

SELECTION OF COMPARISON GROUPS.

INDIVIDUAL-LEVEL CAUSAL EFFECTS: effect of $\begin{cases} \text{trt} \\ \text{exposure} \end{cases}$ on the same individual.

POPULATION-LEVEL CAUSAL EFFECT: average effect of $\begin{cases} \text{trt} \\ \text{exposure} \end{cases}$ on group.

COUNTERFACTUAL:

Simultaneously examine the $\begin{cases} \text{pop} \\ \text{ind} \end{cases}$ with $\begin{cases} \text{trt} \\ \text{without exposure} \end{cases}$

↓
impossible, so impossible to be certain

Choose comparison groups to approximate the counterfactual →

ASSUMPTION OF COMPARABILITY OF GROUPS

ASSUMPTION OF NO RESIDUAL CONFOUNDING.

↓
Assumption probably violated → bias $\begin{cases} \text{point estimates} \\ \text{measures of variability} \end{cases}$

→ ↑ Type I error

C.I don't achieve nominal coverage

↓

C.I can achieve nominal coverage w randomization?

CHOICE OF COMPARISON GROUPS

- INDIVIDUALS : select very similar to the trt / exposed ind (cloned mice, twins)

- POPULATIONS :

• CLINICAL TRIALS :

■ COMPARABILITY OF EFFECT:

- trt effect : placebo

- exposure effect : reference category

■ COMPARABILITY OF POPULATION : randomization.

Avg. incidence of the outcome not with exposure in all groups is =

Does not eliminate confounding → ↑ or ↓ chance of imbalance = confounding.

Imbalance is caused by differential selection into and out of compared groups, w/ the selection being related to risk factors of the outcome.

Ex: healthy worker effect : ↓ disease ind. at workplace compared to general pop.

other pops are comparable → no confounding

■ COMPARABILITY OF INFORMATION : blinding.

STUDY BASE PRINCIPLE:

- PRIMARY DEFINITION: 1st define cohort → 2nd identify cases
sum of the individual-level experience
captured through census / sample
as basis for learning about the
reference population it represents
 - ≠ statistical sample of target pop
= direct referent of the empirical results
- ind. / time sample
" ↓
study results
↓
extrapolate
↓
pop that
is represented
by the sample
" ↓
study base

STUDY BASE

CLINIC TRIAL follow-up experience of
enrolled patients.

- COHORT
person-time experience = incidence rates
of the entire cohort (including censored ind)
- CASE-CONTROL
?
→ need to conceptualize case-control
as a cohort study to understand
study base.

• SECONDARY DEFINITION : 1st identify cases → 2nd define study base

≠ pop. experience from which cases could have come to be included in the study (registry, etc...)

= totality of experience from which each potential case HAD IT OCCURRED would have been included in study.

• Each ind is in a study base at time t_i ; if were it to become a case, it would be included in study

Ex: diabetic dogs sampled from Devis vet hospital
study base = all dogs who, were they to become diabetic, would show up at Devis vet hospital

study base is limited
geographically (hospital compartmented)
temporally (study time)

↓
used in

CASE-CONTROL STUDIES = TRUTHOC

Aim to cohort studies when dis. is rare.

Recruit all cases → draw sample of study base
cases → ↙ controls

constrained by time / place

can yield complete info on the study base.

No of controls: depends on prevalence of exposure in study base and magnitude of OR.

$r = n^b$ of cases

→ Max 4 controls

$$\frac{r}{r+1}$$

- $r = 1 \rightarrow 1/2 \rightarrow 50\% \text{ of info from using whole pop.}$
- $r = 2 \rightarrow 2/3 \rightarrow 60\%$
- $r = 3 \rightarrow 3/4 \rightarrow 75\%$
- $r = 4 \rightarrow 4/5 \rightarrow 80\%$

The case-control approach is particularly valuable when:

- It is expensive, time-consuming, and inefficient to put an entire cohort under surveillance for disease incidence.
- The study base is large enough so that adequate information about exposure distribution can be learned by studying a small fraction of it.

How to sample controls? (no "sampling frame")
From study base → pop. from which potential cases would have been captured. =

ind. who had they developed the illness → would have been in my attention as cases.

Controls should be distinguished from cases only by a last-minute differential dx.

Ex: dx of G.I tumor in FE → study base is FE w/ other tumors dx in the same way (echography).

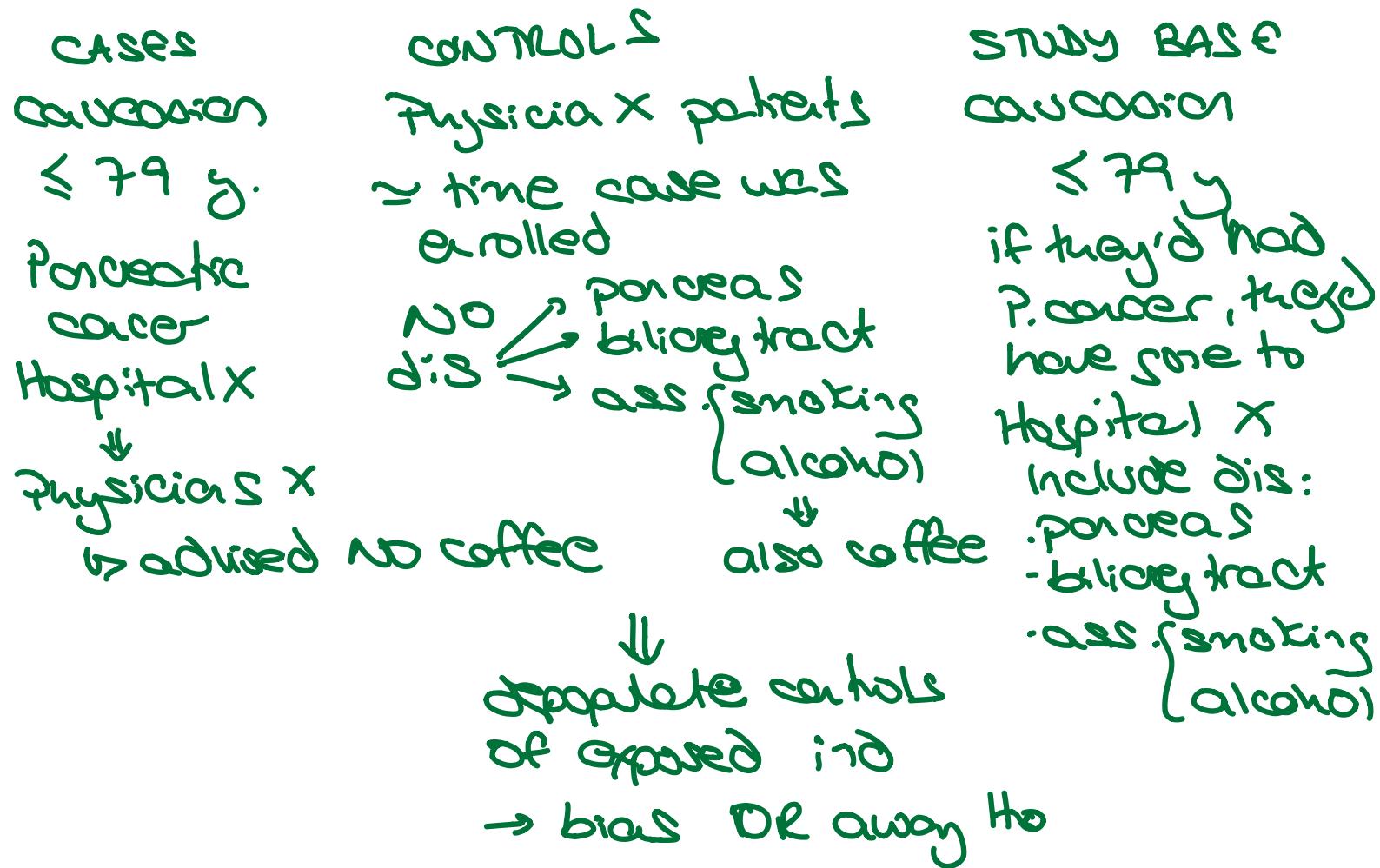
Controls should represent what the case pop. would be if the exposure had not occurred. = Ho → where cases are independent of exposure (conditionally on confounders)

Ideal control: cases of another disease :

- that derive from same registry
- where disease :
 - is unrelated to the exposure
 - is ≈ to dis of interest regarding selection

factors that contribute to it appearing in the registry.

