Intro to Contingency Tables

Author: Nicholas Reich and Anna Liu, based on Agresti Ch 1

Course: Categorical Data Analysis (BIOSTATS 743)

Made available under the Creative Commons Attribution-ShareAlike 4.0 International License.



Distributions of categorical variables: Multinomial

Suppose that each of n independent and identical trials can have outcome in any of c categories. Let

$$y_{ij} = \begin{cases} 1 & \text{if trial } i \text{ has outcome in category } j \\ 0 & \text{otherwise} \end{cases}$$

Then $\mathbf{y}_i = (y_{i1}, ..., y_{ic})$ represents a multinomial trial with $\sum_j y_{ij} = 1$. Let $n_j = \sum_i y_{ij}$ denote the number of trials having outcome in category j. The counts $(n_1, n_2, ..., n_c)$ have the *multinomial distribution*. The multinomial pmf is

$$p(n_1,...,n_{c-1}) = \left(\frac{n!}{n_1!n_2!...n_c!}\right) \pi_1^{n_1} \pi_2^{n_2} ... \pi_c^{n_c},$$

where $\pi_j = P(Y_{ij} = 1)$ $E(n_j) = n\pi_j, \quad Var(n_j) = n\pi_j(1 - \pi_j)$ $Cov(n_i, n_j) = -n\pi_i\pi_i$

Statistical inference for multinomial parameters

Given *n* observations in *c* categories, n_j occur in category *j*, j = 1, ..., c. The multinomial log-likelihood function is

$$L(\pi) = \sum_j n_j \log \pi_j$$

Maximizing this gives MLE

$$\hat{\pi}_j = n_j/n$$

The Chi-Squared distribution

This is not a distribution for the data but rather a sampling distribution for many statistics.

- ► The chi-squared distribution with degrees of freedom by *df* has mean *df*, variance 2(*df*), and skewness √8/*df*. It converges (slowly) to normality as *df* increases, the approximation being reasonably good when *df* is at least about 50.
- Let $Z \sim N(0,1)$, then $Z^2 \sim \chi^2(1)$
- ► The reproductive property: if $X_1^2 \sim \chi^2(\nu_1)$ and $X_2^2 \sim \chi^2(\nu_2)$, then $X^2 = X_1^2 + X_2^2 \sim \chi^2(\nu_1 + \nu_2)$. In particular, $X = Z_1^2 + Z_2^2 + ... + Z_{\nu}^2 \sim \chi^2(\nu)$ with the standard normal Z's.

Chi-square goodness-of-fit test for a specified multinomial

Consider hypothesis $H_0: \pi_j = \pi_{j0}, j = 1, ..., c$, - Chi-square goodness-of-fit statistic (score)

$$X^2 = \sum_j \frac{(n_j - \mu_j)^2}{\mu_j}$$

where $\mu_j = n\pi_{j0}$ is called **expected frequencies under** H_0 .

- Let X²_o denote the observed value of X². The P-value is P(X² > X²_o).
- For large samples, X² has approximately a chi-squared distribution with df = c − 1. The P-value is approximated by P(χ²_{c−1} ≥ X²_o).

LRT test for a specified multinomial

LRT statistic

$$G^2 = -2\log \Lambda = 2\sum_j n_j \log(n_j/n\pi_{j0})$$

For large n, G^2 has a chi-squared null distribution with df = c - 1.

- When H₀ holds, the goodness-of-fit Chi-squiare X² and the likelihood ratio G² both have large-sample chi-squared distributions with df = c − 1.
- ► For fixed c, as n increases the distribution of X² usually converges to chi-squared more quickly than that of G². The chi-squared approximation is often poor for G² when n/c < 5. When c is large, it can be decent for X² for n/c as small as 1 if table does not contain both very small and moderately large expected frequencies.

