

- A retrospective sampling is as effective as a prospective one for estimating Δ (provided (1) the probabilities of inclusion in D and in D^c are homogeneous or their ratio is irrelevant to covariates, and (2) data is reliable, e.g., no problems such as inaccurate or incomplete historical records; or unreliable memory of the subject)

This manipulation is not possible for other links ever mentioned.

Q: When can retrospective sampling work (under logit link)? retrospective

Consider a scenario: a study with response Z and covariates X and

z_j : binary response of j^{th} unit (e.g., disease present/not present)

x_j : covariate values of j^{th} unit in the population

$I_j = 1$ if j^{th} unit is included in the study, 0 if not

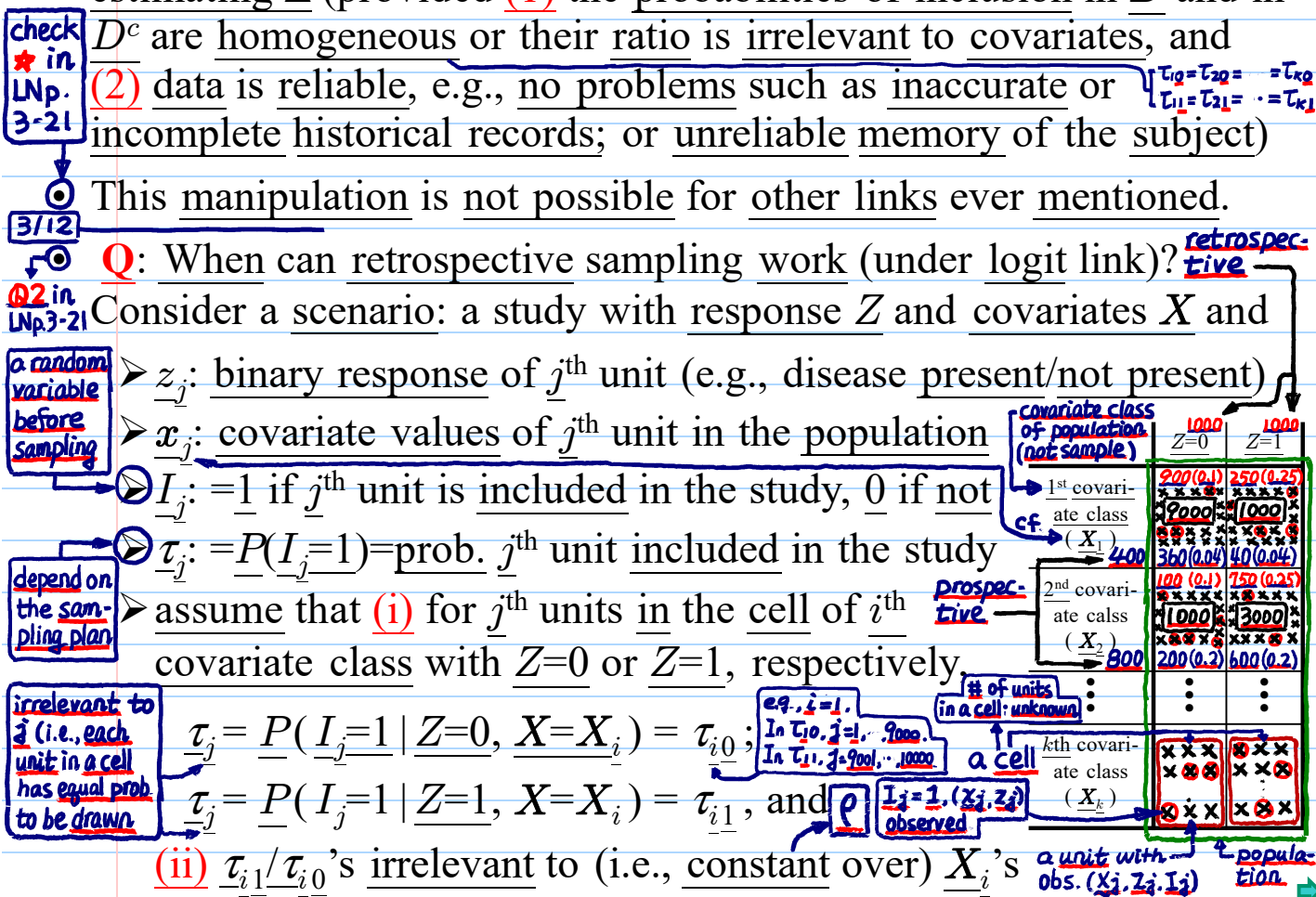
$\tau_j = P(I_j = 1) = \text{prob. } j^{\text{th}} \text{ unit included in the study}$