Ex: MacMahon



■ “However, there is no virtue in using a wide variety of diagnoses. Each diagnosis needs to be justified as to the assumptions, and insofar as they are untenable, the errors have no tendency to cancel out on the basis of variety.”

■ “The need is to have defensible inclusions, and there is no need to defend exclusions.”



controls don't have to come from the study base, they can come from any pop w/ the same distribution of exposure / incidence, conditional on controlled variables in the analysis (it is necessary to control for unknown confounders).

Ex:

Incidence of thromboembolism = OUTCOME

w/ ABO bloodtype = EXPOSURE

in ♀ using oral contraceptives. = BASE POP.

control: you could use ♂ or ♀ w/o oral cont.

bc ABO dist = independent of gender, age, oral cont.

> 1 control group :

"Because there tend to be strong feelings for and against various types of control groups by different investigators, an association will have a much greater likelihood of being accepted as valid if it is found when different types of control groups are used."



Marginally unassociated
conditionally associated:

bc any exposure-dis association
is conditional on being selected
for inclusion in the study.

	A=1		A=0	
	B = 1	B = 0	B = 1	B = 0
C = 1	800	600	400	200
C = 0	200	400	600	800
Total	1000	1000	1000	1000

TOTAL

2000

2000

4000

A and B marginally independent : A=1

Ex: $P(A=1|B=1) = \frac{1000}{2000} = 0.50 = P(A=1) = \frac{2000}{4000} = 0.50$
OR $(A \text{ and } B) = (1000/1000) / (1000/1000) = 1$

	A=1		A=0	
	B = 1	B = 0	B = 1	B = 0
B = 1	800	400		
B = 0	200	600		
Total	1000	1000		2000

A and B NOT conditionally independent.

①

Ex: $P(A=1|B=1,C=1) = \frac{800}{1200} = 0.67 \neq P(A=1|C=1) = \frac{1400}{2000} = 0.70$
Ex: $P(A=1|B=1,C=0) = \frac{200}{800} = 0.25 \neq P(A=1|C=0) = \frac{600}{2000} = 0.30$
Ex: $P(B=1|A=1,C=1) = \frac{800}{1400} = 0.57 \neq P(B=1|C=1) = \frac{1200}{2000} = 0.60$
Ex: $P(B=1|A=1,C=0) = \frac{200}{600} = 0.33 \neq P(B=1|C=0) = \frac{800}{2000} = 0.40$
OR $(A \text{ and } B|C=1) = 0.67$ and OR $(A \text{ and } B|C=0) = 0.67$

①

	A=1		A=0	
	B = 1	B = 0	B = 1	B = 0
C = 1	800		400	1200
Total				

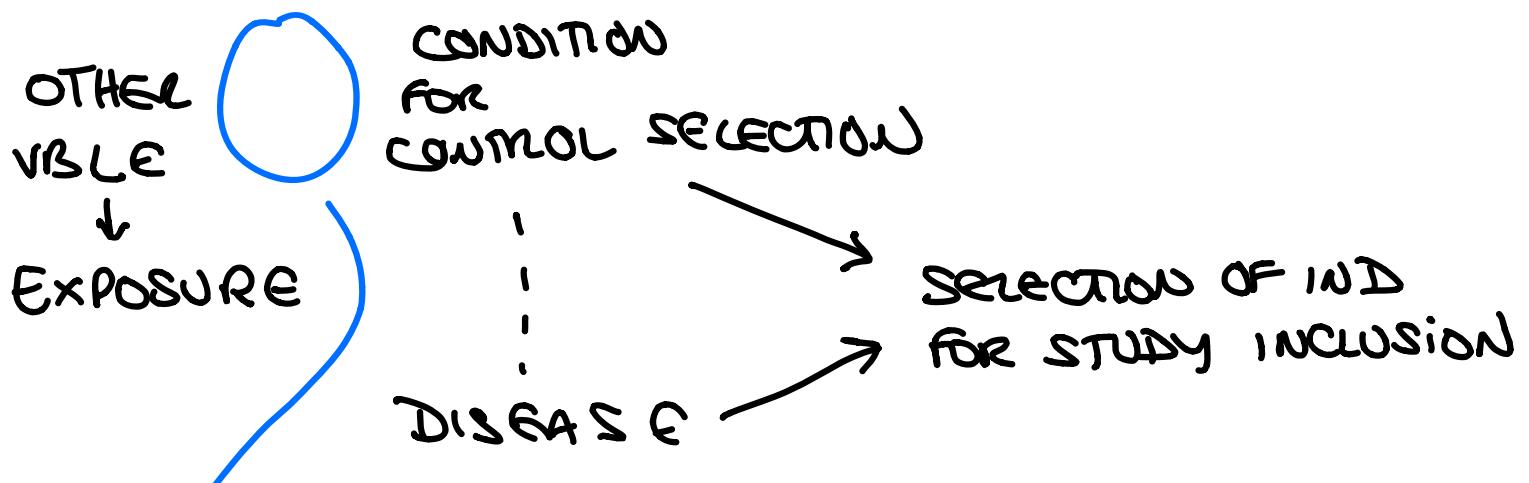
	A=1		A=0	
	B = 1	B = 0	B = 1	B = 0
C = 1	800	600	400	200
Total	1400		600	2000

	A=1		A=0	
	B = 1	B = 0	B = 1	B = 0
C = 1	800	600		
Total				

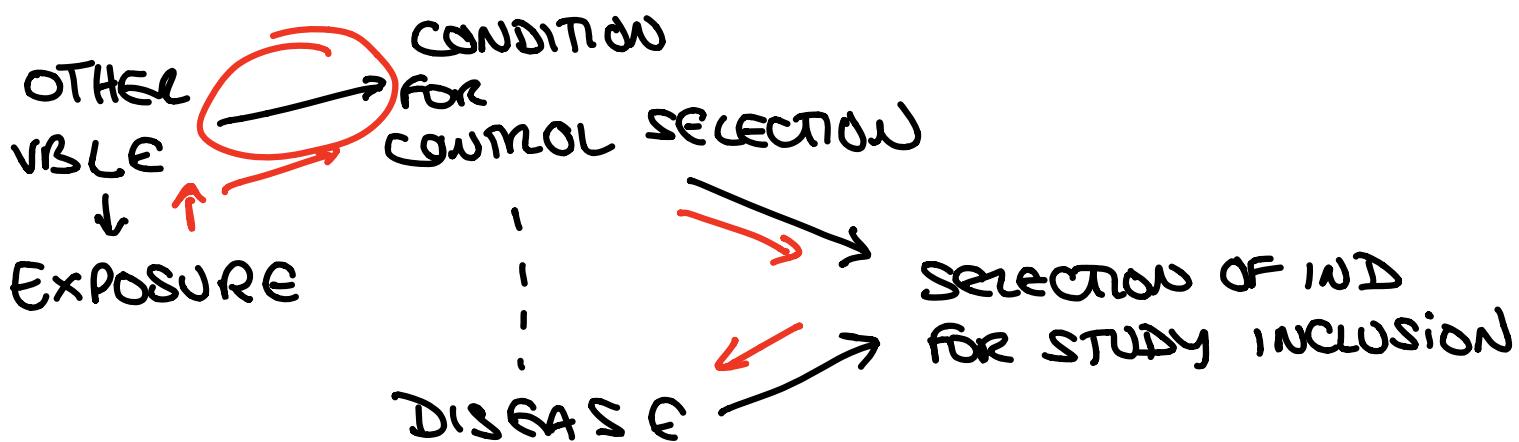
	A=0		A=0	
	B = 1	B = 0	B = 1	B = 0
A = 0				
Total	400	200		

$$\frac{800 \times 200}{400 \times 600}, 0.67 = \text{OR}(A \text{ AND } B | C=1)$$

$$\frac{200 \times 800}{400 \times 600}, 0.67 = \text{OR}(A \text{ AND } B | C=0)$$

