7 days after vaccination.

	Pfizer / BNT162b2	Placebo
Fever $\geq 38.0^{\circ}\text{C}$	331	10
No fever	1,767	2,093
Total	2,098	2,103

* Systemic reactions in persons aged 18–55 years



Are occurrences of fever related to the vaccine/placebo treatment?

Could randomization result in such an imbalance?

H_0 : There is no relation. Any imbalance is due to randomization.

Dependencies between row and column classifications

Example 5. During phase 3 trial, some vaccine recipients were asked to complete diaries of their symptoms during the 7 days after vaccination.

	Treatment	Control	
Symptom	n_{11}	n_{12}	$n_{1.}$
No symptom	n_{21}	n_{22}	$n_{2.}$
	$n_{.1}$	$n_{.2}$	$n_{..}$

$$n_{11} + n_{12} = n_{1.}$$

$$\begin{pmatrix} n_{1.} \\ n_{11} \end{pmatrix} \quad \begin{pmatrix} n_{.2} \\ n_{12} \end{pmatrix} \rightarrow n_{1.} - n_{11}$$

H_1 : There is a relation.

H_0 : There is no relation. Any imbalance is due to randomization.

↓ induces

Hypergeometric distribution :

$$P(N_{11} = n_{11} | H_0) = \frac{\binom{n_{1.}}{n_{11}} \binom{n_{.2}}{n_{1.} - n_{11}}}{\binom{n_{..}}{n_{1.}}}$$

$$R = \left\{ \begin{array}{l} N_{11} \leq C_1 \text{ or} \\ N_{11} \geq C_2 \end{array} \right\}$$

Tea-tasting experiment

H_0 : Bristol has no skills in determining the order.

→

	<u>Milk first</u>	<u>Tea first</u>	
<u>Milk first</u>	4	0	4
<u>Tea first</u>	0	4	4
	4	4	8

← ←

Hypergeometric distribution :

$$P(N_{11} = n_{11} | H_0) = \frac{\binom{4}{n_{11}} \binom{4}{4-n_{11}}}{\binom{8}{4}}$$

N_{11}	0	1	2	3	4
p	0.014	0.229	0.514	0.229	0.014

$$P = \{ N_{11} \leq c_1 \text{ or } N_{11} \geq c_2 \}$$

$$P \left(\underline{N_{11} \leq 0} \text{ or } \underline{N_{11} \geq 4} \mid H_0 \right) = 0.014 + 0.014 = 0.028$$

21

Dependencies between row and column classifications

Example 6. A group of supervisors each examined a personnel file to decide whether to promote the employee or not. The files are identical except for the gender label.

H_0 : There is no gender bias. Any imbalance is due to randomization.

	Male	Female	
Promote	21	14	35
Hold file	3	10	13
	24	24	

* From 13.2 of Rice

Hypergeometric distribution :

$$P(N_{11} = n_{11} | H_0) = \frac{\binom{24}{n_{11}} \binom{24}{35-n_{11}}}{\binom{24+24}{35}}$$



`dhyperv(n11, m=24, n=24, k=35)`

```
p-value = 2*phyper(21-1, m=24, n=24, k=35,  
lower.tail = FALSE)
```

```
[1] 0.04899141
```

Dependencies between row and column classifications


Example 5 cont'd. During phase 3 trial, some vaccine recipients were asked to complete diaries of their symptoms during the 7 days after vaccination.

H_0 : There is no relation. Any imbalance is due to randomization.

	Pfizer / BNT162b2	Placebo	
Fever $\geq 38.0^\circ\text{C}$	331	10	341
No fever	1,767	2,093	3860
	2,098	2,103	

* Systemic reactions in persons aged 18-55 years

Hypergeometric distribution :

$$P(N_{11} = n_{11} \mid H_0) = \frac{\binom{2098}{n_{11}} \binom{2103}{341-n_{11}}}{\binom{2098+2103}{341}}$$


`dhyper(n11, m=2098, n=2103, k=341)`

*Benefits of more trials and repeated tests
⇒ More significant results*

`p-value = 2*phyper(331-1, m=2098, n=2103, k=341,
lower.tail = FALSE)`

`[1] 2.548842e-90`

Dependencies between row and column classifications

Example 7. Phase 3 trial was a large, randomized, double-blind, placebo-controlled clinical trial:

H_0 : Infection rate is not related to vaccine/placebo treatment. \leftrightarrow H_1 : It is related.

	Pfizer / BNT162b2	Placebo	
SARS-CoV-2 infected	9	169	178
No infection	21,711	21,559	43,439
	21,720	21,728	

* Age ≥ 16 , infections observed with onset at least 7 days after the second dose

Hypergeometric distribution :

$$P(N_{11} = n_{11} \mid H_0) = \frac{\binom{21720}{n_{11}} \binom{21728}{178-n_{11}}}{\binom{21720+21728}{178}}$$



`dhyper(n11, m=21720, n=21728, k=178)`

```
p-value = 2*phyper(9, m=21720, n=21728, k=178,  
lower.tail = TRUE)
```

```
[1] 1.702187e-39
```

Dependencies between row and column classifications

Example 7 *cont'd*. Phase 3 trial was a large, randomized, double-blind, placebo-controlled clinical trial:

H_0 : Infection rate is not related to vaccine/placebo treatment. \leftrightarrow H_1 : Infection rate is **lower** in the vaccine group.

	Pfizer / BNT162b2	Placebo	
SARS-CoV-2 infected	9	169	178
No infection	21,711	21,559	43,439
	21,720	21,728	

* Age ≥ 16 , infections observed with onset at least 7 days after the second dose

Hypergeometric distribution :

$$P(N_{11} = n_{11} \mid H_0) = \frac{\binom{21720}{n_{11}} \binom{21728}{178-n_{11}}}{\binom{21720+21728}{178}}$$



`dhyper(n11, m=21720, n=21728, k=178)`

```
p-value = phyper(9, m=21720, n=21728, k=178,  
lower.tail = TRUE)
```

```
[1] 8.510933e-40
```

Tomorrow ...

$I \times J$ table

- Chi-squared test of homogeneity
- Chi-squared test of independence