**● QR Framework ( Quantitative Reasoning)**  — offers precision

1) Frame
2) Specify
3) Collect
4) Analyse
5) communicate

**○ Causal Relationships**
(chpt 1) Controlled experiments
(chpt 2) Observational studies

*assign subj into grps* →

- compare rates
  (eg.) Treatment vs control
- groups possibly having different risks
  → social economic status
  → severity of disease
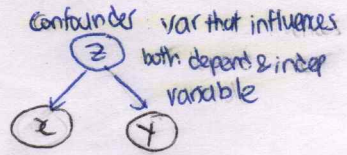- Randomised assignment from group w parental control ⟹ Randomised Controlled **Double Blind** Experiment

Q: Are groups different aside from treatment

⟹ can reduce confounder (using randomise)
∴ can suggest causation

precautions agst bias ! (confounding bias) *

- Placebo
- Blinding Subj
- Blinding Doctors

Diff from experimental study in that groups have been **preassigned subjects**
(eg. can't ask someone smoke 20 yrs)
⟹ ∴ cannot reduce confounders by randomization (but can slice)
∴ cannot suggest/prove causation
(only correlation)

Confounder : var that influences both depend & indep variable


"rate"
○ A, B pop<u>n</u> characteristics <u>w</u> $0 < r(A), r(B) < 1$

| A and B | | Condition | |
|---|---|---|---|
| positively associated | | $r(A\|B) > r(A\|not\ B)$  or  $r(B\|A) > r(B\|not\ A)$ | |
| negatively associated | | $r(A\|B) < r(A\|not\ B)$  or  $r(B\|A) < r(B\|not\ A)$ | |
| not associated | | $r(A\|B) = r(A\|not\ B)$  or  $r(B\|A) = r(B\|not\ A)$ | |

eg. Smoking & lung prob
eg. Polxity & High IQ

B is more common among people w A than among people without A
* always compare $r(A\|B)$ & $r(A\|^\sim B)$
NOT $r(A\|B)$ & $r(\sim A\|B)$

~~Confounding~~

**Confounder**
- Smoking & Heart Disease → Age & Sex confounder (exaggerates effect of smoking)
  *2 common confounders*
- confounder (3rd var) associated <u>w</u> both exposure & disease
* impt to control for confounder in Observational studies
  ↳ "slicing" → splitting grp into heart disease amongst males & h.d. amst females
  [Separating Data sets]

Techniques for controlling confounder ⟨ Slicing (basic but cumbersome)
                                      ⟨ Statistical Technique (eg regression)

thinking of confounders.
Are __ & __ diff in some ways other than ____

**SUMMARY**
○ Association is not causation
○ Observational Studies prone to confounding
→ knowledge & thinking useful for spotting potential confounders.
○ Slicing is effective for controlling confounders.

**Quiz learning**
○ <u>Yule-Simpson Paradox</u> : when there is a disparity btw ⟨ direction of association in <u>most</u> subgroups ⟨ & ~~dir of asso~~ overall association (when subgroups are combined)
(Chpt 1 slide 83)
→ R/nships in subgroups are reversed when rates are combined
⟹ Due to confounding (eg. size of kidney stones)
∴ treat paradox by slicing into grps / confounders

# GER wk 3 & 4

- **Association → Statistical Relationship**
  1. Bivariate Data & Scatter Diagram
  2. Explore Relationship
  3. correlation coefficient
  4. some limitations
  5. Ecological correlation (aggregate measure)
  6. Cautionary Notes
  7. Simple Linear Regression

  → closeness to meanline
  - correlation relation r ──────
  (direction & strength) of **linear** relation
  
  $-1 \quad 0 \quad 1$
  
  - if var indep → r = 0
    r = 0 ↛ var indep ( r only measure linear dependences )
  
  - r close to ±1 → strong association (abt 0.7)   ⎫ positive / negative
    "   ±0.5 → moderate   " (0.3 - 0.7)            ⎬
    "   0 → weak   " (0 - 0.3)                     ⎭

  - computing correlation coefficient (r)
    ① Convert to Standard Unit (SU)   each var   $SU = \dfrac{x_i - \bar{x}}{sd_x}$
    ② Take pdt of SU for father-son pair
    ③ r ⇒ Avg of all products

    $$ r = \frac{1}{n} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{sd_x} \right) \left( \frac{y_i - \bar{y}}{sd_y} \right) \quad \gg = CORREL( \text{rangeX}, \text{rangeY} ) $$