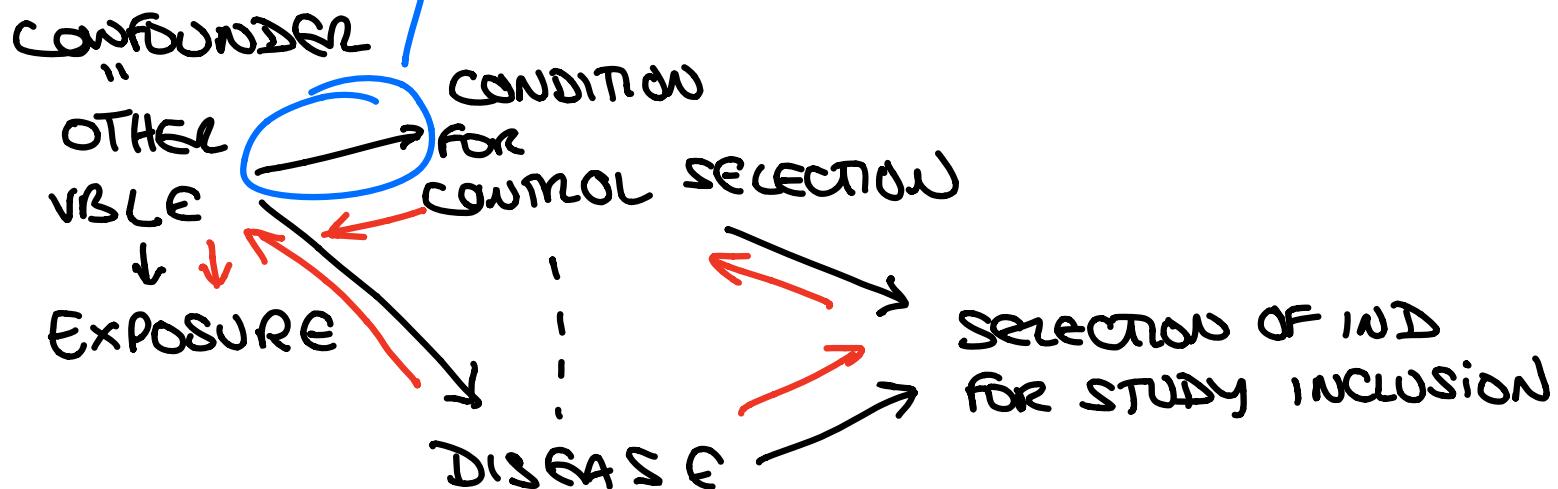


controls properly selected. No ass. between condition for selection and exposure

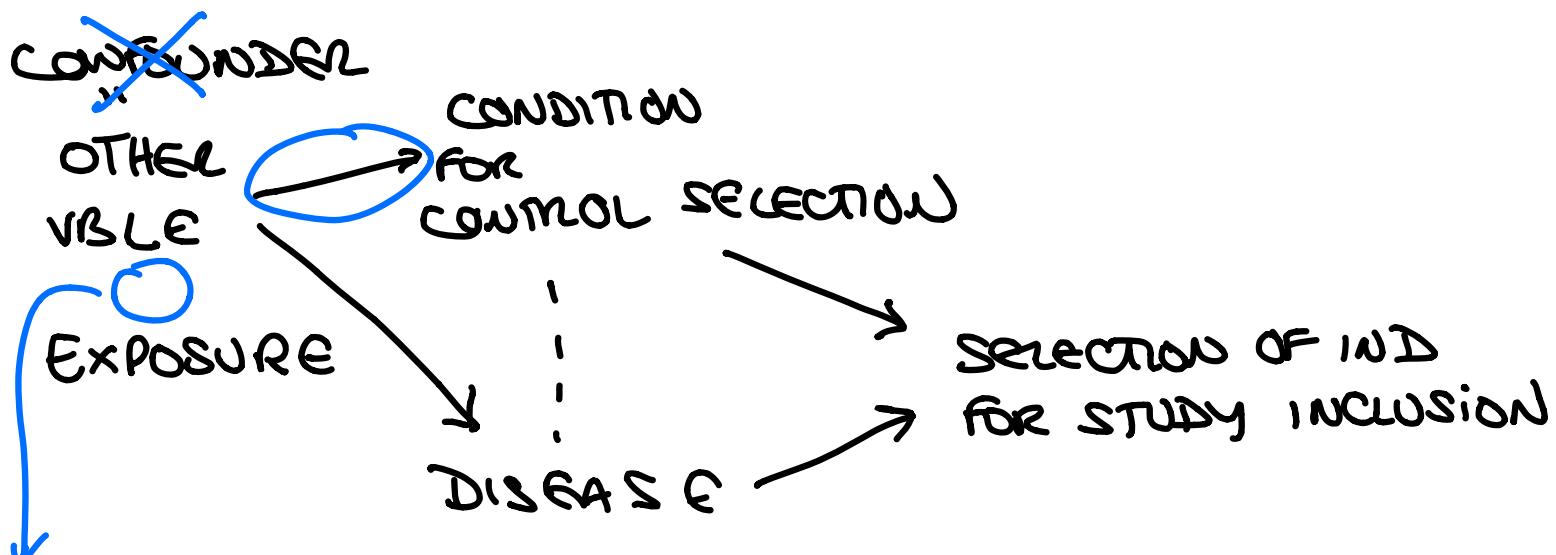


Backdoor path: controls selected for reasons not independent of exposure → exposure → condition → disease or associated → selection bias → shows exposure-disease ass even when there is no causal relationship.

MATCHING: cases are matched to controls on a confounder

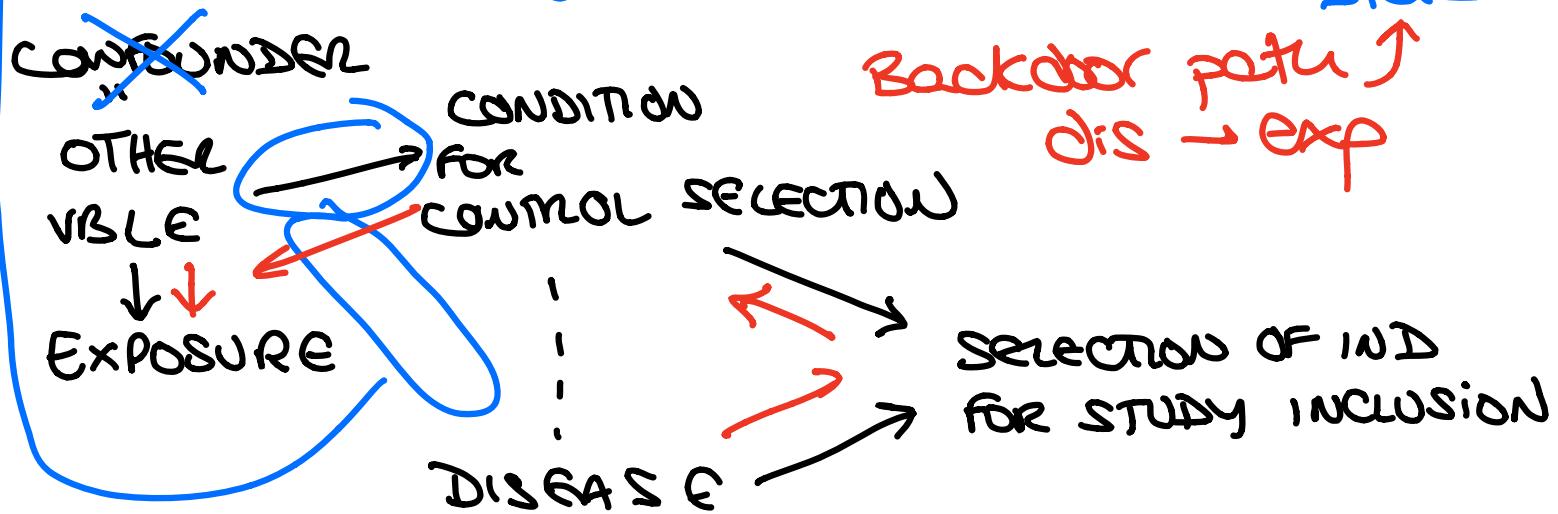


Backdoor path: you still have to control on the confounder, even you match by it, to block the backdoor path $DIS \rightarrow EXP$.



