■ MH statistic combine information of $y_{11k}$'s from $K$ tables:

*[margin left note:]* If some $\Delta_\beta \gg 1$ & some $\Delta_\beta \ll 1$ then some $y_{11k} \gg E(y_{11k})$ & some $y_{11k} \ll E(y_{11k})$ they might cancel out & MH stat. becomes small (not rejected)

*[top annotations:]* $y_{11+} \leftarrow$ Does it actually use only marginal 2×2 table? Ans. No. Check — continuity correction (increasing p-value, more conservative, check LNp.5-7)

$$\frac{\left(\left|\sum_k [y_{11k} - E(y_{11k})]\right| - 1/2\right)^2}{\sum_k Var(y_{11k})} \overset{a}{\sim} \chi_1^2 \leftarrow \dim(H_1) - \dim(H_0)$$

∵ $y_{11k}$'s are independent.

*[right note:]* check (*) in LNp. 5-32

where $E(y_{11k})$ and $Var(y_{11k})$ are calculated under the $H_0$ →

□ can calculate an exact p-value for smaller *[note: # of possible values of $y_{11k}$'s are not large (check LNp 5-12)]* dataset using hypergeometric distribution

*[margin note:]* many small $y_{ijk}$'s

⇒ useful when data is sparse, under which the $\chi^2$ approximations based on asymptotic thm is questionable

➢ MH test is sometimes called Cochran-Mantel-Haenszel test because a version without the 1/2 is published earlier by Cochran (1954).

❖ **Reading**: Faraway (2006, 1st ed.), 4.4

*[box:]* 5/14

# **Ordinal Variables**

*[annotation:]* (LNp. 2-3 ~4) ① continuous (interval) ② discrete interval ③ ordinal ④ nominal

*[annotation:]* ⓐ In analysis, cannot arbitrarily rearrange the order of categories ⓑ the order matters for the objective of analysis Otherwise, we can find order(s) in many (almost any) variables, but some not matter

• Some variables have a nature ordering between categories *[distance = ?]*

➢ e.g., education: HS, BA, MA; political ideology: VL, SL, M, SC, VC

*[annotation:]* cf. • ordinal response 隨 ← LNp.5-28 • ordinal covariate 規

---

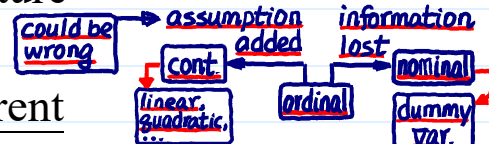*[margin left note:]* no difference in the treatments for ordinal & nominal

➢ The ordinal structure not matter when # of categories = 2 *(□)*

➢ For ordinal variables, can use the methods for nominal variable

■ But, more information can be extracted by *[note: some dummy var., e.g. Helmert, forward difference, ···. more meaningful for ordinal predictor]* taking advantage of the ordinal structure

*[margin note:]* check ★ in LNp. 2-4

➢ Treatments for ordinal response (future *[could be wrong]* lecture) and ordinal covariates are different

*[annotation:]* assumption added → information lost; cont. → nominal; linear, quadratic, ··· ordinal dummy var.

• Treatment for ordinal predictors: assign each category a *score* *[cf. CA in LNp. 5-15]*

➢ It kind of turns an ordinal variable into a continuous variable

*[margin note:]* like adding distance btwn categories

➢ The choice of scores requires some judgment

*[annotation:]* • known values (Δ) ⊙ should increase with order • assignment of scores can be quite flexible (*) ⊙ scores cf. dummy var.

■ If no particular preference, even spacing allows for the simplest interpretation

*[margin box:]* order 1 2 ··· I treat them as quantitative

■ For interval scales, midpoints of the intervals are often used

➢ Should check whether the inference is robust to different assignments of scores *[e.g., sign of β̂]* *[might ignore order & treat ordinal as nominal]*

*[margin note:]* usually defined for continuous variables

■ If *qualitative* conclusions are changed, this is an indication that you cannot make any strong finding based on scores

• Poisson GLM with linear-by-linear association for 2-way tables:

➢ Consider table with ordinal row ($X_1$) and column ($X_2$) variables

- assign scores $u_1 \le u_2 \le \dots \le u_I$ to rows, denoted by $u(X_1)$ | known values
- assign scores $v_1 \le v_2 \le \dots \le v_J$ to columns, denoted by $v(X_2)$

| check LNp.5-5 |
➤ Linear-by-linear association model: | cf. interaction in Poisson GLM

$$\eta_{ij} = \log(\mu_{ij}) = \log(t\,\pi_{ij}) = \log(t) + \log(\pi_{i+}) + \log(\pi_{+j}) \quad |X_1| \quad |X_2|$$
$$+ \gamma \times u_i \times v_j \qquad = u_i(\gamma v_j)$$

where $u_i$'s, $v_j$'s are known scores, | only 1 parameter | interaction | $\gamma v_j \equiv$ ★

and $\gamma$ is an unknown parameter

| use codings for nominal $X_1$: I-1 d.f. $X_2$: J-1 d.f. |

- $Y \sim X_1 + X_2 + u(X_1)v(X_2) \equiv S_{O \times O}$ | $\vec{d}$ $\vec{d}+1$