    - no units
    - not affected by interchange
    - not affected by change of scale ( addition / multiplication ) or units of measurement

## Limitations w Causation   (correlation)
  1. Causality
  2. Outliers in data set ──→ very sensitive, may deflate/inflate correlation
  3. non-linear association

## Ecological correlation
- Data of ~~individuals~~ Aggregate data (not individual data) like groups

→ generally if measured the 2 variables, strength would be overstated (cos same line)
* (similar if var are in same direction)

**Ecological fallacy**: deduce inferences on _corr._ abt individuals based on aggregate data  ⎱ Only causal* correlated relations
**Atomistic Fallacy**: Generalize corr. based on individuals toward aggregate

∴ proceed carefully w scatter diagram

## Cautionary Notes on Correlation
Attenuation **means** → Reduction in Value
- **Attenuation Effect** : ( Phenomenon due to Range Restriction in 1 var )
  → ~~reduction of range~~ Range Restriction
    - bivariate data set formed from restriction of one var
    - data for other variable available for a limited range

* tend to have diminishing effect of correlation (understate)

- Removal of Data (eg.)
  (eg. space shuttle ~~removal~~ removal of non-damaged data obj )

## Linear Regression
- using regression line  ( Y = a + bX )
  a → y-intercept
  b → gradient
  
  excel name ~~best name~~
  "intercept" ⎱ f(x) get
  "slope"     ⎰ - "statistical"
  ...
- specify independent (X) and dependent (Y) var.
- ✓ Determined by least square Method

## Sampling

- How Data is generated
- Generate own Data set / use public ones?
- Terminology : "unit", "population", "sample"
  - census — measurement from every unit in popn
  - sample — measurement from some selected unit in popn

Adv of sampling over taking census : 
- when census not possible (eg. Disease / Blood sample)
- Speed, cost, accuracy
  - (btr task force)
  - ( scale measurement )

Sampling frame → list of sampling units used to identify all units in popn (eg. house address
- Good sample is extendable
- Every unit in popn has possibility to be sampled → people in HDB houses)
- No selection bias
- simplest sampling frame : list of units in popn

### Characteristics of good Frame
- Good coverage
- Up-to-date & complete

→ inclusion of undesired units will not increase cost
→ exclusion of desired units will not have major impacts on outcome of study

eg. Drug Testing (Disease)
- unit — individual & disease
- popn — collection of all individuals defined above
- sample — collection of individuals placed under experiment
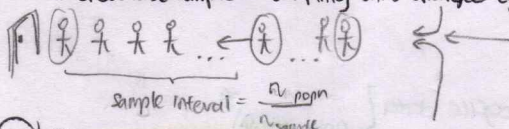
## Probability Sampling Plan

- (eg.) Simple Random Sampling ( every unit same prob of selection )
- (eg.) Systematic Sampling ( application of selection interval K )
  - can apply to situation when exact size of popn not known (rough estimate only)
  - → treated like Simple Random Sampling IF sampling units arranged randomly
  - ✗ Undesirable sample if sampling units arranged cyclically (197197...)

  eg. 110 units in popn,
  Sampling unit desired 10.
  $110/10 = 11$, $rand(1, 11) = 6$
  Sampling unit = [6, 17, 28 ...]



sample interval = $\frac{N_{popn}}{N_{sample}}$

- (eg.) Stratified Sampling Plan
  - → divide popn into grps ("strata")
  - → take prob sample from each group
  - – can involve setting quotas

  (eg. Male & Female)
  (eg. Country as strata
  so local organisations
  help-take results
  → find effectiveness of drug in
  each country (more accurate))

- (eg.) Multistage Sampling Plan
  - In many studies, several stages of sampling req. for reaching popn units.

  (eg.) list of houses →  → everyone interviewed [cluster sampling]
  [Simple Random]
  → one member randomly interviewed [2-stage sampling]

  (eg.) Drug testing of particular disease
  - Sampling frame : list of physicians
  - From each physician, only some of patients included in study.