OVERTMATCHING: match on a vble not ass w/exp.
• last ass w/ DIS → NO selection bias
not a confounder

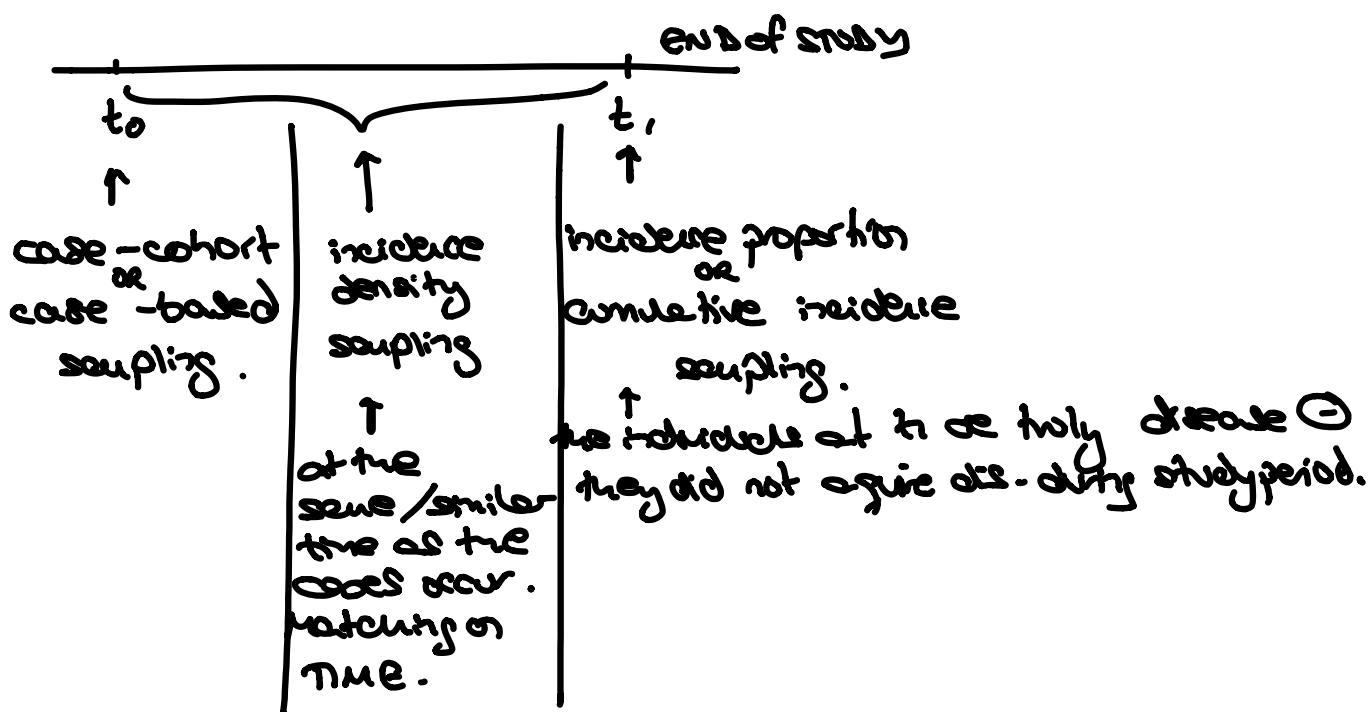
- match on a var not ass w/dise but ass w/ EXP. → not a confounder
- selection bias



To avoid backdoor path you NEED to control for this vble, but if you hadn't matched you wouldn't have needed to in 1st place.

METHODS OF CONTROL SAMPLING:

- CASE - BASED SAMPLING : to
- INCIDENCE SAMPLING : match on time w/ cases
DEENSITY
- INCIDENCE SAMPLING : t_1 , end of study.
PROPORTION
OR
CUMULATIVE
INCIDENCE



ODDS RATIO

Diseased to nondiseased | Exposure \oplus

COTORT STUDY

	Exposed E^+	unexposed E^-	
CASE \exists^+	A_1	A_0	
POP. AT RISK \exists^-	B_1	B_0	

CASE CONTROL STUDY

	Exposed E^+	unexposed E^-	
CASE \exists^+	A_1	A_0	
CONTROL \exists^-	B_1	B_0	
			Exposed to unexposed Disease \ominus

COTORT EXPOSURE OR :

$$\frac{\text{odds Diseased | Exposed}}{\text{odds Diseased | unexposed}} = \frac{A_1 / B_1}{A_0 / B_0}$$



Algebraically similar
but \neq meaning.

CASE-CONTROL EXPOSURE OR :

$$\frac{\text{odds Exposure | case}}{\text{odds Exposure | control}} = \frac{A_1 / A_0}{B_1 / B_0}$$

Exposed to unexposed | Disease \ominus

- CONTROL SELECTION NOT MATCHED ON TIME :
to : CASE - BASED. t_i : INCIDENCE PROP. SAMPLING

ASSUMPTION ① : a constant IDR ; ie $\frac{ID_1(t)}{ID_0(t)} = \overline{IDR}$ over (t_0, t_1)

$$\overline{IDR} : IDR(t) = \frac{ID_1(t)}{ID_0(t)} \Rightarrow ID_1(t) = \overline{IDR} \cdot ID_0(t)$$

ASSUMPTION ② $\pi_i(t) = \pi_i$ stable pop. assumption with respect to exposure

This means that time is not a confounder \Rightarrow we do not have to control for time.

$$OR_{UNMATCHED} = \frac{A_1/A_0}{B_1/B_0} = \frac{\frac{\int_{t_0}^{t_1} a_1(t) dt}{\int_{t_0}^{t_1} a_0(t) dt}}{\frac{\int_{t_0}^{t_1} b_1(t) dt}{\int_{t_0}^{t_1} b_0(t) dt}} =$$

$a_+(t) = a_1(t) + a_0(t)$
 = exposed cases + unexposed cases

CONTROLS

(total n° of cases expected per unit of t) (proportion of $N(t)$ in exposure group i)

$$= \frac{\int_{t_0}^{t_1} N(t) \cdot P_i(t) \cdot ID_i(t) dt}{\int_{t_0}^{t_1} N(t) \cdot P_o(t) \cdot ID_o(t) dt} =$$

CASES

(pop. at risk) (proportion of the pop. n°) (incidence rate) in exposure group i

$$\overline{IDR} = \frac{ID_i}{ID_o} \Rightarrow$$

$$ID_i = \overline{IDR} \times ID_o$$

$$= \frac{\left[P_i \overline{IDR} \int_{t_0}^{t_1} N(t) ID_i(t) dt \right] \left[\int_{t_0}^{t_1} a_+(t) dt \right]}{\left[P_o \int_{t_0}^{t_1} N(t) \cdot ID_o(t) dt \right] \left[P_i \int_{t_0}^{t_1} a_+(t) dt \right]} = \overline{IDR}$$

$$\left[P_o \int_{t_0}^{t_1} N(t) \cdot ID_o(t) dt \right] \left[P_i \int_{t_0}^{t_1} a_+(t) dt \right]$$

we have interpreted the \overline{IDR} from the OR without making the rare disease assumption. Importantly, we assumed that time is NOT a confounder.

