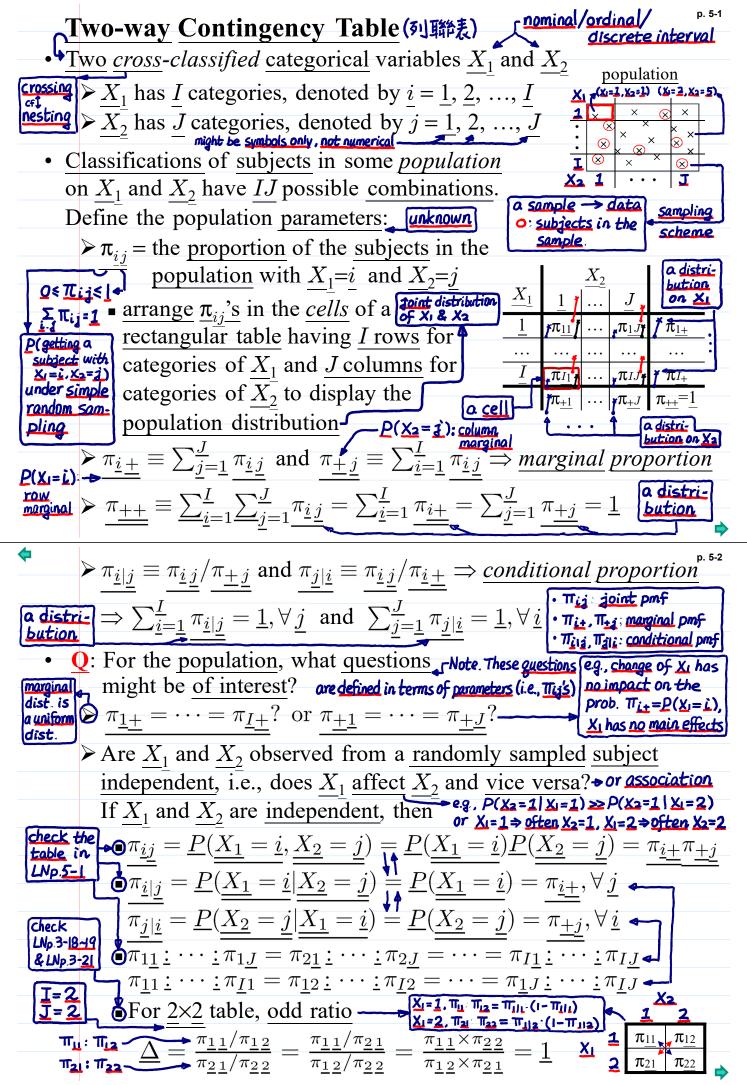
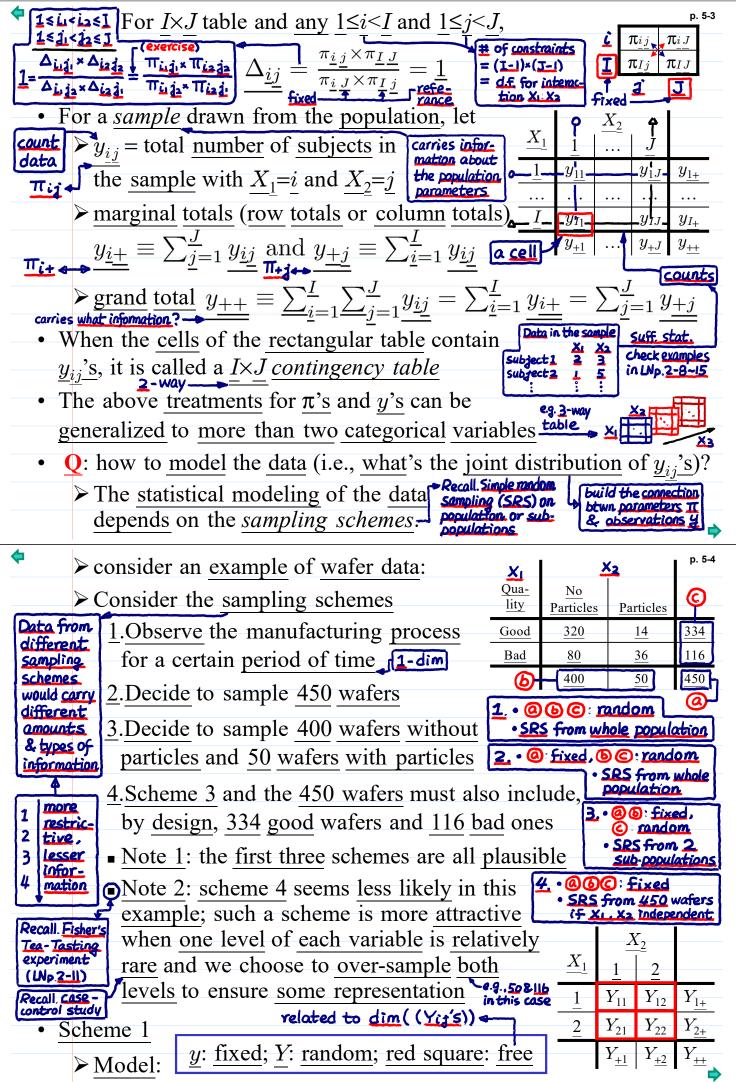
NTHU STAT 5230, 2025



