

Week 4

Tuesday, 3 August 2021 22:23

HANDLING BINARY DATA

QUESTION: YES OR NO?

- DOES A PERSON HAVE A CERTAIN DISEASE?
- DOES THE SUBJECT HAVE A SPECIAL CHARACTERISTIC?
- DOES SUBJECT ENGAGE IN A CERTAIN BEHAVIOUR?

EASIER TO SUMMARIZE DATA

CONTINUOUS DATA (SPREAD, CENTER, LOCATION)

BINARY DATA (SAMPLE PROPORTION)

WE JUST HAVE 1 SUMMARY STATISTIC

THERE ARE DIFFERENT WAY TO COMPARE DATA

BINARY DATA (DEFINITION & SUMMARIZATION)

$$\hat{p} = \frac{\# \text{ POSITIVE OUTCOMES}}{\# \text{ TOTAL OUTCOMES}}$$

"P-HAT"

PROPORTION

$\hat{p} \neq p$
↓ ↓
SAMPLE ESTIMATE TRUE VALUE

VERY SIMILAR TO A MEAN

$$X = \begin{cases} 0 & \text{NEGATIVE OUTCOME} \\ 1 & \text{POSITIVE OUTCOME} \end{cases}$$

$$\hat{p} = \frac{1}{N} \sum_{i=1}^N X_i$$

$$S = \sqrt{N \hat{p} (1 - \hat{p})}$$

STANDARD DEVIATION

- NOT IMPORTANT FOR UNDERSTANDING THE DISTRIBUTION
- USEFUL WHEN COMPARING 2 DIFFERENT DATA

NORMAL DISTRIBUTION

(5) S DOES NOT DEPEND ON \bar{X} (4)

BINOMIAL DISTRIBUTION

(5) S DOES DEPEND ON \bar{X} (4)
(5) (\hat{p})

PERCENTILES WILL EITHER BE 0 OR 1

EXAMPLE 1000 PATIENTS, 208 HIV POSITIVE

$$\hat{p} = \frac{208}{1000} = 0.208 = 20.8\%$$

(NOT VERY
USEFUL)

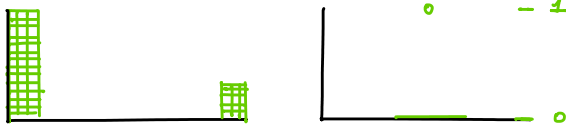
0 0 0 0 0 ... 0 1 1 1 1 ... 1
792 208

$$\begin{aligned} q_{0.25} = Q_1 &= 0 & \chi_{(250)} &= 0 \\ q_{0.5} = Q_2 = M &= 0 & \chi_{(500)} &= 0 \\ q_{0.75} = Q_3 &= 0 & \chi_{(750)} &= 0 \end{aligned}$$

$$\begin{aligned} q_{0.753} &= 1 & \chi_{(753)} &= 1 \\ q_{0.95} &= 1 & \chi_{(950)} &= 1 \end{aligned}$$

\hat{p} CAN BE USE TO KNOW BOTH THE
- CENTER
- VARIABILITY
- PERCENTILES

VISUAL DISPLAY ? NOT USEFUL



UNLIKE CONTINUOUS DATA :

IF WE HAVE \hat{p} WE HAVE THE ENTIRE STORY

PERCENTAGE,
PROPORTION,
PROBABILITY,
RISK

COMPARING BINARY OUTCOME
BETWEEN TWO POPULATIONS

HAVING \hat{p}_1, \hat{p}_2

DIFFERENCE IN PROPORTIONS

$$\hat{p}_1 - \hat{p}_2 = D \quad RD$$

RISK DIFFERENCE
ATTRIBUTABLE RISK) \rightarrow EPIDEMIOLOGY

D% GREATER RESPONSE TO THERAPY
IN GROUP 1 COMPARED TO GROUP 2

(ABSOLUTE RISK) EVEN IF IT'S A GOOD THING

RELATIVE RISK PROPORTIONS RATIO

$$r = \frac{\hat{p}_1}{\hat{p}_2}$$

RR

GROUP 1 HAS r TIMES MORE RISK THAN GROUP 2

$$r = 1.52 \rightarrow 52\% \text{ GREATER RISK}$$

$$OR = \frac{p_1 - p_2}{p_2}$$

RISK DIFFERENCE

CAN BE INTERPRETED AS
IMPACT (ASSUMING CAUSATION)
ON A FIXED # OF PERSONS

THEY WILL ALWAYS
AGREE ON DIRECTION
ASSOCIATIONS

RELATIVE RISK

IMPACT (ASSUMING CAUSATION)
AT THE INDIVIDUAL LEVEL

RELATIVE RISK TENDS TO BE BIGGER
MEDIA OFTEN USES THEM

ODDS OF AN OUTCOME

$$\hat{ODDS} = \frac{\hat{p}}{1 - \hat{p}}$$

$$\frac{P(\text{EVENT OCCURS})}{P(\text{EVENT DOES NOT OCCUR})}$$

HIGHER ODDS
HIGHER PROPORTIONS

$$\text{AS } \hat{p} \rightarrow 1 \quad \hat{ODDS} \rightarrow \infty$$

ODDS RATIO

$$OR = \frac{\hat{ODDS}_1}{\hat{ODDS}_2} = \frac{\frac{p_1}{1 - p_1}}{\frac{p_2}{1 - p_2}}$$