• CONTROL SELECTION MATCHED ON TIME :

INCIDENCE DENSITY SAMPLING

ASSUMPTION ① : a constant IDR; ie $\frac{ID_1(t)}{ID_0(t)} = \overline{IDR}$ over (t_0, t_1)

$$\overline{IDR} : IDR(t) = \frac{ID_1(t)}{ID_0(t)} \Rightarrow ID_1(t) = \overline{IDR} \cdot ID_0(t)$$

CASE CONTROL STUDY

		CONTROLS	
		Exposed E^+	unexposed E^-
Exposed E^+	Exposed E^+	M_{11}	M_{10}
	unexposed E^-	M_{01}	M_{00}

Matched pairs OR

$\frac{M_{10}}{M_{01}} \rightarrow$ proves H
 $\frac{M_{01}}{M_{10}} \rightarrow$ disproves H

C	C_0
E	UnE
UnE	E

The expected number of pairs are $(t_0, t + dt)$ in which the case is in exposure group i and the control B is in exposure group j = $a_i(t) P_j dt$

OR MATCHED

$$= \frac{\int_{t_0}^{t_1} a_{10}(t) dt}{\int_{t_0}^{t_1} a_{01}(t) dt} = \frac{\int_{t_0}^{t_1} a_i(t) P_0(t) dt}{\int_{t_0}^{t_1} a_0(t) P_i(t) dt} =$$

$$= \frac{\int_{t_0}^{t_1} N(t) \cdot P_i(t) \cdot ID_1(t) \cdot P_0(t) dt}{\int_{t_0}^{t_1} N(t) \cdot P_0(t) \cdot ID_0(t) \cdot P_i(t) dt} =$$

$$= \frac{\int_{t_0}^{t_1} N(t) \cdot P_i(t) \cdot \overline{IDR} \cdot ID_0(t) \cdot P_0(t) dt}{\int_{t_0}^{t_1} N(t) \cdot P_0(t) \cdot ID_0(t) \cdot P_i(t) dt} = \overline{IDR}$$

Incidence is a $f(x)$ of time and exposure is a $f(x)$ of time \Rightarrow
time is a confounder. we match on time.
we have 1 less assumption than before and we get that $OR \approx \frac{DR}{IDR}$
still without having to invoke the rare disease assumption.

If you have a rare disease you get a better estimate if you do
incidence density sampling rather than cumulative incidence sampling.
you have the most flexibility in terms of interpreting the OR,
if you match cases on time.

You're better off designing a study that allows you to interpret the OR
on a rate or risk scale.

For each individual i , there is a certain incidence probability
 r_{ii} = probability disease will occur after an individual i 's exposed
 r_{oi} = " " "unexposed

$1 - r_{ii}$ = probability that dis. does not occur when ind i is E⁺ = S_{ii}

$$1 - r_{oi} = S_{oi}$$

Define risk diseased odds : $w_{ii} = r_{ii}/S_{ii}$ } on individual
 $w_{oi} = r_{oi}/S_{oi}$ " " individual level risk odds

Define risk difference : $r_{ii} - r_{oi}$ on individual i

Define risk ratio : r_{ii}/r_{oi} " " "

Define risk (rate) odds difference: $w_{ii} - w_{oi}$ or ind. i

Define incidence/risk/disease odds ratio: w_{ii} / w_{oi} or ind. i .

The expected number of individuals over the risk period:

	Exposed	not exposed	
Disease occurs	$A = \sum_i r_{ii}$	$B = \sum_i r_{oi}$	
Disease not occur	$C = \sum_i s_{ii}$	$D = \sum_i s_{oi}$	
	N_1	N_0	

s = survivors

r = prob. disease will occur

Incidence proportion | exposure: $\frac{A}{N_1} = \frac{\sum_i r_{ii}}{N_1}$

} measures of average risk

Incidence proportion | exposure: $\frac{B}{N_0} = \frac{\sum_i r_{oi}}{N_0}$

Disease odds | exposure: $\frac{A}{C} = \frac{\sum_i r_{ii} / N_1}{\sum_i s_{ii} / N_1}$

} Ratio of average risk to average survival probability.

Disease odds | exposure: $- = \frac{\sum_i r_{oi} / N_0}{\sum_i s_{oi} / N_0}$

Individual level risk odds.

Average disease odds

$$w_{ii} = r_{ii} / s_{ii} \quad \left. \right\} \text{on individual} \Rightarrow \frac{r_{ii}}{s_{ii}} \neq \frac{\sum_i r_{ii}}{\sum_i s_{ii}} \Leftarrow \frac{\sum_i r_{ii}}{N_1} \neq \frac{\sum_i r_{oi}}{N_0}$$

Ratio of average risk to average survival probability \neq Average disease odds

it's NOT equal

Odds of complicated odd need to explain!!

Assume NO confounding (i.e. group comparability)
 then the average risk in the exposed and unexposed subcharts would be the same if the exposed group was unexposed.

$$\frac{\sum_{\text{O}} \text{roi}}{N_0} = \frac{\sum_{\text{E}} \text{roi}}{N_1}$$

This is a critical counterfactual assumption.

The incidence proportion (risk) difference

$$\frac{A}{N_1} - \frac{B}{N_0} = \frac{\sum_i r_{1i}}{N_1} - \frac{\sum_i r_{0i}}{N_0} = \underbrace{\frac{\sum_i r_{1i}}{N_1}}_{\text{Hypothetical unexposed roi}} - \underbrace{\frac{\sum_i r_{0i}}{N_1}}_{\text{and exposed } r_{1i}} = >$$

This only holds if
NO confounding.

Absolute change in the average risk of the exposed subchart produced by exposure

\uparrow
 \uparrow
 \uparrow
 \uparrow

in the same cohort N_1 .

$$= \frac{\sum_i (r_{1i} - r_{0i})}{N_1} = >$$

Average absolute change in risk produced by exposure in exposed individuals.

This is a statement with clear meaning!!

That's why we need to avoid bias, in order to have answers that have causal significance.

Incidence proportion ratio :

$$\frac{A}{N_1} = \frac{\sum_i r_{1i}/N_1}{\sum_i r_{0i}/N_0} = \frac{\sum_i r_{1i}/N_1}{\sum_i r_{0i}/N_1} =$$

\downarrow
 \downarrow
 \downarrow
 \downarrow

proportionate change in the average risk ($\frac{\sum_i r_{1i}}{N_1}$, $\frac{\sum_i r_{0i}}{N_1}$) of the exposed subchart (N_1) produced by exposure (r_0 vs r_1)

\neq average proportionate change in risk produced by exposure

$$\frac{\sum_i (r_{1i}/r_{0i})}{N_1}$$

The causal interpretation looks at what happens to the exposed cohort if they had not been exposed. \Rightarrow this is the point of the unexposed cohort. It's a proxy for the exposed cohort if no exposure had occurred.

The disease odds ratio

$$\frac{A/C}{B/D} = \frac{\sum r_{1i} / \sum s_{1i}}{\sum r_{0i} / \sum s_{0i}}$$

NO confounding

Assumption

$$= \frac{\sum r_{1i} / \sum s_{1i}}{\sum r_{0i} / \sum s_{0i}} = \frac{\sum w_{1i} / N_1}{\sum w_{0i} / N_1}$$

$$= \frac{\sum w_{1i} / N_1}{\sum w_{0i} / N_1} = \frac{\sum w_{1i} / N_1}{\sum w_{0i} / N_1}$$

The proportionate change in the incidence odds (i.e. the ratio of the average risk to the average survival probability) in the exposed population produced by exposure (r_0 vs r_1 , s_0 vs s_1)

This is the correct interpretation of the OR, but it is hard to comprehend.

$$\frac{\sum w_{1i} / N_1}{\sum w_{0i} / N_1} = \frac{\text{proportionate change in the average odds in the exposed population produced by exposure}}{\sum w_{0i} / N_1}$$

$$\frac{\sum \left(\frac{w_{1i}}{w_{0i}} \right)}{N_1} = \frac{\text{average of the individual odds ratios among the exposed}}{N_1}$$

These are more intuitive but incorrect interpretations of the OR

Thus, the incidence odds ratio lacks any simple interpretation in terms of average risk or average odds or an average exposure effect on individual risk or odds.