made by S.-W. Cheng (NTHU, Taiwan)

NTHU STAT 5230, 2025



made by S.-W. Cheng (NTHU, Taiwan)

NTHU STAT 5230, 2025 Lecture Notes p. 5-5 Process (LNp. 4-3) All Yiis SRS e.g., average vield -----Y++~Bisson(t) • For a random sample, can assume $\mu_{ij} = t \times \pi_{ij}$, have same $\pi_{\mu} \pi_{12} \pi_{21} \pi_{22}$ value of the where t is an unknown value of a size variable Y11 Y12 Y21 Y22 <u>size variable</u> $t = \alpha$ parameter. Note. $Y_{\pm\pm} = \sum_{ij} Y_{ij} \sim Poisson(t)$ (cf. LNp. 4-12) • X_1 (=*i*) and X_2 (=*j*) are <u>covariates</u> Yij~Poisson(41) 212 X1=1 -check LNp1-31 2:2 Iog. 🗸 721 722 X1=2 Suppose the <u>data</u> (from an $I \times J$ table) is <u>XB = 20</u> fitted with a Poisson GLM with log link $O = (\mathcal{I}_{11} - \mathcal{I}_{12}) - (\mathcal{I}_{21} - \mathcal{I}_{22})$ Recall. $= \frac{\log \left(\frac{\pi_{11} * \pi_{22}}{\pi_{12} * \pi_{21}}\right)$ • When $\underline{\pi}_{ij} = \underline{\pi}_{i+1} \underline{\pi}_{+j} (\underline{X}_1 \text{ and } \underline{X}_2 \text{ independent}),$ ANOVA model =<u>r</u>+α<u>i</u>+βi $\alpha_{I} = 0$ Bj = 0 · Main effects of X1 • Main effects of X₁ (<u>1-1</u> parameters \therefore $\Sigma_i \pi_{i+=1}$) \longrightarrow corresponds to a main-effect model, i.e., $\underline{Y_{ij}} \sim \underline{X_1} + \underline{X_2} \equiv \underline{S}$ <u>ΞΣjΠ+j=1)</u> • When $\underline{\pi}_{ij} = \underline{\pi}_{i+1} \underline{\pi}_{+j}$ and $\underline{\pi}_{\underline{1}\underline{+}} = \cdots = \underline{\pi}_{\underline{I}\underline{+}}$ (or $\underline{\pi}_{\underline{+}\underline{1}} = \cdots = \underline{\pi}_{\underline{+}J}$) <u>uniform</u> dist. $\underline{\eta_{\underline{ij}}} = \underline{\log(\underline{t})} + \underline{\log(\underline{\pi_{\pm j}})} \quad (\text{or } \eta_{\underline{ij}} = \underline{\log(\underline{t})} + \underline{\log(\pi_{\underline{i\pm}})})$ $\beta_{I} = Q$ $= \Gamma + \alpha_i$ \Rightarrow corresponds to the model $\underline{Y}_{ij} \sim \underline{X}_2$ (or $\underline{Y}_{ij} \sim \underline{X}_1$) $\alpha_1 = 0$ • When $\underline{\pi}_{ij} \neq \underline{\pi}_{i+} \underline{\pi}_{+j} (X_1 \text{ and } X_2 \text{ not independent})$, check LNp.5-3 p. 5-6 \Rightarrow add interaction $\underline{X}_1: \underline{X}_2$ (I-1)(J-1) $\widehat{\mu}_{ij} = Y_{ij}$ $\widehat{\Delta}_{ij} \neq 1$ for a pair of (i.i) $lij = \Upsilon + \alpha_i$ \Rightarrow may consider $\underline{Y}_{ij} \sim \underline{X}_1 + \underline{X}_2 + \underline{X}_1 : \underline{X}_2 \equiv \underline{L}$ (saturated model) + βj+(αβ)ii • Q: what type of π 's corresponds • without or with interactions \Leftrightarrow independent or dependent $\alpha_{I} = 0$ $\beta_3 = 0$ to the following models? • without or with main effects a (αβ)_{ij=0}, <u>γ</u> uniform or non-uniform marginal dist. (αβ)₁j=0. 