

# Recent Development in Epidemiologic Principles

Sept. 1975 - March 1976

(30 - 2 hour classes\*)

Class Notes

\* NB: 1 semester  $\cong$  35 hrs

# Recent Developments in Epidemiologic Principles - Outlines of Classes - 1975

9/3  
9/11

## Introduction - D+I

A. Course in Minn. - work at Harvard - Miettinen

B. Our objectives

1. communicate (reiterate) course at Minn.
2. place in perspective of our learning experiences
3. document our efforts - notes, etc.
4. apply concepts to our own work
5. evaluate their utility - modification

C. Context - "state of the art"

1. many ideas well reviewed are unpublished
2. Minn. course only scratched surface - no set of notes would be "complete"
3. there has been little published criticism of M's work

D. Format/Expectations

1. Background required - memo (little)
2. Format
  - a. meet Wed. 2:30-5:00 - w/ break - not 9/10
  - b. presentation & discussion of material

Time - develop concepts & analytical principles  
D. - apply concepts to quantitative exercises

  - c. readings - reference list - reserve
  - d. my handouts - first drafts

E. <sup>basic course (not complete)</sup> vocab. building - putative, intuit, obviate, quantitative

## II. General Content

A. What won't be covered

1. Descriptive/Experimental/Eval. epid -  
I AM will focus on statistical - ie etiological Ho's

- 2. Ho generation - ie, substantive issues - course is about principles + methodology
- 3. little emphasis on basic issues of study design - much will be assumed
- 4. Statistics, per se - tho we will make use of it

### B. General content

- 1. first part, mostly theoretical
- 2. Second part, mostly computational (applied)
- 3. outline of topics - handout - review
- 4. overall object: to address basic epid. issue of identifying + measuring causality (effect) - "obscured" by statisticians

\* common non-manualist + (w/)

### III. Measures of Disease Frequency - quantify occurrence of dis. (ie, epid.)

#### A. General definitions - other sheet

- 1. Absolute measures
- 2. Relative measures
  - a. rates
  - b. ratios

→ categorical data of non-disease characteristics (as well as disease)

#### B. Prevalence Rates - point / period

- 1. conceptual quantification
- 2. limitations in etiologic research
  - a. not risk - ie, includes new/existing cases
  - b. cause/effect ambiguity

#### C. Incidence Rates

- 1.  $I_c$  = risk - the concern when applying data to individuals
  - a. conceptual, quantif.
  - b. interpretation - referent time period
  - c. is/implies in
    - (1) genl - dis-specific
    - (2) all cause death - (not conditional)
    - (3) attack rates - referent time defined by epidemic
    - (4) case-fatality - den. is cases

\*

d. Selection into den. - PAR p.6

- exclusion of ... (Cervical Ca.)
- problem - PAR is elusive (eg)

Breast cancer

2.  $I_D$  - fundamental of dis. occurrence in pop.

- conceptual - picture "Sea of pop-time"
- quantify - den. = p-years
- compare to  $I_c$ 
  - limits -  $I_D$  may be  $> 1$  ( $< \infty$ )

3.  $I_D$  (as <sup>conceptual</sup> fundamental measure of dis. occur.) may be used to determine  $I_c$

a. exponential decay <sup>model</sup> - constant  $I_D$  over period

$$N = N_0 e^{-Rt} \quad R = I_D$$

b.  $I_c = 1 - N_t/N_0 = 1 - e^{-I_D \cdot t}$

c. if  $I_D$  not constant, subdivide t

$$I_c = 1 - e^{-\sum I_{Dj} t_j} \quad - t_j = \text{age range}$$

d. interpretation of  $I_c$  - over all  $t_j$

e. why compute  $I_c$ ? - conceptually

Simple & easy to communicate

4.  $I_D$  as misunderstood concept

a. common criticisms

- $I_D$  not good estimator of  $I_c$  ∴ inferior
- $I_D$  is assumed constant over obser. - but  $I_D$  is  $\neq$  if not, it can be stratified into subperiods

5. next time - relate P and  $I_D$

9/17 Outline

I. Review and Elaboration - Rates + Risks + Prevalence

- A. handouts - Miettinen's manuscript \$4.00
  1. Rates & Measures of Effect in Epid. Designs
  2. Adj. and Standardization of Rates and Measures of Effect
  3. Notations and Symbols
  4. Epid. Measures used to Quantify Illness (outline)

- B. last time defined freq. measure used to quantify disease occurrence - 3 important ones
  1.  $P = N_D / N$  - proportion [0-1]
  2.  $I_c = I / PAR$  - proportion [0-1] - also conditional  $P_t$  - time referent
  3.  $I_D = I / T$  - "force of mortality/morbidity" ("sort of instantaneous risk") [0-∞] to be considered fundamental measure of disease occurrence - ∴ can use  $I_D$  to derive  $I_c$  or  $P$

- C. first change in notations
  1.  $t$  = point in time
  - $p$  = period of time
  2.  $P =$  prev. rate (genl)
  - $P_D =$  " " of disease D

D. we had:

$$I_c = 1 - e^{-\sum I_{Dj} \cdot P_j}$$

where:  $P_j$  = hypothetical follow-up period of  $j$ -th (age) category - ie, age range of the category

$\sum p_j = p =$  total (hypothetical) period of follow-up - total age span  
 $I_{0j} = I_0$  for  $j$ -th stratum

(01)  $I_c = 1 - \prod (1 - I_{c_j})$

where:  $I_{c_j} = 1 - e^{-I_{0j} \cdot p_j}$

E. for example:  $I_j =$  # events observed in  $j$ -th category during follow-up

age $\sigma$	$I_j$	$T_j$ (p-yr)	$I_{0j}$	$I_{c_j}$
30-39	6	2000	.003	.030
40-49	15	3000	.005	.049
50-59	10	1000	.010	.095
Total	31	6000	.0052	X $\rightarrow$ [NB $I'_c = .144$ ]

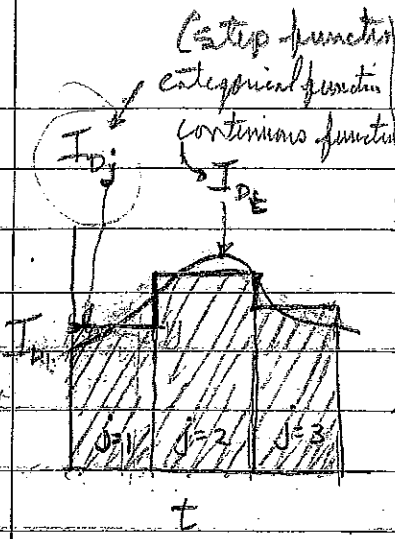
$I_c(30-59) = 1 - \prod (1 - I_{c_j}) = 1 - (1 - .030)(1 - .049)(1 - .095) = .165$

$I_c(30-59) = 1 - e^{-\sum I_{0j} p_j} = .165$  also

interpretation of  $I_c$ : condit. Pr.

\* F. 2 very important assumptions:

- $I_{0j}$  remains constant throughout each age interval - NB. this implies a categorical analysis (ie, categorical assumption)
  - Since we're assuming that  $I_0$  jumps suddenly at ages 40 + 50
  - $I_{0j}$  is defined at  $t=0$ , so we aren't considering a person of (eg) 39 changing  $I_{0j}$  if he is followed for over a year.



there is no way out of this unless we avoid the assumed consistency of  $I_D$  over time periods by: 
$$I_D = 1 - e^{-\int_{t_0}^{t_0+P} I_{DE} dt}$$
 (this is momentary  $I_D$ ) but this requires that we mathematically define the function  $I_{DE}$  in practice this is rather difficult. I'm not sure whether it is justified. We categorize into small enough intervals such that  $I_D$  remains constant.

2. the period-specific experience of our subjects must represent their full longitudinal experience - eg,  $I_D(30-39)$  must be the same as  $I_D(50-59)$  group was 20 yrs prior to the study.

(or) if the source pop were stable (ie, hypothetically) distrib. by age remains constant over time,  $I_{Dj}$  would also remain constant over time - this is quite analogous to

3. this assumption needed  $R_x$ : annual death rates for a pop. in order to derive a life expectancy at birth (in demographic life table analysis) - X-S experience must represent longitudinal experience.

G. this last assumption points out the dual nature of time -

1. calendar (chronological) time - eg, 1960-1970
2. age (eg) 30-59 years old

3. this feature allows us to convert an  $I_D$  (w/ no period referent) to a  $I_C$  (w/ a particular period referent) - this is done by constructing a hypothetical follow-up period for the  $I_C$  as defined by the age span of the subjects at the onset of the study

H.  $\frac{1}{2}$   $I_D$  is fundamental measure in epid. why compute  $I_C$  at all? 2 reasons

1.  $I_C$  is conceptually simpler & therefore easier to communicate

2. this provides an alternative to characterizing the  $x$ -year experience of a group w/ a disease when the ~~incidence rate varies~~ <sup>incidence rate varies</sup> considerably w/ age - otherwise we would compute an age-standardized rate which (like any "arranging" process) does not define any identifiable age category but instead suppresses info. Re: the variability of rates among age categories ( $\rightarrow$  limitation of using the mean) - While we will see the utility of standardization to control for conf. effects; it does not provide any "scientific" info. about the groups experience (as a single SR) i.e. SR must be compared to other SR's to be of value



## II. Relationship between $I_D$ and $P$ (NB. not a rigorous proof)

A. In order to use  $I_D$  to derive  $P$ , we

1. need to know something about the duration of the disease (ie, from clinical onset to termination) -  
let  $\{d_i\}$  = distrib. of durations for a given disease  $D$  -  $d_i$  is in years -  $d_i = 1$  etc.

Categorical are 1-yr intervals

2. let  $I_i$  = # new cases per year (same units as  $d$ ) of a given duration  $i$  -  
well assume:  $I_i$  remains constant over time (i.e. proof not rigorous)

B. we can make following relationships (see my paper for graphical illustration)

$$N_{(D)}i = I_i d_i = \# \text{ existing cases at 1 pt. in time who have a duration } i$$

$$\text{then: } \bar{d} = \frac{\sum I_i d_i}{\sum I_i}$$

$$\therefore \boxed{N_{(D)}} = \sum N_{(D)}i = \sum I_i d_i = \bar{d} \sum I_i = \boxed{\bar{d} I}$$

Note: relationship of freq. - not rates

C. In order to get relationship of rates:

$$I_D = \frac{I}{N - N_{(D)}}$$

(hard to intuit)  
 $I_i$  remains constant  
 $\therefore I_D$  is constant for each year  
 (ie, categorical assumption)

$$P = \frac{N(t)}{N} = \frac{\bar{d}I}{N} = \frac{\bar{d}I_0(N - N(t))}{N}$$

$$= I_0\bar{d} - \frac{I_0\bar{d}N(t)}{N} = I_0\bar{d} - I_0\bar{d}P$$

∴  $P(1 + I_0\bar{d}) = I_0\bar{d}$

$$P = \frac{I_0\bar{d}}{1 + I_0\bar{d}}$$

$$\frac{P(1 + I_0\bar{d})}{P + P I_0\bar{d}} = \frac{I_0\bar{d}}{I_0\bar{d}}$$

$$I_0(\bar{d} - P\bar{d}) = P$$

$$I_0 = \frac{P}{\bar{d} - P\bar{d}}$$

Note:  $P \neq I_0\bar{d}$  (as is sometimes assumed; McMahon p 66)

D. Miettinen expresses this relationship as an interval (with a more general formulation) & points out the following limits not immediately apparent from the above formula:

$\lim_{\bar{d} \rightarrow 0} P = 0$       disease is extremely short in duration

$\lim_{\bar{d} \rightarrow \infty} P = I_0$  (from birth to age  $t_0$ )      other extreme: disease is irreversible & non-fatal

### III. Adjustment and Standardization of Rates -

A. As mentioned before that (SR) is a specific transformation (or averaging) of a stratum-specific rates, to form an overall rate that is used to control for the confounding effects of known factors such as age, race, sex, etc. - & to that extent, forms the heart of epid. methodology

B. Miettinen has actually reconceptualized standardization to carefully define its utility and limitations

c. But we will begin with the common formulation of 2 separate methods of rate standardization:

a. Direct Method - DR

category-specific rates observed in the study pop. are applied to, (ie, multiplied by) the corresponding category sizes of the "standard" pop (which is an independently observed distribution)

$$DR = \frac{\sum n_{sj} R_j}{\sum n_{sj}} = \frac{\sum n_{sj} R_j}{n_s}$$

b. Indirect Method - IR

category-specific rates of the "standard" pop. are applied to the corresponding subgroup sizes of the study group

	D	D	
study	a	b <sub>j</sub>	n <sub>j</sub>
standard	a <sub>s</sub>	b <sub>s</sub>	n <sub>s</sub>

$R_j = a_j / n_j$

(c)  $R = \frac{\sum a_j}{\sum n_j} = \frac{a}{n}$

$R_{sj} = a_{sj} / n_{sj}$

(e)  $R_s = a_s / n_s$

$$IR = \frac{\sum n_j R_j (R_s)}{\sum n_j R_{sj}} = \frac{a (R_s)}{E[a]}$$

SMR · R<sub>s</sub> ← crude rate in Standard

D. Problems w/ common methods

1. concept of "stand. pop" is ambiguous - in DR, it is where the n's come from ; in IR, it is where the R's come from
2. it appears that standardiz. is actually 2 different methods, that are in one sense opposites - but it will be maintained that a standardiz. of (a study) distribution according to a correlate of rate (eg, age) has only a single conceptualization - i.e. there is only one type of standardiz. (DR) which we will see has several forms, both for rates and measures of effect (RR)

E. We will first define 2 terms

1. referent group (or pop) - the series used to compare a given study (test-exposed group - ie, a basis of comparison - eg) E
2. Standard pop. - the series from which we get category-specific "weights" that we will apply to the corresponding rates of the test group, in order to derive an overall (stand.) rate

F. Now define 3 types of rates:

1. CR = no stratification by a potential confounding var. is modifying factor - we can consider num. + den. to be total numbers observed

$$CR = \frac{\sum a_j}{\sum n_j} = \frac{a}{n}$$

this also might be seen as having "latent weights" = the stratum specific sample sizes of the study group. - i.e

$$CR = \frac{\sum n_{j1} R_j}{\sum n_j} = \frac{\sum n_j \cdot \frac{a_j}{n_j}}{\sum n_j} = \frac{a}{n}$$

2. AR = includes stratification by potential conf. var. - and incorporates arbitrary weight

$$AR = \frac{\sum w_j R_j}{\sum w_j} \quad (\text{general formula})$$

eg,  $w_j = 1000$

3. SR = special type of AR in which the weights are derived from the experience of a real group (ie, from the data) -  $\therefore$  standardization is an adjustment to a common distrib. (ie, stand. pop.) w/ respect to a potential confounding factor - same formula as AR

G. So Stand. is a form of weighting rates -

Being more specific re: SR, we may define 2 major classes according to the source of the weights

- SR(I) 1. internal stand. → weights come from stratum-specific sample sizes of the study group -  $w_j = n_j$   
 but note: this is equal to the CR, which has  $n_j$  as implicit rates
- SR(E) 2. external stand → weights come from a group other than the study group (E)
- a. Standard = referent  
 $\therefore w_j = n_{Ej}$  (as  $d_j$  in  $e/\bar{e}$ )
  - b. Standard = another category of E (E')
 $\therefore w_j = n'_{j}$  (as  $d_j$  in  $c/\bar{e}$ )
  - c. Standard = a general pop. (eg) U.S. (S)
 $\therefore w_j = n_{sj}$

H. For example: Same data as before + (referent group)

	Study group (E)			referent group (E')		
	$T_{(E)}$	$I_{D(E)}$	$T_{(E)} I_{D(E)} = I_{E'}$	$T_{(E')}$	$I_{E'}$	$I_{D(E')}$
30-39	2000	.003	6	3000	5	.0017
40-49	3000	.005	15	2000	5	.0025
50-59	1000	.010	10	1000	5	.0050
total	6000	.0052	31	6000	15	.0025

$$DR = \frac{\sum T_{(E)j} I_{D(E)j}}{\sum T_{(E)j}} = \frac{9 + 10 + 10}{6000} = \boxed{.0048}$$

NB. • referent group is considered "Standard",  
 • weights are similar for  $I_c$  as  $I_D$

$$IR = \frac{\sum T(E)_j I_{D(E)_j}}{\sum T(E)_j T_{D(E)_j}} \times I_{D(E)} = SMR \cdot I_{D(E)}$$

crude rate

$$= \frac{31}{(3.33+7.5+5)} (.0025) = \boxed{.0049}$$

\* might mention that in IR, it is customary to pool all exposure categories to get what is commonly considered a "standard" pop. - but we see here that the standard is actually the E group (where the weights come from) - thus the common pooling technique is actually forming a pooled referent, which represents a rather ambiguous category since it is dependent on the specific distrib. of the pop. - not a known category

$$SR(E) = \frac{\sum T(E)_j I_{D(E)_j}}{\sum T(E)_j} = \boxed{DR} = \boxed{.0048}$$

(where standard = referent)

$$SR(I) = \frac{\sum T(E)_j I_{D(E)_j}}{\sum T(E)_j} = \frac{31}{6000} = \boxed{.0052} = GR$$

∴  $SR(I) \neq IR$   
 $SR(E) = DR$   
st = ref.

While it appears that we have simply eliminated the IR from our bag of analytical procedures, we will see next week that, using the method of SR(I), IR is highly related to the Standardization of effect measures (not of rates)

I. An SR (of any type) is a hypothetical measure and has no meaning unless compared to other SR's (for other groups) that incorporate the same standard. But first we must discuss the general nature of rate comparisons to see how we measure the association between 2 variables - (E → D)

#### IV. Measures of Association - 3 general types (see other outline)

- A. Non-effect Measures of Assoc.
  - 1. Correlations (concurrent validity)
  - 2. Causation measures (instrumental validity)
- B. Measures of Effect
  - 1. Absolute Measures of Effect
  - 2. Relative " " Effect
- C. Potential Impact Measures
  - 1. for causal factors (EF)
  - 2. for protective factors (PF)



V. To be discussed Next time

A. Measures of Assoc. in more detail  
(esp. B + C)

1. ~~it~~ might be obvious now why "relative risk" is not used by Miettinen since it does not refer to the ratio of 2 I<sub>0</sub>'s - similarly w/ "attributable risk"
2. in C/E study, we estimate RR by calculating an OR - in the past, this approx. was considered valid only if we assumed a rare disease - in fact Miettinen makes this assumption in his manuscript but it is not necessary (→ my paper + his) - based on the distinction between I<sub>0</sub> and I<sub>c</sub> - measures of assoc. are compared for various study designs in latter portion of paper (I<sub>0</sub>R, I<sub>c</sub>R, OR,  $\phi$ )

B. Standardizing of measures of effect:  
(RD and RR)

eg) 
$$\frac{SRR(I)}{SRR(E)} = SMR \quad (ie, SMR \neq SR)$$

9/24 - 10/1 - 10/7

# Outline

## I. Introduction

- A. discuss class schedule, pace, etc. - 2:30 - 4:30
- B. Debbie - exercises (later in class if time)
- C. handouts: Quality of Data & Findings (Study Design) before 10/8
- D. topics today
  1. measures of association - concentrating on 3 measures of effect
  2. standardization of these 3 measures

## II. Types of measures of association (ie, indications of criterion validity) (see other outline #4)

### A. Non-effect measures of assoc.

1. Correlations (concurrent validity → criterion is another factor, besides the predictor was. at issue; and both are measured at same point in time <sup>+ independently</sup>) - used in prevalence studies -
  - a. of interest to us: nominal data - generally when analyzing a contingency table, a  $\chi^2$  statistic (of indep. or homogeneity is calculated) - but  $\chi^2$  only indicates the existence of an assoc., not its strength (or magnitude). in 2x2 table

predictor → criterion →

type of correl. coef depends on kind of data:

(genl) 
$$\chi^2 = \sum \sum \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

$$\hookrightarrow \chi^2 = \frac{n(ad-bc)^2}{m_D m_B n_E n_{\bar{E}}}$$

but notice that the statistic is dependent on n; the larger the n, the larger the  $\chi^2$  thus the smaller the p-value

b. We can form a correl. coef ( $\phi$ ) which measures the strength of the assoc. -  $\phi = [-1, +1]$

$$\phi = \sqrt{\frac{\chi^2}{n}} = \frac{ad - bc}{\sqrt{m_D m_B n_E n_{\bar{E}}}}$$

c. actually this is a Pearson correl. coef ( $r$ ) applied to nominal data, assuming any numbers for the categories

d. Significance level is found w/ p-value of the corresponding  $\chi^2_{(1)}$

e. for  $m \times n$  table: use Contingency coef. ( $c$ )

$$C = \sqrt{\frac{\chi^2}{n + \chi^2}} \quad \text{range: } \left(0 - \sqrt{\frac{k-1}{k}}\right)$$

where:  $k = m = n$

$\phi$  &  $\Phi$  are only stable when the researcher has not arbitrarily fixed any of the marginals:

	D	$\bar{D}$	
E	a	b	$n_E$
$\bar{E}$	c	d	$n_{\bar{E}}$
	$m_D$	$m_{\bar{D}}$	n

thus,  $\phi$  can only be used w/ type II sampling (ie X-S selection of a pop. -  $\therefore$  we know the sampling fraction(s))

with:  $(m-1)(n-1)$  degrees of freedom for the  $\chi^2$  statistic

2.  $\Phi$  Correctness of Measurement - (instrumental validity  $\rightarrow$  criterion is other known measure of the same factor, measured at same point in time) - also for X-S sampling; but generally info. in correl. coef. is separated into  $\geq 1$  measure:

- a. Sensitivity
  - b. Specificity
  - c. Predictive Value
- } screening tests

B. Measures of effect - (predictive validity  $\rightarrow$  criterion is another factor measured after the predictor var. (in calendar time))

B.11. Absolute measures of effect

a. rate difference = RD = R<sub>E</sub> - R<sub>E'</sub>

b. excess mortality = a - E[a]

2. Relative Measures of Effect

a. relative rate (or rate ratio) = RR =  $\frac{R_E}{R_{E'}}$

b. excess relative rate ("relative effect") = ERR = RR - 1

C. Potential Impact Measures

1. for causal factors - ARP = % of cases that could have been prevented if persons were unexposed

2. for protective factors - (PF) = % of potential cases that might have been prevented

III. Measures of of Effect

A. RD

1. I<sub>C</sub>D = I<sub>C(E)</sub> - I<sub>C(E')</sub>; when period of follow-up is same for both rates

2. I<sub>D</sub>D = I<sub>DCE</sub> - I<sub>DCE'</sub>; when follow-up is comparable

3. limits (of both) for causal exposure are:

0 ≤ I<sub>C</sub>D ≤ ∞

0 ≤ I<sub>D</sub>D ≤ 1

(10/1)

B. RR

1. I<sub>C</sub>R = I<sub>C(E)</sub> / I<sub>C(E')</sub>

2. I<sub>D</sub>R = I<sub>DCE</sub> / I<sub>DCE'</sub>

3. limits of both for causal factors:

3. a.  $1 \leq I_c R < \infty$   
 $1 \leq I_b R < \infty$

b. also:  $\lim_{p \rightarrow \infty} I_c R = 1$  ie, if we could follow everyone for an infinite period of time & they didn't die of any other cause, eventually, everyone would die of the disease at issue —  $\therefore 1$  is lower limit of  $I_c R$  for chronical exposure

c. also:  $\lim_{p \rightarrow 0} I_c R = I_b R$

ie, assuming  $I_{bc}$  is a continuous function, as  $p$  approaches 0, the  $I_c$  for each group becomes  $I_b \cdot p$  — thus:

$$\lim_{p \rightarrow 0} \frac{I_c(E)}{I_c(E)} = \frac{I_b(E) p}{I_b(E) p} = I_b R$$