CONFOUNDING, BIAS

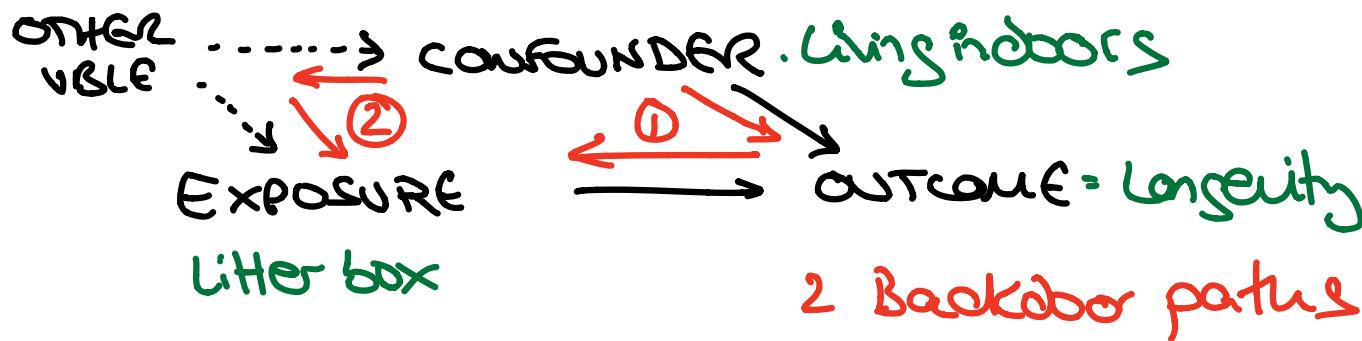
Error of "composition validity" → cohorts not comparable → biased effect measure estimates.

confounder: factor responsible for confounding.

CONFOUNDING present if EXP assoc w/ DIS.

as if all EXP effects were blocked/reversed.

Ex: cats w/ litterbox live longer. Actually die to living indoors?



- IN COHORT STUDIES:

① CASES w/o
OR_{cohort} ≠ 1

Living indoors Longevity
confounder ass. w/outcome
regardless of exposure (presence
litterbox use absence)

Confounder predictive of outcome
in unexposed ind.

② CASES w/e
OR_{CE} ≠ 1

Living indoors Litter box
confounder ass w/ exposure
(d-connected)

This can only be tested w/ data; if the assumptions of
RANDOM SAMPLING = COMPARABILITY of cohorts = True

• IN CASE-CONTROL STUDIES

① CASE w/o

$$OR_{co1} \bar{E} \neq 1$$

Living indoor

longevity

confounds ass. w/outcome

regardless of exposure (presence
littlebox use) / absence

confounder predictive of outcome
in unexposed ind.

② CASE w/e

$$OR_{CE} \neq 1$$

Living indoors

Litter box

confounds ass w/ exposure
(d-connected)



only applies inc-c studies if we know about
these associations in the source population;
we cannot look at this ass. in the data bc
we don't have a census of pop. cases and controls
but selectively sample non-cases

EXPOSURE

CONFOUNDER

CASES FE die young →

Litter box

Living indoors

CONTROLS FE die old →

we expect association



conditional on being selected in study

To examine the data for confounder-exposure association



② ALT C ASS w/ E

OR $CE | \bar{\delta} \neq 1$

living indoor

confounder ass. w/ exposure
regardless of outcome (present
longevity) [assence]

The C-E association (statistical dependence between the putative confounder and exposure) is not always identifiable from the data, because it requires comparing the distribution of the "confounder" in unexposed individuals with the distribution it would have had in the exposed individuals if exposure had been absent or blocked.

③ confounder not on causal path EXPOSURE → OUTCOME.

can't verify from the data

STATISTICAL TESTING OF CONFOUNDING.

It reflects :

- magnitude of association

- n° of observations → sufficient n° = everything is significant

should be avoided

The magnitude of confounding should have nothing to do with the size of the study.

Instead, look at changes in the estimates with confounder variable → if the estimate without changes > cut-off (5-10%) → significant association

COMPETING DEFINITIONS OF CONFOUNDING.

• COLLAPSIBILITY:

A measure of effect is confounded if the stratum-specific measures \neq crude ^{measures} collapsed

congve \neq results depending on effect measure of choice in the same pop!

(RR might show collapsibility but OR might not)

• COMPARABILITY:

confounding = invalid comparison of groups of ind.

No confounding = groups comparable \rightarrow absence of exp. \Rightarrow same incidence

Equivalent results between these two definitions may be obtained when using measures that compare averages of individual responses, such as incidence proportion ratios and differences, and rate ratios.

Odds ratios do not compare average responses in populations!

| This implies that collapsibility-based definitions of confounding will fail when using odds ratios unless the study outcome is rare; i.e., the odds ratio is interpretable as an incidence proportion or rate ratio.

| This criterion must hold in all strata of controlled variables.

Crude SRR = O/E = 656 / 258 = 2.54 = Adjusted SRR, = O/E = 656 / 258 = 2.54
 O: observed cases : 656
 E: expected cases : $\frac{258}{4000} \times 4000 \rightarrow$ exposed pop
 $\frac{258}{4000} \rightarrow$ pop. expected cases

O: 124 + 532 = 656
 E: $\frac{43}{2000} \times 2000 + \frac{215}{2000} \times 2000 = 258$
 NO confounding by gender

Table I—Example showing the association of corticosteroid (Co) treatment and canine urinary tract infection (UTI) by gender

Disease	Male		Female		Crude	
	Co	No Co	Co	No Co	Co	No Co
UTI	124	43	532	215	656	258
No UTI	1,876	1,957	1,468	1,785	3,344	3,742
Total	2,000	2,000	2,000	2,000	4,000	4,000
Incidence	0.062	0.022	0.266	0.108	0.184	0.065
Risk difference	...	0.041	—	0.159	—	0.100
Risk ratio	...	2.88	—	2.47	—	2.54
Odds ratio	...	3.01	—	3.01	—	2.85

Heterogeneity RD and RR } collapsibility = confounding
 Non-collapsibility OR } ignorability ≠

Crude SMR = $O/E = 1.775 \neq$ Adjusted SMR = $5120/8167.162 \approx 1.62 \Rightarrow$ confounding by gender

$O: 5120$
 $E: \frac{2,075}{6,000} \times 8,850 = 2,887.7$

$O: 5120 = 4,588 + 532$
 $E: \frac{1,859}{4,000} \times 6,350 + \frac{216}{2,000} \times 2,000 = 8,167.162$

Table 4—Numerical example of odds ratio collapsibility, but risk difference and risk ratio noncollapsibility

Disease	Male		Female		Crude	
	Co	No Co	Co	No Co	Co	No Co
utI	4,588	1,859	532	216	5,120	2,075
No utI	1,782	2,141	1,468	1,784	3,230	3,925
Total	6,350	4,000	2,000	2,000	8,350	6,000
Incidence	0.723	0.465	0.266	0.108	0.613	0.348
Risk difference	0.258	—	—	0.158	—	0.267
Risk ratio	1.55	—	—	2.46	—	1.77
Odds ratio	3.0	—	—	3.0	—	3.0

Co = corticosteroid; utI = urinary tract infection.

Heterogeneity RD and RR } collapsibility = confounding
 Collapsibility OR } comparability = "

Gender ass. w/ exposure
 disease incidence
 not \Rightarrow causal pathway

①
 ②
 ③
 ④
 ⑤
 requires
 confounder
 def.