assume that (1) for j^{th} units in the cell of i^{th} covariate class with $Z=0$ or $Z=1$, respectively.

$$\tau_j = P(I_j = 1 | Z=0, X=X_i) = \tau_{i0}$$

$$\tau_j = P(I_j = 1 | Z=1, X=X_i) = \tau_{i1}, \text{ and}$$

(ii) τ_{i1}/τ_{i0} 's irrelevant to (i.e., constant over) X_i 's



For fair prospective study, $\tau_{i0} = \tau_{i1}, i=1, \dots, k$. $\tau = 1$ check example in LNp. 3-22

For retrospective study, often $\tau_{i0} \neq \tau_{i1}$ and τ_{i1} much larger than τ_{i0} $\tau > 1$

The probability of interested is

$$p_i \equiv P(Z=1 | X=X_i) \leftarrow \text{describe how covariates } X \text{ affect the response } Z$$

We can use data (either from prospective or retrospective sampling) to study the probability:

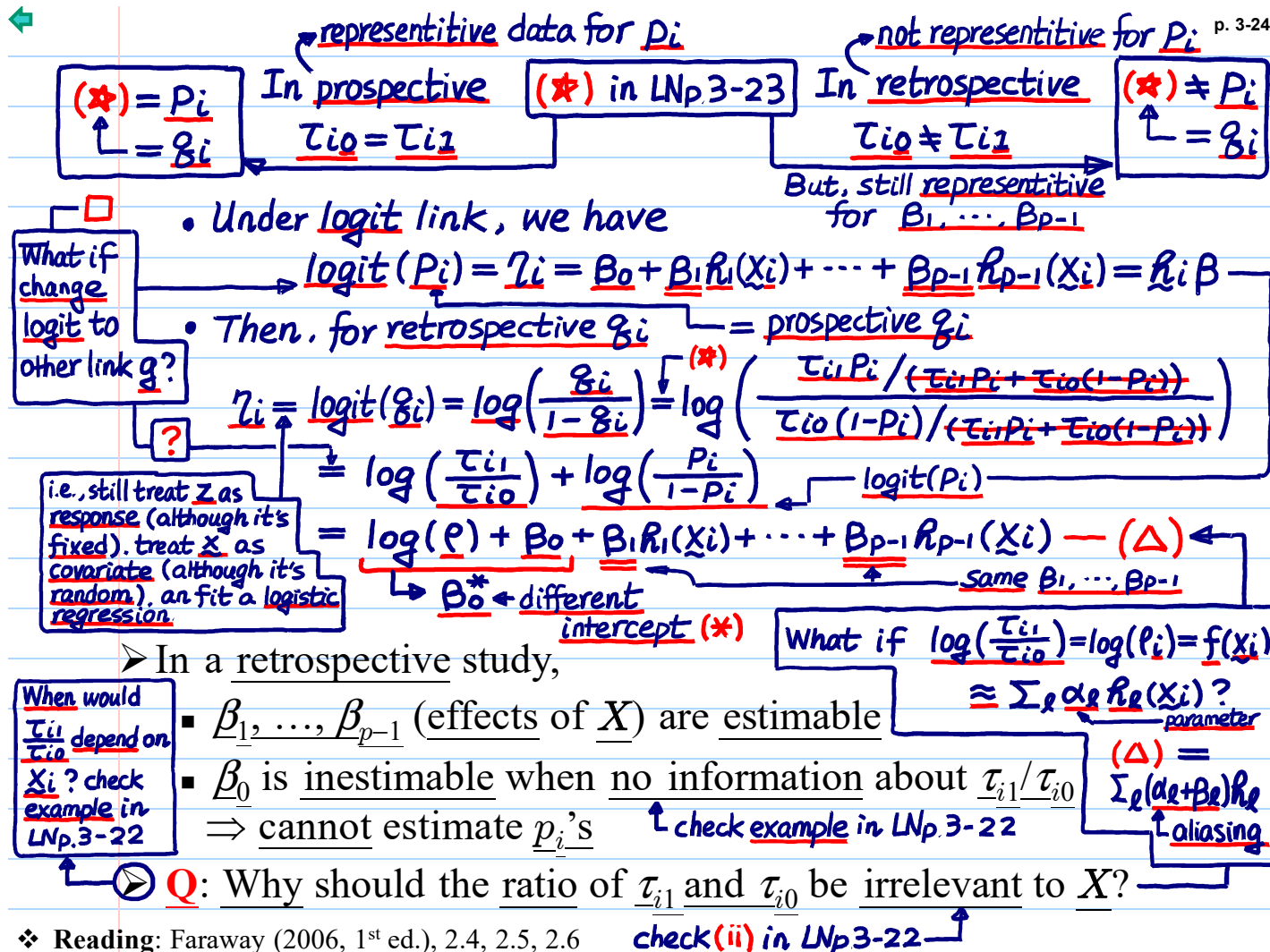
$$q_i \equiv P(Z=1 | I=1, X=X_i)$$

in observed data

Q1: $p_i \neq q_i$, check this

Q2: the way X_i affect p_i \neq the way X_i affect q_i (Recall ANCOVA model)

$$\begin{aligned}
 q_i &= \frac{P(Z=1, I=1 | X=X_i)}{P(I=1 | X=X_i)} \\
 &= \frac{\tau_{i1} \cdot p_i}{\tau_{i1} \cdot p_i + \tau_{i0} \cdot (1-p_i)} \quad (*) \\
 &= \frac{P(Z=0, I=1 | X=X_i) + P(Z=1, I=1 | X=X_i)}{P(I=1 | Z=0, X=X_i) \cdot P(Z=0 | X=X_i) + P(I=1 | Z=1, X=X_i) \cdot P(Z=1 | X=X_i)} \\
