$P_x$ = average of all $P_{ij}$'s at $X=x$, ∵ random sampling clusters LM $\underset{cf.}{\overset{s.f.}{\leftrightarrow}}$ what if not?

(*) In LM, if assume
$Y \sim N(X\beta, I)$
t-test → Z-test
F-test → $\chi^2$-test
if assume
$Y \sim N(X\beta, \sigma^2 I)$ cf.

◉ $E(y) = \sum_i E(s_i) = \sum_i \{E[E(s_i|\pi_i)]\} = l \times mp = np$   cf. p. 3-40  (→ $m\pi_i$)

◆ $Var(y) = \sum_i Var(s_i) = \sum_i \{E[Var(s_i|\pi_i)] + Var[E(s_i|\pi_i)]\}$

$s_1,\ldots,s_\ell$ independent   $= \sum_i \{E[m\pi_i(1-\pi_i)] + Var(m\pi_i)\}$    $Var(y)$ if $y \sim$ binomial

$y$ not binomial if $\sigma^2>1$   $= l \times \{mp - m[\tau^2 p(1-p) + p^2] + m^2 \tau^2 p(1-p)\}$

$\boxed{\sigma^2 \equiv}$ ★ $= [1+(m-1)\tau^2]\,np(1-p) \geq np(1-p)$    $1 \leq \sigma^2 \leq m \leq n$

◆ Overdispersion cannot arise when $n=1$ (sparse case) & $m=1$

$Var(y)$ ③/26
$= Var(\sum_u z_u)$   $\sim B(1,p)$
$= \sum_u Var(z_u)$   $\sim np(1-p)$
$+ 2 \sum_{u<v} cov(z_u, z_v)$

■ Violation of independence assumption can cause over-dispersion, e.g., response has a common cause, say a disease is influenced by genes, the responses will tend to be positively correlated

Check 2 Z's from same cluster in LN p. 3-39 cf. LM
$E(z_u z_v) = E[E(z_u z_v|\pi_i)]$
$= E(\pi_i^2) = \tau^2 p(1-p) + p^2$
$cov(z_u, z_v)$ ←

assumption:
$\sigma^2 = \dfrac{Var(y_x)}{n_x P_x(1-P_x)}$
is a constant over $X$ & $n_x$ (△)

□ under-dispersion, e.g., when food supply is limited, survival probability of an animal may be increased by the death of others, i.e., negatively correlated   cf.

$P(z_u=1|z_v=0)$
$> P(z_u=1|z_v=1)$
survive →

Note, not assign a likelihood

• **Q**: how to model overdispersion and do analysis?   why? check ★   in addition to mean structure $\eta = \sum_\ell \beta_\ell \hat{R}_\ell$

➤ Introduce one additional dispersion parameter $\sigma^2$, i.e.,   check (*)

$Var(y_x) = \sigma^2 \times n_x p_x(1-p_x) \Leftarrow$ notice its similarity to linear model

defined by $g^{-1}(\eta_x) = E(y_x)/n_x$   often include under-dispersion ⟶ $\sigma^2<1$ ⟶ $0<\sigma^2<\infty$

(standard binomial case $\Rightarrow \sigma^2=1$; over-dispersion $\Rightarrow \sigma^2>1$)

---

➤ For a model $S$, $\sigma^2$ can be estimated using   why? check ★ in LN p.34   p. 3-41

$$\hat{\sigma}_S^2 = X_S^2/(k-p)$$

When $n_i$'s large, using $D_S$ is OK.