Table 2—Numerical example demonstrating that control of a factor meeting the three criteria for confounding can produce bias

Type of individual	Male		Female		Crude	
	EXP	UNEXP	EXP	UNEXP	EXP	UNEXP
Disease* inevitable	100	200	67	133	118	139
Exposure causal	100	0	18	0	118	0
Immune to disease	100	33	164	328	264	361
Total	300	100	200	400	500	500
Incidence	0.67	= 0.67	0.18	0.18	0.47	0.28
Incidence ratio	exposure effect $\neq 0$	1.00	exposure effect $= 0$	1.00	— \neq —	1.70

* = urinary tract infection. EXP = exposed to corticosteroids; UNEXP = unexposed to corticosteroids.

Incidence in exposed:
 $100 + 100 + 18 + 13 = 0.47$

$200 + 200$

CAUSAL INCIDENCE

$$\frac{D^+ E^+ + D^- E^+}{500} = \frac{285}{500} = 0.57$$

$$\frac{D^+ E^-}{500} = \frac{118}{500} = 0.24$$

If you remove exposure, 118 no. would still be decreased.

118 D⁺ causal Exp

118 D⁺ indep Exp

264 D⁻

NON-COLLAPSIBILITY IR

Incidence proportion if there had been no exposure in the exposed group? $\frac{(0.995 \times 5000 + 0.0002 \times 1000)}{6000} = 0.829$

Confounding? ↴ Non-comparability

Incidence prop. in the unexposed group: 0.1660

effect modification operates on the measurement scale (multiplicative vs additive)

TABLE 1

Hypothetical cohort showing inadequacy of odds-ratio collapsibility as a criterion for ignoring a covariate

	Stratum 1		Stratum 2		Combined (crude)	
	Exposed	Unexposed	Exposed	Unexposed	Exposed	Unexposed
Diseased	4,999	995	5	1	5,004	996
Nondiseased	1	5	995	4,999	996	5,004
Total	5,000	1,000	1,000	5,000	6,000	6,000
Incidence proportion	0.9998	0.9950	0.0050	0.0002	0.8340	0.1660
Incidence ratio	1.00	\neq	25.0	\neq	5.02	
Incidence difference	0.0048	\neq	0.0048	\neq	0.6680	
Incidence odds ratio	25.1	=	25.1	=	25.2	

when you don't have disease the OR is not a good approximation of incidence ratio.

Greenland et al. 1988

COLLAPSIBILITY OR
HETEROGENEITY IR
NON-COLLAPSIBILITY ID

OR
collapsibility ≠ confounding
comparability P = "

Incidence proportion if exposure had been absent: $\frac{\frac{995}{1000} \times 1000 + \frac{1}{1000} \times 1000}{2000} = 0.498$

Odds Ratio Noncollapsibility

Incidence proportion in unexposed group = $\frac{996}{2000} = 0.498$

TABLE 2

Hypothetical cohort data, no effect modification; parameter is odds ratio; exposure and strata independent in joint source population

	Stratum 1		Stratum 2		Summary	
	Exposed	Nonexposed	Exposed	Nonexposed	Exposed	Nonexposed
Diseased	999	995	5	1	1,004	996
Nondiseased	1	5	995	999	996	1,004
Total	1,000	1,000	1,000	1,000	2,000	2,000
	$\Theta_1 = 5.02$		$\Theta_2 = 5.02$		$\Theta_s = \Theta^* = 1.02$	

IPR $\frac{999/1000}{995/1000} = 1.004 \neq \frac{5/1000}{1/1000} = 5 \neq \frac{1004/2000}{996/2000} = 1.008$

Grayson, 1987

IPD $\frac{999}{1000} - \frac{995}{1000} = 0.004 = \frac{5}{1000} - \frac{1}{1000} = 0.004 = \frac{1004}{2000} - \frac{996}{2000} = 0.004$

NON-COLAPSIBILITY OR HETEROGENEITY IPR

HOMOGENEITY IPD
NO CONFOUNDING!

Confounding? Effect modification?

TABLE 3

Hypothetical cohort data; effect modification present; exposure and strata independent; parameter is odds ratio

	Stratum 1		Stratum 2		Summary	
	Exposed	Nonexposed	Exposed	Nonexposed	Exposed	Nonexposed
Diseased	999	996	6	1	1,005	997
Nondiseased	1	4	994	999	995	1,003
Total	1,000	1,000	1,000	1,000	2,000	2,000
	$\Theta_1 = 4.01$		$\Theta_2 = 6.03$		$\Theta_s = \Theta^* = 1.02$	

OR: noncollapsibility, heterogeneity (note disease incidence)

RR: heterogeneity

RD: heterogeneity

SMR crude = SMR adjusted = 1.008

STANDARDIZATION OF RATES.

▪ **Incidence** tells us about the occurrence of new disease in a population (remember: numerator, denominator and specification of time are required).

▪ **Incidence rate (sometimes just called incidence density rate, hazard rate, or hazard for short)**: the denominator consists of the sum of different times each individual was at risk: that is, the “person-time” in the denominator. The numerator same as for cumulative incidence for same time period.

Age	Age proportion of both populations	Population 1 mortality rate per 100,000	Population 2 mortality rate per 100,000
Young	33.3%	7.75	26.25
Middle aged	33.3%	18.5	18.5
Elderly	33.3%	26.25	7.75
Overall	100%	17.5	17.5

- The age group-specific mortality rates vary enormously by population.
- But because although age is not associated with the populations (they have the same age distribution/structure), the differences cancel out.
- In this case, it would make more sense to report age-specific mortality rates for each population than to report overall mortality rates.

Age	Population 1		Population 2	
	Age proportion	Mortality rate per 100,000	Age proportion	Mortality rate per 100,000
Young	10%	10	70%	15
Middle aged	20%	20	20%	25
Elderly	70%	65	10%	70
Overall	100%	50.5	100%	22.5

- Note that in each of the three age groups, the mortality rate in Population 2 is greater than Population 1.
 - Young: 15 vs. 10 deaths per 100,000
 - Middle aged: 25 vs. 20 deaths per 100,000
 - Elderly: 70 vs. 65 deaths per 100,000
- Nevertheless, the overall mortality rate in Population 1 (50.5 deaths per 100,000) is greater than in Population 2 (22.5 deaths per 100,000).

AGE IS A CONFOUNDER

▪ Notice that AGE is associated with population: Population 1 is much older (70% elderly) than Population 2 (70% young).

▪ Notice that AGE is predictive of mortality: the rates in young individuals are much less than older individuals.

Age	Population 1		Population 2	
	Age proportion	Mortality rate per 100,000	Age proportion	Mortality rate per 100,000
Young	10%	10	10%	15
Middle aged	20%	20	20%	25
Elderly	70%	65	70%	70
Overall	100%	50.5	100%	55.5

▪ Now note the mortality ratios comparing Population 2 to Population 1:

$$\begin{aligned} 15/10 &= 1.500 \\ 25/20 &= 1.250 \\ 70/65 &= 1.0769 \end{aligned}$$

▪ The crude mortality ratio is $55.5/50.5 = 1.099$, correctly reflecting the higher mortality in Population 2 compared to Population 1.

AGE IS A CONFOUNDER

ELIMINATION OF CONFOUNDING.

- ELIMINATE ASSOCIATION BETWEEN ^{CONFOUNDER} RATE
"
- STANDARDIZATION

• DIRECT STANDARDIZATION :

- The reason for standardization is so one can compare populations using an overall summary (synoptic) rate measurement that are inherently "non-comparable" because of their inherent differences that influence the outcome rates.
- Standardization makes non-comparable populations "comparable".

- The directly standardized rate is motivated by the idea of determining what the crude disease or death rate would have been in a cohort if it had an age distribution like that of a standard population
- That is, we apply age-specific cohort rates to the standard age distribution.
- w_j represents the proportion of individuals in the standard population in the j^{th} age group.

