

SIMPSON'S PARADOX

The value of a measure of association btw 2 vars, may be identical w/in the levels of a third variable but can take an entirely \neq value when this 3rd variable is disaggregated and the associated is computed from pooled data.

POLYTOMOUS EXPOSURE AND OUTCOME guarda slide

• n° of parameters in saturated models $(IJ) - 1$

• n° of parameters if we assume independence $(I-1) + (J-1)$

we test the hypothesis of independence w/ χ^2 test or chi squared

2 approaches: - Test indep. : btw exposure and outcome by comparing observ vs exp freq

DRIDIT SCORES = compares 1 or + gr to a stand. gr

RIDIT = ① Frequency A

RIVEDERE

② $\frac{1}{2}$ Frequency

③ cumulative frequency starting from zero

④ 2 + 3

$\bar{R} = \frac{\text{sum of prod}}{N_B}$ $\left\{ \begin{array}{l} > 0.5 \text{ more likely} \\ < 0.5 \text{ less likely} \end{array} \right.$

⑤ $4 / N_A$

$\frac{\bar{R}}{1-\bar{R}} = \text{ODDS}$

GENERALIZED LINEAR MODELS

linear
- Binomial
- Poisson

① LOGISTIC MODEL

$$\log\left(\frac{\pi}{1-\pi}\right) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$\beta_0 = \log \text{ODDS OF SUCCESS}$

$\text{EXP}(\beta_0) = \text{ODDS OF SUCCESS}$

$\beta_1 = \text{slope, is the change in the log odds of success for } +1 \text{ unit increase in } X \rightarrow \log \text{ODDS RATIO}$

$\text{EXP } \beta_1 = \text{ODDS RATIO}$

② LINEAR MODEL

$$\mu = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

$\beta_0 = \text{average for baseline, or expected } m$

$\beta_1 = \neq \text{ btw 2 averages change in mean ratio btw 2 expected } m \text{ w/ } +1 \text{ unit increase}$

MEASURING PREDICTIVE POWER

AUC = area under the ROC curve that describes (y) sensitivity and 1-specificity (x)
Also the proportion of concordant pair

$$\begin{array}{l} \frac{C}{C+T+D} \\ C_{Stat} \frac{C+OST}{C+T+D} \\ \text{gini} \\ \text{Coeff} \end{array} d = \frac{C-D}{C+T+D} \left. \begin{array}{l} C: y_1=1 \quad y_2=0 \\ D: y_1=1 \quad y_2=0 \\ T: y_1=1 \quad y_2=0 \end{array} \right\} \begin{array}{l} \pi_1 > \pi_2 \\ \pi_1 < \pi_2 \\ \pi_1 = \pi_2 \end{array}$$

TEST OF AGREEMENT (For matched data)

① INTER OBSERVERS = In the diagonal we observe perfect agreement, outside = disagreement

sum of pr in the diagonal = 1

COHEN'S KAPPA = $\frac{\text{sum of the pr in the diag} - \text{sum pr under indep}}{1 - \text{sum pr independence}}$

$H_0: K=0$ pr in the diag = pr ind
 $K=1$ num and den are the same, sum = 1

TEST OF HOMOGENEITY

① BD TEST: TEST OF COLLAPSIBILITY, if all stratum specific have same OR or RR

$H_0: \varphi_1 = \varphi_2 = \dots = \varphi_k$ if so, they can be reported as weighted averages
if we reject BD test, the stratum specific estimates must be reported separately

② CMH TEST: used to test specific hypothesis

$H_0: \varphi_1 = \varphi_2 = \dots = \varphi_k = 1$

③ MC NEMAR TEST: For matched data, we test marginal homogeneity

$H_0: \pi_{1+} = \pi_{+1}$

(glm) $\rightarrow T = \frac{m_{21} - m_{12}}{\sqrt{m_{21} + m_{12}}}$

(clogit) $\rightarrow T = \frac{m_{121}}{m_{212}}$

if $T > p_{val}$ we reject H_0 and accept $H_1: \pi_{1+} \neq \pi_{+1}$

FACTORIAL NUMBERS

$$Pr(y) = \frac{m!}{y!(m-y)!} \cdot p^y (1-p)^{m-y}$$

STATISTICS FOR HEALTH

• **Linear relationship** — COVARIANCE $\text{COV}(X, Y) = E\{(X - E(X))(Y - E(Y))\}$
 ↳ if = 0 no linear rel