$\eta = \sum_{\ell=1}^p \beta_\ell R_\ell \leftarrow D_S \cong X_S^2 = \sum_i (r_{i,S})^2 \overset{cf.}{\leftrightarrow}$ RSS in LM

(using deviance $D$ in place of $X^2$ is not very recommended as $D$ may be inconsistent for sparse data)

# of covariate classes   # of parameters in $\eta$ ($S$)

check LN p.3-14 → $D/(k-p)$ might not converge to $\sigma^2=1$ as $k\to\infty$

$n_i$'s $=1$, $\sigma^2=1$

MLE ➤ Estimation of $\beta$ is unaffected since $E(y_x) = n_x P_x$ is not changed (Why? Note that $y_x$ is not $= n_x g^{-1}(\eta_x)$ $\sim$ binomial so that likelihood is different) $= n_x g^{-1}(\sum_\ell \beta_\ell R_\ell)$ → quasi-likelihood

rationale
• IRWLS only need mean & var functions
• weights $\propto [Var(y_x)]^{-1}$ $= 1/[\sigma^2 n_x P_x(1-P_x)]$ $\propto 1/[n_x P_x(1-P_x)]$ ← binomial case

Intuition In LM.
if $\varepsilon \to c\varepsilon$
• $Var(\varepsilon) \to c^2 Var(\varepsilon)$
• $\hat{\varepsilon} \to c\hat{\varepsilon}$
• RSS → $c^2$RSS
• $E(\hat{\beta})$ not change
• $Var(\hat{\beta}) \to c^2 Var(\hat{\beta})$

⊘ But, $Var(\hat{\beta}) \approx \sigma^2(X^T W X)^{-1}$ and $\hat{Var}(\hat{\beta}) = \hat{\sigma}^2(X^T \hat{W} X)^{-1}$

parameter   from IRWLS under binomial | estimate

⊘ For $S$ nested in $L$, difference in their deviances   cannot use it as a null dist.

$$D_S - D_L \overset{a}{\approx} \sigma^2 \chi^2_{df_S - df_L} \quad \text{(under } S\text{)}$$

parameter   check LN p.3-10 cf.

Let $L$ be the saturated model. Then, under $S$
$D_S, X_S^2 \overset{a}{\approx} \sigma^2 \chi^2_{df_S}$

➤ When comparing models, e.g., testing $H_0:S$ vs. $H_1:L\backslash S$, can use

$D_{sat}$   $D_L$   $D_S$

$$F = \frac{(D_S - D_L)/(df_S - df_L)}{\hat{\sigma}_L^2} \overset{a}{\approx} F_{df_S - df_L, df_L} \quad \text{(under } S\text{)}$$

$X_L^2/df_L$   LM cf.

i.e., adding a $\sigma^2$ multiplier on binomial variance

check graph ★ in LN p.3-37   ∵ no information about true $\sigma^2$ and $X_L^2=0$ if $L$ is saturated

more possible that $Var(y_x) / [n_x P_x(1-P_x)]$ is a constant (check ★ in LN p. 3-40)

➤ No goodness-of-fit test is possible

➤ This dispersion parameter method is more appropriate when the covariate classes are roughly equal in size (i.e., $n_1 \approx n_2 \approx \ldots \approx n_k$)

➤ Alternative approaches to over-dispersion  *more flexible than dispersion parameter method*  **p. 3-42**

- beta-binomial method (Williams, 1982; Crowder, 1978)

- quasi-likelihood: specify only how the mean and variance of the response are connected to covariates. → *1st moment* → *2nd moment*

*likelihood is known (cf. dispersion parameter method)*

$P_x \sim beta(\alpha_x, \beta_x)$
$y_x | P_x \sim B(n_x, P_x)$
$y_x \sim beta\text{-}binomial$

$E(y_x) = n_x \left(\frac{\alpha_x}{\alpha_x + \beta_x}\right) \underset{P_x^*}{\underbrace{\quad}}$   *a linear function of $n_x$ (cf. (▲) in LNp. 3-40)*

$Var(y_x) = n_x \left(\frac{\alpha_x}{\alpha_x + \beta_x}\right)\left(\frac{\beta_x}{\alpha_x + \beta_x}\right)\left(\frac{\alpha_x + \beta_x + n_x}{\alpha_x + \beta_x + 1}\right)$
$\underset{P_x^*}{\underbrace{\quad}} \underset{1 - P_x^*}{\underbrace{\quad}} \underset{\geq 1}{\underbrace{\quad}}$