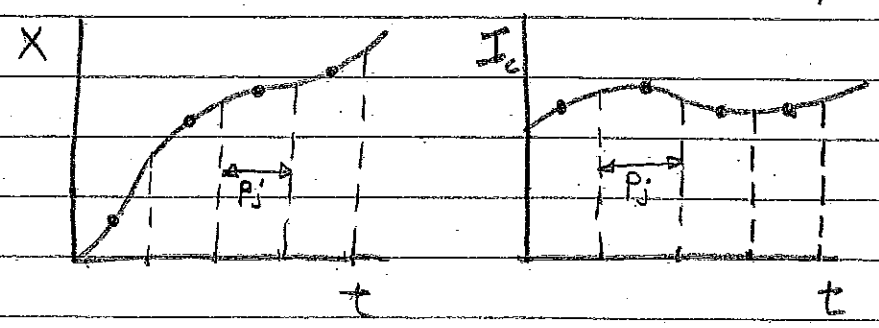
d. putting these 2 limits together:

$$1 \leq I_c R \leq I_b R < \infty$$

if  $I_b$  is constant over time

e. But the constancy of  $I_b$  (for any  $p$ ) has important implications to epid. research — we can think of  $I_{bc}$  as analogous to velocity ( $V_c$ ); and

categorical  $I_D$  (ie,  $I_{Dj}$ ) as analogous to average speed ( $\bar{S}$ ) for time period  $p$  (NB:  $\Delta t = p$ ) - distance =  $X$  -  $j = \text{time period}$



$$\bar{S}_j = \frac{\Delta X_j}{P_j}$$

$$I_{D,E} = \lim_{P_j \rightarrow 0} I_{Dj}$$

$$V_t = \lim_{P_j \rightarrow 0} \bar{S}_j = \text{slope}$$

$$I_{D,E} = \lim_{P_j \rightarrow 0} \frac{I_{D,j}}{P_j}$$

$$V_t = \lim_{P_j \rightarrow 0} \frac{\Delta X_j}{P_j}$$

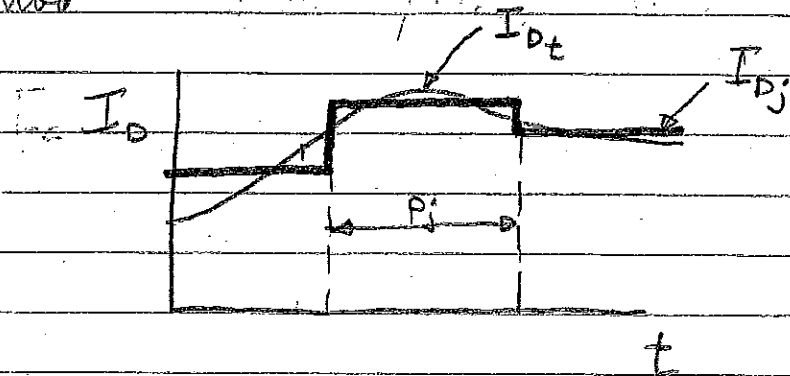
$$I_{D,E} = \lim_{P_j \rightarrow 0} \frac{\Delta I_{D,j} / N_j}{P_j}$$

Despite this analogy, the actual situation in physics is quite different from that of Epid. In physics, we have laws (thanks to Newton, etc.) which describe the mathematical relationship between  $V$  and other variables - eg,

$$V_t = \sqrt{2gX} \quad (\text{for a free-falling body encountering no friction})$$

But in epid, we have no such laws of nature describing  $I_{D,E}$ . So we approximate it with a categorical measure

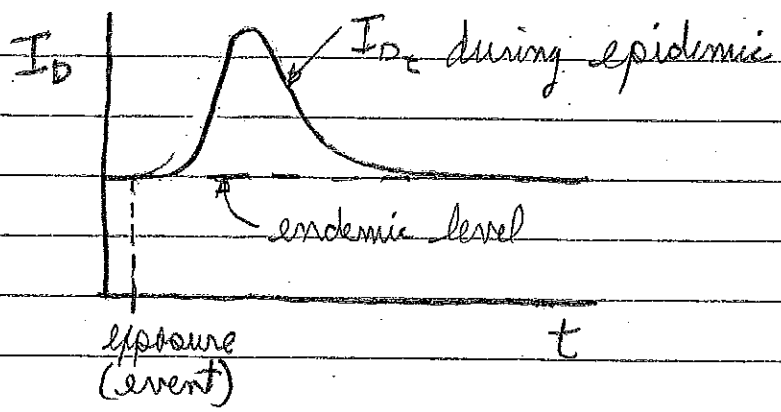
(using person-years of follow-up) - just like we use  $\bar{S}$  to characterize an assumed constant velocity over a time period -



$$\text{where: } I_{c(p)} = 1 - e^{-\int_{t_0}^{t_0+p} I_{D,t} dt}$$

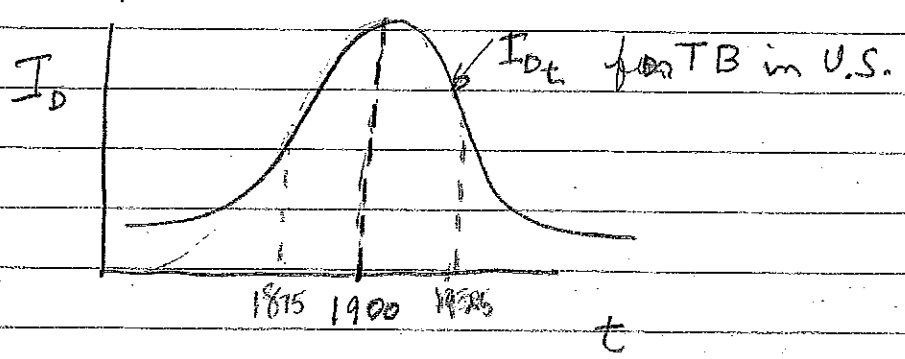
$$\text{or: } I_{c(p)} = 1 - e^{-\sum I_{D,j} P_j}$$

f. Thus Epid. needs a Newton; but is there any evidence that  $I_0$  follows any laws of nature. The key to using (categorical)  $I_0$  is the assumption of its constancy for  $p_j$  and its constancy over time w/in age categories. We have already seen that  $I_0$  should not remain ~~so~~ uniform over wide age spans, but is there always reason to believe that  $I_0$  remains constant over calendar time. For example, imagine a short term epidemic in which exposure to contaminated food at one meal causes a sudden outbreak.



Thus, there would be a short-term time variation within age bands.

But we know that even <sup>deaths from</sup> chronic diseases follow long-term epidemic trends, that pose a similar phenomenon (very crudely)



In fact, we even know (according to Dr. C) that this curve in any pop. shows a basic consistency for many diseases: the ratio of F/M and young/old is greater on the rise (when slope  $\geq +$ ) than on the decline. It is possible that the curve may be described by sociological/environmental variables in the society. —



One researcher has actually claimed that he can predict (or determine) when a pop. is chronologically in relation to its industrial revolution by the ratio of M/F TB deaths per year. Perhaps there are possibilities for describing some general principles of disease occurrence that would enable us to predict  $I_{D_2}$  on a basis independent of observing disease occurrence (post hoc)

4.  $OR = O_{(D)} / O_{(\bar{D})}$  (usually computed in c/c study, but may compute anytime)

$$Pr[\text{event}] = \frac{[\# \text{ favorable outcomes possible}]}{[\# \text{ total outcomes possible}]} = \frac{1}{2}$$

$$O[\text{event}] = \frac{[\# \text{ fav. outcomes possible}]}{[\# \text{ unfav. " " }]} = 1$$

for cases: Odds of exposure =  $O_{(D)}$

$$O_{(D)} = Pr[E|D] / Pr[\bar{E}|D] = \frac{a}{a+c} = \frac{c}{c+d}$$

$$= \boxed{e_D / 1 - e_D}$$

for non-cases:

$$O_{(\bar{D})} = Pr[E|\bar{D}] / Pr[\bar{E}|\bar{D}] = \frac{b}{b+d} = \frac{d}{b+d}$$

$$\boxed{e_{\bar{D}} / 1 - e_{\bar{D}}}$$

$$\text{So, } OR = \frac{O_D}{O_{\bar{D}}} = \frac{e_D(1-e_{\bar{D}})}{e_{\bar{D}}(1-e_D)} = \frac{ad}{bc}$$

(which is a ratio of 2 ratios)

It has also been shown using Bayes' theorem that in a  $10/\bar{0}$  study

$$OR \cong IER$$

if  $P(D)$  is small ( $< .20$ )

5. In a  $X-S$  study (w/ type II

Sampling)

$$OR_p = \frac{P(E)(1-P(\bar{E}))}{P(\bar{E})(1-P(E))} = \frac{e_D(1-e_{\bar{D}})}{e_{\bar{D}}(1-e_D)} = OR$$

if Cases +  $\bar{C}$  of both types of studies are from same pop.

Note:  $OR_p$  (in  $X$ -sectional) cannot technically be used as a measure of effect since we don't know the temporal relationship

Note:  $OR_p \neq PR = \frac{P_D}{P_S}$

PR is often used to estimate RR - eg, Slomes manual

insert  
V  
p.9

b. Derive: (see next insert)

$I_D R$  in c/e study

- a. for existing cases
- b. for new cases (+ advantages)

General conclusions from my paper #1

- c/e study
- a.  $OR = I_D R$  (when d's are equal regardless of  $P_D$  or  $P_S$ )
  - b.  $OR = I_C R$  (only when disease is rare, ie  $P_D < .20$ )

c. But despite the lack of inherent quantitative flaws in estimating effect in c/e study, we must remember the logical flaws that accompany the specific study design - (eg)

- (1) Biased recall of exposure
- (2) Selective survival of all cases

d.  $I_C R < I_D R$  + this difference widens w/ incr.  $P_D$

e. OR is quite ineff. in estimating RR in: CH (when:  $.10 > P_D > .20$ )  
X-S (when:  $.25 > P_S > .50$ )

(ie, CI are very wide)  
also: OR is unstable (invalid) for CH study (Type I sampling) for  $P_D > .20$

f. OR (ingent.) is always ineff at.  
 $P_{(E)} < .05$

$$P_{(D)} < .10$$

g.  $\phi$  is unstable (invalid) for:

$$CH: .10 < \phi < .70$$

$$CE: \phi > .10$$

but:  $X-S$  — fairly invariant (assuming  
 type II sampling)

#### IV. Potential Impact Measures

A. ARP — % of the disease that could potentially  
 have been prevented if no one had been exposed

1.  $ARP_E$  (for exposed cases only)

$$ARP_E = \frac{RR-1}{RR} = \frac{a - E[a]}{a}$$

(potentially)  
 = % of exposed cases that would have  
 been prevented if no one had been  
 exposed →

$$2. ARP_N (\text{pop. A.R.P.}) = e_D ARP_E = \frac{a - E[a]}{a + c}$$

= % of all cases that would have been  
 prevented if no one had been exposed

3. ARP requires that cases be representative  
 of the pop. in order for it to apply to  
 the pop. as an "expected prevention rate"  
 (ie) ARP applies to a specific pop. at a  
 period in time (unlike  $RR + RD$  which  
 refers to the nature of a disease)

4. ARP pertains only to observ. studies — as  
 an expected % of preventable cases, it is  
 like a  $R_r$ , not a measure of prevention

or treatment effect that one would estimate in an experiment - ∴ ARP is a measure of assoc. (like RD or RR)

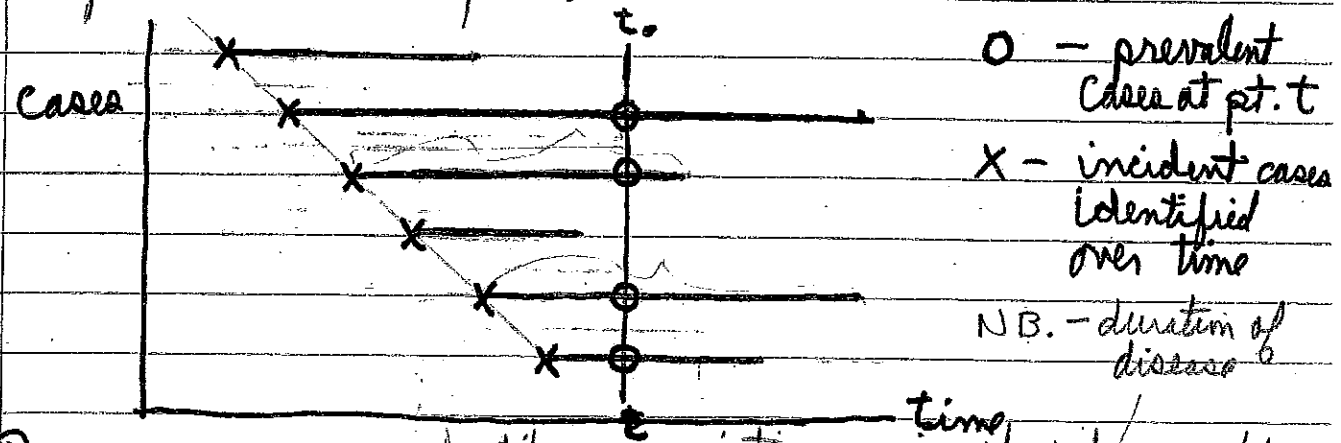
- 5. 2 additional assumptions needed to regard ARP as a prevention measure;
  - a. E can be Δ'ed in source pop. - ideally all E → Ē
  - b. changing a person from category E to Ē will place him at the same risk (rate) that previously unexposed persons had

VI. (next page)

last time discussed Measures of Assoc. moved too quickly over estimation of effect in e/E study, so I'll go over it in more detail

# IV. Estimating the I<sub>D</sub>R in a c/e study - Review in more detail

a. define incid & prev. cases -



- ① prev. cases - identifying existing cases that ~~to~~ might have had the disease for any # of years (NB. dis not any duration of cases at time t)
- ② incid. cases - new (fresh) cases identified as they occur over time -

to get  $\hat{I}_D R$

①  $N(t) = I \bar{d}$  note what  $\bar{d}$  is : from onset to termination

$$\bar{I}_D = \frac{I}{N - N(0)} = \frac{I}{N_B}$$

$$P = \frac{I_D \bar{d}}{1 + I_D \bar{d}}$$

$$I_D = \frac{P}{\bar{d}(1-P)}$$

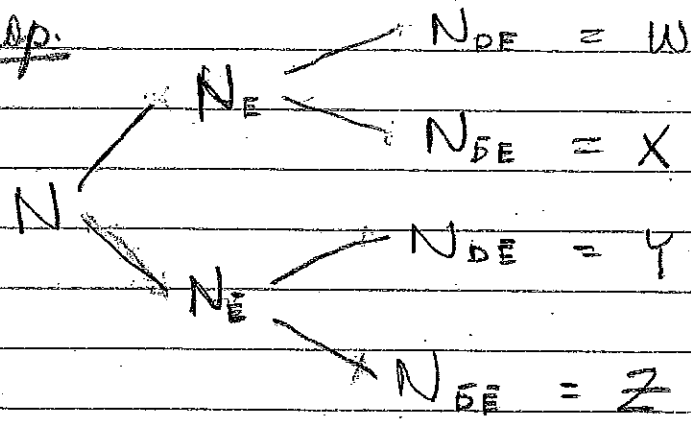
$$I_D R = \frac{I_{D(E)}}{I_{D(\bar{E})}} = \frac{P_{(E)}(1-P_{(E)})}{P_{(\bar{E})}(1-P_{(\bar{E})})} \cdot \frac{\bar{d}_{(\bar{E})}}{\bar{d}_{(E)}} = OR_P \cdot \frac{\bar{d}_{\bar{E}}}{\bar{d}_E}$$

X-S of pop.

(10)

	P	D
E	W	X
E	Y	Z

in a pop.



$$P_{(E)} = \frac{W}{W+X}$$

$$P_{(\bar{E})} = \frac{Y}{Y+Z}$$

X-S of pop.

	P	D

$$OR_{D|E} = \frac{\frac{W}{W+X} \left(1 - \frac{Y}{Y+Z}\right)}{\frac{Y}{Y+Z} \left(1 - \frac{W}{W+X}\right)} = \frac{\frac{W}{W+X} \frac{W+Z}{Y+Z}}{\frac{Y}{Y+Z} \frac{Y+X}{W+X}} = \frac{W}{W+X} \frac{W+Z}{Y+Z} \frac{(W+X)(Y+Z)}{Y(W+X)}$$

$$= \frac{W(Y+Z) - WY}{Y(W+X) - YW} = \frac{WY + WZ - WY}{YW + YX - YW}$$

$$= \boxed{\frac{WZ}{YX}}$$

$$\therefore I_{DR} = \frac{WZ}{YX} \frac{d_{DE}}{d_{\bar{D}E}}$$

E values in c/c study: where we have sampled cases + noncases w/ certain sampling fractions

note both  $F$ 's are equal in each column if we assume no selective bias

	D	$\bar{D}$
E	$F_D W$	$F_{\bar{D}} X$
$\bar{E}$	$F_D Y$	$F_{\bar{D}} Z$

$$F_D = \frac{\# \text{ cases in study}}{\# \text{ cases in source pop.}}$$

$$F_{\bar{D}} = \frac{\# \text{ noncases in study}}{\# \text{ noncases in source pop.}}$$

$$OR = \frac{ad}{bc} = \frac{F_D W \cdot F_{\bar{D}} Z}{F_{\bar{D}} X \cdot F_D Y} = \frac{WZ}{XY} = OR_P$$

$$I_{DR} = \left( \frac{ad}{bc} \right) \frac{\bar{d}_E}{d_E}$$

Thus,  $OR = I_{DR}$  if  $\bar{d}_E \approx d_E$   
(w/ prev. cases)

② for incident cases - don't need to know about durations - but NB:  $\bar{d} \neq 0$  (it's the same as before)  
 nor can we say:  $d_E = \bar{d}_E$ , simply because we have incident cases  
 but we can say:

	D	$\bar{D}$
E	$F_D I_E$	$F_{\bar{D}} N_{DE}$
$\bar{E}$	$F_D I_{\bar{E}}$	$F_{\bar{D}} N_{\bar{D}\bar{E}}$

same as before

⊗ if we assume  $I + N_B$  remain stable over time of study



from before:  $I_D = \frac{I}{N - N_0} = \frac{I}{N_D}$

for exposed:  $I_{D(E)} = \frac{F_D I_E}{F_E N_{DE}}$

for unexposed:  $I_{D(\bar{E})} = \frac{F_{\bar{D}} I_{\bar{E}}}{F_{\bar{E}} N_{D\bar{E}}}$

$\therefore I_{DR} = \frac{F_D I_E / F_E N_{DE}}{F_{\bar{D}} I_{\bar{E}} / F_{\bar{E}} N_{D\bar{E}}} = \frac{I_E N_{D\bar{E}}}{I_{\bar{E}} N_{DE}}$

OR =  $\frac{ad}{bc} = \frac{F_D I_E F_{\bar{E}} N_{D\bar{E}}}{F_{\bar{D}} N_{DE} F_D I_{\bar{E}}} = \frac{I_E N_{D\bar{E}}}{I_{\bar{E}} N_{DE}} = I_{DR}$

$\therefore$   $OR = I_{DR}$  w/o any assumption about d

So advantages to using incident cases

- (1) (1) meth - don't need to know durations to estimate I<sub>DR</sub>
- (2) ~~methodological~~ methodological adv. (inference)
  - (a) having new cases <sup>helps</sup> eliminates possibility that D predated E - ie, temporal relationships
  - (b) also reduce selection bias. Since cases have had not had a <sup>much of</sup> chance to die from the dis. ie, reduce selective survival

So what are practical adv. of this:

- (1) don't need to make assumption R:  
rare dis. — but usually no big deal
- (2) imp. adv. to using new (incid.) cases in  
C/E study

Repeat: C/E study

- (1) no inherent meth. flaw in estimating IR
- (2) but several important methol. problems  
that are of particular relevance in  
C/E studies (differential bias)

- (a) biased recall of exposure history
- (b) selective survival of all cases  
(esp. prev. cases)

- (c) often difficult to attribute cause +  
to effect (esp. w/ prev. cases)  
& with chronic diseases w/ long  
induction time

So, practical advantages of this theory:

- (1) We don't need to make assumption of rare disease - but in most situations, this isn't an added benefit because most conditions we study are rare ( $< .20$ )
- (2) Using incidence cases in a C/E study allows us to gain math. + method. advantages

To repeat - how this affects the utility of a C/E study in practice

- (1) no inherent math flaw in estimating I<sub>D</sub>R (esp. w/ incident cases)
- (2) but several important methodological problems that are of particular relevance to C/E designs

- a. might be biased recall of exposure history by cases + non-cases
- b. selective survival of all cases (esp. w/ prev. cases)
- c. often difficult to attribute cause + effect relationship (esp. w/ prev. cases and/or chronic diseases that have long induction periods) - this is more of a quantitative problem since the I<sub>D</sub>R might vary according to the time that the exposure occurred

## VI Standardization

A. Review:  $SR_{\text{app}}$

1. adjusted rate = AR =

$$\frac{\sum w_j R_j}{\sum w_j}$$

2. 2 types of AR

a. general adj. - when  $w_j$  are chosen arbitrarily

b. standardized - when wts. are taken from empirical data source - 2 types:

i.  $SR(I) = \frac{\sum n_j R_j}{\sum n_j} = \frac{a}{n} = CR$

(ie, wts come from study group)

ii.  $SR(E) = \frac{\sum w_{sj} R_j}{\sum w_{sj}}$

in which  $w_{sj}$  come from:

- (1) referent  $\rightarrow$   $SR(E) = DR$
- (2) another category(s) of E - i.e.  $E'$  (besides  $E + \bar{E}$ )
- (3) general pop. - (actually a general case of  $SR(E)$ )

3.  $IR = SMR \cdot R_E \neq SR$   
 (Only 1 way to standardizing a rate - the "direct method")

B. Standardization of RR

1.  $SR$  is meaningless by itself - i.e. must be compared to another  $SR$  - when both  $SR$ 's have the same standard pop. (i.e. a common distribution) we call the ratio of the 2  $SR$ 's an SRR.
2. At this point, mention difference in terminology between my use of the terms "internal" & "external" vs. M.'s - I refer to (I) and (E) as an indication of the source of weights for each  $SR$  - for an SRR, the terms internal & external refer to the source of weights relative to the test (exposed) group - M. only uses (I) and (E) to refer to the type of RR standardization (not SR)  
 Thus: (M)
  - a. internal <sup>standardized</sup> both component rates ( $R_E + R_{E'}$ ) have been adj. to the same distrib. (in my terms, the defn. of SRR)
  - b. external <sup>stand.</sup> in addition, the standard is not the study pop. - i.e.  $SR_E$  is not crude