CORRELATION COEFFICIENT -1 0 +1
 neg lin rel no corr pos lin rel
 ↑ ↓ ↓ ↑ ↑ ↑ ↓ ↓

• **ESTIMATION - SAMPLE VARIANCE** $S^2 = \frac{\sum (y_i - \bar{y})^2}{m}$ correct $S^2 = \frac{\sum (y_i - \bar{y})^2}{m-1}$

↳ **STANDARD ERROR** SE: $\sqrt{\frac{\sigma^2}{m}} = \frac{\sigma}{\sqrt{m}}$

↳ **MAXIMUM LIKELY ESTIMATOR** $L: (p) = [p^{y_1} (1-p)^{1-y_1}] \cdot \dots \cdot [p^{y_m} (1-p)^{1-y_m}]$

↳ **POINT ESTIMATION - SAMPLE MEAN**: $\bar{\mu} = \frac{\sum \mu_i}{m}$

↳ **SAMPLE VARIANCE**: σ^2 or $S^2 = \frac{\sum (y_i - \bar{y})^2}{m-1}$

↳ **INTERVAL ESTIMATION** $CI = \bar{\mu} \pm SE \cdot z_{1-\frac{\alpha}{2}}$

TEST ON THE MEANS

we have 1 population w/ mean μ , $H_0 = \mu = \mu_0$ and $H_1: \mu \neq \mu_0$
 H_1 is a two sided hypothesis bc μ can be both $>$ or $<$ of μ_0

if we don't know σ^2 we can use the estimate of $S^2 = \frac{\sum (y_i - \bar{y})^2}{m-1}$

① **TWO SIDED TEST ON THE MEAN**: $\frac{|\bar{y} - \mu_0|}{S/\sqrt{m}} \geq t_{1-\frac{\alpha}{2}, (m-1)}$
 ↳ critical value
 if so we REJECT H_0

② **ONE SIDED TEST ON THE MEAN**: $\frac{\bar{y} - \mu_0}{S/\sqrt{m}} > t_{1-\alpha, (m-1)}$

if we have 2 pop or 2 independent samples y_1 w/ σ_1^2 / y_2 w/ σ_2^2
 we want to test the equality of the 2 means
 $H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$

we assume that $\sigma_1^2 = \sigma_2^2$

③ **TEST OF EQUALITY IN 2 MEANS TWO SIDED**: $\left| \frac{\bar{y}_1 - \bar{y}_2}{Sp \cdot \sqrt{\frac{1}{m_1} + \frac{1}{m_2}}} \right| \geq t_{1-\frac{\alpha}{2}, (m_1 + m_2 - 2)}$
 $Sp = \frac{(m_1 - 1)S_1^2 + (m_2 - 1)S_2^2}{m_1 + m_2 - 2}$ — pooled var

if we assume $\sigma_1^2 \neq \sigma_2^2$ $\left| \frac{\bar{y}_1 - \bar{y}_2}{\bar{S}} \right| \geq t_{1-\frac{\alpha}{2}, (m)}$
 $\bar{S} = \frac{S_1^2}{m_1} + \frac{S_2^2}{m_2}$