Distributions of categorical variables: Poisson

One simple distribution for count data that do not result from a fixed number of trials. The Poisson pmf is

$$p(y) = \frac{e^{-\mu}\mu^y}{y!}, y = 0, 1, 2, \dots E(Y) = Var(Y) = \mu$$

For adult residents of Britain who visit France this year, let

- Y_1 = number who fly there
- ► Y₂=number who travel there by train without a car
- ► Y₃=number who travel there by ferry without a car
- Y₄=number who take a car

A poisson model for (Y_1, Y_2, Y_3, Y_4) treats these as independent Poisson random variables, with parameters $(\mu_1, \mu_2, \mu_3, \mu_4)$. The total $n = \sum_i Y_i$ also has a Possion distribution, with parameter $\sum_i \mu_i$. Distributions of categorical variables: Poisson

The conditional distribution of (Y_1, Y_2, Y_3, Y_4) given $\sum_i Y_i = n$ is *multinomial* $(n, \pi_i = \mu_i / \sum_j \mu_j)$

Example: A survey of student characteristics

In the R data set survey, the Smoke column records the survey response about the student's smoking habit. As there are exactly four proper response in the survey: "Heavy", "Regul" (regularly), "Occas" (occasionally) and "Never", the Smoke data is multinomial.

library(MASS) # load the MASS package
levels(survey\$Smoke)

[1] "Heavy" "Never" "Occas" "Regul"
(smoke.freq = table(survey\$Smoke))

Heavy Never Occas Regul ## 11 189 19 17

Example: A survey of student characteristics

Suppose the campus smoking data are as shown above. You wish to test null hypothesis of whether the frequency of smoking is the same in all of the groups on campus, or $H_0: \pi_j = \pi_{j0}, j = 1, ..., 4$.

```
##
## Chi-squared test for given probabilities
##
## data: smoke.freq
## X-squared = 382.51, df = 3, p-value < 2.2e-16</pre>
```

Thus, there is strong evidence against the null hypothesis that all groups are equally represented on campus (p<.0001).

Example (continued): expected and observed counts

x2.test\$expected

Heavy Never Occas Regul ## 59 59 59 59

x2.test\$observed

##

Heavy Never Occas Regul ## 11 189 19 17

Testing with estimated expected frequencies

In some applications, the hypothesized $\pi_{j0} = \pi_{j0}(\theta)$ are functions of a smaller set of unknown parameters θ .

For example, consider a scenario (Table 1.1 in *CDA*) in which we are studying the rates of infection in dairy calves. Some calves become infected with pneumonia. A subset of those calves also develop a secondary infection within two weeks of the first infection clearing up. The goal of the study was to test whether the probability of primary infection was the same as the conditional probability of secondary infection, given that the calf got the primary infection. Let π be the probability of primary infection. Fill in the following 2x2 table with the associated probabilities under the null hypothesis:

	Secondary	Infection	
Primary Infection Yes	Yes	No	Total
No			

Example continued

Let n_{ab} denote the number of observations in row a and column b.

	Secondary	Infection	
Primary Infection	Yes	No	Total
Yes	<i>n</i> ₁₁	<i>n</i> ₁₂	
No	<i>n</i> ₂₁	n ₂₂	

The ML estimate of π is the value maximizing the kernel of the multinomial likelihood

$$(\pi^2)^{n_{11}}(\pi-\pi^2)^{n_{12}}(1-\pi)^{n_{22}}$$

The MLE is

$$\hat{\pi} = (2n_{11} + n_{12})/(2n_{11} + 2n_{12} + n_{22})$$

One process for drawing inference in this setting would be the following:

- ► Obtain the ML estimates of expected frequencies: μ̂_j = nπ_{j0}(θ̂) by plugging in the ML estimates θ̂ of θ
- Replace μ_j by $\hat{\mu}_j$ in the definition of X^2 and G^2
- Use the approximate distributions of X^2 and G^2 are χ^2_{df} with $df = (c-1) dim(\theta)$.

Example continued

A sample of 156 dairy calves born in Okeechobee County, Florida, were classified according to whether they caught pneumonia within 60 days of birth. Calves that got a pneumonia infection were also classified according to whether they got a secondary infection within 2 weeks after the first infection cleared up.

	Secondary	Infection
Primary Infection	Yes	No
Yes	30(38.1)	63(39.0)
No	0	63(78.9)

The MLE is

$$\hat{\pi} = (2n_{11} + n_{12})/(2n_{11} + 2n_{12} + n_{22}) = 0.494$$

The score statistic is $X^2 = 19.7$. It follows a Chi-square distribution with df = c - p - 1 = (3 - 1) - 1 = 1.

Example continued

The p-value is

$$P(\chi_1^2 > 19.7) =$$

1-pchisq(19.7, df=1)

[1] 9.060137e-06

Therefore, the evidence suggests that the probability of primary and secondary infections being the same is not supported by the data. Under H_0 , we would anticipate that many more calves would have secondary infections than did end up being infected. "The researchers concluded that primary infection had an immunizing effect tht reduced the likelihood of a secondary infection."