### 3. Comparison of terminology & defns.

type	mine	Miettinen's
A.	<p><u>Internal Stand.</u> (wts. come from E)</p> $= SRR(I) = \frac{SR_E(I)}{SR_E(E)} = \frac{CR_E}{SR_E(E)}$ $= \frac{a}{E[a]} = SMR$ <p>(in stand. study)</p>	<p>Internally standardized/ Externally crude = = semi standardized</p>
B	<p><u>External Stand.</u> (wts. don't come from E)</p> <p><math>SRR(E)</math></p>	<p>Internally standardized/ Externally standardized</p>
1.	<p>Standard = referent</p> $\frac{SR_E(E)}{SR_E(I)} = \frac{SR_E(E)}{CR_E} = \frac{E[E]}{c}$ <p>(in ref. = stand)</p>	<p>= mutually standardized</p>
2.	<p>standard <math>\neq</math> referent</p> $\frac{SR_E(E)}{SR_E(F)} = \frac{E[a]}{E[c]} \text{ (in stand.)}$	
etc.	<p><u>Standardized excess RR</u> = <math>SERR</math></p> <p>= <math>SMR - 1</math> (null value)</p>	<p>"crude" <math>\checkmark</math> relative effect = "Crude Abs. effect" <math>CR_E \checkmark</math> null ratio</p>
	<p><u>Standardized Rate difference</u> = = <math>SRD</math></p>	<p>crude <math>\checkmark</math> absolute effect</p>
	<p><u>Relative Rate</u> = Rate Ratio = <math>RR</math></p>	<p>Risk Ratio</p>

(10/15) B. 3. take closer look at: (generally)

$$a. \text{SRR} = \frac{SR_E}{SR_{E^*}} = \frac{\sum w_j R_{Ej} / \sum w_j}{\sum w_j R_{E^*j} / \sum w_j} = \frac{\sum w_j R_{Ej}}{\sum w_j R_{E^*j}}$$

recast as:

$$= \frac{\sum w_j RR_j}{\sum w_j}$$

where:  $w_j = w_j (R_{E^*j})$

$$RR_j = R_{Ej} / R_{E^*j}$$

	D den.*	
E	$a_j$	$T_{Ej}$
E	$c_j$	$T_{E^*j}$

I will compute SRR in a few cases but paper has several others

\* = # persons in  $E_c$  study, or # p yrs in  $I_b$  study

b. cohort study: computations straightforward

$$\text{SRR}(I) = \frac{\sum T_{Ej} R_{Ej}}{\sum T_{E^*j} R_{E^*j}} = \frac{\sum \frac{T_{Ej} a_j}{T_{Ej}}}{\sum \frac{T_{Ej} c_j}{T_{Ej}}} =$$

$$= \frac{\sum a_j}{\sum \frac{T_{Ej} c_j}{T_{Ej}}} = \frac{a}{\frac{\sum T_{Ej} c_j}{T_{Ej}}} = \frac{a}{E[a]} = \text{SMR}$$

$$\text{SRR}(E) = \frac{\sum \frac{T_{Ej} a_j}{T_{Ej}}}{\sum \frac{T_{Ej} c_j}{T_{Ej}}} = \frac{\sum \frac{T_{Ej} a_j}{T_{Ej}}}{c} = \frac{E[c]}{c} = \text{SRR}$$

stand = ref. (SRR)

	D	$\bar{D}$
E	$a_j$	$b_j$
	$c_j$	$d_j$

c. Case-Control Study: derivations are more complicated, but results are understood intuitively. - see paper for derivations - NB. can't compute  $R_j$ 's

$$SRR(I) = \frac{a}{\sum \frac{b_j c_j}{d_j}} = \frac{a}{E[a]} = SMR$$

the trick comes, however, in seeing how the denominator is  $E[a]$

ie,  $E[a]$  = what "a" would be (in the standard) under the  $H_0$  that all strata had  $RR=1$

$$\therefore OR_j = I_D R_j = \frac{E[a] d_j}{b_j c_j} = 1$$

N.B. wts. ( $w_j$ ) are implicit in formula + are not  $n_j$  but  $b_j / E[b_j]$  = # exposed controls in the source pop.

$$E[a] = \frac{b_j c_j}{d_j}$$

$$SRR(E) = \frac{\sum \frac{a_j d_j}{b_j}}{e} = \frac{E[e]}{c}$$

THESE WTS. ARE NOT THE SAME AS IN CASE-CONTROL STUDIES  $\rightarrow$  (over)

\* 4. SMR's - 4 common mistakes

- a. note correction in paper #2, p. 19, 3rd line "rate" should be ratio -
- b. often in epid. research, SMR's for c/e studies are calculated incorrectly

$$"SMR" = \frac{a}{\sum \left( \frac{b_j}{b_j + d_j} \right) (a_j + c_j)} = \frac{\sum e_o}{\sum e_b} \neq SRR$$



(4') in the general case w/ standard pop.  $S$ , which is different from  $E, \bar{E}$  or any combination of exposure groups, it is possible for weights to be in # persons and exposure group denominator to be in person-years - in this case: the  $SRR(E)$  will not actually be a ratio of observed to expected cases

(eg)

$$SRR(E) = \frac{\sum \frac{N_{sj} a_j}{T_{Ej}}}{\sum \frac{N_{sj} c_j}{T_{Ej}}} \neq \frac{E[a]}{E[c]}$$

It should be noted that this does not apply to the situations when the standard is the referent, etc., because in order to compare  $\bar{E}$  and  $\bar{E}$ , they must have the same kind of incidence rates: either  $I_0$  or  $I_0$

(ans)

$SRR \neq \frac{a/c}{b/c}$

c. often SMR's are computed w/ the stated intention of "indirectly" standardizing the rate - what is called a stand. pop. is actually a ref. - but it is often recommended that we pool the exposure categories to form this so-called "stand." - done for 2 reasons:

- (1) get stable rates in the "stand." by creating a group w/ larger #'s
- (2) to make the "stand." similar to the indiv. exposure groups w/ respect to ~~ex.~~ stratum-specific rates - ~~if~~ it is recognized that IR adj. can ~~give~~ lead to weird results if the stand. is quite different from any of the compared groups

but forming a pooled referent (not stand) is actually being done in this makes the context of the ref. group variable according to the distrib. of the pop. or the selection procedures - thus any measure of effect has an ambiguous interpretation (what actually is being compared?) - it is thus best to carefully select the referent group - preferable a single category of exposure - N.B. Col's Bladder Ca. Study uses several jobs categories as ref. since they are hypothesized as unexposed.

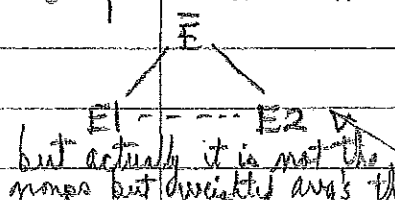
d. often SMR's are used as a form of IR adj. (esp. when ratios in the exposed groups are unstable) - but there is a problem when we attempt to compare 2 or + SMR's (of various exposure categories) even using the same referent

(eg) The Study

		D	$\bar{D}$	
> 1ph	E1	a <sub>1</sub>	b <sub>1</sub>	n <sub>E1</sub>
< 1ph	E2	a <sub>2</sub>	b <sub>2</sub>	n <sub>E2</sub>
0 cig.	E	c	d	n <sub>E</sub>
		m <sub>D</sub>	m <sub><math>\bar{D}</math></sub>	n

Suppose we make two comparisons to estimate RR:  $E1/\bar{E} + E2/\bar{E}$ , controlling for some confounding factors - Since SMR's are actually internally stand. RR's, the weights for each comparison come from the test group: E1 in the first comparison and E2 in the second. - thus O.M. has maintained that the 2 SMR's would not be comparable - ie, they would give us the wrong idea of the relative effect between E1 and E2; if the variable (E) is ordinal, this means that a comparison of SMR's would not be appropriate in describing a dose-response relationship - thus weighted averages (SMR's) cannot be compared unless the weights (ie, the standard pop.) is identical for all comparisons

apparently common belief: that things compared to equal things can be compared to each other



(eg) - p. 20 #2

I<sub>e</sub>

	S <sub>1</sub>			S <sub>2</sub>		
	D	$\bar{D}$		D	$\bar{D}$	
E1	10	0	10	90	0	90
E2	90	0	90	10	0	10
$\bar{E}$	5	45	50	50	0	50

Note that in each stratum of E1 and E2, the rate of ~~dis~~ D is the same (1) - i.e., everyone gets the disease - we should therefore expect that a comparison of effects due to E1 and E2 should be equal - however, the SMR's are not equal; but SRR's are equal because it involves ~~an~~ the same stand. ( $\bar{E}$ ) in both comparisons

NB. I would like to question the universality of this argument - i.e., that SRR's are always appropriate (+ not SMR's) in comparing several groups with the same referent. - 3 questions

- (1) what is meant by comparability?
- (2) what happens if the "weights" for both E1 and E2 are equal - could we then use SMR's?
- (3) ~~what~~ does it make any difference if the conf. factor also modifies the effect of E on D?

Set up 4 conditions to study -  
Conclusions

5. Why use SMR's at all? - to determine whether conf. is present -  $RR^* = \frac{CRR}{SMR}$

### c. Standardization of RD

1. (general)

$$SRD = SR_E - SR_{\bar{E}} = \frac{\sum w_j R_{E_j}}{\sum w_j} - \frac{\sum w_j R_{\bar{E}_j}}{\sum w_j}$$

only done for I rates =  $a - E_0(A)$

$$SRR = \frac{a}{E_0(A)}$$

2. 2 types:

a.  $SRD(I) \Rightarrow$  wts. come from  $\bar{E}$

b.  $SRD(E) \Rightarrow$  " " "  $E$

c. same principle of compatibility as for SRR

### 3. relationship between SRD + SRR

$$RD = R_E - R_{\bar{E}} = R_{\bar{E}} RR - R_{\bar{E}} = R_{\bar{E}} (RR - 1)$$

$$SRD(I) = SR_{\bar{E}}(I) - SR_{\bar{E}}(E) \leftarrow \text{stand. = ref.}$$

since,  $R_{\bar{E}} = SR_{\bar{E}}(I) (= CR_{\bar{E}})$ ,

$$\therefore SRD(I) = \frac{R_{\bar{E}}}{SR_{\bar{E}}(I)} [SR_{\bar{E}}(I) - SR_{\bar{E}}(E)] =$$

$$= \frac{R_{\bar{E}} SR_{\bar{E}}(I)}{SR_{\bar{E}}(I)} - \frac{R_{\bar{E}} SR_{\bar{E}}(E)}{SR_{\bar{E}}(I)} =$$

$$= R_{\bar{E}} - \frac{R_{\bar{E}}}{SMR} = R_{\bar{E}} \left(1 - \frac{1}{SMR}\right) =$$

might use C/E study + other info. to obtain SRD

$$= \left( \frac{SMR - 1}{SMR} \right) R_{\bar{E}} \quad \left( \text{note difference in form from above (RD)} \right)$$

### D. Standardization of Measures of Potential Impact

1. ARP - causal factors (p. 7)

$$(S)ARR = EF$$

a. for exposed <sup>pop</sup>:  $EF_E = \frac{SMR-1}{SMR} = \frac{a-E[a]}{a}$

$$\frac{SMR-1}{SMR} = \frac{\frac{a}{E[a]} - 1}{\frac{a}{E[a]}}$$

where:  $SMR \geq 1 \Rightarrow$  (EF would be negative)

b. for total pop.:  $\frac{a-E[a]}{a} = EF_E$

$$a+c \left[ \frac{a-E[a]}{a} \right] = EF_N = \frac{a-E[a]}{a+c}$$

$$\frac{a-E[a]}{a+c}$$

c. EF for several categories of exposure var. (ie,  $E_i$ )

$$EF_N = 1 - \sum \frac{e_{Di}}{SMR_i}$$

where:  $EF_N \geq 0$  (not  $SMR_i \geq 1$ )

$$e_{Di} = \frac{a_i}{\sum a_i} = \text{\% of cases in the } i\text{-th category of } E \text{ (including ref.)}$$

$SMR_i = SMR$  for  $i$ th category of  $E$

NB. correction ① p. 23 # 2, bottom formula: put  $\sum$  at beginning of expression after "="

③ # 5, p. 11 - cross out "note" under 3. b.

② # 2, p. 11, bottom line  $\frac{E[a]}{E[a]}$  should be  $\frac{E[e]}{E[a]}$

Harvey's Questions - back of p. 18

d. for several exposure vars;  $E_a, E_b$  etc.

if effect of  $E_a + E_b$  are indep. <sup>(no interactive effect)</sup>

$$EF_{a,b} = 1 - (1 - EF_a)(1 - EF_b) =$$

$$EF_a + EF_b - EF_a EF_b$$

where:  $EF_a$  is EF, given non-exposed state of  $E_b$ .

$EF_b \rightarrow$  " "

if  $EF_{a,b} > \dots \rightarrow$  Synergism

$EF_{a,b} < \dots \rightarrow$  antagonism

2. preventive factors - <sup>preventive fraction</sup> analogue to EF but it is not quantitatively an "inverse" of EF - (as was the case for RR)

because  $\rightarrow$  PF is % of potential cases (actual + prevented) that ~~would~~ have been prevented ~~if~~ because of exposure to protective factor

\* PF  $\rightarrow$  exposed category is considered the low risk group (eg. heavy physical <sup>exposure</sup>)

a. for exposed cases:  $PF_E = 1 - SMR =$

$$PF_E = 1 - SMR = 1 - \frac{a}{E[a]} = \frac{E[a] - a}{E[a]} \quad \text{where: } SMR \leq 1 \quad = \frac{E[a] - a}{E[a]}$$

b. for total cases:  $PF_N = e_D^* PF_E$

where:  $e_D^* = \% \text{ of potential cases that are exposed (ie, in low-risk category)} = \frac{E[a]}{E[a]+c}$  (if E is dichotomized)  
ie,  $e_{D0} = 1 - e_D$

(over)

$$PF_N = \frac{e_D (1 - SMR)}{e_D + (1 + e_D) SMR} = \frac{E[a] - a}{E[a] + c} \text{ (over)}$$

where:  $SMR \leq 1$

c. for several categories of exp. var. -  $E_i$

see notes

$$PF_N = \frac{\sum e_{Di} (1 - SMR_i)}{e_{D0} + \sum e_{Di} SMR_i}$$

where:  $e_{D0} = \% \text{ of cases falling in ref. group}$   
 $SMR_i = \text{doesn't have to be } \leq 1 \text{ for all } i$

$$PF_N \geq 0$$

d. for several var.:  $E_a, E_b$   
if effects of  $E_a + E_b$  are indep.

$$PF_{a,b} = 1 - (1 - PF_a)(1 - PF_b) = PF_a + PF_b - PF_a PF_b$$

where:  $PF_a, PF_b$  for  $E_i$ ; given non-exposed state (hi-risk) of  $E_b$

if:  $PF_a > \dots \rightarrow$  synergism  
 $PF_b < \dots \rightarrow$  antagonism

VI Next time: Study Design



$$e_D^* = \frac{e_D}{e_D + e_{D_0} \text{SMR}} = \frac{a}{a+c} + \left(1 - \frac{a}{a+c}\right) \frac{R}{E[a]}$$

$$+ [a] \frac{a}{a+c} + \frac{a^2}{E[a](a+c)} = \frac{\alpha E[a]}{\alpha E[a] + a(a+c)}$$

$$\frac{E[a]}{E[a]+c}$$

$$PF_N = e_D^* PF_E = \left( \frac{E[a]}{E[a]+c} \right) \left( \frac{E[a]-a}{E[a]} \right) = \frac{E[a]-a}{E[a]+c}$$

11/5 Outline

next week: end 4:00

next sem: time

I. Intro.

A. Outline for remaining wks. of Semester

1. criteria for evaluating quality of one's study design and data
2. dimensions and types of epid. study designs + methodological implications
3. clarification + elaboration of 3 concepts: conf., EM, and I.

B. Today - mostly based on handout #5, which was intended to give some feel for the different perspectives of viewing the issues dealing w/ the quality of one's design: 1. conceptual (outlined #1); statistical (#2); and methodological (#3) -> focus of today's discussion + based on O.M. (Ch. 2) - p. 9

II. 2 bases or general criteria on which we can evaluate studies

A. Feasibility of implementation - practical considerations, including economic, social, moral, + political issues -

P. 9

B. Informativeness - accuracy in estimating the object of study (i.e. measure of assoc.) -

1. in statistical terms, this means studying the variability in one's data, separating 2 components
  - a. true var. - reflects etiologically based differences in experience w/ D
  - b. Error variability - reflects inaccurate estimation of experience w/ D
    - i. sampling (usually don't sample the entire source pop.)
    - ii non-sampling (measurement error)

2. from epid. perspective (O.M.), there are 2 components of I.

a. internal validity - the lack of systematic (non-random) error or lack of epid. bias

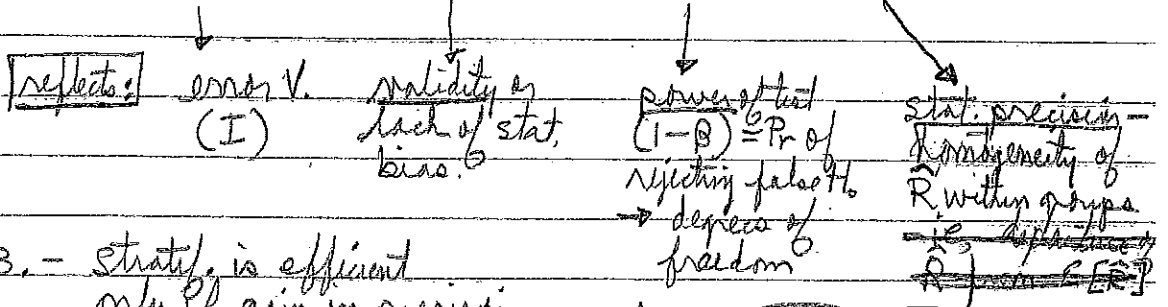
b. precision - lack of random errors; which

is reliability - is the departure of any given estimate from the mean value (if study were repeated)

III. there are subtle differences between statistical bias (& validity) and epid. bias (& internal validity)

A. look at stat. bias - consider in terms of CI of pt. est.

$$CI[\hat{R}] = \hat{R} \pm CV_r SE[\hat{R}]$$



NB. - strat. is efficient only if gain in precision is not outweighed by loss in df.

stat. efficiency random error - departure of  $\hat{R}$  from  $E[\hat{R}]$  - reflected in size (width) of CI

See addendum (A)

B.  $\rightarrow$  lack of statistical bias  $\Rightarrow$  if study were done an  $\infty$  # times, the est. of the effect (or any parameter) would produce an accurate assessment of the pop. parameter (R) - ie, "true" value of effect

$$ie, E[\hat{R}] - R = 0$$

NB: St. B is dept. on the b

C.  $\rightarrow$  lack of epid. bias  $\Rightarrow$  if study were expanded to  $\infty$  size (n) it would provide accurate est. of R

$$ie, \lim_{n \rightarrow \infty} \hat{R} = R$$

NB. doesn't depend on n ( $\equiv$  stat. consistency)

(me) if we sample entire pop, bias must be 0 since every estimate of parameter (w/ n = N) will be same - but this does not address issue of conf. factors that may be operating - even in total pop. → want stat. consistency

D → OM chooses this defn. to use for internal val. because:

- a. epid. bias is conceptually cleaner
- b. its not dep't. on n & ∴ reflects the type of bias that is most prevalent & important in epid. studies - he notes that st. bias can often be corrected in analysis by approp. transformations

(HIA) ? E. to stat. efficiency is measured by

$$1 / \text{Var}[\hat{R}] = 1 / \text{SE}[\hat{R}]^2$$

where: SE[ $\hat{R}$ ] is a function of the "deviation"  $\hat{R} - E[\hat{R}]$

F. to epid. precision is basically same thing except, unlike the other, it ~~doesn't~~ isn't defined in the context of unbiasedness - epid. precision is taken measured indep. of validity

#### IV. Epid. Precision - 2 components

- A. Size-efficiency ⇒ I / subject
- B. Size - is, # subjects - gen'l principle: the larger the n (everything else being =), the more precise an estimate is
  1. yet in practice, we often think in terms of an optimum size - decide a priori how large an effect is worth detecting & the degree of precision desired ( $\alpha + \beta$ ) - then compute n according to stat. model
  2. OM argues this is not optimum size because we don't know the benefit of I that we might get from study - this would involve putting \$ value on study results which is impossible
  3. ∴ only optimum sizes are 2 extremes:

3. a.  $n = 0$  - if no info is afforded by the study - implies that Val. + Prec. are not to be traded-off for one another -  $\therefore$  unwarranted to pursue high Prec. in study that may be quite invalid
- b.  $n = \infty$  - very high prec. is necessary to "prove a H<sub>0</sub>" (ie, effect = 0) - so w/ any  $n < \infty$ , all we can do is disprove H<sub>0</sub>s, not prove them

IV. Internal Validity - Systematic error - 3 components + 3 types of B

A. Selection (often undefined) - when response var. of study influences the selection of subjects in such a way as to bias R - in:

(spirit type of conf.)

1. CH<sup>o</sup> study - since D cannot influence selection of subjects, selection is not an issue w/ prospective study
2. C/E study - E (response) influences Pr of selection of subjects, but differentially among cases + controls - otherwise no bias - eg, Berkson's fallacy, where differential admission rates into study (crosswalkings from different diseases) may bias results

(one) except in certain circumstances: we select subjects w/o spitting dis. + follow-up for dis. occurrence - selection bias exists if large % (>20%) of source pop. have dis. (thus ineligible) + dis. prev. is assoc w/ dis. incidence - ie, sick people are most likely to get sick in future (Lawrence, Public H Reports 63:1507-21, 1948)

3. X-S study - influence of D on selection of subjects  $\rightarrow$  need not be differential - only that  $Pr[\text{selecting } C] \neq Pr[\text{sel. } \bar{C}]$
- $\frac{\# C \text{ in study}}{\# C \text{ in source pop}}$

if  $Pr[C] > Pr[\bar{C}] \Rightarrow$  bias toward null - of PR OR P  
 if " < "  $\Rightarrow$  bias away from "  
 NB. 2 Pr's would be equal under conditions of perfect simple random sampling (w/ strat. B)

B. information

1. non-differential bias - error occurs identically in all compared groups - net effect is to bias estimate toward null - not spurious assoc. - (eg) in  $c/\bar{c}$ , we overestimate % of exposed persons, but equally in  $C+E \therefore RR \rightarrow 1^0$

N.B. this occurs whenever:  $Se \geq \%T - Sp$   
(almost always)

- a. mislabelling or misclass. (of E or D)
- b. non-particip. of selected subjects (w/o obs.)
- c. loss to follow-up (migration, non-cooperativeness, death)

2. differential bias - bias is different between compared groups - results in spurious assoc. - eg over-estimate of E is different in C than controls - eg, living C + Smoking

a. group comparison - 2 issues

- 1. appropriateness of referent group (planning stage)
- 2. confounding factors - (planning + analysis) - conf is when the effects of 2 D covariates get mixed so that the apparent (crude) effect of one actually includes the other -  $\therefore$  assoc. of E + D ~~is~~ is spurious owing to the effects of other var's that have not been separated thru statistical analysis -  
~~Crude effect~~ crude effect ~~is not~~ true (unconf.) effect -  
 if CE is closer to null than TE  $\rightarrow$  " - " conf.  
 " " " further from " " "  $\rightarrow$  + conf.

note  
change  
RCE

We will go into, to more detail later & ~~also~~ measure the magnitude of conf. effect in order to know when to control.

VI Preferability - combining Val. + Prec. into overall measure of I (both rand + non-random error)

$$I = \frac{1}{B^2 + V[\epsilon]} \quad \epsilon = \text{random error}$$

but, B +  $\epsilon$  are not to be traded-off

VII Optimality - combining feasibility + I.

