

# Chapter 1

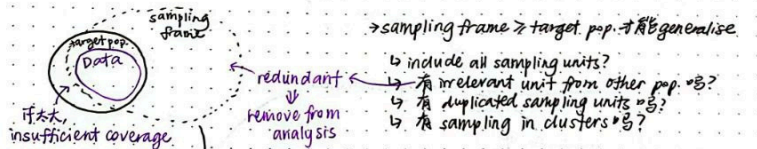
population = grp of individual / subject we want know smtg abt

population parameter. 想研究的指标

sample

pop. 呈 proportion to 参加 study, sample size ↑, 误差 ↓, random error ↓  
③ 用大 sample size

estimate



在做 sampling frame 时, 可能有 bias

- selection bias: ① 选 sampling frame 时, imperfect selection frame. ② 选 probability sampling, know & randomised self choice (计算出来选中概率).
- non-response bias: ③ minimise non-response non-response bias (non-disclosure of info, will happen naturally 无它 method) \* non-response bias = random error

non-probability sampling

- convenience s.: proximity & availability → selection & non-response bias
- volunteer s.: researchers seek volunteers 但有愿意参与的人更能参加

probability sampling

- simple random s.: all units equal chance → select randomly two replacement
- stratified random s.: 将 similar characteristic 分组, 再在每个组 simple random (按比例决定每个组选多少人)
- cluster sampling: 把 unit 分成 similar cluster, 再选一个 cluster 做代表
- systematic s.: 在 n 个 unit 里, 选间隔 k, 再从 1 到 k 这个范围选点, 开始! ③ 选点, 开始! ③ 选点, 开始!

Fallacies

- ecological fallacy: 用 aggregate (group) 的 correlation 判断 individual 的
- atomistic fallacy: 用 individual 的 correlation 判断 group 的

ecological correlation: consider grp characteristic → compute grp average 之后把数据代入计算 r

aggregate points: 看数据是点还是一个组 / w slicing (≠ ecological) 还是真的有几个 aggregate point

→ 可能和 individual 反方向, 也可能 overstate assoc. of individuals

Studying non-linear 2 variables

\* r 不准, 看 scatter plot 的形状才是关键!

log →  $y = cb^x$  → 两边加 ln

$ln y = ln(cb^x)$   
 $ln y = ln c + x ln b$

Y X m

① compute (t, lny)  
② plot lny against t  
③ find linear regression line  
④ 有 line 就有  $m \times x$ , 可以反推出  $ln c$  &  $ln b$  里  $c$  &  $b$  是什么 ④

\* r ≠ gradient of regression line  
 $m = \frac{S_y}{S_x} r$

Linear Regression

obtaining regression line method of least squares 找让  $\sum e_i^2$  最小的线

一定 pass through ave ( $\bar{x}, \bar{y}$ )

只能看 X 判断 Y (因为线是 minimise  $\sum b_i e_i$ , 要判 Y 判断 X 就要找  $e_i = Y - \hat{Y}$  最小)

超出已有的 X 的范围去做 prediction 会不准

$Y = mx + b$  → 可以 predict data

基一条直线进 graph

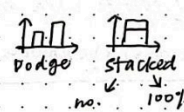
observed value of Y for  $X_i$

predicted value of Y for  $X_i$

$e_i$  = residual of i-th observation (error of i-th)

## Chapter 2: Categorical Variable

Analyse 1 variable → table & plots → rate



Analyse 2 variables → 2x2 table

Outcome treatment	Success	Fail	Row T.
X			
Y			
Column T.			

Rates

- marginal r. =  $\frac{\text{category no.}}{\text{grand total}}$  → Rate (success)
- joint r. =  $\frac{\text{交集 no.}}{\text{grand total}}$  → Rate (Y & failure)
- conditional r. (among / given) =  $\frac{\text{交集 no.}}{\text{subgroup after 1}}$

Rules

Symmetry Rule:  $R(A|B) > R(A|NB) \Leftrightarrow R(B|A) > R(B|NA)$   
 $R(A|B) < R(A|NB) \Leftrightarrow R(B|A) < R(B|NA)$

Basic Rule of R: overall R(A) lies btw  $R(A|B)$  &  $R(A|NB)$

Subgroup rate & overall rate 的关联

Association ≠ causation

association present

$R(A|B) > R(A|NB)$  → +ve btw A & B  
presence of A is stronger when B is present than when B is absent  
 $R(A|B) > R(A|NB)$   
 $R(B|A) > R(B|NA)$   
 $R(NA|NB) > R(NA|B)$   
 $R(NB|NA) > R(NB|A)$

$R(A|B) < R(A|NB)$  → -ve btw A & B  
他可以说 +ve btw A & NB

association absent

$R(A|B) = R(A|NB)$

confounders → third variable assoc. to both dependent & independent variable. check confounder → 看和每个 variable 有没有 assoc.