4j $\underline{IJ}-\underline{I}-\underline{J}+\underline{I} \underbrace{\underline{H}_{21}}_{\underline{I}_{21}} \sim \underline{1} \qquad \underline{Y}_{\underline{ij}} \sim \underline{X}_{\underline{1}} + \underline{X}_{\underline{1}}: \underline{X}_{\underline{2}}$ $\underline{Y}_{ij} \sim \underline{X}_2 + \underline{X}_1 : \underline{X}_2 \text{ exp} \hat{\mathbf{1}} = \underline{\mathbf{X}} \hat{\mathbf{1}}$ \blacktriangleright <u>Recall</u>. For a <u>Poisson GLM</u> with log link, $\underline{X}^T \underline{Y} = \underline{X}^T$ $=Y_{1+}+Y_{2+}$ For models without interactions, canonical reg. 2*2 table $= Y_{\pm 1} + Y_{\pm 2}$ $\Rightarrow X^T Y \text{ is only related to marginal totals} \qquad X = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{11} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot X^T Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{12} \\ Y_{12} \\ Y_{12} \end{bmatrix} \cdot Y = \begin{bmatrix} Y_{12} \\ Y_{1$ sufficient statistics of B YE+'S.Y+a's Y ++ Yit's Y+i's \Rightarrow the fitted values $\hat{\mu}$ is a In general, consider <u>2-way</u> function of marginal totals $\rightarrow \hat{\mu}_{ij} = \exp(\hat{i}_{ij}) = \exp(\hat{i}_{ij} + \hat{a}_{i} + \hat{b}_{j})$ ANOVA mode $= t \pi_{i+} \pi_{i+} \Rightarrow$ for example, for main-effect model $Y_{ij} \sim X_1 + X_2$ in DOE $= \underbrace{\underbrace{\mathcal{U}_{i\pm}}_{\mathcal{U}_{\pm\pm}}, \underbrace{\mathcal{U}_{\pm\pm}}_{\mathcal{U}_{\pm\pm}} = \underbrace{\mathcal{U}_{i\pm} \times \mathcal{U}_{\pm\pm}}_{\mathcal{U}_{\pm\pm}} \Rightarrow \underbrace{\hat{\mu}_{\underline{i}\underline{j}}}_{\underline{i}\underline{j}} = \underbrace{Y_{\pm\pm}}_{Y_{\pm\pm}} \underbrace{\hat{\pi}_{\underline{i}\underline{i}}}_{\underline{i}\underline{i}\underline{j}} \stackrel{\frown}{=} \underbrace{Y_{\underline{i}\underline{i}\underline{j}}}_{\underline{i}\underline{j}} \underbrace{Y_{\pm\underline{j}}}_{\underline{i}\underline{j}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm}}_{\mathcal{U}_{\pm}} \underbrace{Y_{\pm}} \underbrace{Y_{\pm$ using Poisson

made by S.-W. Cheng (NTHU, Taiwan)