A. might assume common units - eg, cost everything being  $\Rightarrow$  a larger study provides more info but is costlier -

but optimal design is not identifiable, because we have no basis for determining whether a more inform. study is worth the added cost  $\rightarrow$  simply because we can never place \$ figures on the future consequences of a study that has not been conducted - but there are some ("rational planners") who devise complex methods for doing benefit-cost analysis under these premises  $\Rightarrow$  "social benefit functions"

B. in practical terms - we should max. I, per cost - i.e., consider cost effic  $\rightarrow$  I/\$ ; not trade-off th. 2

VIII External validity - extending empirical results to persons outside range of source pop. - i.e., generalizability - moving from particular experience to scientific generalization (research to theory) - involves judgmental inferences about observed results

A. external vs. internal - p. 17

1. types of inference - p. 13-14

2. 3 pts - p. 14-16

B. Again, no trade-off between internal + ext val. same terms is premisses for th. latter.

IX

Next chore is to see how these concepts of it are considered methodologically - ie, in context of specific study design - but first we must find some systematic ~~classific~~ basis for classifying research designs

A. general perspective of epid. pursuits which includes quantification of disease occurrence in groups of indiv. - p. 2, 3 # 6 - 3 types of studies

B. look at analytical (obs.) studies in more detail - division of analyt. studies into 3 types: Ctt, C/E, X-S, is an ~~over~~ over-simplification that might be useful when taking 160 but certainly presents limitations when studying study design in more depth -

C. next time, we'll look at study design in terms of 3 dimensions & combine them to form a typology of study types. - NB: flow diagrams, (#7) - 1 dimension: directionality



11/12 Outline

I. Introd.

A. Handouts:

1. 2 page replacement for pp 10-11 (#6)

Types of Study designs

2. EM & I (#9)

B. today: dimensions & types of epid. study designs  
& their methodological implications

II. Underlying task of today's discussion is the classification of design methodologies

A. broad view of epid.: includes quantification of D occurrence in pop's - see #6, pp. 2-3 -

3 types of studies, differ by objectives

- 1. Observational
- 2. Experimental
- 3. Evaluative

B. focus on Analytical (Obs.) to test etiological H<sub>0</sub>

1. traditional division into 3 types (CH, C/E, X's) may be fine for an introd. to Epid, but presents limitations when studying design in more depth - eg, what effect measure to use

2. will try another approach that begins by defining study design (Obs.) into 5 dimensions

- a. directionality - flow diagrams
- b. timing
- c. subject selection + sampling
- d. temporal arrangement of observations
- e. units of measurement

see other outline p. 5 (#6)

III. Typology of Study Designs - based on the 5 dimensions - derive categories of study design that may be treated as methodologically ~~superior~~ distinct - see other sheet; (just handed out)

I. Longitudinal

A. Forward

- 1. Cohort
- 2. Panel

B. Ambidiv.

C. Non-direct.

II. Case-referent

A. Backward (case-history)

B. Ambidiv.

C. Non-direct. (X-S)

III. Cross-sectional - nondirect,

A.

IV. Major methodological considerations for selecting a particular obs. study design (based on O.P.M., Ch. 3) - divide issues into 3 categories

A. Type of Inf. Sought

B. Internal Validity

C. Feasibility of Study

(over) somewhat different than the outline already handed out (#6, pp. 13-15)

A Typology of Study Designs in Observational Research  
Involving Etiologic Hypotheses

I. Longitudinal study - observations of D are made at more than one point in time

A. Forward directionality (follow-up)

1. Cohort study - calculation of incidence rates

a. Type IA (selective) sampling of exposure groups - either prospective, retrospective, or ambispective

i. fixed cohort -  $I_C$  (or  $I_D$ )

ii. dynamic cohort -  $I_D$

b. Type II (nonselective) sampling of a source population - usually prospective

i. fixed cohort -  $I_C$  (or  $I_D$ )

ii. dynamic cohort -  $I_D$

2. Panel study - calculation of a change in prevalent status (of the condition - D) - usually prospective

a. Type IA (selective) sampling of exposure groups - (see: I.A.1.a.)

b. Type II (nonselective) sampling of a source population - (see: I.A.1.b.)

B. Ambidirectional study - Type IB (selective) sampling of disease groups within a cohort - retrospective or ambispective

1. Fixed cohort

a. "density type" - subjects are "observed" within the risk period (of D) - for chronic disease -  $I_D$

b. "cumulative type" - subjects are "observed" at the end of the risk period (of D) - for acute diseases, recurrence, complications -  $I_C$  or  $I_D$

2. Dynamic cohort - "density type" -  $I_D$

I. C. Non-directionality (momentary)

1. Multiple, cross-sectional study - two or more cross-sectional surveys (see: III) on a dynamic population, over time - calculation of changes in prevalence rates - usually Type II sampling

(case-control)

II. Case-referent study - observations of E and D are made concurrently; with Type IB (selective) sampling of disease groups - retrospective or ambispective

A. Backward directionality ("case-history" study)

1. use of prevalent (i.e., existing) cases
2. use of incident (i.e., new) cases
  - a. "density type" - subjects are "observed" within the risk period (of D)
  - b. "cumulative type" - subjects are "observed" at the end of the risk period (of D)

B. Ambidirectionality - (see: I.B.) - i.e., a case-referent study within a cohort that is followed up

C. Non-directionality (momentary)

1. use of prevalent (i.e., existing) cases
2. use of incident (i.e., new) cases - actually a "sequential cross-sectional study" ~ "density type"

? cum. type?

III. Cross-sectional study - observations of E and D are made at one point in time - non-directionality

A. Type I (selective) sampling

1. Type IA sampling of exposure groups
2. Type IB sampling of disease groups - (see: II.C.)

B. Type II (nonselective) sampling of a source population

Start

# IV Considerations for Selecting a Study <sup>Decision</sup> in Observational Research

- 3 criteria:
- I. Type of Info. Sought
  - II. Internal Validity
  - III. Feasibility of Study (cost)

## I. Type of Info. Sought

A. Substantive focus - # (+ type) of primary var's to be studied (E+D)

1. specific focus - one E and 1 D - *all of the study designs; espec. ambider. - main interest*

2. Semi-specific

*philos* a. several E's, 1 D - all designs, esp. those w/ Type II sampling - also appropriate to use model fitting technique - (fishing expedition)

b. 1 E and several D's - follow-up study or X-S study w/ Type IA sampling (fishing expedition) - also ambidirectional

3. non-specific - several E's and D's - cohort and X-S w/ Type II sampling - the main idea should be to narrow focus

*if we don't have to resample for each new disease by sampling the entire cohort, not just noncases*

B. Type of Ho <sup>unit of measurement</sup>

1. • ecologic fallacy - can't attribute to indiv. what we know about groups - also

• individualistic fallacy - can't attribute to groups what we know about indiv. (bias against this concept in epid.) - in paper

*paper #6*

• conclusion: serve different purposes (ecol. vs. non-ecol., indiv. analysis → etiology of a dis., ecol. analysis → prognostic/policy decisions re: the use of preventive action in a given pop.)

I. B. 2. O.M (Ch. 3) alludes to this division of objectives by distinguishing:

NB: he doesn't define Scientific by the "methods" used - of the Scientific method physicians use etc. method

Scientific Study Orientation - we attempt to generalize Re: causal determinants of a D.   
 - the findings are based on a source pop. they are not made dept. on the distrib. of the pop w/r: other D correlates - (RR)   
 Particularistic Study Orient. - aim is to improve the Health status of a particular pop. - empirical estimates are made dept. on distrib. of pop. (eg. by age) - EF (PF) are partic. questions - best to use Type II Sampling (Simple or stratified)

NB: refers to external validity, not statistical inference

3. I've extended this dichotomy (#6, p. 20)   
 6 domains of the generation - ie not only need to distinguish 6 domains   
 b. as PH professionals, we might be more concerned w/ intervention, which requires different H's + correspondingly study designs   
 c. in order to do this, we must put more emphasis on evaluation of public programs and complete record systems

views of Remwick, Tervis

c. Parameter estimation (of assoc.) - (over) by type of study

# C. C. Parameter estimation by type of study

## Appropriate Measures of Association by Type of Study Design

### I. Longitudinal

#### A. Forward directionality

##### 1. Cohort study (a or b.)

a/b. Type IA or II sampling - (forward)

i. fixed cohort -  $I_cR, I_cD (I_c)$   
especially in acute diseases, complications, recurrences & diseases w/ short induction periods; because  $I_0$  does not remain constant over follow-up (eg, infectious dis, epidemic, congenital disorders, or maternal complications during or following pregnancy) - but often can derive  $I_D R, I_D D (I_0)$  if  $I_0 = const$

- ① call it: "dynamic pop."
- ② or "running cohort" (where all subjects who enter are followed to death or some desired event)

ii. dynamic "cohort" -  $I_D R, I_D D (I_0)$  - can derive  $I_cR, I_cD$  - also EF, PF  
Since dym. cohort provides representative cases; while fixed in fixed cohort are not repres. Since the cohort ages (NB. original cohort is repres.) (see ch. 3) gives a method for making appropriate adjustments to estimate EF or PF in a fixed cohort

Chronic diseases where  $I_0$  remains constant in follow-up

problem (eg) Lung Ca. - age-adj. Fract. increased by 250% (1950-70)

2. Panel Study - non-categorical procedures (eg, difference in means, model fitting) - can be used to derive certain measures of effect depending on the statistical model, -  
N.B. statistical nature of analysis is different (eg, Myerson's study: ABP → CHB)  
or Smoking → ΔBP

I. B. Ambidirectional study - type IB sampling

1. fixed cohort

- a. density type -  $I_cR, I_cD (I_c)$  - if ref. group (PD) is a random sample of the cohort - but really want:  $I_0R, I_0D (I_0), EF$ , which we can get by adjusting # aggregated pers. (OM, ch. 3)
- b. cumulative type -  $I_cR, I_cD (I_c)$  - if ref. series is random sample of cohort - for acute dis.

2. dynamic cohort -  $I_0R, I_0D (I_0), EF, PF$   
 - if ref. series is random sample of cohort (NB: simpler computations than w/ fixed cohort)

show ex. fixed CH. (1-5)

C. Non-directionality

- 1. Multiple X-S -  $\Delta P (PD)$  by subgroup of E not very useful etiologically because pop. changes between X-S studies - better for secular trends of D

II. Case-referent study - ("Case-control") - See App. C

A/C. Backward ("case-history") or nondirectional (X-S)

	D	$\bar{D}$
E	a	b
$\bar{E}$	c	d

1. prevalent cases -  $OR = \frac{ad}{bc} = OR_p = \frac{P_E \cdot \frac{a}{c}}{P_E \cdot \frac{b}{d}} = \frac{a}{b} \cdot \frac{d}{c}$

$I_{DR} = OR_p \cdot \frac{\bar{d}_E}{d_E} = OR \cdot \frac{\bar{d}_E}{d_E}$

so:  $I_{DR} = OR$  if  $\bar{d}_E = d_E$   
 Cornfield (1951) - if D is rare so that  $OR_p \approx I_{DR} (\approx I_{DR})$

Note: assumption is not good if same factor (E) is a causal of disease onset and prognostic indicator of case fatality - e.g. alcohol consumption  $\rightarrow$  angina



Ambid. Study - Example, Calculation of  $I_{cR} + I_{cE}$   
 fixed cohort - cum. type  
 (non-observed quantities in parentheses)

(over-flow diagram)

	D	$\bar{D}$
E	$I_E$	$(N_{ED})$
$\bar{E}$	$I_{\bar{E}}$	$(N_{\bar{E}\bar{D}})$
	I	$N_D = N_{\bar{D}} = N$

take sample from  $N_D (= N)$  of size  $n$

estimate  $Pr[E] = \frac{N_E}{N}$  from this sample

$\hat{Pr}[E] = n_E/n$  where  $n_E = \#$  in sample who are exposed

$Pr[\text{being } E+D] = Pr[D|E] Pr[E] = Pr[E|D] Pr[D]$

$Pr[D|E] = \frac{Pr[E|D] Pr[D]}{\hat{Pr}[E]} = I_{cE}$

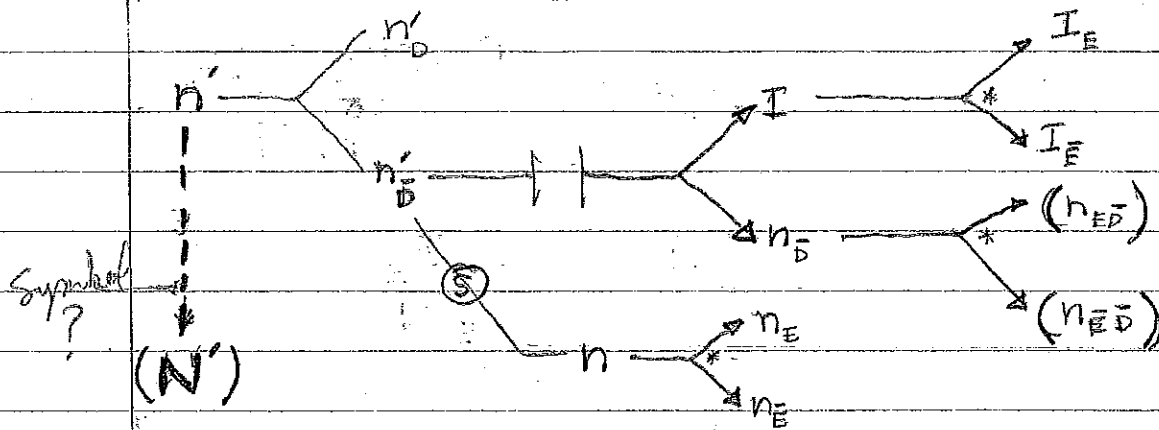
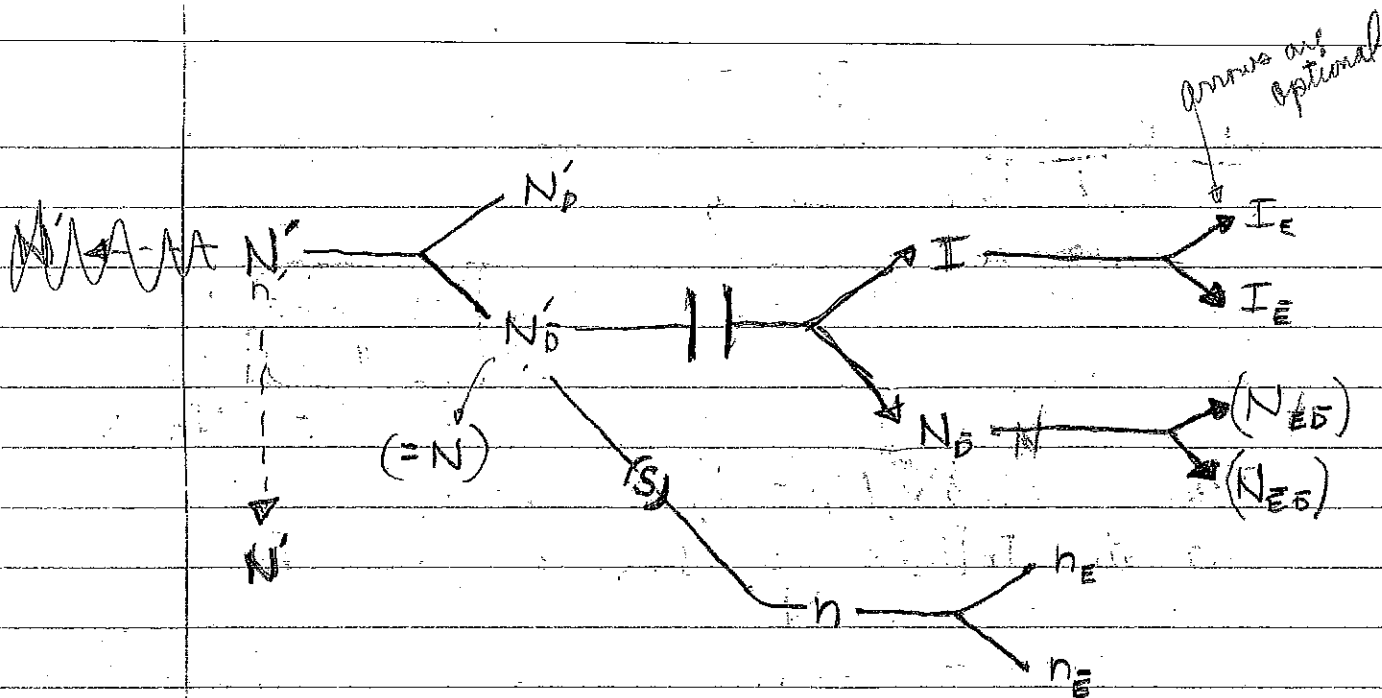
$I_{cE} = \frac{(\frac{I_E}{I} \cdot \frac{I}{N})}{\frac{n_E}{n}} = \frac{n_E I_E}{n_E N}$

Similarly:  $I_{c\bar{E}} = \frac{Pr[\bar{E}|D] Pr[D]}{\hat{Pr}[\bar{E}]} = \frac{\frac{I_{\bar{E}}}{I} \cdot \frac{I}{N}}{\frac{n_{\bar{E}}}{n}} = \frac{n_{\bar{E}} I_{\bar{E}}}{n_{\bar{E}} N}$

$\hat{I}_{cR} = \frac{Pr[E|D] \hat{Pr}[E]}{Pr[\bar{E}|D] \hat{Pr}[\bar{E}]} = \frac{n I_E / n_E N}{n I_{\bar{E}} / n_{\bar{E}} N} = \frac{n_{\bar{E}} I_E}{n_E I_{\bar{E}}}$

Since  $\hat{Pr}[E]$  is an estimate, can get CI of  $I_{cE} + I_{cR}$   
 see Kupper et al., 1975 JASA (70(351):524  
 (also can do stratified analysis)

See: II/c for dyn. cohort



\* - not nec. arrows  
 ( ) - not actually obs.

C-R study

II. A/c. 2. Incident cases

a. density type -

$I_D R = OR_P = OR = \frac{ad}{bc}$  (already shown)

if pop. is dynamically stable  
also can get  $I_D$  fr E and  $\bar{E}$  if we  
have an estimate of  $I_D$  (for entire source  
pop.)

$EF_N = \left( \frac{RR - I_D}{RR} \right) e_D = \frac{I_D - I_{D(E)}}{I_D}$   
 can be standardized

can get w/o  $I_D$

prove, when I have time

$\frac{I_{D(E)}}{I_D} = 1 - EF_N$

$I_{D(E)} = I_D (1 - EF_N)$

$I_{D(E)} (I_{DR}) = I_D (I_{DR}) (1 - EF_N)$

$I_{D(E)} = I_D (I_{DR}) (1 - EF_N)$

(or)  $I_{D(E)} = I_D / (1 - PF_N)$

$I_{D(E)} = I_D (I_{DR}) / (1 - PF_N)$

prev. of source pop. (t)

	D	$\bar{D}$
E	A	B
$\bar{E}$	C	D

but, sometimes problem in estimating  $I_D$  because denominator is taken as both persons w/o disease at risk + prevalent cases

$I_D^* = \frac{(a+c)}{(A+B+C+D)(t_1 - t_0)} = \frac{\text{total cases}}{p \cdot yrs}$   
 prev. cases

from other sources

II. Aft. 2. a. it should be:

$$I_D = \frac{a+c}{(c+d)(t_1-t_0)}$$

A susceptible

also, non-specific assessment of E (did you ever smoke) which not good for getting RR (since it will tend toward 1), but since E are also increased, this type of exposure history is good for getting EF (esp. if you divide smokers into groups - by # pack-years)

So:  $I_D \approx \frac{I_D^*}{1 - \hat{P}(D)}$

prev. rate ( $\hat{P}$ ) may be known from another source or even from refer. series if it was drawn w/o inclusion of prevalent cases

incid. in source pop.

	D	D
E	W	X
$\bar{E}$	Y	Z

b. Cumulative type: "risk odds ratio"

$$OR_R = \frac{I_{CE}(1 - I_{CE})}{I_{\bar{C}\bar{E}}(1 - I_{\bar{C}\bar{E}})} = \frac{W}{W+X} \frac{(1 - \frac{Y}{Y+Z})}{(1 - \frac{W}{W+X})} = \frac{W \cdot Z}{Y \cdot (W+X)}$$

never, cases & never, noncases

$$\approx \frac{W/Y}{X/Z} \approx \frac{a/c}{b/d} = \frac{ad}{bc} = OR$$

↑ can get this

if D is rare so that  $\frac{(1 - I_{CE})}{(1 - I_{\bar{C}\bar{E}})} \approx 1$

$$OR = I_{CE} \hat{R} (= I_{DR})$$

in order to get risks for exposure groups, we need  $I_{CE}$  of source pop. from other sources in the study itself - OM (Ch. 3) describes procedure for estimating  $I_{CE}$  and  $I_{\bar{C}\bar{E}}$

See addendum for different CR design

### III. Cross-sectional study

#### A. Type I sampling

1. Type IA (E groups) - PD (or PR - not good)

$$OR = OR_{PR} \hat{R} I_{E|R} \hat{R} \quad \text{if } \bar{d}_E = \bar{d}_{\bar{E}}$$

+ distributions are same

but  $OR \neq I_{D|R}$ , unless  $d$  is very short for both groups

2. Type IB (D groups) - eD

Similar to case-referent study (see: II, A/c)

if E status for each person is stable over time also good for testing validity of a screening detection test -  $SE$ ,  $Sp$ ,  $Fv$  - in fact, any case-history <sup>(backward)</sup> may be used to test validity of screening test based on risk factors

#### B. Type II sampling -

correlation coefficients -  $\Phi$  w/ categorical data, or PD or PR or  $OR_p$  - (P or T if ordinal)

#### Summary

1. Conceptually, methodologically simplest study is <sup>in chronic dis</sup> density type ambidirectional study - w/ complete ascertainment of cases over a period of time
2. Other studies affording complete estimability of all parameters; follow-up studies - except Type II sampling tends to be less efficient, the representative of a specific pop.

usually we say:  
COP → can't determine antec./conseq.  
CH - can " " " "

# I. D. Temporality Issues

1. antecedent/consequent determination - theoretically possible from backwards, forward or ambidix. - not from X-S - but even in a follow-up study, determination of temporal relationship may not be established - eg E (as a Δ in status) is measured over the same period of follow-up as when new cases of D are detected → ie, longitudinal nature of pop.

eg, ASES → CHD  
if both from 1960-67,  
2 problems:  
① temporality ambiguity  
② reduction of effect because of lag time

on the other hand, temporal ambiguity can be largely avoided in (retrospective) history study (backward direct.) by using new (incident) cases

2. period of follow-up must be selected carefully  
eg, if  $p \approx \infty$  (∞ years) - RR → 1 - chronic dis.  
 $p \approx 0$  - RR → 1 - for acute dis. or clinical trials (through ambidix)

eg, if we attempt to observe thromboembolism w/in 6-mo. after taking oral contraceptives, might show no effect because not enough time for D to develop - but early cases (≈ 1 mo. after) even in E group do not bias final results only reduce efficiency.

- 1. principle for etiologic (scientific) research; be highly selective of comparison groups
- 2. matching - to be covered later

## E. Subject Selection (restriction)

# II. Internal Validity - lack of bias in estimation of effect

A. Selection bias - response var. influences selection of subjects in a biased way

3. hard to follow-up a lack of participation by part of study pop. (exceptions):  
1. remittance dis. P is aware w/I & P is large - sick people are more likely to get sick in future

1. follow-up study - selection is not an issue if follow-up starts after admissibility criteria are met - NB: maybe a problem in retros. CHD<sup>st</sup> subject  
2. case-ref. study - E influence Pr of selection but differentially among cases + controls - eg hospital based study of pill → thromboembolism - it's possible that the use of the pill causes physicians to be more likely to suspect a patient of

E is both 2. E is both 2. E is both 2. E is both 2. E is both 2.