 &= \frac{\tau_{i0} \cdot P(Z=0 | X=X_i) + \tau_{i1} \cdot P(Z=1 | X=X_i)}{\tau_{i0} \cdot (1-p_i) + \tau_{i1} \cdot p_i}
 \end{aligned}$$



Recall tolerance distribution (LN 3-16-17) → logit, probit, CLL → Choice of link function → Note: different links, different interpretation of β (cf., different interpretation of β under different dummy variables)

• Usually not possible to make the choice based on the data alone (i.e., based on fits) because

\hat{y}_i or \hat{p}_i ← e.g. use deviance of the fit ($y_i \leftrightarrow \hat{y}_i$)

① the fits based on the 3 links would be quite similar for moderate p (why?)

usually, do not have large n_i data on x_i with $p_i \approx 0$ or $p_i \approx 1$

different p small or large → LN p.3-5

② for p_i that is close to 0 or 1, a very large amount of n_i would be necessary to distinguish them (why?)

based on $\hat{p}_i = p_i$

Q: But, Var(\hat{p}_i) usually small when $p_i \approx 0$ or $p_i \approx 1$, e.g., if $\hat{p}_i = y_i/n_i$, ≈ 0 when $p_i \approx 0$ or 1

$\text{Var}(\hat{p}_i) = \frac{1}{n_i} p_i(1-p_i)$

If \hat{p}_i is an accurate est. or, why need large n_i ?

