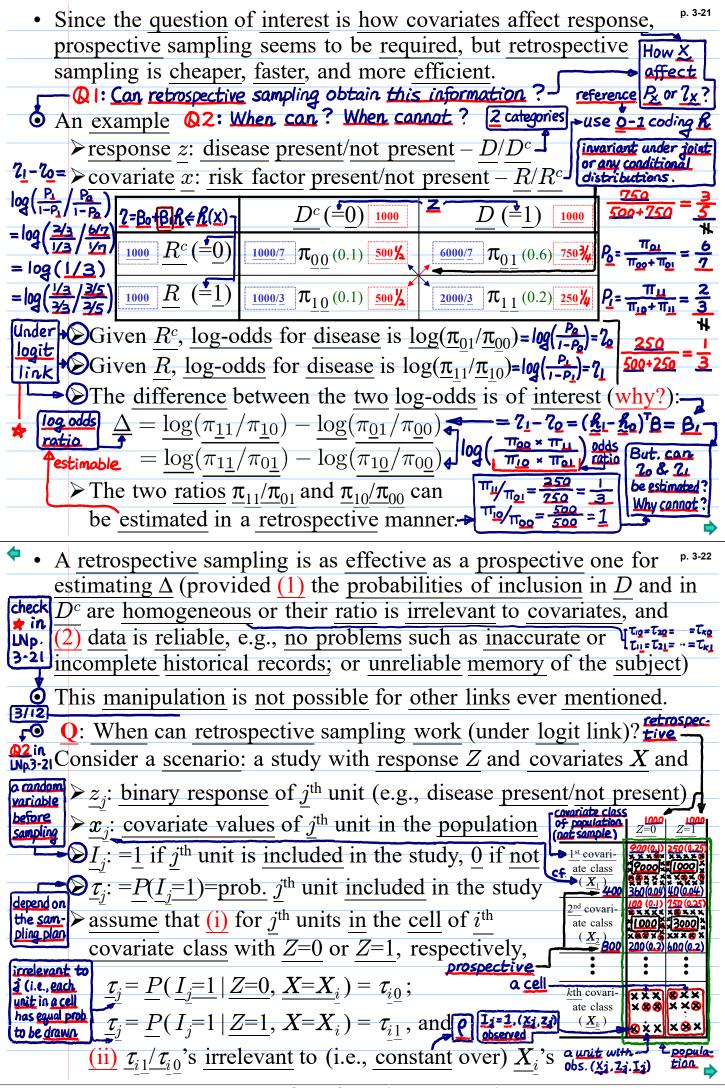
NTHU STAT 5230, 2025

Q1: Which one should be used to characterize the effect of p. 3-19
• Odds ratio and relative risk - Changing from X1 to X2? (always) C - C - Changing from X1 to X2? (always) C - C - C - C - C - C - C - C - C - C
$r_{atia} > Suppose the probability of successes at \underline{x}_1 (say, in the presence$
(cf., in LM. of some condition) is \underline{p}_1 and \underline{p}_2 at \underline{x}_2 (say, in its <u>absence</u>)
rence μ_1 μ_2 $Relative risk = p_1/p_2$ difference $\left[P_1 \right] = \log(\frac{p_1}{1-p_1}) - \log(\frac{p_2}{1-p_2}) \right] = 11 - 12 = \pi_1 B - \pi_2 B$
$\frac{P_2}{(1-P_2)} = \frac{P_1}{P_2} = \frac{P_1}{P_2} = \frac{P_1}{P_2} = \frac{P_1}{P_2} = \frac{P_1}{P_2} = \frac{P_2}{P_2} = \frac{P_1}{P_2} = \frac{P_1}{P_2$
$= (\underbrace{I-P_2}_{I-P_1})(\underbrace{P_1}_{P_2}) = \underbrace{\text{Log odds ratio}}_{\text{log}} = \underbrace{\log(o_1/o_2)}_{\text{(for 1 x)}} $
For rare outcomes, relative risk \approx odds ratio, but for larger $= 41-42$
$P_{i} \approx 0$ probabilities they may be substantial differences model $\eta_{i-R} = 1 + P_{i}(x_{i})$
$-\underline{r_2}$
do not medal I here is some debate over which is the more intuitive = $R_i^T B_i^T $
$\begin{array}{c} \underline{P_{\approx}P_{\ast}} \\ \hline \underline{P_{\ast}P_{\ast}} \\ \hline \underline{P_{\ast}} \\ \hline \underline{P_{\ast} \\ \hline \underline{P_{\ast}} $
be large Prospective and retrospective sampling - issue in survey
InGLM, convariates (fixed), InGLM, response (random), $\in \{0,1\}, \sim Bernoulli(P_i) \neq true? sampling, not in$
• Data: $(x_{j1}, x_{j2}, \dots, x_{jm}, z_j), j = 1, 2, \dots, \underline{K},$
$(x_{i1}, x_{i2}, \dots, x_{im}, y_i), i = 1, 2, \dots, \underline{k}$. It of units
Q: how is the data collected? $\sim binomial(n_i, p_i) \Leftrightarrow true?$
\underline{Q} : <u>how</u> is the <u>data collected</u> ? ~ <u>binomial</u> (<u>Ni.pi</u>) \triangleleft <u>true?</u>
• Sampling methods: choose a sub-population. different Recall In DOE p. 3-20
• <u>Sampling</u> methods: then draw a sample from it <u>offerent</u> Recall in DOC
Prospective sampling: the covariates x are fixed and then the
Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called <i>cohort</i> study.
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study.
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data analysis z:like An infant respiratory disease example:
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data analysis Z: like An infant respiratory disease example:
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data called in data analysis Z:like An infant respiratory disease example: Z (z=1, present; z=0, not covariate x and then the present) Select a sample of newborn boy/girl whose parents had present
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariates x are observed, called case-control study. In data chen the covariate the
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data callection X: like An infant respiratory disease example: Select a sample of newborn boy/girl whose parents had present; z=0, not chosen a particular method of feeding, and then monitor there first year. What are the whether disease present or not present for their first year.