II

A. (2)

having thrombocytopenia (ie, pill is used as diagnostic evidence of T) -> thus "cases" are more likely to have used the pill - also selective survival of cases in case-referent study may cause bias if the surviving cases are more (or less) likely to be "exposed" than those who died - NB: [this kind of selective bias can be sharply reduced by using new cases]

possibility of selection bias can often be ruled out judgmentally  
eg. Coffee -> MI  
OM: found RR > 1, for selective bias to be inflating effect, MI patients who died must have drunk less coffee than cases who were in study -> ? not likely

3. X-S study - influence of D on selection of subjects but not necessarily differentially  
eg. if  $Pr[\text{selecting case}] \neq Pr[\text{selecting D}]$  there will be bias - also selective survival  
NB: bias enhanced because P rates are dept. on incidence of dis. as well as duration, which also might vary by E groups  
 $d_E = d_{\bar{E}}$  more likely  $Pr[D] \neq Pr[\bar{D}]$   
ie, exposed (or E) cases are more likely to die & thus not get in stud  
(must be differential) in cases & non cases

See additional

if getting P or PD

B. Information Bias

1. non-differential

- a. mislabelling - can be a source of error in any type of study, esp. retros. because of incomplete or inaccurate records - if it
- b. lack of participation - esp. w/ non-retrospective data, esp. when using mailed questionnaires or compliance with researchers is an inconvenience - ~~usually~~ can't be more of a problem w/ Type II sampling
- c. loss to follow-up - in any follow-up study, esp. prospective - if too hi (as in b) & reasons are unknown, study is invalid

but if getting PR, must be differential

influence of E on selection of cases (must be differential)  $Pr[E] \neq Pr[\bar{E}]$  but different for D

- II. B. 2. differential measurement bias - can be source of error in any study, but esp.
- case-history study (backwards) - eg, biased recall of exposure as in Lung Ca & Smoking.
  - ~~non-retrospective~~ retrospective study (prospective, ambispective or X-S) - ie primary data observations of recollections are influenced by what results they expect
- C. Comparison bias ~~Factor~~ confounding
- choice of referent (Comparison) series
    - in follow-up study w/ distinct E, referent group can be genl. pop. - \*not likely in case-ref study <sup>to which D is rare unless E were known in pop.</sup>
    - generally not good to ~~form~~ form a pooled referent unless a logical "unexposed" group does not exist - in which case it is best to assign ref. group a priori, before data is collected
  - confounding factors - may be a problem in any type of study design - esp. in Type II sampling of a source pop. -  $\therefore$  better to restrict study pop.

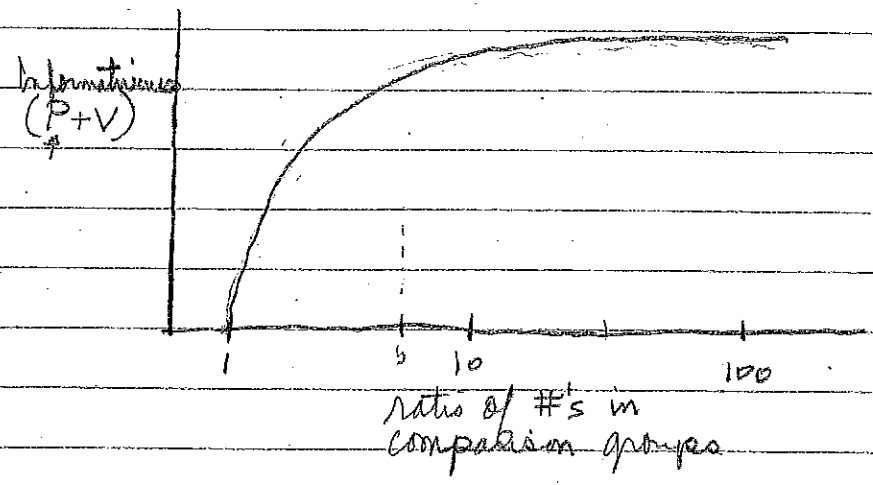
### III. Feasibility of the Study (cost-efficiency)

- Freq. of E in ~~pop~~ source pop. - extremely common or rare exposures (Prate) lend themselves best to follow-up or X-S studies w/ Type IA sampling -
- freq. of D in study pop. - extremely common or rare dis. lend themselves best to case-ref. studies



III

- c. length of induction period of D - chronic dis.  
w/ long induction periods ~~are best~~ lend themselves best to retrosp. cohort studies or case-ref. studies (not prospective follow-up studies)
- D. cost of getting E data - advantage of ambidirectional study (relative to follow-up) is that we don't need the exposure history on as many subjects (ie, only a sample of the non-cases in the cohort) - N.B. never gain much efficiency by having groups in ratio of  $> 4$  or  $5$  (controls on E) to 1 (case on E).



CR: normally done:-

	D	$\bar{D}$
E <sub>1</sub>	a	b
E <sub>0</sub>	c	d
	m <sub>1</sub>	m <sub>0</sub>

$\hat{OR} = \frac{ad}{bc} \approx \hat{I}_cR$  if D is rare  
(because of Bayes theorem)

$\hat{OR} = \hat{I}_dR$  if  $d_1 \approx d_0$  if prev. cases  
no assumpt. if new "

CR: alternative design:-

	D	PAR - ie, sample total pop, not new cases
E <sub>1</sub>	a	n <sub>1</sub>
E <sub>0</sub>	c	n <sub>0</sub>
	m <sub>1</sub>	n

$\hat{OR} = \frac{a/n_0}{c/n_1} = \frac{a/n_1}{c/n_0} = \hat{I}_cR$

So might use this design when:

- ① disease is not rare
- ② especially, \*cumulative type study (acute conditions)  
ie, sampling takes place after risk period  
so don't want I<sub>c</sub>R but risk ratio

Otherwise, it would be best to get I<sub>c</sub>R from I<sub>d</sub>  
(j = age categories) using I<sub>d</sub> (cohort) study  
or ambid. study

NB: Sampling PAR may be:

- ① easier than sampling  $\bar{D}$  +
- ②  $\therefore$  might not introduce same kind of selective bias as common in CR study

# 12/3 Outline

I. Introd. - same time next semester; begin 1/12; same room;  
change place

A. Handouts:

- 2. #11, Summary of Quant. Procedures - Conf, EM, J
  - 1. #10, Elements of Statistical Analysis:  $\chi$ , CI
- "meat of Minn. Course" - note: Table 1, pp.

B. today: Methodological implications of study designs

II. Outline (on board) - go to outline: 11/12

## Considerations for Selecting a Study Design in Obs. Research

I. (3 criteria)

I. Type/nature of info. sought

- A. Substantive focus
- B. Type of Hypothesis
- C. Parameter Estimation
- D. Temporality
- E. Subject Selection (restriction)

II. Internal Validity (lack of bias)

- A. Selective bias
- B. Information bias
- C. Comparison bias - confounding

III. Feasibility of study (cost-efficiency)

- A. Frequency of E in source pop.
- B. Freq. of D in study pop.
- C. Lag period of D
- D. Cost of getting E data

1/14/76 Outline + 1/21

## I. Intro.

A. would like to finish course (lectures) by March before OM visits (3/24)

### B. handouts

#10. Basic Principles of Categorical Analysis:  $H_0$  testing + interval estimation (CI) - heart of Minn. course

#11. A Summary of Quantitative Procedures Used to Identify & Measure Conf., EM, + I

#12. Matching and Matched Analysis in Obs. Studies

#8 (brief revision) Methods of controlling for conf. var.

### C. today's topics

1. principles + vocabulary of causal inference (epid. Standpt.)
2. confounding

## II. Causal Inference

A. purpose of this discussion is to propose a common language (vocab.) for establishing etiological relationships in epid. research

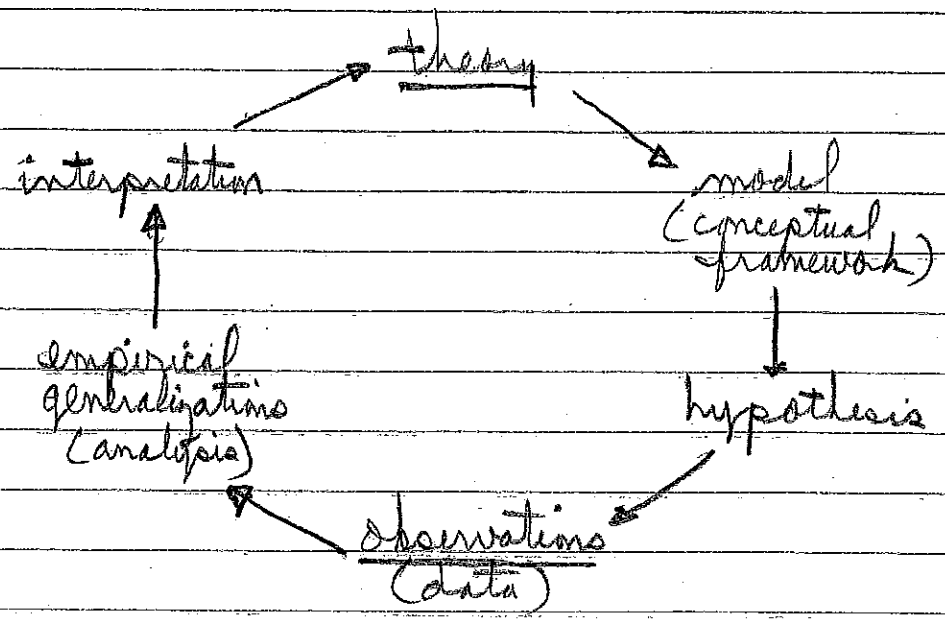
B. must distinguish 2 basic levels of <sup>scientific</sup> inquiry

1. theoretical level - <sup>involving</sup> biological/behavioral mechanism of causation - think in terms of expected changes

2. Empirical level - statistical associations among variables in one's data - think in terms of observed changes (covariation)

C. any scientific process must bridge the gap between these 2 levels - often simplified as the "scientific method" (over)  
often causes conceptual + semantic confusion

Scientific Process - overly simplified cyclic conceptualization



Epid. H's are generated on the basis of models or conceptual frameworks which are based on existing knowledge of the subject area as well as "intuitive leaps" or hunches.

A model is a translation of postulated causal mechanisms; while emp. genl. are a translation of analysis of statistical relationships in the data. Thus, study results are analyzed and interpreted thereby leading to a revised theoretical understanding.

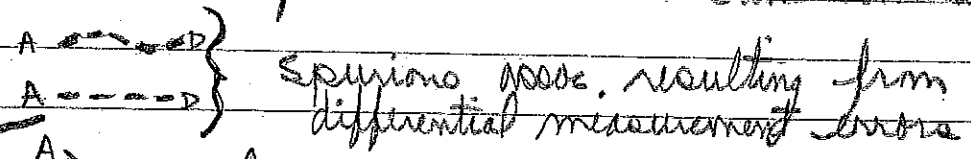
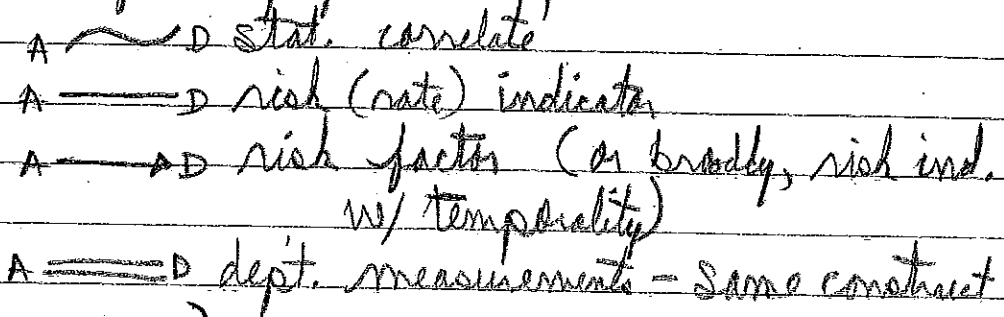
Some confusion exists between modelling & emp. genl., largely because of the lack of specificity in denoting theoretical postulates and statistical associations - thus a need to adopt common vocabulary -

II

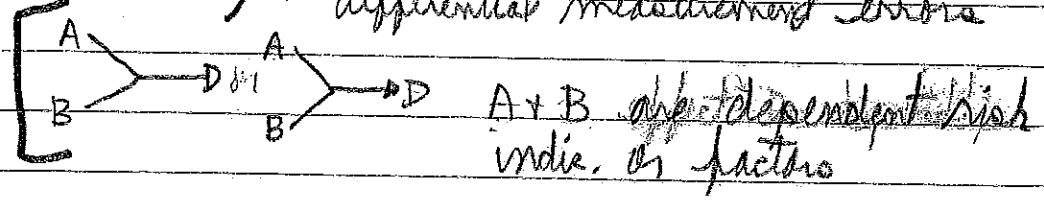
D. Begin by defining 4 conceptual levels of assoc. between 2 var. (E + D) 2

see p. 3 (#9) - hierarchy of assoc.; each level includes evidence for the levels above it. - i.e., graduation from data to theory

1. define 4 levels
2. define 3 types of inference needed to bridge gap between 4 levels
3. define graphical symbols -



no good



from empirical point of view perhaps theoretically meaningful

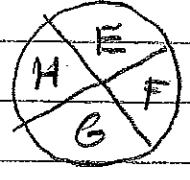
4. From analytical viewpoint, philosophical inferences are irrelevant. We can only hope to establish that statistical correl. are unbiased (ie risk indic.) and build an argument that they are risk factors of D. - still probabilistic statements - ie, stochastic view of causation

E. One procedure for modelling etiologic theory was presented at Minn. - involving concepts of "sufficient" conditions. [Modelling is the methodologic link for generating epid. H<sub>0</sub>s

II

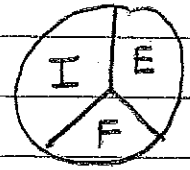
(eg.)  
E

I.



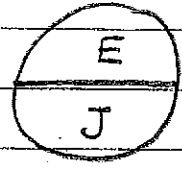
or

II.



or

III.



→ D

any  
model

3 sets of sufficient conditions  
 each circle consists of a set of dependent conditions  
 (ie. D results from interplay of all factors)  
 each set of suff. cond. is indep. of other sets  
 any factor (eg. E) which is present in all  
 suff. sets, is a necessary condition

F. Crucial point is that we recognize disparity  
 between theory & data; & between models  
 and emp. gen. - there is no consistent  
 1:1 correspondence between biol./beh.  
 mechanisms of causation + statistical assoc. -  
 made more evident in coming weeks -

(eg.) we must be very specific w/ regards to  
 graphical symbols  $\Rightarrow E \rightarrow D$   
 What does this mean?  $\rightarrow$

Does it refer to a postulated causal mech.,  
 or one or more sets of empirical evidence

(N.B.) path analysis  
 establishing "causal" associations is the essence of  
 analytic epid. & is what distinguishes us from  
 biostat.; et al.

defn. of indep. risk indic. - E is an indep. risk ind. of D  
 if it is assoc. w/ D in every category of C  
 E is a risk indic. if it is assoc. w/ D in at least  
 1 category of C - E is a dept. risk indic. if there  
 is statistical interaction

not same thing as "main effect?"  
 but a qualitative defn. (categorical analysis)

perhaps  
 define  
 later  
 FM

III. Confounding - most of # 8 isn't worth reading except pp. 1-15; revised outline

A. defn. - the heart of epid. methodology  
A defn. - conf. is the mixing of two or more <sup>effect</sup> risk indicators of a given disease - ie, the effect of E on D is confounded by C when the 2 effects have not been separated thru ~~the~~ control methods involved in subject selection or statistical analysis - assuming we are primarily interested in the effect of E, we would like to know whether the observed (crude) effect is distorted by C.

objects:

- (1) identify conf. factors
- (2) measure their degree of conf.

Conceptual

B. requirements of a single confounder C:

1. assoc. must exist in the data

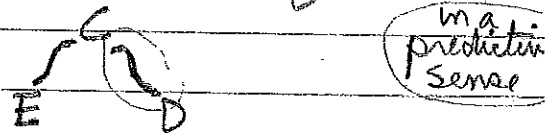
a.  $C \sim E$       w/ E

b.  $C \sim D$

[to primary] 2. + further specifications - (outline on board)  
[2 Secondary]

a. both assoc. must exist in the data (study pop.) - ie, conf. is an empirical condition not an inherent ~~own~~ property of certain variables

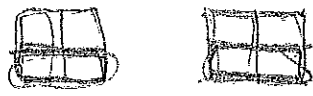
b. assoc. between C + D cannot be spurious  
C must be a risk indicator of D -  $\downarrow$   
 $\therefore$  not conf. if:



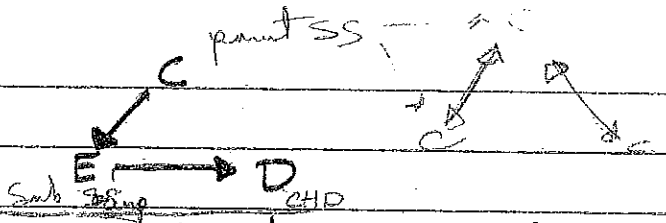
relates to  
SMP (or SSP)  
PP#

c. the prediction of D by C must occur without mediation by E - ~~ie, C must be a risk indicator~~ <sup>must be</sup> assoc. w/ D conditional on ~~an~~ E  
 $\therefore$  not conf. if.





III B. 2. c. (cont.)



C can be a ~~dept.~~ risk indicator  
 NB: an indep. risk indicator would mean that C was assoc. w/ D for every category of E (not just E)

~~thus C must also be indep. (not dept.) risk indicator of D~~

d. C may be conf. even if assoc. between E + C is spurious resulting from selection bias, differential measurement bias, or another conf. factor - also C may not be a true risk factor of D but only an "estensible risk factor" or risk indic. might occur because of differential errors in case detection - actually C might be only assoc. w/ a risk factor; but this ~~theoretically~~ is often difficult to distinguish - eg, is age or immunologic competence a risk factor of D?

secondary (e)

C must not be intervening link between E + D - ∴ not conf. if:



I think this is still comb. math. →

See notes from Black tabs

NB: not often possible to tell whether E predicts C since they are generally measured simultaneously - yet 1 piece of info. would suggest this pathway (i.e. C not conf.); no assoc. between E + D conditional on the lowest risk category of C (if C is causal) - also this pathway can often be ruled out a priori - eg, smoking (E) doesn't predict age (C).

Why not highest or any level?

but also true if

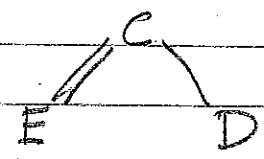


III

B. 2. f. when dealing with several potential conf. ~~they~~ their conf. effects must be considered simultaneously; this viewpoint is often violated in practice by examining pot. conf. sequentially -

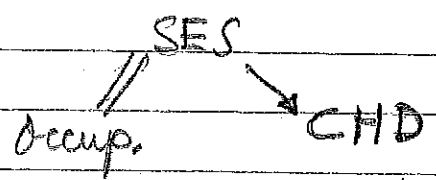
- when looking at the conf. effect of a single variable, it must be examined conditionally on all other conf. - thus the importance of a one conf. is reduced

(me) g. C must not be measuring one aspect or part of E - i.e. the measurement of C should be independent of E (or) C must not be measuring the same underlying theoretical construct as does E - ∴ not conf. if:



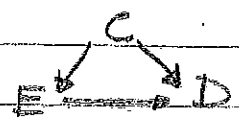
tho C would be conf in a math. sense, to control for it would lead to mis-interpretation - eg,

Similar to e.



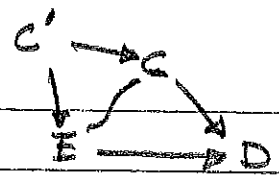
C. how conf. comes about - results from selection of subjects and measurement bias - [might define types of conf. by their causal configurations - causal config. may differ

1. "pure" conf - C predicts both E + D -  $\phi_{EC}$  is not spurious:

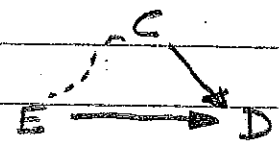


2. "multiple" conf. - spurious assoc. between E + C results from other (exogenous) factors which are themselves confounders - eg.

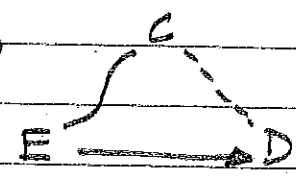
III C.C. 2.



3. "artificial" conf. - spurious assoc. between C + E results from correlation of their errors (ie, diff. meas. bias)

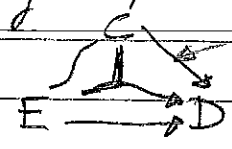


4. in a C-R study,



4. "complex" conf. - same factor (C) is both a conf. + a dept. risk factor of D (being dept. on the effect of E) -

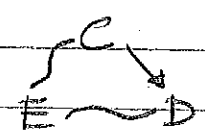
eg,



not necessary

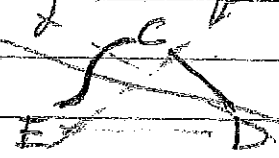
NB. strength of arrows may vary

5. "secondary assoc." (Special case of 1-3) - assoc. between E + D is entirely spurious - eg,



NB: sort of ultimate case of "conf." - yet there is no real assoc. between E + D to conf. - ∴ not mixing actually "mixing" of effect

~~D. with some reservation we can adopt the following <sup>general</sup> symbolization for conf.~~



needed

more sections  
v

III

E. but at the same time we should ~~not~~ be very careful in identifying conf. factors from statistical correlations - this <sup>fact</sup> may be elaborated w/ 3 principles.

NB: 1 set of symbols is inadequate to describe conf. - even for CH study (might be dept. effect)

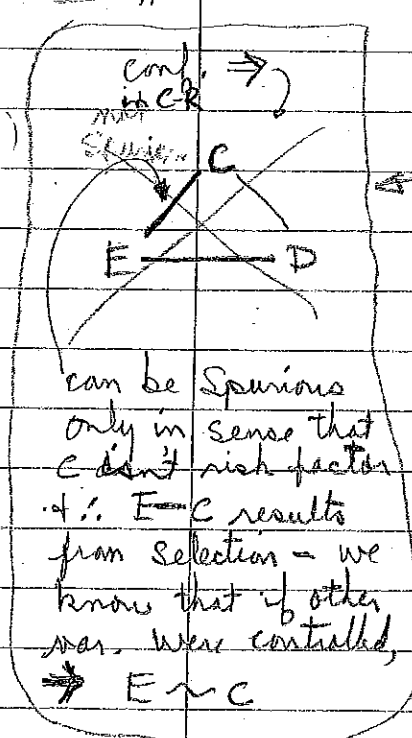
1. conf. cannot be identified by observing only crude assoc. - eg, conf. of car still occur even if  $\phi_{CD} = 0$  -  $C - D \Rightarrow \phi_{CD.E} \neq 0$

2. pertinent correlations (which reflect conf.) are different for CH + C-R studies -  $\phi_{C:conf.}$  to exist;

a. CH:  $\phi_{EC} \neq 0$  } if either is 0, no conf.  
 $\phi_{CD.E} \neq 0$  }

b. C-R  $\phi_{EC.D} \neq 0$  } if either is 0, no conf.  
 $\phi_{CD.E} \neq 0$  }

3. degree of conf. (which we have not yet specified how to measure) does not depend on the significance levels of the previous  $\phi$ 's because conf. does not involve <sup>statistical</sup> inferences to a source pop. remember: conf. only characterizes one's data - a large enough n would make even the smallest  $\phi$  correlation significant.