*But, not the whole dist. family (i.e., likelihood unknown) use them to define a function working as likelihood for estimation & testing ⇒ same estimator & test stat. for all dist. with same mean & var functions*

❖ **Reading**: Faraway (2006, 1st ed.), 2.11

*LNp 3-19 ~24*

# Matched Case-Control Studies —— *Recall 1. blocking in DOE*

*Recall 2. paired $t$ (within block comparison) vs. 2-sample $t$*

- **Q**: In a case-control study, how should we choose the controls if there exist some *confounding variables* $W$, say age and sex, that may affect the outcome in addition to the *risk factors* $X$? *covariates of main interest*

*block factor in DOE, i.e., covariates of no interest, but should be considered in design and analysis.*

*deal with $W$ in data analysis*

➤ Approach 1: record and include confounding variables as covariates in GLM analysis (however, we may not be interested in the effects of the confounding variables)

*cf. ANCOVA*

*Why include covariates of no interest in model? Recall. In LM, true: $Y = X_1\beta_1 + X_2\beta_2 + \varepsilon$, fitted: $Y = X_1\beta_1 + \varepsilon$*

*→ may cause a large number of covariate classes (R) & possibly sparse data*

*deal with $W$ in data collection*

➤ Approach 2: confounding variables are explicitly adjusted for in the design

*cf. blocked design*

*→ Recall. In paired $t$ comparison block effect removed.*

---

- *Matched case-control design (MCCD)*: match each case with one **p. 3-43** or more controls that have the same or similar values of some set of potential confounding variables. A group of a case and its corresponding controls is called a *matched set*, e.g.,

*having same value of $W$*
*allows within block comparison in this contingency table*

*units in a matched set are regarded as homogeneous (∵ they have same values of $W$, check ★ in LNp. 3-39, $W$ = other covariates)*

| | 1st case | | | 2st case | | | $n$st case | |
|---|---|---|---|---|---|---|---|---|
| | age=20; sex=male | *a block* | | age=20; sex=female | *a block* | | age=70; sex=female | |

*a matched set ↓cf.*

| controls → | $D^c$ | $D$ | | $D^c$ | $D$ | | $D^c$ | $D$ |
|---|---|---|---|---|---|---|---|---|
| *values of risk factor* $X^c$ | $y_{100}$ | 1 | $X^c$ | $y_{200}$ | 0 | $X^c$ | $y_{n00}$ | 1 |
| $X$ | $y_{110}$ | 0 | $X$ | $y_{210}$ | 1 | $X$ | $y_{n10}$ | 0 |

⊕ = 1
⊕ = M
● ● ●

*case → controls →*

➤ 1:$M$ MCCD: $M$ controls for each case   *Recall. s.e.$(\bar{y}) = \sigma/\sqrt{n} \downarrow$ as $n \uparrow$ But,*

- $M$ typically small, can vary in size in every matched set

*more controls higher efficiency*

- Each additional control yields a diminished return in terms of increased efficiency in estimating the effects of risk factors

- It is usually not worth exceeding $M=5$

*$Var(\bar{y}_1 - \bar{y}_2) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$ if $\sigma_1^2/n_1 \ll \sigma_2^2/n_2$, $\sigma_2^2/n_2$ dominates. Then, it's of not much use to increase $n_1$*

- Some disadvantages of MCCD

*cf. Approach 1 in LNp. 3-42 can estimate effects of $W$*

➤ Lose the possibility of discovering the effects of the confounding variable $W$

*→ cannot estimate true block effects of $W$, e.g., in each matched set, estimate of $P(Z=1|W=w_j)$ is $\frac{1}{(M+1)}$ due to 1:M setting ⇒ irrelevant to $w_j$*

➤ The data will likely be far from a random sample of the population of interest

*→ might not able to gain some information of the population, e.g.*

*check example in LNp. 3-21*