Consider \hat{p}_i/p_i . $E(\hat{p}_i/p_i) = 1$

$\text{Var}(\hat{p}_i/p_i) = \frac{1}{n_i} \left(\frac{1-p_i}{p_i} \right)$

large when $p_i \approx 0$

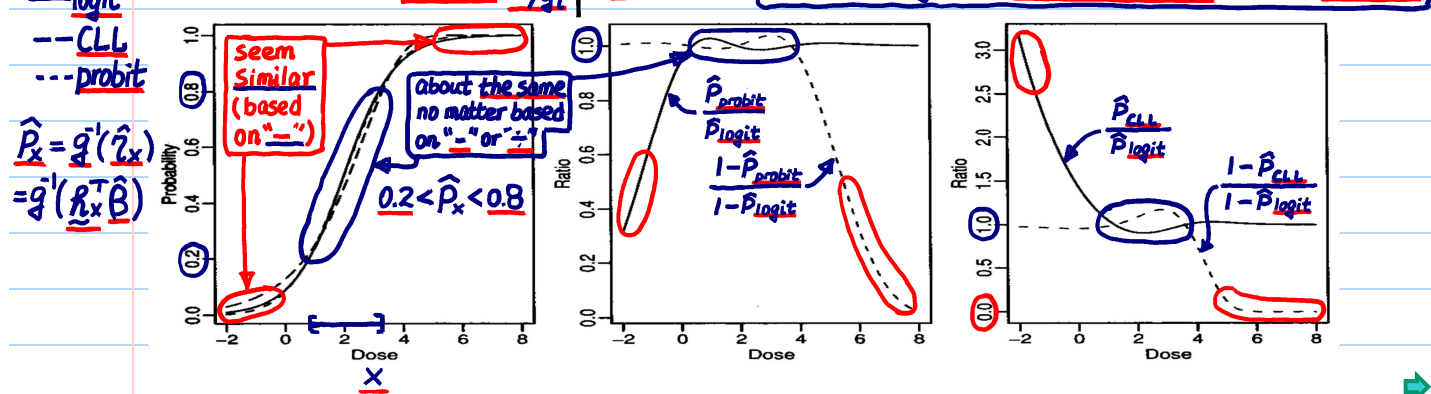
Recall $y_i \leftrightarrow \hat{y}_i$
RSS: $y_i - \hat{y}_i$
deviance: y_i/\hat{p}_i

e.g., $0.01 \leftrightarrow 0.0001$
diff. = 0.0099
ratio = 100

same thing for $(1-\hat{p}_i)/(1-p_i)$ when $p_i \approx 1$

• A lab example

limit



- The choice of link function is usually made based on assumptions derived from physical knowledge or simple convenience e.g.

→ e.g., the idea of tolerance distribution (LNp3-16) can be helpful in some cases.

- Logit link (logistic regression) is popular. Reasons to use logit:

canonical link

➤ theoretically/mathematically simpler due the intractability of Φ

➤ easier to interpret – logit is related to log-odds

➤ retrospective sampling

★ in LNp.3-18 & 19

★ in LNp.3-21, □ in LNp.3-24

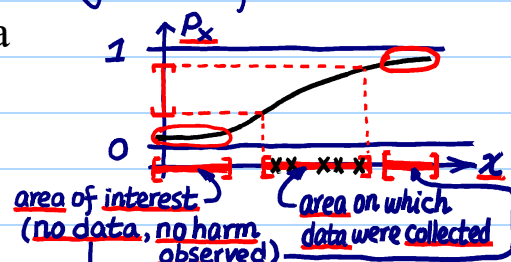
e.g., Normal cdf (probit link) has no explicit form

- In some cases, harmful effects only become apparent at large dosages

more extreme cases (observed data), which make us start to realize something is harmful.

regular condition, connected with rare event

➤ In order to estimate the probability of a harmful effect at low dose, it would be necessary to select an appropriate link function. check graphs in LNp.3-25



➤ However, the data for high dosages will be of little help in selecting an appropriate link

extrapolation: it becomes more difficult (compared to LM) because of the lower (0) & upper (1) bounds in p.