E. to quantify degree of conf. - bids - determine what crude effect (CRR) would be if exposure were to have no effect on the development of D. (ie, under  $H_0$ ) - this <sup>conf</sup> component of CRR (label it  $RR^*$ ) is then the crude effect that

III. F. would be expected in the absence of any true (unconf.) effect.

$$\therefore RR^* = E[CRR] = \frac{E[a]}{c/n_E}$$

	D	$\bar{D}$	
E	a	b	$n_E$
$\bar{E}$	c	d	$n_{\bar{E}}$

in CH study

1. CH study:

$$RR^* = E[CRR] = \frac{E[a]/n_E}{c/n_E}$$

$$SMR = \frac{a}{E[a]}$$

$$E[a] = a/SMR$$

$$\therefore RR^* = \frac{a/SMR/n_E}{c/n_E} = \frac{a/n_E}{c/n_E} \left( \frac{1}{SMR} \right) =$$

$$RR^* = \frac{CRR}{SMR}$$

2. C-R study:

$$RR^* = E[CRR] = \frac{E[a]d}{bc}$$

$$= \frac{ad/SMR}{SMRbc} = \frac{CRR}{SMR} \quad (\text{same})$$

$\therefore$  CRR has 3 components:  $RR^* \neq SMR$   
 G. if  $RR^* \neq 1$ , there is confounding  
 G. but in order to measure the amount of confounding, it is necessary to compare the conf. effect ( $RR^* - 1$ ) with the unconf. effect ( $SMR - 1$ ) -  $\therefore$  index

III.

G. Thus a useful index of conf. may be derived in order to decide whether stratification is necessary in analysis

RCE = (RR\* - 1) / (SMR - 1)

if RCE > 0 => "positive conf."
RCE < 0 => "negative conf."
RCE = 0 => no conf.

pos. conf => CRR is further from null value than is SMR + both are > 1

neg. conf => CRR is closer to null value than is SMR and both are

> 1 < 1 (or) CRR >
CRR > 1 / SMR < 1 (or)
CRR < 1 / SMR > 1

still no single cutoff for substantial conf. w/ RCE esp. for any single factor

H. this procedure implies that conf. is a quant. (not qual.) phenomenon - it takes the first 4 specifications into consideration but not the 3 (I believe) - i.e., just because C is math. conf. does not mean that it should be controlled in the interest of internal validity.

IV.

Illustrations and general principles of conf.

A. computations in a CH study.

Table with 3 columns: D, D-bar, and a total value. Rows include E and E-bar for C1.

Table with 3 columns: D, D-bar, and a total value. Rows include E and E-bar for C2.

IV A.  $CRR = \frac{14/80}{11/120} = 1.909$

$SMR = \frac{6+8}{\frac{9(20)}{100} + \frac{2(60)}{20}} = 1.795$

$RR^* = \frac{1.909}{1.795} = 1.064$

$RCF = \frac{1.064-1}{1.795-1} = .081$  ~~pos. conf.~~

would stratify  
if only conf.  
but don't know  
w/ more

Note: there was conf. even tho  $\phi_{ED} = 0$

	D	$\bar{D}$	
E1	15	105	120
$\bar{E}2$	10	70	80
	25	175	200

$\phi_{ED} = \frac{15(70) - 105(10)}{\sqrt{120(80)(25)(175)}} = 0$

but E was still a risk indicator of D because its assoc. w/ D conditional on  $\bar{E}D$

B. computations in a C-R study

e1

	D	$\bar{D}$	
E	40	10	50
$\bar{E}$	10	40	50
	50	50	100

e2

	D	$\bar{D}$	
E	25	25	50
$\bar{E}$	25	25	50
	50	50	100

$CRR = \frac{65(65)}{35(35)} = 3.449$

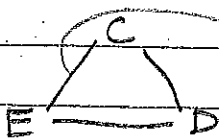
$SMR = \frac{40+25}{\frac{10(10)}{40} + \frac{25(25)}{25}} = 2.364$

$RR^* = \frac{3.449}{2.364} = 1.459$

IV. B. (cont.)  $RCE = \frac{1.459-1}{2.364-1} = .337 \Rightarrow$  pos. conf.

NB. there is conf. even tho  $\Phi_{EC} = \Phi_{CD} = 0$   
(ie, crude assoc are 0)

note: that criteria for conf. in C-R study:



\* can be "spurious" only in the sense that C doesn't cause E, but if we controlled for other relevant vars.,  $E \sim C$  - (p.9)

implications  
C. procedure often used in epid. study to eliminate possibility that C (eg, age) was conf. after eliminating ~~the~~ sign. assoc. between E + D

1. C-R study: compare cases + noncases with respect to age distrib.  $\Rightarrow \Phi_{CD} \neq 0$   
no good since criteria is  $\Phi_{CD.E} = 0$

eg, age	D	$\bar{D}$
Y	-	-
0	-	-
$\vdots$	-	-

(Kolmogorov-Smirnov test)

2. CH study: compare  $E + \bar{E}$  w/ respect to age distrib.  $\Rightarrow \Phi_{EC} = 0$   
this is ok since it is a criterion for conf.



CH study

IV D' -

		D	$\bar{D}$	
E		10	10	20
$\bar{E}$		6	24	30
		16	34	50

		D	$\bar{D}$	
E		30	70	100
$\bar{E}$		10	40	50
		40	110	150

$$CRR = \frac{140/120}{16/80} = 1.667$$

$$SMR = \frac{40}{4+20} = 1.667$$

because:

$$\phi_{CD, \bar{E}} = 0$$

$$RR^* = 1 \quad RCE = 0 \Rightarrow \text{no conf. due to C.}$$

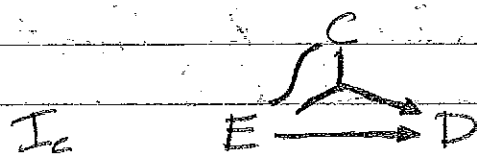
		D	$\bar{D}$	
E	C	10	10	20
$\bar{E}$	C	30	70	100
		40	80	120

		D	$\bar{D}$	
E	C	6	24	30
$\bar{E}$	C	10	40	50
		16	64	80

$$CRR_c = \frac{16/50}{40/150} = 1.2$$

$$SMR_c = \frac{16}{6+6} = 1.333$$

$$RR^* = .900 \quad RCE = -.300 \Rightarrow \text{neg. conf. due to E}$$

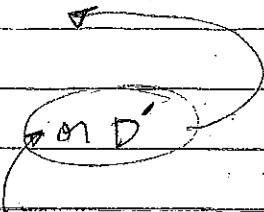


both C+E are risk indicators  
but C is not indep. risk indicator

	C	$\bar{C}$
E	.5	.5
$\bar{E}$	.3	.1

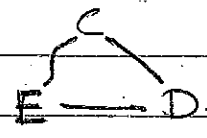
∴ better example because:  
 $SMR_c \neq 0$   
 $SMR_E \neq 0$

NB: Same results if analyzed as C-R study  
 $RBE_c = 0 \quad RCE_E = -.300$



IV.

B. Sometimes interested in <sup>causal</sup> effect of 2 var. (E+C)  
 conf by c (of E-D assoc.) is not necessarily = to conf. b (of C-D assoc.)



eg consider IV, A. - RCE = .081 (pos conf.)  
 B imagine E as a potential conf of C-D assoc.

E	D	$\bar{D}$		E	D	$\bar{D}$	
c1	6	14	20	c1	9	2	11
c2	8	52	60	c2	91	18	109

$$CRR_c = \frac{15/31}{99/169} = .826$$

$$SMR_c = \frac{6+9}{\frac{8(20)}{60} + \frac{91(11)}{109}} = 1.266$$

$$RR_c^* = \frac{.826}{1.266} = .652$$

$$RCE_c = \frac{.652-1}{1.266-1} = -1.307 \text{ neg. conf.}$$

Thus magnitude + sign of conf. is different for E+C.  
 "conf. is not symmetrical" in (CHA C-B study)

IV. E. control tables should not be used to identify or measure conf. - tho they often are, under the mistaken view that conf. is a qualitative phen. - ie, seeing conf. as the presence of secondary assoc. - eg, CH Study

C1			C2				
	D	$\bar{D}$		D	$\bar{D}$		
E	30	70	100	E	4	16	20
$\bar{E}$	4	16	20	$\bar{E}$	6	54	60
	34	86	120		10	70	80

Control table:  $T_c$

	C1	C2
E	.30	.20
$\bar{E}$	.20	.10

⇒ could say still E effect, after controlling for C, but is C conf.?

CRR =  $\frac{34/120}{10/80} = 2.267$

SMR =  $\frac{30+4}{\frac{30(100)}{20} + \frac{4(20)}{60}} = 1.545$

RR\* =  $\frac{2.267}{1.545} = 1.467$

RCE =  $\frac{1.467-1}{1.545-1} = .856$  pro. conf.

conceptual

(p. 15)

3 reasons why we need to know if conf. exists + how much

1. know whether simple analysis (main eff) can be used
2. When considering several part. conf. <sup>simult.</sup> we would like to select the combin. which accounts for the most bias (most eff.)
3. generally, conf. is quantitative phen. + thus we should like to know the magn. of unconf. effect (SMR)

but suppose data were:

C1			C2				
	D	$\bar{D}$		D	$\bar{D}$		
E	15	35	50	E	10	40	50
$\bar{E}$	10	40	50	$\bar{E}$	5	45	50
	25	75	100		15	85	100

IV. E (cont) Same control table,

$$CRR = \frac{25/100}{15/100} = 1.667$$

$$SMR = \frac{15+10}{\frac{10(50)}{50} + \frac{5(50)}{50}} = 1.667$$

RR\* = 1      RCE = 0      no conf.

3 NAMA  
P. 14

However, control tables do allow us to identify secondary associations, tho they seldom exist in practice - (ie) SMR = 1; but doesn't tell whether conf. exists which is the criterion for using simple or stratified analysis (simple analysis is more efficient) - eg, Ic

	E1	E2
E	.3	.3
Ē	.1	.1

⇒ no E effect ∴ SMR = 1  
but what is CRR?

\* NB: (I) can't be used in C-R studies (II) but do allow us to look for CM + I

(note: can do as C-R study)

F. Consider 2 potential conf. - Sex + race (CH)

	WM			BM			WF			BF		
	D	D̄		D	D̄		D	D̄		D	D̄	
E	56	104	160	2	18	20	1	19	20	46	69	115
Ē	14	26	40	6	54	60	3	57	60	6	9	15
	70	130	200	8	72	80	4	76	80	52	78	130
	RR = 1.00			1.00			1.00			1.00		

Note: IR<sub>i</sub> = 1 for all 4 strata = RR<sub>i</sub>

$$CRR = \frac{105/315}{29/175} = 2.01 \quad n = 490$$

34 356  
27

IV. F. (cont.): control for race only

	W			B		
	D	$\bar{D}$		D	$\bar{D}$	
E	57	123	180	48	87	135
$\bar{E}$	17	83	100	12	63	75
	74	206	280	60	150	210

$$SMR = \frac{57 + 48}{\frac{17(180)}{100} + \frac{12(135)}{75}} = 2.01$$

RR\* = 1      RCE = 0      no conf.

control for sex only

	M			F		
	D	$\bar{D}$		D	$\bar{D}$	
E	58	122	180	47	88	135
$\bar{E}$	20	80	100	9	66	75
	78	202	280	56	154	210

$$SMR = \frac{58 + 47}{\frac{20(180)}{100} + \frac{9(135)}{75}} = 2.01$$

RR\* = 1      RCE = 0      no conf.

control for sex + race simultaneously

$$SMR = \frac{105}{\frac{14(160)}{40} + \frac{6(20)}{60} + \frac{3(20)}{60} + \frac{6(115)}{15}} = 1$$

RR\* = 2.01      RCE =  $\frac{2.01-1}{1-1} = ?$       no pos conf.

IV E.(cont.) Thus when controlling for both simult., there is no E effect - yet this fact is not discovered by controlling for Sex + race sequentially.

Also, Control table analysis of each var. fails to show no E effect.

		I <sub>c</sub>				
		M	F			
E	.322	.348		E	.317	.356
E	.200	.120		E	.170	.160

In order to determine whether one var (eg, sex) is conf., it must be viewed conditionally on the other - thus,

for whites

$$CRR = \frac{57/180}{177/100} = 1.863$$

$$SMR = \frac{56 + 1}{\frac{14(160)}{40} + \frac{3(20)}{60}} = 1$$

$$RR^* = 1.863 \quad RCE = \frac{.863}{0} = ? \quad \infty \text{ pos. conf.}$$

for Blacks

$$CRR = \frac{48/135}{12/75} = 2.222$$

$$SMR = \frac{2 + 46}{\frac{6(20)}{60} + \frac{6(115)}{15}} = 1$$

$$RR^* = 2.222 \quad RCE = \frac{1.222}{0} = ? \quad \infty \text{ pos. conf.}$$

IV

G. It is conceivable that when considering more than one pot. conf. variable, more bias will be discovered when stratification is done by combinations of factors - eg,  $C \times C'$  as an interactive term - this was done by Kleinbaum et al. in Evans county when they derived a Multiple logistic <sup>risk</sup> function predicting CHD. eg,  $age \times chol.$ ;  $age \times BP$  - only problem might come if we insist on interpreting the central variables; what do they mean - yet with several pot. conf. & a limited  $n$ , considering interactive terms systematically can make the estimation more precise - OM, however has formulated an alternative method that we will discuss later in the term; it's designed to deal with several conf. factors by combining them into a single risk function thru multivar. analysis - it appears that such a methods makes the consideration of interactive conf. terms a secondary issue

not needed if  
that want to  
central for  
conf.

V

methods of central - see revised outline  
2 basic approaches (#8)

H. (p.19)

addendum: also relates to Interaction

IV, H. recall spec. "e"; math. conf. but not to be controlled if.



eg, CH. study

	C			$\bar{C}$ (low risk)			
	D	$\bar{D}$		D	$\bar{D}$		
E	30	120	150.	E	2	48	50
$\bar{E}$	5	45	50.	$\bar{E}$	16	144	150
$I_{c(c)} = .175$	35	165	200	8	192	200	$I_{c(c)} = .04$
	$\hat{RR}_E = 2$			$\hat{RR}_E = 1$			

$$\hat{CRR} = \frac{32}{11} = 2.909$$

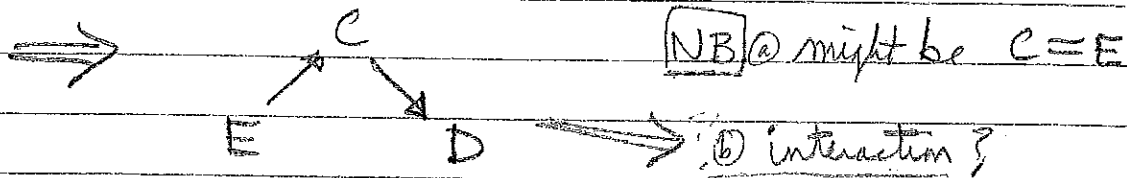
$$\hat{SMR} = \frac{32}{15+2} = 1.892$$

$$RR^* = 1.545 \quad RCE = .618 \text{ (pos. conf.)}$$

but  $\hat{RR}_E = 1$  (no effect)

$I_c$

	P	$\bar{C}$
E	.2	.04
$\bar{E}$	.1	.04



? what about if  $\hat{RR}_E = 1$ ? (over)



IV. H<sub>0</sub>: CH study -  $\hat{RR}_c = 1$   $\hat{RR}_c > 1$

	D	$\bar{D}$	
E	15	135	150
$\bar{E}$	5	45	50
	20	180	200

	D	$\bar{D}$	
E	5	45	50
$\bar{E}$	5	145	150
	10	190	200

I<sub>c</sub>

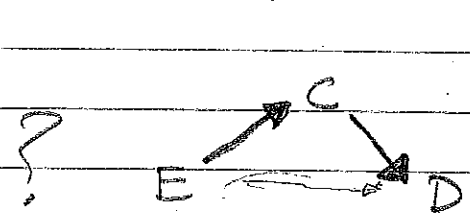
	c	$\bar{c}$
E	.1	.1
$\bar{E}$	.1	.0333

$CRR = 20/10 = 2$

$SMR = \frac{20}{15 + \frac{5}{3}} = 1.2$

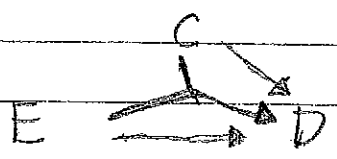
$RR^* = 1.667$   $RCE = 3.333 \Rightarrow$   
*(pos. conf.)*

$\hat{RR}_c = 1$  (hi risk cat. of c)



not as much  
 no? because "c" subjects (ie, hi risk) may have gotten that way because of other causes of c besides E

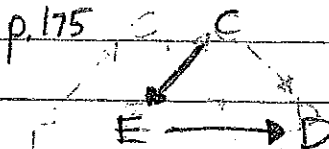
\* Can this config. also manifest in the data as statistical (and/or "causal") interaction? yes



$\therefore$  underlies difference in theoretical & empirical levels: (eg)  $\text{smk} \rightarrow \text{age} \rightarrow \text{CHD}$

\* same problem w/  $\hat{RR}_c = 1$  (p. 19)


Notes on Causal Inference

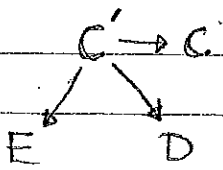
1. p.175  (Blalock, 1964) [E as background variable]

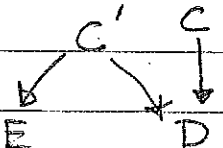
observe:  $r_{ED}$  vs  $r_{ED.E}$  (partial correl.)  
 $b_{ED}$  vs  $b_{ED.E}$  (partial regression)  
 [C is not a conf.]

prove w/ SMR

If E is intervening (as drawn), controlling for E will reduce r (even if  $E \rightarrow D$ ), but will not reduce  $b$  (similarly w/  $RR_{E.C} = SMR_E$ )  
 the stronger  $r_{EC}$ , the more reduction in  $r_{ED.E}$

  $\Rightarrow$  both r and b are reduced by controlling for E (see # 4)

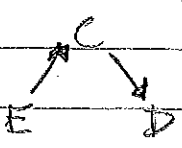
2. p.177   $\Rightarrow r_{ED.C} < r_{ED}$

  $\Rightarrow r_{ED.C} \geq r_{ED}$  ( $r_{C'C} = 0$ )

3. p.190  $r_{xy} = b_{xy} \left[ \frac{S_x}{S_y} \right]$  [in brackets]  $\rightarrow$  peculiar to pop.  
 $b$ 's  $\rightarrow$  peculiar to variable

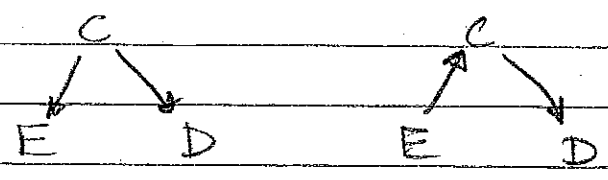
$$r_{xy.z} = b_{yx.z} \left[ \frac{\sqrt{1-r_{xz}^2} \cdot S_x}{\sqrt{1-r_{yz}^2} \cdot S_y} \right] = \frac{r_{xy} - r_{xz} r_{yz}}{\sqrt{1-r_{xz}^2} \sqrt{1-r_{yz}^2}}$$

(me) (4)

  $RR_{E.C} (=SMR) \neq CRR_E$  tho  $\leq$  should not be controlled for in interest of interpretation

p.175

4.



both cases:  $r_{ED,C} \cong 0$

$r_{ED,C} < r_{ED}$

same w/ b

1/21/76 Outline 1/28 3/4

I. Introd.

- A. handouts - revised outline of Epid. Measures - #4
  - B. today - finish conf.
- EM + I

change -  
to 3/4 divide  
of EM.

II. Confounding (cont.) see: 11/14/76 p. 11 RR\* RCE

III. Effect Modif. - defn.

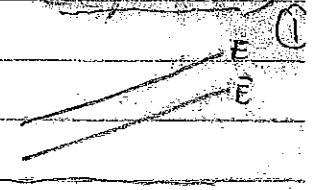
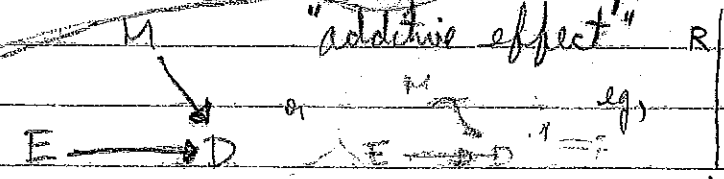
- A. multiple risk indicators (factors) in all dis.
  1. Stat. correl. of dist. = factor is assoc. w/ rate of dis.
  - (EM) 2. correlate of effect = factor (M) is assoc. with the effect of E on D  $\Rightarrow$  M is an effect modifier if the effect (RR or RD) between E + D is not uniform among all categories of M.

risk indicators  
Stat. correlates

- B. EM requires that 2 (or) variables (E + M) are risk indicators w/ rate of dis. - doesn't
- C. EM  $\therefore$  doesn't represent a singular biosocial mechanism of causation, but may result from other classes of conceptual configurations - ie, 4 types of EM

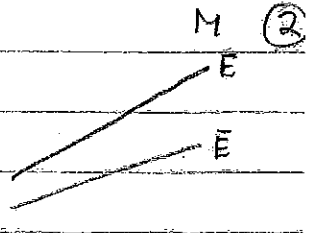
NO  
RCE

1. E + M are indep. risk factors (indicators) "additive effect" R

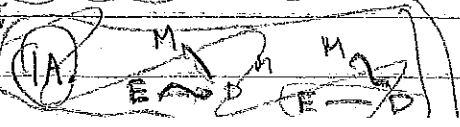


ie, E + M is assoc. w/ D for every category of M + subcategory

must use RR to show EM " " RD " " " "

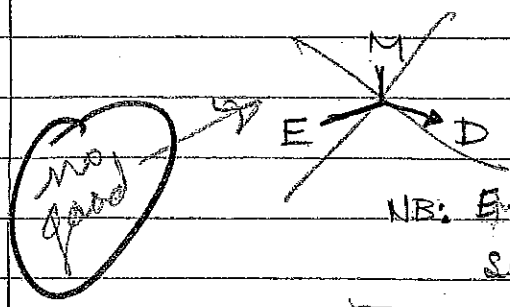


needn't be causal

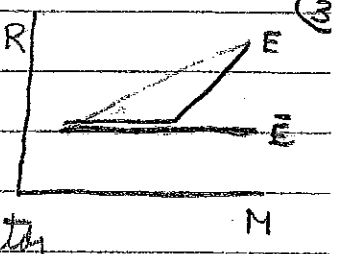


but not additive effect

III C. 2. M + E are dept. risk factors



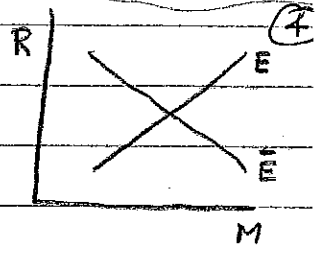
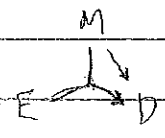
"interactive effect"  
(when curves are non-parallel)



NB: E + M are not indep. risk factors

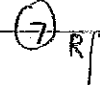
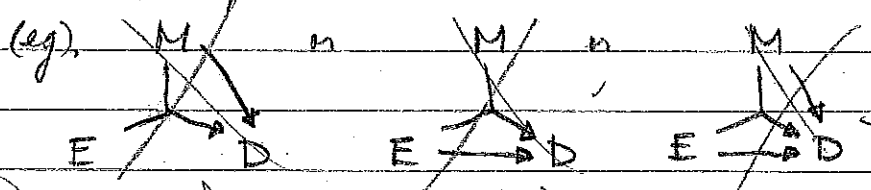
see (2) - using RD, no EM, yet interaction

yet sometimes I does necessitate EM



3. indep. + dept. risk factors aren't mutually exclusive categories - a var. can be both