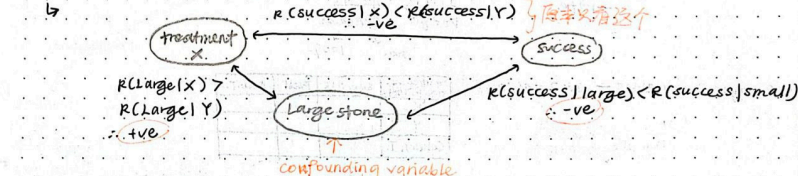
- draw table of 2 variables we 调查
- calc. conditional rates btw variables
- compare assoc. btw variables
- 重复 1 和 2 都写出 assoc. → 第 3 个 variable 和 1 和 2 有联系! = confounder

在 study 时要考虑 background info

\* Simpson → confounder  
没 Simpson ≠ 没 confounder

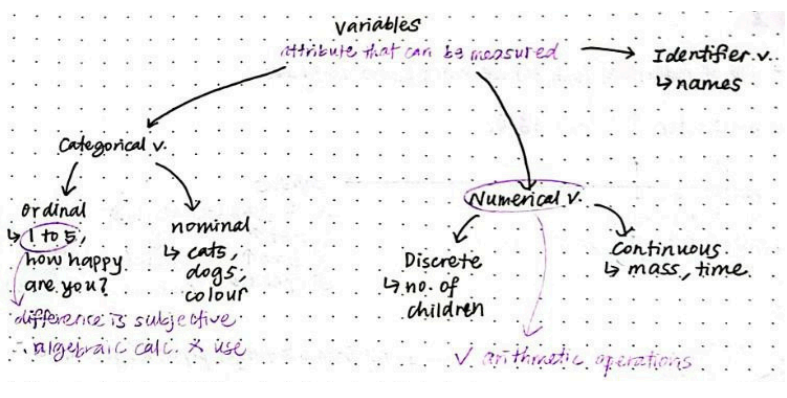
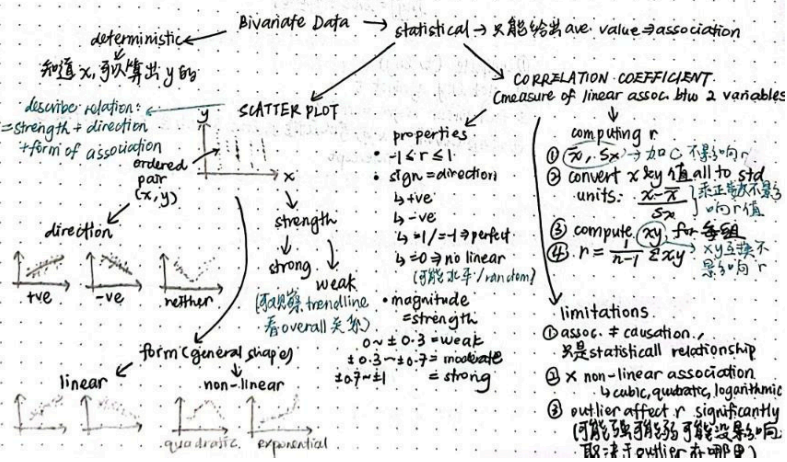
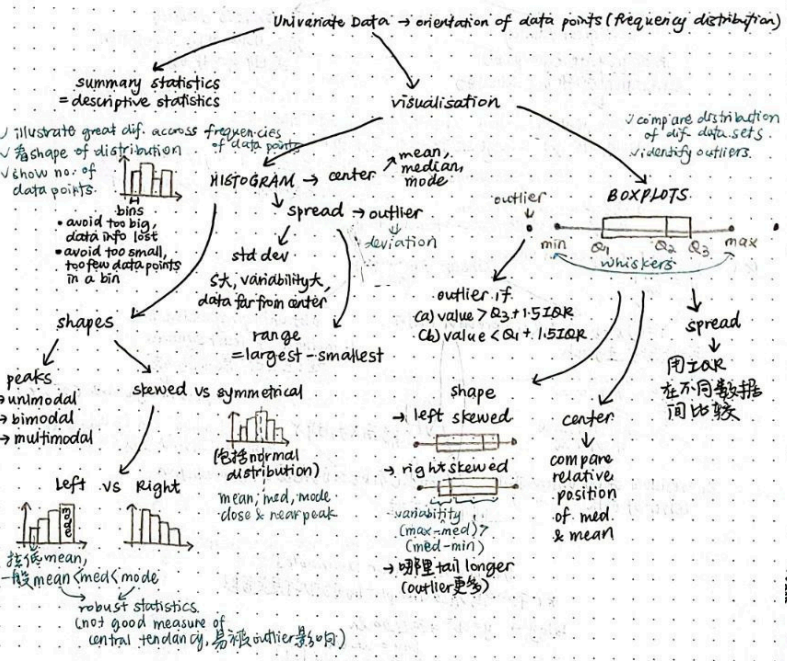
B: is positively associated with y  
A: Rate (y|x) vs. Rate (y|mx)

Simpson's Paradox → trend appears in majority of several grps of data but disappears / reverses when groups combined





# Chapter 3: Numerical Data



# Chapter 4: Statistical Inference