➤ Some notes about $\gamma$: | only 1 d.f. | larger $|\gamma|$, stronger association (under same $u_i$'s, $v_j$'s)

| $S_{O \times O}$ reduced to $Y \sim X_1 + X_2$ main-effect model ⟷ independent |
- values of $\gamma$ represents the amount of association
- ⊙ $\gamma = 0 \Leftrightarrow$ independence | $\gamma > 0$: $X_1 \uparrow$, $X_2 \uparrow$ / $X_1 \downarrow$, $X_2 \downarrow$ → check (✳) in LNp 5-34 / $\gamma < 0$: $X_1 \uparrow$, $X_2 \downarrow$ / $X_1 \downarrow$, $X_2 \uparrow$

| e.g., $u_{i+1} - u_i = v_{j+1} - v_j = 1$ or $u_{i+2} - u_i = \beta$ $v_{j+2} - v_j = \ell$ |
- positive and negative $\gamma$
- Interpretation of $\gamma$ by log-odds-ratio:

$$\log\left(\frac{\pi_{i,j}\,\pi_{i+1,j+1}}{\pi_{i,j+1}\,\pi_{i+1,j}}\right) = \log\left(\frac{\mu_{i,j}\,\mu_{i+1,j+1}}{\mu_{i,j+1}\,\mu_{i+1,j}}\right)$$

| $\log\left(\dfrac{\pi_{i,\vec{d}}\,\pi_{i+\beta,\vec{d}+\ell}}{\pi_{i,\vec{d}+\ell}\,\pi_{i+\beta,\vec{d}}}\right) = \gamma(u_{i+\beta} - u_i) \times (v_{\vec{d}+\ell} - v_{\vec{d}})$ | (exercise) |

$$= (\eta_{i,j} + \eta_{i+1,j+1}) - (\eta_{i,j+1} + \eta_{i+1,j}) = \gamma(u_{i+1} - u_i)(v_{j+1} - v_j)$$

| $= \gamma u_i v_j + \gamma u_{i+1}v_{j+1} - \gamma u_i v_{j+1} - \gamma u_{i+1}v_j$ | scores affect the interpretation of $\gamma$ |

| uniform association for 3-way table (LNp.5-25~26) cf. |
- ◻ for evenly spaced scores, these log-odds-ratios are equal ⟹ called uniform association in Goodman (1979) | fixed $\beta$, $\ell$ any $i$, $j$

---

- Latent (continuous) variable $Z$ motivation for $\gamma$: | SS1, Poisson (LNp.5-4~5) or SS2, multinomial (LNp.5-8)

| Z can not be directly observed, but, a function of Z, $X = f(Z)$, can be observed to gain information of Z |
- ◻ Assume $\pi_{ij}$'s are obtained by putting a grid on an approximately bi-variate Normal $(Z_1, Z_2)$ for latent variables and $u_i$'s and $v_j$'s are cutpoints | This explains under what conditions a model like $S_{O \times O}$ is appropriate for a 2-way table

| $N\left(\binom{u_L}{u_2}, \binom{\sigma_1^2 \ \rho\sigma_1\sigma_2}{\rho\sigma_1\sigma_2 \ \sigma_2^2}\right)$ |

| $\gamma$ / $\rho$ : positive ↔ positive, negative ↔ negative |
- ⊙ $\gamma$ can then be identified with the correlation coefficient $\rho$ of the latent variables (cf., positive and negative $\rho$) | $\gamma \approx \rho/(1-\rho^2)$ if $u_i$'s & $v_j$'s are standardized (Agresti, 2013, 10.4.1)

| Note. indep. $\subset S_{O \times O}$ $\subset S_{N \times N}$ |
⊘ Q: for the tests of independence or goodness-of-fit, what is the benefit of using $S_{O \times O}$ over the nominal approach, i.e., fitting a nominal-by-nominal model $S_{N \times N}$: $Y \sim X_1 + X_2 + X_1{:}X_2$? | $X_1 = i$, $X_2 = \vec{d}$ ⟺ $Z_1 \in (u_{i-1}, u_i)$ & $Z_2 \in (v_{\vec{d}-1}, v_{\vec{d}})$

| $X_1$: 7 levels $X_2$: 7 levels |
As shown in a lab example, | (I-1)(J-1) = 6×6 = 36 d.f. | saturated model

| $H_0$: indep $H_1$: $S_{N \times N}$ not significant |
- in the $N \times N$ approach, interaction effects reduce a deviance of 40.743 on 36 degrees of freedom, but | SS3, product multinomial (LNp 5-10) regression line

| cf. $H_0$: indep $H_1$: $S_{O \times O}$ significant |
- the $O \times O$ interaction effect reduces a deviance of 10.175 on *one* degrees of freedom, i.e., the other 35 interaction effects only reduce a deviance of 30.568 | $Z_2 | Z_1 = 3_1 \sim$ simple regression model | ordinal | deviance & its d.f. of $S_{O \times O}$