LESS INTUITIVE
USED IN CASE/CONTROL STUDIES

CANNOT DIRECTLY ASSESS THE RISK OF INTEREST

NOT A DIRECT COMPARISON RISK BUT A FUNCTION OF THE RISKS CALL ODDS

ALL 3 RISKS AGREE ON DIRECTION BUT NOT ALWAYS ON MAGNITUDE

$$\begin{aligned}\hat{RD} &= p_1 - p_2 \\ \hat{RL} &= \frac{p_1}{p_2} \\ \hat{OR} &= \frac{p_1 / (1 - p_1)}{p_2 / (1 - p_2)}\end{aligned}$$

EXAMPLE $\hat{p}_1 = 0.07$ $\hat{p}_2 = 0.22$

	$\hat{RD} = p_1 - p_2 = -0.15$	(-15%)	SAME MAGNITUDE CHANGES SIGN
INV	$\hat{RD} = p_2 - p_1 = +0.15$	$(+15\%)$	
	$\hat{RL} = \frac{p_1}{p_2} = 0.32$	(-68%)	VERY DIFFERENT
INV	$\hat{RL} = \frac{p_2}{p_1} = 3.1$	$(+210\%)$	
	$\hat{OR} = \frac{p_1 / (1 - p_1)}{p_2 / (1 - p_2)} = 0.27$	$(-73\% \text{ IN ODDS})$	
INV	$\hat{OR} = \frac{p_2 / (1 - p_2)}{p_1 / (1 - p_1)} = 3.7$	$(+270\% \text{ IN ODDS})$	

WHY DIFFERENT MAGNITUDE IF
DIRECTION IS REVERSED ?

SCALES OF RATIO IS NOT SYMMETRIC

$\ln X$ EQUALIZES THE VALUES (REGARDLESS OF DIRECTION)

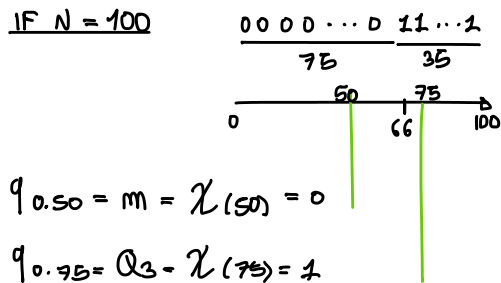
$$\ln\left(\frac{x}{y}\right) = \ln x - \ln y$$

$$\ln(\hat{RL}) = \ln\left(\frac{p_1}{p_2}\right) = \ln p_1 - \ln p_2 = K$$

$$\ln(\hat{RL}_{INV}) = \ln\left(\frac{p_2}{p_1}\right) = \ln p_2 - \ln p_1 = -K$$

EXERCISES

- $\hat{p} = \bar{x}$ SINCE $\hat{p} = \frac{1}{N} \sum_{i=1}^N x_i$
- IF $\hat{p} = 0.35$ 35% YES 65% NO



- $\hat{RL} = 0.78$ DRUG
PLACEBO

- 22% LOWER RELATIVE RISK IF YOU TAKE THE DRUG

- $\hat{RL} = 3$ 3x MORE RISK FOR WHO HAS THAT GENETIC MUTATION
- ABSOLUTE RISK DIFFERENCE $\hat{RD} = p_{\text{GENETIC MUTATION}} - p_{\text{NO MUTATION}} = 0.002$

(ASSUMING CAUSALITY) IF $N = 10,000$
WE WILL HAVE 20 MORE CANCER IF THE WHOLE
POPULATION HAS THE MUTATION

- IF $\hat{ODS} = 1$, $p = ?$

$$ODS(p) = \frac{p}{1-p} = 1 \quad 1-p = p$$

$$2p = 1 \quad \hat{p} = 0.5 = 50\%$$

- STUDY

$$X \text{ BLOOD LEAD LEVELS (MALES)} \sim N(1.97, 1.61^2)$$

$$Y \text{ BLOOD LEAD LEVELS (FEMALES)} \sim N(1.27, 1.16^2)$$

$$N_{Y1} = 100 \quad S_N \leq S_{100} \quad ? \quad \text{CANNOT KNOW}$$

$$N_{Y2} = 2,128$$

$$\delta = \bar{y} - \bar{x} = -0.7 \text{ mg/dL}$$

"HIGH BLOOD LEAD LEVELS: $X, Y > 3.10$ "

$$B \sim \begin{cases} 0 & \text{if } BL < 3.10 \\ 1 & \text{if } BL > 3.10 \end{cases}$$

$$\hat{p}_{F, \text{HIGH BL}} = 5\% \quad \hat{RR} = \frac{\text{FEMALES}}{\text{MALES}} \quad \hat{p}_{M, \text{HIGH BL}} ?$$

$$\hat{p}_{M, \text{HIGH BL}} = 3 \hat{p}_{F, \text{HIGH BL}} = 15\%$$