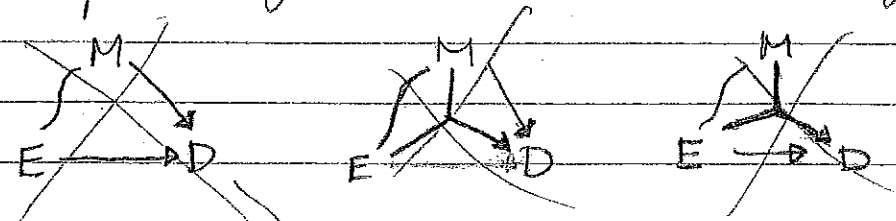
3. M + E as combination



NB. from a statistical standpoint, a var. (E + M) is not an indep. risk indicator doesn't mean it does not have a significant "main effect" on D

(in qualit. sense)

4. Complex conf. - M is both EM + confounder



D. confusing pt. C.5 does not automatically imply that conf. automatically suggests EM - altho the diagram seems to indicate such - will be illustrated later

III. E. Note difference in terminology

\* mention main statistical nature of dept/underp.

me	AM.	Fleiss
EM	"trends in effect"	"homogeneity of assoc."
I.	"EM"	-

begin

IV. Comparison of Conf., EM, + I as separate issues - [p. 14]  
 A. I. Conf.  $\rightarrow$  purpose in identifying - how internal validity is threatened  
 Conf. is a f. of (determinants)

pooling estimates no stratification p. 5

B. 2. EM - " " how to control for conf  
 C. 3. I - " " causal dept.  
 (D) note disagreement w/ Rothman - Conf. is f. (data); I = f. (causal assoc.)  
but based on empirical relationships

V. Identifying EM + I - control table analysis (CH study)

A. EM -

		I <sub>c</sub>				RR	
		M	$\bar{M}$			M	$\bar{M}$
E		R <sub>11</sub>	R <sub>01</sub>	E		RR <sub>11</sub>	RR <sub>01</sub>
$\bar{E}$		R <sub>10</sub>	R <sub>00</sub>	$\bar{E}$		RR <sub>10</sub>	RR <sub>00</sub>

ref.  $\rightarrow$

RR<sub>11</sub> = RR<sub>01</sub>  $\Rightarrow$  no EM  
 RD<sub>11</sub> = RD<sub>01</sub>  $\Rightarrow$  " "

		RD	
		M	$\bar{M}$
E		RD <sub>11</sub>	RD <sub>01</sub>
$\bar{E}$		0	0

B. I. control table

		I <sub>c</sub>				RD	
		M	$\bar{M}$			M	$\bar{M}$
E		R <sub>11</sub>	R <sub>01</sub>	E		RD <sub>11</sub>	RD <sub>01</sub>
$\bar{E}$		R <sub>10</sub>	R <sub>00</sub>	$\bar{E}$		RD <sub>10</sub>	0

RD<sub>11</sub> = RD<sub>01</sub> + RD<sub>10</sub>  $\Rightarrow$  no I. (ie, additive effects)  
 " > " "  $\Rightarrow$  I - synergistic  
 " < " "  $\Rightarrow$  I - antagonistic

\* poor words because they suggest bully, precease, etc

ERR = RR - 1

IV. B. also, contrast control table via CR study

$ERR_{11} = ERR_{01} + ERR_{10} \Rightarrow \text{no I}$

	M	$\bar{M}$
E	ERR <sub>11</sub>	ERR <sub>01</sub>
$\bar{E}$	ERR <sub>10</sub>	0

C. Examples.

1. GH study:

	M	D	$\bar{D}$	
E	60	90	150	
$\bar{E}$	10	40	50	
		70	130	200

	M	D	$\bar{D}$	
E	10	40	50	
$\bar{E}$	15	135	150	
		25	175	200

	M	$\bar{M}$
E	.4	.2
$\bar{E}$	.2	.1

$RR_M = RR_{\bar{M}} = 2 \Rightarrow \text{no EM (RR)}$

NBM:  $RD_M = .2$   
 $RD_{\bar{M}} = .1 \Rightarrow \text{EM (RD)}$

	M	$\bar{M}$
E	3	1
$\bar{E}$	1	0

$3 > 1 + 1 \Rightarrow \text{I (Synergy)}$

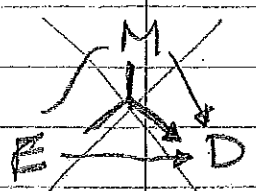
EM + I can be deduced from control table

$CRR = \frac{70/200}{25/200} = 2.8$

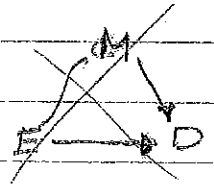
$SMR = \frac{10 + 60}{\frac{10(150) + 15(50)}{50}} = 2.0$

$RR^* = \frac{2.8}{2} = 1.4$

$RCE = \frac{1.4 - 1}{2 - 1} = .4 \Rightarrow \text{pos. conf.}$



$\therefore \text{conf.}$



doesn't imply EM (RR)

(don't know I<sub>c</sub>)

II. c. Ex. - G-R study - no control table, but similar principle

→

M	D	$\bar{D}$	$\bar{M}$	D	$\bar{D}$
E	50	10	60	60	20
$\bar{E}$	30	10	40	20	20
	80	20	100	80	40

ref.

$$OR_{11} = \frac{I_{DR_M}}{I_{DM_M}} = \frac{50(10)}{10(30)} = 1.6$$

$$OR_{01} = \frac{I_{DR_{\bar{M}}}}{I_{DM_{\bar{M}}}} = \frac{60(20)}{20(20)} = 3.0$$

⇒ EM (OR)

stat. test.

∴ overall RR derived from pooling strata of M would be misleading - tho, recall we can standardize to control for conf.

\* 2 basic ways to form weighted avg. of strata & spec. estimates

$$SRR = \frac{\sum w_j R_{Ej}}{\sum w_j R_{\bar{E}j}}$$

↓

① Adjustment

$$A_0 = \frac{\sum W_j A_j}{\sum W_j}$$

$A_0$  = overall MoB assoc. (eg. RR)  
 $A_j$  = measure of Assoc.  
 $w_j$  = based on stratum's stand.  
 $W_j$  = based on  $Var[A_j]$

② genl. stratification

ERR

	M	$\bar{M}$
E	4	2
$\bar{E}$	2	0

3-1 = 2

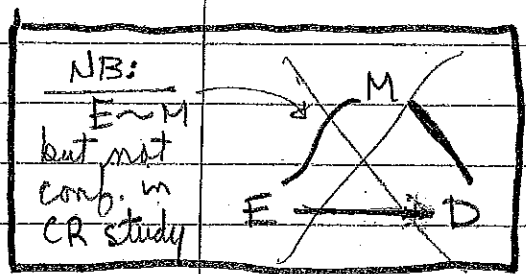
30(20) - 1 = 2  
10(20)

50(20) - 1 = 4  
10(20)

4 = 2 + 2 ⇒ no I.

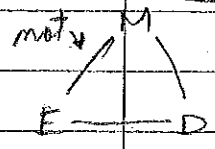


II. c. (cont.)  $\hat{CRR} = \frac{110(30)}{130(50)} = 3.2$



$\hat{SMR} = \frac{50 + 60}{\frac{10(30)}{10} + \frac{20(20)}{20}} = 2.2$

$RR^* = 1$      $RCE = 0 \Rightarrow$  no conf.



i. I. does not necessarily imply conf. actually conf, EM, + I are separate issues

III. Quantifying EM + I - statistical procedures

- A. EM - qualitatively - ie, is EM statistically significant?
1. qualitatively - ie, is EM stat. sign? partitioning of  $\chi^2 = f[A]$  q strata of

$$\chi^2_{u(q-1)} = \sum \chi^2_j(q) - \chi^2_0(1)$$

$\chi^2_j$  = usual  $\chi^2$  for 1 strata - several ways to compute, most common =  $\sum_{all\ cells} \frac{(O-E)^2}{E}$

will be discussed in more detail later (#10)

$\chi^2_0$  = test of overall assoc. between E and D controlling for M (if it is conf.)

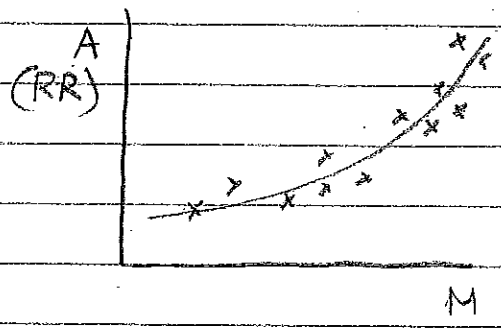
=  $A_0^2 \left( \frac{\sum w_j}{\sum w_j} \right)^2$      $A_0$  = overall M of assoc.     $w_j = 1/N[A_j]$

will also be covered later

$\chi^2_u$  = test of uniformity of effect (EM); ie, if  $p_0$  is signif (< .10) there is signif EM.  $A_0$  is not a good estimate of overall effect

VI

A. 2.  $\chi^2$  test is not very sensitive test since no order is assumed in the categories of M -  $\therefore$  one highly weighted stratum could distort the overall picture of EM -  
 a more sensitive test might be made if M is at least ordinal - we could fit the data to a specific statistical model



B. I - normally measured qualitatively thru statistical models ("multivariate") - eg,

linear:  $\hat{R} = \beta_0 + \beta_E E + \beta_M M + \beta_{ME} ME = \beta_i X_i$

not good fit for  $R < .24 > .8$

log linear:  $\ln \hat{R} = \beta_i X_i$   
 $\hat{R} = \exp[\beta_i X_i]$

often gives E values  $< 0$  or  $> 1$

logistic model:  $\ln \left[ \frac{\hat{R}}{1 - \hat{R}} \right] = \beta_i X_i$

2 or 3 disadvantages  $\rightarrow$

$\hat{R} = \frac{1}{1 + \exp[-\beta_i X_i]}$

(1) don't get degree of interaction (except  $\beta$ )

disadv of cont interaction since natural trends of interaction remain hidden

while this model is fine (esp. for many vars) there is another easier, more direct way to not only detect interaction but to measure it quantitatively (based on cont tables)  
 there is also another disadvantage in using model fitting which we'll discuss later (ie. when we'd like to control for conf. effects + study I, simultaneously)

VI. B. Rothman's index of interaction (S) = RIE = <sup>relative</sup> <sub>inter-</sub> <sub>effect</sub>

= ratio of observed joint effect from indep. var. (E+M) to the effect expected if each var. were acting indep.

$$RIE = \frac{ERR_{EM}}{\sum [ERR]_{indep.}}$$

how to compute:

	ERR	
	M	M
E	ERR <sub>11</sub>	ERR <sub>01</sub>
E	ERR <sub>10</sub>	0

could have more categories

$$RIE = \frac{ERR_{11}}{ERR_{01} + ERR_{10}}$$

- RIE = 1 ⇒ no I
- RIE > 1 ⇒ "Synergy"
- RIE < 1 ⇒ "antagonism"

Rothman also gives SE & CI of RIE

Criticisms of RIE:

NB. Stat. interaction does not have to imply some biologic interaction (of any factors)

1. if there are many categories of E and/or M, the effect for each cat. of E (or M) controlling for the other must be increasing (or b) monotonically, - if not, computation of RIE is ambiguous & misleading

(eg)

ERR	M1	M2	M3
E1	4	3	2
E2	3	6	3
E3	2	3	0

RIE = ?

interaction is an artifact

both refer to same etiologic process

not really a criticism of RIE but how it is interpreted (theoretically)

when RIE < 1, interpretation of antagonism may be <sup>especially</sup> deceptive, esp. if M or E is an abstract var. eg SES. - if M = E (measuring something), RIE would be < 1, so this measure of I may not relate to biological or soc. processes - ie, the statistical assoc. may have no etiologic significance

8A

A. ~~there is no~~

A. stat. interaction does not imply there is any biological/behavioral interaction (of any sort) because, might have:

- ? not sp. ① M + E are measuring same thing  $\Rightarrow RIE < 1$
  - ' spurious ② differential measuring errors (eg, correlated errors)
  - spurious ③ M<sup>ME</sup>: assoc. w/ lag period of Dis.
  - not sp. ④ non-diff. meas. errors  $\Rightarrow RIE \rightarrow 1$  ?
- 3 types of artifacts bias

- ⑤ selection ie, 3 factors that prod. bias between E + D
- ⑥ conf. 2 factors above (1+3)

B. ~~there~~ in general, there is no 1:1 correspondence between stat. ~~in~~ and bio-beh. mechanisms  $\rightarrow$

VI

e. this underscores the lack of direct correspondence between biosocial mechanisms of causation (theoretical level) and statistical assoc. (empirical)

consider the example of infectious hepatitis (IH study) -

		SES	
		hi	low
age	old	hi	↑
	young	↓	hi

opposite trends for each stratum of SES

∴ Age + SES interact. But so what? does it indicate "synergy" in a biomedical sense? no

the hypothesis put forth & later validated was that SES was actually reflecting the difference in hygienic practices. young persons of low SES families w/ poor hygiene have a hi risk of developing dis (the symptoms are minimal); as they get older, immunity is developed & fewer low SES persons are susceptible - w/ high SES, young children are protected from virus until they get older, at which time the risk of D increases

VII

or stress and social supports

thus: relationships between statistical interaction + biosocial interaction is not necessarily direct nor is the connection always obvious, esp. when presented w/ statistical info. & attempting to infer biosocial mechanisms - yet biological or bios. interaction is extremely important to pathogenesis of chronic disease. acknowledge the importance of (at least) Env + genetic factors

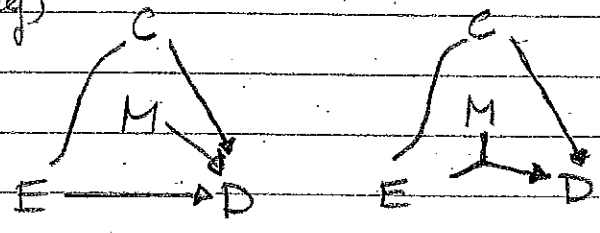
also control for conf.

VII

EM and I

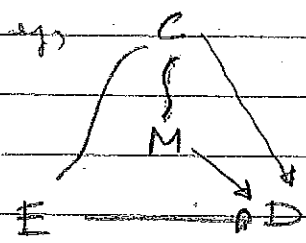
Consider the study of EM or I, when having to control for part. conf. factors simultaneously

(eg)



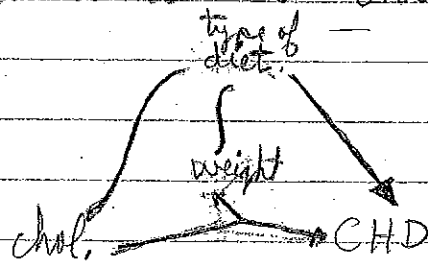
etc.

usual procedure of studying interaction of E+M & controlling for C is thru multivar. analysis specifically covariance - OM has pointed out that multivariate would be deceptive if C were highly assoc. w/ M (which they usually are, even spuriously)



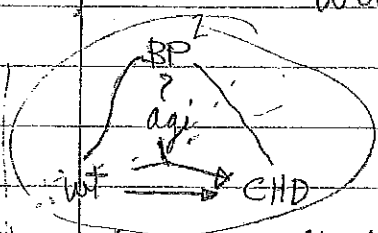
the model controls for C by keeping it "fixed", i.e. any correlate of C (eg, M) will be made statistically meaningless with regards to its interactive effect on D

(eg)



if we control for dieting, wght will be rendered a hollow concept (to a certain extent)

turn expression



Ho: wt. is more critical for young than old

OM's solutions: involve correct choice of effect measure

- 1. study EM - use SMR (internal stand.) for each category of M & compare - N.B. that a different "standard" is used for each stratum of M, but not a problem

begin → (2/4)  
 view 2 disadvantages of model-fitting to study interaction  
 @ interp. & conf. p.7

- VII. A. 1. (cont) comparability because each stand. measure is computed within 1 stratum of M, comparing E and  $\bar{E}$  (ie, only 2 categories of E)  
 $\therefore$  the importance C as a conf. is within each stratum of M
2. Study I - can find RIE by externally standardizing for C within each stratum of M -  $\therefore$

$$RIE = \frac{SERR(E)_{\bar{E}}}{E[SERR]_{indep.}}$$

where:  $SERR = SRR - 1$  (ext. stand.)  
 must be ext stand. since stand. for each comparison must be identical.

- stand. 1. ref. categories (cont)  
 2. ref. gen. pop.  
 3. other gen.

eg. #9 p. 42 Oral Ca.  
 3. do a problem \* numerically - CH study  
 look for interaction of  $\bar{E}$  and  $\bar{M}$  controlling for sex (C)

Exam

	WM			WF			BM			BF		
	D	$\bar{D}$		D	$\bar{D}$		D	$\bar{D}$		D	$\bar{D}$	
E	25	75	100	15	35	50	10	40	50	50	50	100
$\bar{E}$	5	45	50	15	85	100	5	95	100	5	45	50
	30	120	150	30	120	150	15	135	150	55	95	150
	$I_c R = 2.5$			2.0			4.0			5.0		

a. Suppose we ignored M (race) & determined if C (sex)  
 $CRR = (I_c R) = \frac{100/300}{30/300} = 3.33$  where conf.

$\Delta$   $SMR = \frac{100}{\frac{10(150) + 20(150)}{150}} = 3.33$

$RR^* = 1$   $RCE = 0 \Rightarrow$  no conf. due to sex but wrong

VII, A. 3. (cont.) - we might then proceed to calculate the interactive effect between E+M (race), ignoring C (sex); if we did,

		D	$\bar{D}$	
E	40	110	150	
$\bar{E}$	20	130	150	
	60	240	300	

		D	$\bar{D}$	
E	60	90	150	
$\bar{E}$	10	140	150	
	70	230	300	

		W	B
E	.267	.400	
$\bar{E}$	.133	.067	

 $\Rightarrow$ 

		W	B
E	3	5	
$\bar{E}$	1	0	

Since  $3 < 5+1 \Rightarrow \underline{I}$  ("antagonism")

$$RIE = \frac{3}{1+5} = \underline{.5}$$

b. Suppose we now consider C (sex) in the above analysis thru stratif. standardizing.

		SERR(E)		
		W	B	
E	3.00	3.50	$\xrightarrow{SRR=1} SERR = \frac{10(100) + 50(50)}{10} = 17.5$	
$\bar{E}$	.75	0		
			$\xrightarrow{SRR=1} SERR = \frac{5(100) + 15(50)}{10} = 1 = .75$	
			$\xrightarrow{SRR=1} SERR = \frac{5(100) + 15(50)}{10} = 1 = 3.00$	

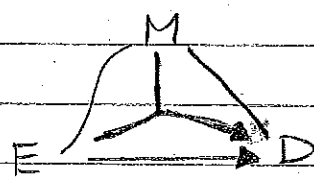
$$RIE = \frac{3}{3.5+.75} = \underline{.71} \text{ (antag.)}$$

\* but:  $[\underline{.71} \neq .50]$  - why?  
 because sex was conf. - RR\* (above) was calculated incorrectly



VII. 1 A. 3. RR\* was incorrect because it ignored race (M) which is itself a conf. - when considering the conf. effect of a variable, it must be viewed conditionally on all other confounders - thus in getting RIE controlling for sex, we must standardize separately for each comparison made; and stand. must be done externally to ensure that the same stand. is used for all comparisons

4. another possibility is that C = M, ie, "complex conf" since C is same as M, it is impossible to consider it simult. as both a conf. + E Mod. -



perhaps we could control for M w/in each category of M?  
 can't

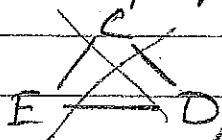
in the analysis, M as a E mod. becomes irrelevant since we must control for conf. in the interest of internal validity - use SMR w/ 2 categories of E or SRR w/ > 2 categories if looking for dose-response gradient

5. ~~planning a study on the basis of expected conf. in the data (or lack of it) can then be complicated by the presence of EM or E. - eg, C-R study:~~

~~ABO Blood group (E) → → Cervical Ca. (D)~~

III

A. 5. can we have a control series of male non-cases, altho this ~~paradoxical~~ design seems paradoxical we might defend it on the basis that sex is not assoc. w/ ABO Bl group in the source pop. & thus conf. won't be expected in the data ~~if no selection is introduced~~ but recall that  $\phi_{CE} \neq 0$  is not a condition for conf. - instead  $\phi_{CE, D}$ ; ie age can't be assoc. w/ ABO Bl group, conditional on being a non-case;



this condition will not be met, since the risk of D is greater for unexposed females than unexp. Males (where it is 0) - this fact suggests an interactive effect between sex + ABO Bl groups. to illustrate, consider a hypoth. C-R study of all females (assume all other ~~vars.~~ have been controlled for)

		D	$\bar{D}$	
ABO	F	20	10	30
	M	30	40	70
		50	50	100

$$I_{DR}^{\hat{}} = \frac{20(40)}{10(30)} = 2.67$$

if done w/ F cases + M controls - would expect following e's

$$e_D (F) = \frac{20}{50} = .4 \text{ (from F above)}$$

$$e_D (M) = \frac{30}{100} = .3 \text{ (from total)}$$

ie,  $\phi_{sex, ABO} = 0$

VII. A. 5. ∴ if did study w/ 50 F cases + 50 M controls → e

	D	$\bar{D}$	
F	20	15	35
$\bar{F}$	30	35	65
	50	50	100

14 × 50  
3 × 50 = 15

$$I_{DR} = \frac{20(35)}{15(30)}$$

$$1.56 < 2.67$$

∴ bias

in order to use M controls, it would have been necessary for either:

$$\phi_{\text{sep } ABO, \bar{D}} = 0$$

⇒ no conf. difficult to engineer data - assumes no selection

$$\text{or } \phi_{ABO, D, \bar{E}} = 0$$

which is impossible for chr. Ca. because of interaction

We may conclude that there was conf. (indirectly) but notice that it can't be computed because an SMR is undefinable

F	D	$\bar{D}$	
F	20	0	20
$\bar{F}$	30	0	30
	50	0	50

M	D	$\bar{D}$	
$\bar{E}$	0	15	15
$\bar{E}$	0	35	35
	0	50	50

$$SMR = \frac{20+0}{\frac{0(30)}{0} + \frac{15(0)}{35}} = ?$$

8.10

### VII. B. EM and I - the relationship

- 1. as noted before, EM depends on the meas. of A. chosen as well as the type of design selected (i.e.  $I_D$  vs  $I_C$ ) - See: p. 48 for illustration -

\* there are 4 ways of determining EM

all give different results

- 1.  $I_D D$ ,  $I_C D$ ,  $I_D R$ ,  $I_C R$
- 2. there exists an inherent relationship between EM + I such that computation of ~~EM~~ EM, using  $I_D D$  amounts to a test for interaction - the proof of this is given in #9, pp. 47-53 - eq,

$$I_D$$

	M	$\bar{M}$
E	.04	.02
$\bar{E}$	.02	.01

$(.04 - .02) \neq (.02 - .01) \Rightarrow$  interaction

the other way:

$$ERR$$

	M	$\bar{M}$
E	3	1
$\bar{E}$	1	0

$RIE = \frac{3}{2} = 1.5 \Rightarrow$  interaction

can use central table ( $I_D$ ) to look for interaction

$\therefore$  observing trends in  $I_D D$  is tantamount to identifying an interactive effect

NB. - this also implies that  $I_D$  (not  $I_C$ ) should be used in computing RIE - something we have not done thus far

VII. B. 3. These principles only apply, however, to causal factors - when analyzing the interaction of protective factors, the appropriate measure is PF (instead of I<sub>→</sub>D) = % of all pat. cases th. have been prev. because of exp.