## Difficulties in Sampling
① Using imperfect sampling frame
  - exclude desired & include unwanted
  - ① redefine target popn ↳ increase cost of study
  - ② Assess impact of excluding these units in own study
② Non-Response
  - Not anyone willing take part in study
  - ↳ incentive increase response rate (eg. Medicine study < payment / free checkup)
③ Getting Volunteer / Self-selected sampling [Aka non-prob sampling]
  - (eg.) News / Media conduct polls on website
    - ↳ bias as only people w strong views willing to give answers
④ Using a Convenience or Haphazard Sample
⑤ Taking a Judgement Sample : interviewer choose units by discretion
⑥ Selecting a Quota Sample (typical in Market Research) (distinguish from Stratified)
  - ↳ compare prop of quota sample to category units / census units — use convenience sample to hit Quota

Statistician John Tukey " I would trade all you 18,000 case histories for 400 in a prob sample "

GERI000 wk6

## Ch3 Estimating Parameter

- Parameter : numerical fact bout popn (usually unknown)
  → ~~estimate~~ estimated from sample

  { Estimate = Param + Random Error + Bias }
  
  Easy to Quantify    Hard to Quantify

→ Random Error

Larger sample size ——→ likely smaller random error

### Confidence Interval ← confidence level
          ← range

- range of values we are reasonably certain
  unknown param lies in

|———|———|———|
0.18   0.10   0.22

95% CI : 0.20 ± 0.02

Confidence level x% → x% ~~chance~~ likelihood
that ~~the~~ actual param value in range → resp & unique
(aka x% of researches will have intervals containing popn param)

→ Large sample size → likely → smaller range of CI
        (width of CI) also affected
            by variation
             in popn

- We are 95% confident that [0.18, 0.22] contains the popn param

## Chapter 7 Uncertainty <sub>not 'chance'</sub>
    unit : measuring uncertainty

### Probability — the measure of likelihood

- Interpretation ( Relative Freq   vs   Personal Prob )

| Relative Freq | Personal Prob |
|---|---|
| can be quantified | can't be exactly quantified |
| Based on repeated observation of outcomes | Based on own personal belief |

→ Proportion of times over the long run

---

### Types of Bias
1) Selection Bias
2) Non-response Bias
  ... other types

Systematic tendency on
the part of sampling procedure
to exclude one kind of person
from sample

Caused By? → imperfect sampling frame
         → Non probability sampling methods

Systematic tendency from subjects
who do not respond to survey / questionaire

Caused By → Diff btw non-respondent & respondents

** non-response rate ↑ → non-response bias significance ↑
       likely

### Other types of Bias

- Phrasing of qn / tone / attitude of interviewer ( including order)

- When subj have tendency to understate
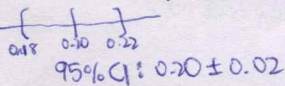  undesirable social habits (like smoking)

→ Conclusion ( to minimize bias in survey)
1. Include every popn unit in frame
2. Use prob sampling frame
3. Get 100% response rate

(Chpt4) more on observational studies
PARTE 1: RISKS (RATE IN POPD)

Two kinds of observational studies
> RISKS : rate in population
> Odds

Risk Ratio (RR)
(aka Relative Risk)

(eg.) risk(Diabetes | Female) = $\frac{D \cap F}{F}$ = $\frac{72,000}{216,000}$ ≈ 0.33

risk(D | M) = 0.25

Risk Ratio = $\frac{risk(D|F)}{risk(D|M)}$ ≈ $\frac{0.33}{0.25}$ ≈ 1.33

o Use Probability Samples to estimate
(SRS)
2 samples strategies : Simple Random Sample
From
> o each exposure group (female & males)
> o each disease group (diabetic & healthy)

|  | Diabetic | Healthy | Row Total |
|---|---|---|---|
| Females | 72k | 144k | 216k |
| Males | 52k | 156k | 208k |
| Col Total | 124k | 300k | 424k |

Cohort

Prospective

"Relatives are the same"
(Relative Risk)

Exposed Cohort     OUTCOMES
                    A
                    B
UNEXPOSED           C

Case Control
Retrospective
A          O

"at Odds w one another"
(odd ratio)

EXPOSURE     CASES
A
B
C
          CONTROL
  t

### Cohort & Case Control studies

| Study | Samples from | Advantage |
|---|---|---|
| Cohort | Exposure grp | Risks & RR can be estimated from sample table |
| Case-Control | Disease grp | Good for rare diseases |
|  |  | ↳ estimation not possible |