⊙ Ordinal-by-nominal model (or nominal-by-ordinal model) → for interactions in Poisson GLM association    p. 5-37

*Note row & column are ordinal variables*

— row — column (check (□) in LNp 5-34) — row — column

➤ Rows (or columns) assigned scores, but column (known) (or row) variable treated as a nominal variable   use dummy variables

➤ called *column* (or *row*) *effects model* because the columns (or rows) are not assigned scores; instead, their effects are estimated

*Recall. coding for an interaction: componentwise product of a coding of $X_1$ & a coding of $X_2$*

■ alternative viewpoint: the scores of the ordinal columns (or rows) regarded as parameters (unknown)   *i.e., coefficients of dummy var*

*Under Poisson GLM. Check LNp 5-5*

⊙➤ Column effects model:

Under ✱ in LNp 5-35   proportion unknown (estimated by data) in O×N
(△) $\sigma_1 : \sigma_2 : \cdots \sigma_J$ = $v_1 : v_2 : \cdots v_J$ — proportion known in O×O

$$\eta_{ij} = \log(\mu_{ij}) = \log(t\,\pi_{ij})$$

$X_2$ only   interactions

$$= \log(t) + \log(\pi_{i+}) + \log(\pi_{+j}) + u_i \times \gamma_j$$

$X_1$ only     $X_1$ as continuous (scores)

where   $u_i$'s, $i=1,\ldots, I$, are known scores, and
$X_2$ as nominal (coefficients of dummy var.)

*each column (jth) a parameter $\gamma_j$* → $\gamma_j$'s, $j=1,\ldots, J$, are unknown parameters (over-

*adding 1 constraint. e.g., $\gamma_1 = 0$ or $\Sigma_j \gamma_j = 0$*   (*) parameterized; only requires $J-1$ parameters),

dummy var   reference level

∵ [$v_j$] in LNp 5-35 ∈ the space of [$\gamma_j$]

*use codings for nominal $X_1$: I-1 d.f., $X_2$: J-1 d.f.*

■ $Y \sim X_1 + X_2 + u(X_1){:}X_2 \equiv S_{O\times N} \,(\supseteq S_{O\times O})$

— (J−1) d.f.

— denoted by $S_{N\times O}$ for row effects model

➤ Some notes about $\gamma_j$'s, called the *column effects*:

---

*Under the constraint (✱) in LNp 5-37, it becomes $\gamma_1 = \cdots = \gamma_J = 0$ ⇒ no interaction model ⇔ independent*

■ Equality of the $\gamma_j$'s (then, $u_i \times \gamma_j = u_i \times \gamma$) corresponds to the hypothesis of independence between $X_1$ and $X_2$   p. 5-38

It becomes a main effect of $X_1$ & is merged into the $X_1$ term in $S_{O\times N}$ in LNp 5-37

*check (△) in LNp 5-37*

■ For ordinal column variable, if the model $S_{O\times O}$ were a good fit, we would expect the estimates of the $\gamma_j$'s in $S_{O\times N}$ to be roughly proportional to $v_j$'s (e.g., for evenly spaced $v_j$'s, estimates of $\gamma_j$'s should follow a linear trend)

*Note. $u_i$'s are fixed, i.e., under the assumption that $u_i$'s are suitable scores for rows*

data determine →

■ We can use the estimates of $\gamma_j$'s in $S_{O\times N}$ (1) to   of order (by (△) in LNp 5-37) examine whether the chosen scores for columns in $S_{O\times O}$ (i.e., $v_j$'s) are appropriate, or (2) to possibly suggest better scores (see an example in lab)

*Note. In $S_{O\times N}$, the estimated $\hat\gamma_j$'s might not increase or decrease with the order of $X_2$ (cf., (△) in LNp 5-34)*

• Some advantages of using scores for ordinal variables

*identifying few but meaningful significant interactions*

➤ helpful in reducing the complexity of models for categorical data with ordinal variables

cf. nominal →

➤ especially useful in higher dimensional table where a reduction in the # of parameters is particularly welcome

| | d.f. |
|---|---|
| main effect, | $X_1$: I−1, $X_2$: J−1, $X_3$: K−1 |
| 2 f.i., | $X_1 X_2$: (I−1)(J−1) |
| | $X_1 X_3$: (I−1)(K−1) |
| | $X_2 X_3$: (J−1)(K−1) |
| 3 f.i., | $X_1 X_2 X_3$: (I−1)(J−1)(K−1) |

➤ can also sharpen our ability to detect associations

❖ **Reading**: Faraway (2006, 1st ed.), 4.5   CA(LNp.5-13~17) ← cf. → interactions in Poisson GLM