STUDY DESIGN

- ① **COHORT STUDIES**: Healthy population divided in E^+ and E^-
 calculate the OCCURRENCE of disease to see if it's
 the same btw E^+ and E^-
 cons: time consuming, exp, no good for rare dis
 Longitudinal
- ② **CASE-CONTROL STUDIES**: unhealthy (cases) vs healthy (controls) ppl
 they are already ill
 cons: can not estimate prevalence or independence
 pro: good for rare dis, multiple exposure
 RETROSPECTIVE
- ③ **CROSS-SECTIONAL STUDIES**: target population, strict criteria
 pop divided in: D^+ D^- E^+ E^-
 cons: only prevalence, no good for rare dis
 pro: not expensive

RISK

$$RR = \frac{a/m_1}{c/m_0} = \frac{\lambda_1}{\lambda_0}$$

$$\text{EXCESS RISK} = \lambda_1 - \lambda_0$$

$$\text{Attributable risk} = \frac{\lambda_1 - \lambda_0}{\lambda_1}$$

$$\text{Crude risk} : \pi = \frac{m \text{ cases}}{m \text{ pop}}$$

$$\text{Incidence rate} : \lambda = \frac{m \text{ cases}}{p \cdot \text{years}}$$

$$\text{RISK} = R : 1 - \exp(-\lambda t)$$

$$\text{ODDS RATIO} : OR = \frac{ad}{bc} \quad OR = RR \left(\frac{1 - \lambda_1}{1 - \lambda_0} \right)$$

$$\text{SENSITIVITY} = \frac{TP}{TP + FN}$$

$$\text{SPECIFICITY} = \frac{TN}{TN + FP}$$

$$\text{PPV} = \frac{TP}{TP + FP}$$

$$\text{NPV} = \frac{TN}{TN + FN}$$

CI w/ WOOLF'S METHOD

$$RR : \log(RR) \pm z \cdot \sqrt{\frac{1}{a} - \frac{1}{a+b} + \frac{1}{c} - \frac{1}{c+d}}$$

$$OR : \log(OR) \pm z \cdot \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$$

then exp
 LB and UB

CONFOUNDING

When we are talking about the association btw outcome and exposure, we have to take into account the impact of other variables

- ① Is a risk factor
- ② Is outside of the causal chain w/ exposure
- ③ The association must occur in the absence of exposure

③ POISSON

$$\log(\lambda) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

$\beta_0 = \log$ of the intercept

$\exp(\beta_0) = \lambda$ is the mean m of y when all $x = 0$
expected value

$\exp(\beta_i) =$ RATE RATIO of the expected m of event

OFFSET = log of the population, to make comparison of rates
 \rightarrow transform counts in rates

QUASI-POISSON = if residual deviance is very large compared to n means that the model

Needs additional var, and if we don't have those we use quasi poisson, that generates a new coeff estimated, has a $> SE$ and larger CI

$$\frac{\text{residual dev / new coeff est}}{\text{residual df}}$$

ZERO INFLATED MODEL = zero inflation is a common issue w/ Poisson
 It's an excess of zero counts that can alter outcome and results ($\downarrow \lambda$)

The zero inflated model removes the bad effect of those zeros

\rightarrow or effect modification

INTERACTIONS = is an additional variable that can change the relation btw exposure and outcome

$$\log \frac{\pi}{1-\pi} : \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \gamma_{12} x_1 x_2$$

MEASURING GOF \rightarrow goodness of fit

④ LRT : For nested models

We start from model 1: $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$

We add to mod 1 some : $y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \beta_{k+1} x_{k+1} \dots \beta_{k+h} x_{k+h}$
 \xrightarrow{n}

Our $H_0 : \beta_{k+1} = \beta_{k+2} = \dots = \beta_{k+h} = 0$

$D = -2 \{ \log L(\text{smaller}) - \log L(\text{larger}) \}$ if $D > \chi^2$ or p-value we reject H_0

② RESIDUAL DEV
DEG OF FREED

The ideal value is 1
 > 1 means that the model can't explain all

res dev = is the variability that the model can't compute

③ AIC = if we don't have nested models

$$AIC = \underbrace{2 \log L}_{\text{GOF}} + \underbrace{2m}_{\text{complexity}} \quad \text{THE SMALLER, THE BEST}$$