### Remarks on Sampling
(risk estimates too high)
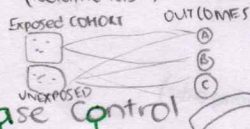(depends on sampling % from each)
> Randomised Exp : subj need not resemble popn group
  Extrapolation to popn is an issue
> Observational studies : Both extrapolation &
  confounders are impt issues
  o Cohort studies rely on random samples for
    accurate estimation of risks

### Summary (Risks)
o Risks like rates are affected by confounding
o Risk Ratio /Relative Risk measures association
o Cohort Study has accurate estimation of
  popn risks & RR w random samples
o Case-Control Study does not.
· EXPOSURE → DISEASE
· R(Disease | Exposure) = P(getting Disease if you
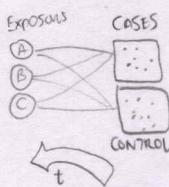                                are exposure grp)

## PARTE 2: Odds
### Risk & Odds

$$odds = \frac{risk}{1-risk}$$

Odds (Diabetes) among females = $\frac{72k}{144k}$ ≈ 0.50 ↔ $\frac{0.33}{1-0.33}$

Odds (Diabetes) among males = $\frac{52k}{156k}$ ≈ 0.33 ↔ $\frac{0.25}{1-0.25}$

o odds > risk
o if risk is small, odds is small

* 
OR=1 → No diff. in disease risk btw 2 groups: RR=1
OR>1 → Higher Risk in first group : RR >1
OR<1 → Lower Risk in first group : RR <1

~ODDS RATIO~
baseline
OR (Diabetes) btw Female & Male
≈ $\frac{0.50}{0.33}$ ≈ 1.50

only situation
der the value
can be deduced
from OR

CROSS-PRODUCT RATIO → to calculate OR

|  | Diabetic | Healthy |
|---|---|---|
| Female | 378 | 702 |
| Male | 53 | 155 |

$\frac{378 \times 155}{702 \times 53}$ ≈ 1.57

|  | Diabetic | Healthy |
|---|---|---|
| Female | 364 | 142 |
| Male | 256 | 158 |

$\frac{364 \times 158}{142 \times 256}$ ≈ 1.58

Table must be set proper
o Event of interest (1st col)
o 1st Group (1st row)

✓ Risk Ratio
→ Only Cohort Studies

✓ Odds Ratio
→ Both Cohort & Case-Control studies

✓ Cohort Study
① All subjects disease-free @ beginning
② Sometimes, OR is used by researchers
→ Compare odds of developing disease between 2 exposure groups

### Multi-Level Contingency Table
→ Form 2x2 table (choose interest & baselines)

⚡ error for observational studies,
correlation → causation

## Chapter : Uncertainty (measuring it)

Probability _ measure of likelihood

(Diabetes eg.)
Risk_female = 0.33 = $\frac{72k}{216k}$ = P(female has Diabetes)
Risk_male = 0.25 = P(Male has Diabetes)

### Simple Probability Rules
o Mutually Exclusive A,B → P(A∩B) = P(A) + P(B)
o Complement Rule

| Interpretations | Relative Freq. | Subjective Prob Personal Prob |
|---|---|---|
|  | can be quantified exactly | can't be quantified exactly |
|  | Based on repeated observation of outcomes | Based on own personal beliefs |

### odds revisited
Odds = $\frac{\text{# of event}}{\text{# of non-event}}$ = $\frac{P(event)}{P(event\ not\ occurring)}$

o Independent Event A,B → P(A∩B) = P(A)P(B)
(ocurrence of one does not affect
ocurrence of the other)

### Relative Freq vs Prob

| Tossing coin | Weather forecast |
|---|---|
| Precise Prob | Imprecise Prob |
| Assumption abt physical reality | circumstances repeat & outcome observed |
| Prob used to predict relative frequency | Relative Freq used to predict Prob |

# Geri000 wk8

Unit 3) Expected value
Unit 4) Uncertainty p-values

Overview:
p-values
Null Hypothesis   Hypothesis Test   Statistical Significance

p-value:
- The probability of obtaining an outcome **equivalent to** or **more extreme** than the observed.