Directly standardized rate $\hat{r} = \sum_{j=1}^J w_j * \hat{r}_j$

\hat{r}_j : no deaths / dis in j^{th} category

j : variable we standardize on (age)

w_j : no persons / time in j^{th} cat.

w_j : no or proportion of ind. in the
Standard pop in the j^{th} category

j^{th} category
stratum - spp
incidence
rate in
cohort.

- The result of direct standardization is an actual rate.
- The weights are supplied by the standard population.
- The rates are supplied by the populations (or cohorts) that are being compared.

- It is important to recognize that direct standardization involves taking *weighted* averages of age (and/or other stratum)-specific rates.

- If one is interested in the effect of some exposure in a particular population, then the "exposed" population should provide the weighting for standardization.

EXAMPLE

standard pop.

Age (yrs)	Dairy 1			Dairy 2		
	Population	Number of mastitis cases	Age-specific mastitis rate	Population	Number of mastitis cases	Age-specific mastitis rate
2-3	200	14	0.070	1,000	20	0.020
4-5	1,200	36	0.030	100	14	0.140
6-7	800	12	0.015	400	12	0.030
All ages	2,200	62	0.028 (crude)	1,500	46	0.031 (crude)

Age (yrs)	Dairy 1	Dairy 2	Combined	Percent (w_j)
2-3	200	1,000	1,200	0.324
4-5	1,200	100	1,300	0.351
6-7	800	400	1,200	0.324
	2,200	1,500	3,700	1

- Note that by using the SAME age distribution (weights) to get a weighted average of each population's age-specific mastitis rates, we remove the association between age and population, and hence remove confounding by age.

- The populations can now be compared.

Age (yrs)	Dairy 1			Dairy 2		
	Standard weight	Age-specific mastitis rate	$w_j * \lambda_j$	Standard weight	Age-specific mastitis rate	$w_j * \lambda_j$
2-3	0.324	0.070	0.0227	0.324	0.020	0.0065
4-5	0.351	0.030	0.0105	0.351	0.140	0.0491
6-7	0.324	0.015	0.0049	0.324	0.030	0.0097
Total			0.0381			0.0653

- Note that the ratio of the two age-adjusted (standardized) rates comparing Dairy 2 to Dairy 1 = $0.065 / 0.038 = 1.71$.
- In contrast, the ratio of the crude mortality rates = 1.11.
- The difference in these ratios is a consequence of confounding by age: Dairy 2 has younger cows, and the incidence of mastitis varies by age.

- INDIRECT STANDARDIZATION: when you want to compare the incidence rate in a cohort to the incidence rate in a standard pop, controlling for confounders.
- $\lambda_j^* = \frac{\lambda_j}{\eta_j^*}$
- CRUDE MORTALITY RATE \star :
- $\underline{\Sigma} \eta_j^* = \frac{D}{T} = \frac{\text{# deaths}}{\text{time}}$
- $\Sigma \eta_j^* = \frac{\text{ind}}{\text{time of all ind in } j\text{th stratum}}$

EXPECTED N° DEATHS/DIS: n° of ind we would expect to die in the cohort if it had the same death rates as \star ?

OBSERVED N° DEATHS

$$D = \sum d_j$$

$$\frac{\text{observed}}{\text{Expected}} = \frac{D}{E^*} = \frac{\sum d_j}{\sum \lambda_{j*}} = \frac{\sum d_j / E_{nj}}{\sum \lambda_{j*} / E_{nj}} = \frac{\text{crude rate in cohort}}{\text{weighted average of category-specific rates}}$$

weights are supplied by the cohort

E_{nj} → you can't compare across cohorts

rate supplied by the standard pop &

$$\lambda_{j*}$$

"Of course, there is no guarantee that this process will identify those diseases or causes of death that are most closely associated with the exposures. Cause-of-death specific rates for unexposed cohort members may be less than those in the general population, whereas rates for the exposed members are higher, and the two effects may cancel each other out when averaged over the entire cohort."

Example : healthy worker effect

EXAMPLE : compare crude mortality rates
and adjust for age.

Age	Davis			Walnut Creek		
	Pop ⁿ	Deaths	Rate	Pop ⁿ	Deaths	Rate
0-9	6,000	54	0.009	9,000	45	0.005
10-19	10,500	33	0.003	16,500	36	0.002
20-29	13,500	60	0.004	3,000		
30-39	18,000	36	0.002	9,000		
40-49	4,500	36	0.008	7,500	54	0.007
50-59	3,000	75	0.025	4,500	42	0.009
60+	4,500	126	0.028	7,500	45	0.006
All ages	60,000	420	0.007	57,000	315	0.0055

Walnut Creek = Exposed
Davis = Standard pop.

CRUDE MORTALITY RATE ratio : $\frac{0.0055}{0.007} = 0.79$

Observed = 315

Expected = $\sum n_j \lambda_j^*$

SMR = $\frac{315}{543} = 0.58 \neq 0.79$

	Davis	Walnut Creek	Age-specific expected values
Age	Rate (λ_j^*)	Pop ⁿ (n_j)	$\lambda_j^* n_j$
0-9	0.009	9,000	81
10-19	0.003	16,500	49.5
20-29	0.004	3,000	12
30-39	0.002	9,000	18
40-49	0.008	7,500	60
50-59	0.025	4,500	112.5
60+	0.028	7,500	210
All ages	0.007	57,000	543

Confounding by age.

	Zimbabwe			United States		
Age (yrs)	Deaths	Pop ⁿ	Rate	Deaths	Pop ⁿ	Rate
0-4	?	1,899,204	?	44,000	19,204,000	0.00229
5-24	?	5,537,992	?	45,000	72,244,000	0.00062
25-44	?	2,386,079	?	147,700	82,197,000	0.00180
45-64	?	974,235	?	368,800	46,751,000	0.00789
65-74	?	216,387	?	478,600	18,280,000	0.02618
75+	?	136,109	?	1,084,900	13,484,000	0.08046
Total	98,808	11,150,006	0.00886	2,169,000	252,160,000	0.00860

standard pop.

Age (yrs)	US rate (λ_j^*)	Zimbabwe pop ⁿ (n_j)	Expected = $\sum n_j \lambda_j^*$
0-4	0.00229	1,899,204	4,349
5-24	0.00062	5,537,992	3,434
25-44	0.00180	2,386,079	4,295
45-64	0.00789	974,235	7,687
65-74	0.02618	216,387	5,665
75+	0.08046	136,109	10,951
Total			36,381

$$\frac{\text{Observed}}{\text{expected}} = \frac{98,808}{36,381} = 2.72$$

The overall mortality rate in Zimbabwe would be 2.72 times greater than the mortality rate in the US if Zimbabwe had the same age-specific mortality rates that the US does.

The number of deaths observed in Zimbabwe is 2.72 times greater than the number of deaths that would be expected in Zimbabwe if Zimbabwe had the same age-specific mortality rates that the US does.

Not directly comparable \Rightarrow
 \neq age dist.
 Can't do direct
 standardization
 bc we don't
 have the age-specific
 rates in Z.

The effects of exposure are identical in both pop
 the difference in the overall SMR is that their age
Exercise 3 (Rothman) distribution weights are \neq

Age	General pop ⁿ	Exposed Cohort 1	Exposed Cohort 2
Young			
cases	50	50	5
person-time	100,000	10,000	1,000
incidence	0.0005/yr	0.005/yr	= 0.005/yr
SMR		10	= 10
rate ratio			1
Old			
cases	400	4	40
person-time	200,000	1,000	10,000
incidence	0.002/yr	0.004/yr	= 0.004/yr
SMR		2	= 2
rate ratio	$\frac{50+4}{(0.005 \cdot 10,000) + (0.002 \cdot 1000)} =$	1	
Overall SMR	$\frac{50+4}{(0.005 \cdot 1000) + (0.002 \cdot 10,000)} =$	7.7	$\neq 2.2$ Why!