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data analysis Z: like An infant respiratory disease example: Select a sample of newborn boy/girl whose parents had present; z=0, not chosen a particular method of feeding, and then monitor whether disease present or not present for their first year. What are the sub-populations? Find infants coming to a doctor with a respiratory (fixed)
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, then the covariates x are observed, called case-control study. In data callestian, callestian, then the covariates x are observed, called case-control study. In data callestian, callestian, callestian, then the covariates x are observed, called case-control study. In data callestian, ca
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. An infant respiratory disease example: Select a sample of newborn boy/girl whose parents had chosen a particular method of feeding, and then monitor whether disease present or not present for their first year. What are the sub-populations? Find infants coming to a doctor with a respiratory disease in the first year and then record their sex and method of feeding; also obtain a sample of respiratory (random)
Prospective sampling: the <u>covariates x</u> are fixed and <u>then</u> the response z (or y) is <u>observed</u> , called <u>cohort</u> study. Retrospective sampling: the response z (but not y) is fixed and In data then the <u>covariates x</u> are observed, called <u>case-control</u> study. Leg. $z = 1$ eg. $z = 0$ Covariate z : like An infant <u>respiratory disease</u> example: z : like Select a <u>sample</u> of newborn <u>boy/girl</u> whose parents had <u>present</u> chosen a <u>particular</u> <u>method of feeding</u> , and <u>then</u> monitor whether disease present or not present for their first year. z = 0 (random) z = 0 (r
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. An infant respiratory disease example: 2 (z=1, present; z=0, not covariate z) (z=1, present; z=0, not covaria
Prospective sampling: the <u>covariates x</u> are <u>fixed</u> and <u>then</u> the response z (or y) is <u>observed</u> , called <u>cohort</u> study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called <u>case-control</u> study. Z:like An <u>infant</u> respiratory disease example: Covariate X:like Select a <u>sample</u> of newborn <u>boy/girl</u> whose parents had chosen a particular <u>method of feeding</u> , and <u>then</u> monitor whether disease present or not present for their first year. What are the sub-populations? What are the sub-populations? Control Case Z = 1 Case Ca
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. In data analysis In d
 Prospective sampling: the covariates x are fixed and then the response z (or y) is observed, called cohort study. Retrospective sampling: the response z (but not y) is fixed and then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariates x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x are observed, called case-control study. Z then the covariate x and then monitor the covariate values are the study to be irrelevant to the covariate values are the check INp3-22-24

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

NTHU STAT 5230, 2025



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