0.05%
small ←——————→ Large

- unlikely for the observed to occur by chance
- unlikely that null hyp is true

- More likely that null hyp is true
- more likely that observed occured by chance

Null hypothesis rejected at 5% level of statistical significance

## Hypothesis Testing

1. Identify the question (Frame)
2. State the null hypothesis (Specify)
3. Conduct the experiment (Collect)
4. Compute P-value (Analyse)
5. Make Conclusion abt the null hypothesis (Communicate)

(summary)

small p-value
- Prob too small for observed data to be purely by chance
- Statistically significant to reject null hypothesis
- Observed effect in sample is likely to reflect effect in population

Large p-val
- observed data may be due to chance
- Not statistically significant to rej $H_0$
- observed effect may or may not reflect effect in pop $n$

Unit 5) Conditional Prob
- independent: $P(A|B) = P(A)$  / $P(A \cap B) = P(A) P(B)$
   A, B
- $P(A \cap B) = P(A|B) P(B)$

Unit 6) Uncertainty Testing Rare Events

eg. can tossing, is coin fair?

HHHTH  } outcome
THHHH  } equivalent to outcome
HTHHH
HHTHH
HHHHT
HHHHH  } more extreme than outcome

Assume
(assumes that observation is due to chance)
Null Hypothesis: the coin is fair

Evidence that coin is biased in favor of head
P-value = P(TH·IHH) + P(HTHHH) +...
= 0.18

P-value = 0.18 > 0.05
→ Do not reject null hypothesis at the 5% significance level
→ cannot conclude that the coin is not fair

conclusion: case 1) reject $H_0$ (accept ~~~~ $H_1$ at x% level of significance)
case 2) fail to reject $H_0$ (cannot ≠ accept $H_0$)

* $H_0$ is never accepted

### Drug Testing eg.

1) Is drug effective in some pop $n$?
2) $H_0$: drug has no effect in pop $n$
3) Give drug to three patients and observe the number of patients who survived
4) Compute p-value
5) Conclude (rej $H_0$?)

Disease D, 40% fatal
- 40% fatal → 60% survive → P(survive) = 0.6
- P(all 3 patients survive) = $0.6^3$ = 0.216
- P-val = 0.216 > 0.05 ⟹ Do not reject null hypothesis at 5% significance level

If 4 out of 6 survived
P-val = P(4 survive) + P(5 survive) + P(6 survive)
$= \binom{6}{4}(0.6)^4(0.4)^2 + \binom{6}{5}(0.6)^5(0.4)^2 + \binom{6}{6}(0.6)^6$

interesting eg.
D: double cot death
E: Sally is innocent

| P(D|E) | vs | P(E|D) |
|--------|----|----|
| - Prob quoted by expert | | - Prob required by prosecutor |

eg. Framing an: [ tested positive ⟹ Treatment? ]  P(disease|positive)
Specifying what to measure:

| | | |
|---|---|---|
| Base Rate | P(disease) | = A/N |
| Sensitivity of Test | P(positive|disease) | = C/A |
| Specificity of Test | P(negative|no disease) | = D/B |

Collecting Data:
1. Carry out study on random sample (N) from pop $n$
2. Record no. of pep $w$ disease (A) & $wo$ disease (B)
3. Among those $w$ disease, record no. of pep tested positive (C)
4. Among those $wo$ disease, record no of people tested negative (D).

### Communicate the Findings

pos | dis
P(pos|pos) = 0.95
P(neg|nodis) = 0.9

- Test has high sensitivity & specificity

- Less than 1% tested positive have the disease: P(disease|positive) = 0.0094

- More than 99% tested positive have no disease:
  P(no disease|positive) = 0.9906

→ if tested positive, not confident that test is correct
→ ALWAYS happens, when disease is rare.

SUMMARY
- P(event happening | event suspected)
- Base rate, sensitivity, specificity REQUIRED!
- Contigency Table set up
- True positive & False positive

Analyze Tables Data:
Contigency Table: assume pop $n$ of 100,000

Base rate: 0.001
Sensitivity: 0.95
Specificity: 0.9

| | Test positive | Test negative | Row sum |
|---|---|---|---|
| Have Disease | 95 | 5 | 100 |
| Do not have Disease | 9990 | 89910 | 99900 |
| Col sum | 10085 | 89915 | 100 000 |

fill in table

True positive
False positive
False negative
True negative

$P(disease|positive) = \frac{No. of true positive}{No. of positives} = 0.0094$ & $P(no disease|positive) = \frac{No. of false positive}{No. of pep tested pos} = 0.9906$