❖ Reading: Faraway (1st ed.), 2.7 other example: highly reliable products

an important objective in regression-type analysis

Prediction and Effective Doses

- Recall: At covariate values $\underline{x}_0 = (x_{01}, x_{02}, \dots, x_{0m})^T$,

$$\underline{p}_{\underline{x}_0} \xrightarrow[g]{g} \underline{\eta}_{\underline{x}_0} = \underline{h}_0^T \underline{\beta}, \text{ where } \underline{h}_0 = (1, h_1(\underline{x}_0), \dots, h_{p-1}(\underline{x}_0))^T$$

e.g., $g = \text{logit} \Rightarrow \log(\frac{\eta_{x_0}}{1-\eta_{x_0}}) = \eta_{x_0}$

denoted by $\hat{\eta}_{x_0}$ & $[\underline{\eta}_{x_0}, \bar{\eta}_{x_0}]$

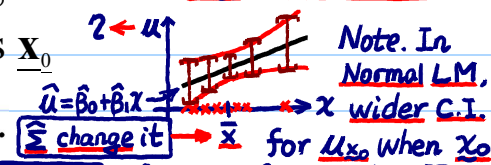
Prediction of the probability p_{x_0} (or odds η_{x_0})

similar to \hat{p}_{x_0} & $[\underline{p}_{x_0}, \bar{p}_{x_0}]$

for parameters in model

or $\underline{\eta}_{x_0}$) of success for given covariate values \underline{x}_0

➤ For the MLE $\hat{\beta}$, we have $\hat{\beta} \stackrel{a}{\sim} N(\underline{\beta}, \underline{\Sigma})$.



consistent

Denote the estimate of $\underline{\Sigma}$ by $\hat{\underline{\Sigma}}$.

$= (\underline{X}^T \hat{\underline{W}} \underline{X})^{-1}$ from IRWLS

lik $\hat{\eta}_{x_0}$ in LM

➤ Predict $\underline{\eta}_{x_0}$ by $\underline{h}_0^T \hat{\underline{\beta}}$, denoted by $\hat{\underline{\eta}}_{x_0}$.

➤ Predict \underline{p}_{x_0} by $g^{-1}(\hat{\underline{\eta}}_{x_0})$, denoted by $\hat{\underline{p}}_{x_0}$.

➤ 100(1- α)% confidence interval for $\underline{\eta}_{x_0}$:

Note. In Normal LM, wider C.I. for \underline{u}_{x_0} when \underline{x}_0 is away from center \bar{x} (i.e., has larger leverage), same for binomial GLM? Ans. Yes for $\underline{\eta}_{x_0}$. No for \underline{p}_{x_0} . But, $\log(\frac{\eta_{x_0}}{1-\eta_{x_0}})$ would get larger as $\underline{\eta}_{x_0} - \bar{\eta}_{x_0}$ can be regarded as a log-odds-ratio

se($\hat{\eta}_{x_0}$)

$$\hat{\underline{\eta}}_{x_0} \pm z(\alpha/2) \times \sqrt{\underline{h}_0^T \hat{\underline{\Sigma}} \underline{h}_0} \equiv [\underline{\eta}_{x_0}, \bar{\eta}_{x_0}] \Leftarrow P\left(\left|\frac{\hat{\underline{\eta}}_{x_0} - \underline{\eta}_{x_0}}{\sqrt{\underline{h}_0^T \hat{\underline{\Sigma}} \underline{h}_0}}\right| < z\left(\frac{\alpha}{2}\right)\right) = 1 - \alpha$$

because $\hat{\underline{\eta}}_{x_0} = \underline{h}_0^T \hat{\underline{\beta}} \stackrel{a}{\sim} N(\underline{h}_0^T \underline{\beta} = \underline{\eta}_{x_0}, \underline{h}_0^T \underline{\Sigma} \underline{h}_0)$.

➤ 100(1- α)% confidence interval for \underline{p}_{x_0} :

$$[g^{-1}(\underline{\eta}_{x_0}), g^{-1}(\bar{\eta}_{x_0})] \equiv [\underline{p}_{x_0}, \bar{p}_{x_0}]$$

$P(\underline{\eta}_{x_0} \in [\underline{\eta}_{x_0}, \bar{\eta}_{x_0}])$ random \downarrow $= P(\underline{p}_{x_0} = g^{-1}(\underline{\eta}_{x_0}) \in [g^{-1}(\underline{\eta}_{x_0}), g^{-1}(\bar{\eta}_{x_0})])$ monotone