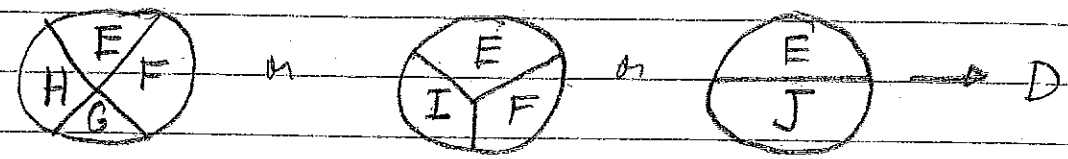
$$(PF_{11}) = (PF_{a,b,\dots} = 1 - (1 - PF_a)(1 - PF_b) \dots)$$

> ⇒ Synergism (+ interaction)  
 < ⇒ antagonism (- int.)

4. Similar relationship <sup>h.o.</sup> epistatic causal factors

$$EF_{a,b,\dots} = 1 - (1 - EF_a)(1 - EF_b) \dots$$

Suppose we had computed EF's for the following model.



I. EF = .20      II. EF = .30      III. EF = .40

? ∴ EF for (H) = .20 (ie, if indep. effect of H were 0)  
 $EF_{12} = 1 - (1 - .2)(1 - .3) = .44 = EF_F$   
 ∴ 44% of all cases could have been prevented if factor E had been absent in source pop. (ie, no one had been in hi risk category of F)  
 $EF_E = 1 - (1 - .2)(1 - .3)(1 - .4) = .664$

d F had no indep. effect

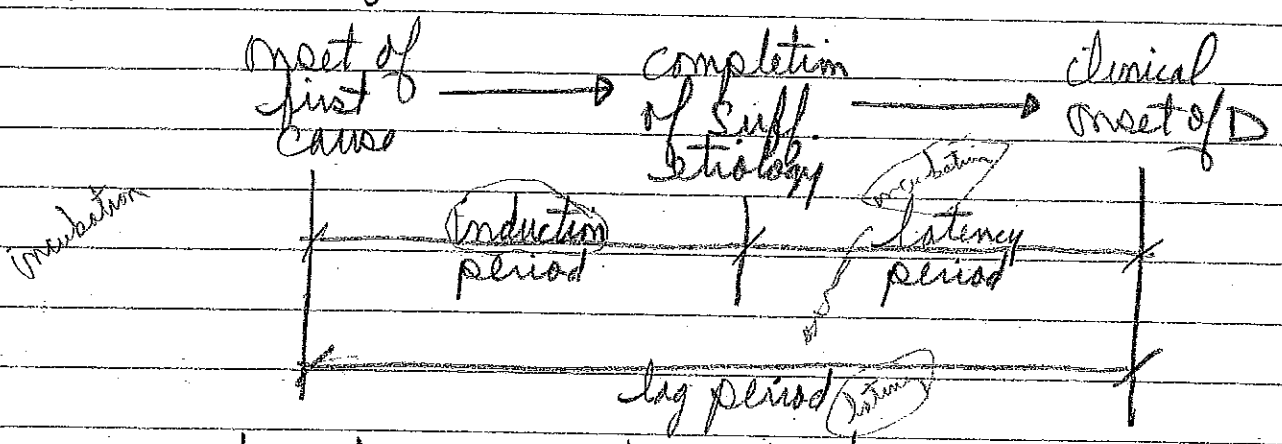
VII, B. 4. (cont.) thus all measured factors (A-J) did not account for  $1 - .664 = 33.6\%$  of all cases of disease D - *new model or factors?*

VIII, \* Biased estimate of statistical I:  $\frac{3}{2}$  factors (Selection, meas., conf.)  
Lag period - (M = E, lag period)

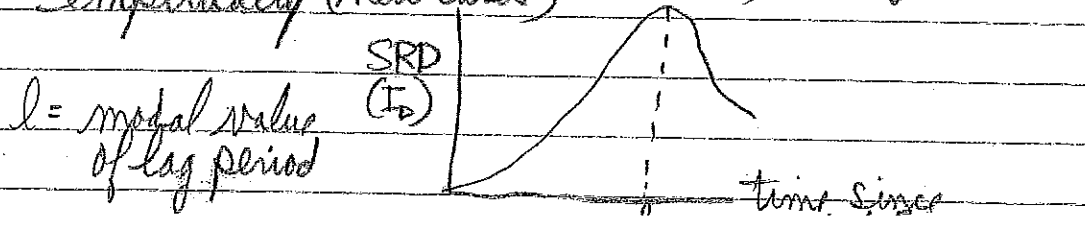
A. ~~it is~~ it is unlikely that any single factor is a "sufficient" cause of <sup>chronic</sup> VD, since it is bound to statistically interact w/ a no. of other factors

Stat. interaction need not imply any bio-beh. interaction

B. in studying diseases w/ multiple risk factors, the lag time between onset of first cause, & clinical onset of dis. may create complex analytical problems - consider the following conceptualization



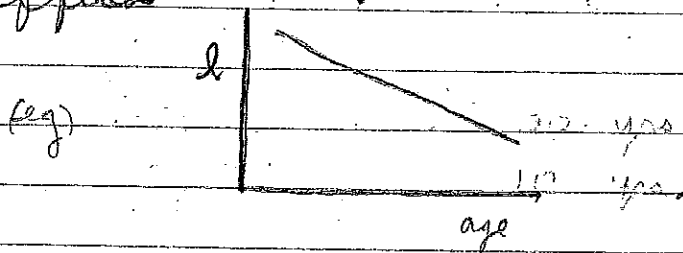
C. but from practical standpoint, it is only possible to estimate the "lag" period - if the E's begins at a certain age & remain stable thereafter - if this is possible we could estimate lag period empirically (new cases)



VIII

D. (with) the fact that the lag period is not constant but, in fact, has Variance  $> 0$ , ~~and~~ complicates the estimation of EM + I - ~~this~~: even without a true EM effect, there may appear to be a EM (Spurious) - i.e. the effect of E on D may be observed to vary along categories of M, even if M is not a risk indicator (i.e. not true E mod).

\* [ any var. which is assoc. with the lag time in study pop. (eg age) may appear to modify the effect of E or interact w/ it, thereby distorting the true nature of effects



modifying effect of age would be distorted (the not entirely spurious)

E. #11 -

2/4/76 Outline 3/11

I. Introd.

A. topic - principles of stat. analysis of categorical data - today; Simple (unstratified) analysis

B. analysis of etiologic research question involving assoc. of  $E \rightarrow D$  - 2 components

- II. (is there an assoc?) 1. hypothesis testing - qualitative  
 (how much?) 2. estimation of effect - quantitative (more important)
- a. pt. estimation - eg.  $\hat{RR}$
  - b. interval estimation - eg.  $CI[\hat{RR}]$

II. H. testing: nominal (categorical data) - simple analysis

(A) common procedure:  $\chi^2$  test of homogeneity (CR or CH indep. (X-S))

	D	$\bar{D}$	
E	a	b	$n_E$
$\bar{E}$	c	d	$n_{\bar{E}}$
	$m_D$	$m_{\bar{D}}$	n

in genl.  $\chi^2 = \sum_h \sum_i \frac{(O_{ih} - E_{ih})^2}{E_{ih}} = \frac{n(ad-bc)^2}{n_E n_{\bar{E}} m_D m_{\bar{D}}}$   
 (2x2 table)

(B)  $Z = \sqrt{\chi^2} = \chi$  ("chi-statistic") (more exact tables)

suggesting that a test of homog. is also testing the difference between proportions - same test

c. actually  $\chi^2$  or  $\chi$  can be derived from any measure of assoc. (A) that is N distrib.

A no variable  $\chi = Z = \frac{A - \bar{A}}{\sqrt{V[A]}} = \frac{A - E[A]}{SE[A]} = \frac{A}{E[A]}$

(("diff in prop.))  
 if  $A \equiv RD$  (I, D, eD, PD)  $\Rightarrow \chi_{RD}^2 = \frac{A}{SE[A]}$   
 $\rightarrow = \sum \sum \frac{(O-E)^2}{E}$  as above

$\chi^2 = \frac{A}{SE[A]}$



I. or C-R  
hypergeometric model

	D	$\bar{D}$	
$E_1$	a	b	$n_1$
$E_0$	c	d	$n_0$
	$m_1$	$m_0$	$N$

So we sample  $n$ , exposed from pop of  $N$ , = # of exp cases (a) is distrib. as hypergeometric distrib.

$$N = n$$

$$n = n_1, \quad p = \frac{m_1}{N}$$

$$D = m_1$$

$$E[a] = \frac{nD}{N} = np = \boxed{\frac{n_1 m_1}{N}}$$

$$V[a] = \frac{nD(N-D)(N-n)}{N^2(N-1)} = np(1-p) \frac{(N-n)}{(N-1)} = \boxed{\frac{n_1 m_1 m_0 n_0}{N^2(n-1)}}$$

II.  $\chi^2$  (cont.) but could use any A that is N distrib. tho the value of  $\chi$  might be slightly different  
 begin  $\rightarrow$  D. consider: RR  $\sim$  logN distrib  $\therefore \ln RR \sim N$

$$\therefore \chi = \frac{\ln RR - E[\ln RR]}{SE[\ln RR]} = \frac{\ln RR}{SE[\ln RR]}$$

could then estimate SE  $\xrightarrow{\text{eg.}}$   $SE[\ln RR] \approx \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$   
 (lge. sample approx.)

E, another A: we will use for remaining discussion on analysis -  $A \equiv \hat{M}H$

$$A \equiv \hat{M}H = \frac{n(n-1)(RD)}{m_D m_B} = \frac{n(n-1)(ad-bc)}{n_E n_{\bar{E}} m_D m_B}$$

$\rightarrow$  not usually computed directly because hard to interpret - NB. null value = 0

but we can use to get  $\chi \rightarrow$  M-H test

$$\chi_{MH} = \frac{\hat{M}H - E[\hat{M}H]}{SE[\hat{M}H]} \quad \hat{M}H \sim N$$

in  $I_e$  study:  $E[\hat{M}H] = 0$

$$(M+H) SE[\hat{M}H] = \sqrt{\frac{n^2(n-1)}{n_E n_{\bar{E}} m_D m_B}}$$

$$\therefore \chi_{MH} = \frac{n(n-1)(ad-bc)}{n_E n_{\bar{E}} m_D m_B} \sqrt{\frac{n^2(n-1)}{n_E n_{\bar{E}} m_D m_B}} =$$

consider distrib of "a"

$$= \frac{a - \frac{(m_D n_E)}{n}}{\sqrt{\frac{n_E n_{\bar{E}} m_D m_B}{n^2(n-1)}}}$$

$$= \frac{a - E[a]}{\sqrt{V[a]}}$$

(1) (in this context) using hypergeometric model:  $E[a] + \sqrt{V[a]}$   
 Sample w/ replacement

II E. (cont.) NB. not even necessary to compute  $\hat{MH}$

F.  $\chi_{MH}$  is same for C-R study =  $E[a] + V[a]$

	D	p-w.
E	a	$T_E$
$\bar{E}$	c	$T_{\bar{E}}$
	$m_D$	T

but somewhat diff. for  $T_D$  study: use binomial model because  $N$  &  $n$  is unknown

$$E[a] = \frac{m_D T_E}{(nD)} = np \quad V[a] = \frac{m_D T_E T_{\bar{E}}}{T^2} = np(1-p)$$

$n = m_D$   
 $p = T_E/T$

G. example:  $T_D$  study

Sample  $m_D$  cases from pop. w/p =  $T_E/T$ ; What is # cases, a?

	D	$\bar{D}$	
E	40	50	90
$\bar{E}$	30	80	110
	70	130	200

$$\chi_{MH} = \frac{a - E[a]}{\sqrt{V[a]}} = \frac{40 - \frac{70(90)}{200}}{\sqrt{\frac{90(110)(70)(130)}{200^2(199)}}} = 2.527 \quad (2p = .012)$$

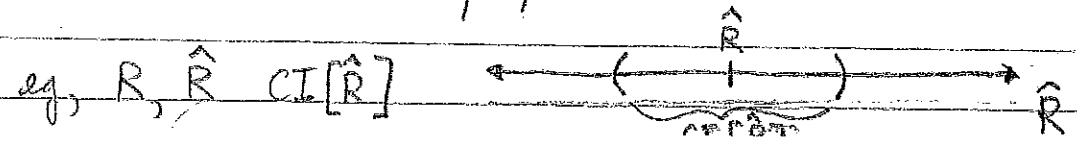
we will see next week that this technique is easily adapted to stratified analysis, but first CI

III. Estimation of Effect

A. Point Est. - covered previously:  $\hat{R}_R, \hat{R}_D$  w/o conf. (crude)

B. Interval Est. - CI

1. CI  $\equiv$  empirically derived range of values (or interval) expressing the  $P_r$  that the interval (constructed around  $\hat{p}$ ) includes the true pop. measure



III

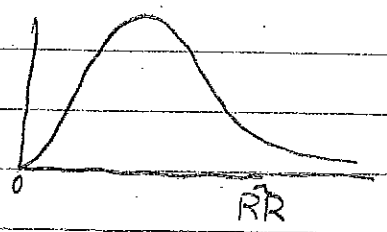
B. 2. we express CI as  $CI_\gamma[\hat{R}]$   
where:  $\gamma = \text{conf. level}$ ,  $\gamma$  represents designated precision of the CI

3. in genl., if  $\hat{R}$  is N distrib.

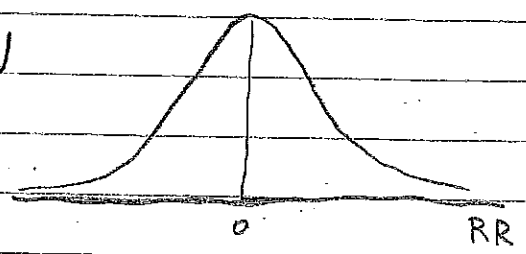
$$CI_\gamma[\hat{R}] = \hat{R} \pm CV_\gamma SE[\hat{R}]$$

where:  $CV_\gamma = \text{crit. value}$ , from known distrib.  $Z$  or  $t_{(n-1)}$  if  $\sigma^2$  is not known

C. CI around ratio, eg  $RR \sim \text{not N distrib}$   
but  $RR \sim \ln N$  distrib



so,  $\ln RR \sim N$



thus,

$$CI_\gamma[\ln \hat{RR}] = \ln \hat{RR} \pm Z_\gamma SE[\ln \hat{RR}]$$

$$CI_\gamma[\hat{RR}] = \exp[\ln \hat{RR} \pm Z_\gamma SE[\ln \hat{RR}]]$$

but how to estimate SE?

if  $RR = OR \Rightarrow SE[\ln OR] \cong \sqrt{\frac{1}{a} + \frac{1}{b} + \frac{1}{c} + \frac{1}{d}}$  large sample

could also use approx. from Taylor Series (genl)

$$SE[\ln \hat{RR}] \cong \frac{SE[\hat{RR}]}{\hat{RR}} \rightarrow \text{but what is this}$$

CI[ $\hat{EF}$ ] - proof:

if:  $(1 - \hat{EF}) \sim \ln N$

$\ln(1 - \hat{EF}) \sim N$

$\therefore CI[\ln(1 - \hat{EF})] = \ln(1 - \hat{EF}) \pm Z_{\alpha} SE[\ln(1 - \hat{EF})]$

$SE[\ln(1 - \hat{EF})] = \ln(1 - \hat{EF}) / \chi$

$\therefore CI[1 - \hat{EF}] = \exp \left[ \ln(1 - \hat{EF}) \pm Z_{\alpha} \frac{\ln(1 - \hat{EF})}{\chi} \right]$

$= \exp \left[ \ln(1 - \hat{EF}) \left( 1 \pm \frac{Z_{\alpha}}{\chi} \right) \right]$

$= (1 - \hat{EF})^{(1 \pm Z_{\alpha}/\chi)}$

$CI[\hat{EF}] = 1 - (1 - \hat{EF})^{(1 \pm Z_{\alpha}/\chi)}$

but is this good for the study?

if:  $\hat{EF} \sim N$  (or  $(1 - \hat{EF}) \sim N$ )

$\therefore CI[\hat{EF}] = \hat{EF} (1 \pm Z_{\alpha}/\chi)$

CI

CI[ $\hat{I}_a$ ] - CI[ $\hat{I}_c$ ] - #12 on ref. list

III. D. new method for computing CI - method of "test-based intervals" - avoids direct computation of SE

recall:  $\chi = \frac{A - E[A]}{\sqrt{V[A]}} = \frac{\ln \hat{RR}}{SE[\ln \hat{RR}]}$

$SE[\ln \hat{RR}] = \frac{\ln \hat{RR}}{\chi}$  as computed before =  $\frac{a - E[A]}{\sqrt{V[A]}}$   
an exact expression

Subst. into formula for CI:

$CI_{\gamma}[\hat{RR}] = \exp \left[ \ln \hat{RR} \pm z_{\gamma} \cdot \frac{\ln \hat{RR}}{\chi} \right] = \exp \left[ \ln \hat{RR} \left( 1 \pm \frac{z_{\gamma}}{\chi} \right) \right] =$

antilog of log =  $\hat{RR} \left( 1 \pm \frac{z_{\gamma}}{\chi} \right)$  where:  $\chi$

E. can also use this method w/ RD ~ N distrib.

$\chi_{RD} = \frac{\hat{RD}}{SE[\hat{RD}]}$

$SE[\hat{RD}] = \frac{\hat{RD}}{\chi}$

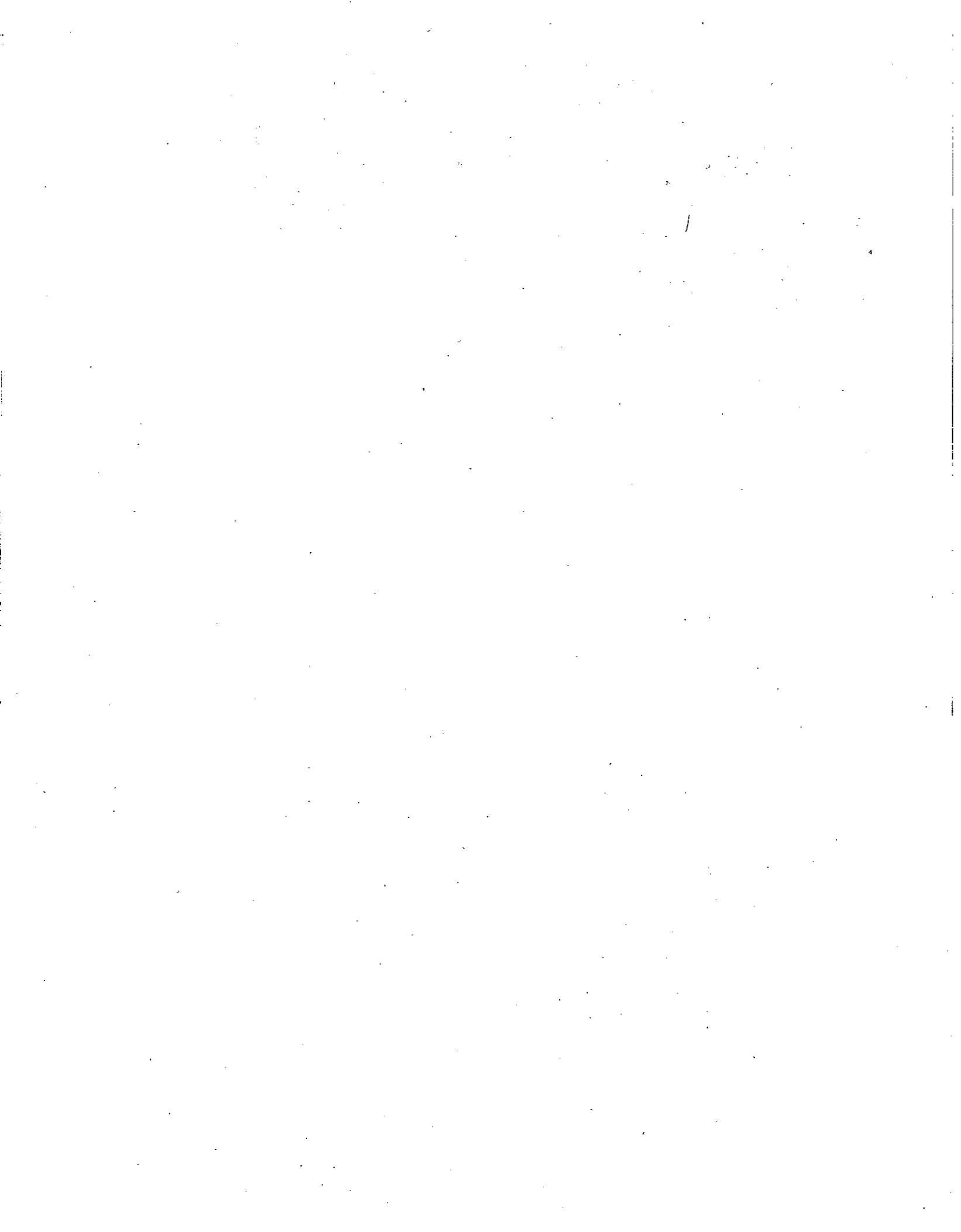
$\therefore CI_{\gamma}[\hat{RD}] = \hat{RD} \pm z_{\gamma} \frac{\hat{RD}}{\chi}$

=  $\hat{RD} \left( 1 \pm \frac{z_{\gamma}}{\chi} \right)$

F. + for EF (MPE)

$(1 - \hat{EF}) \sim \ln N \Rightarrow CI_{\gamma}[\hat{EF}] = 1 - \hat{EF} \left( 1 \pm \frac{z_{\gamma}}{\chi} \right)$  (w/ EF ≥ 0)  
note CI, must ≥ 0 + PE < 1

proof ↻



new

III

F'

$i_j$  (exp. or unexp.)

$$CI_j [\hat{I}_{d_j}] = \exp \left[ \ln \hat{I}_{d_j} \pm z_j SE[\ln \hat{I}_{d_j}] \right]$$

$$= \hat{I}_{d_j} e^{\pm z_j SE[\ln \hat{I}_{d_j}]}$$

could remove  $i_j$ 's to get  $CI[\hat{I}_{d_j}] \Rightarrow$  only need

$e^{\ln \hat{I}_{d_j}}$   
 $e^{\pm z_j SE}$

where:  $SE^2[\ln \hat{I}_{d_j}] = V[\ln \hat{I}_{d_j}] =$

$$V[\ln \hat{I}_{d_j}] = \frac{(\ln \hat{I}_{d_j} - \ln \hat{I}_d)^2}{\chi^2}$$

where:  $V[\ln \hat{I}_d] = \frac{1}{a + c}$  inverse of total # cases in the estimate

F''

$CI_j [\hat{I}_{c_j}] =$  similar to above except

duration

$$V[\ln \hat{I}_{c_j}] = \frac{1}{\hat{I}_{c_j}^2} \sum \left[ \frac{I_{d_j}^2}{(a_j + c_j)} \right] (p_j^2)$$

NB.  $I_{d_j}$  is for both  $E + \bar{E}$

When  $I_{c_j}$  is computed from  $I_{d_j}$  rate

$$\hat{I}_{c_j} = 1 - e^{-\sum I_{d_j} p_j}$$

but what about in  $T_c$  study?  
I guess, you can consider  $\hat{I}_{c_j} \sim N$ ?



III

G. NB: "I" indicates calculation of CI<sub>L</sub> + CI<sub>u</sub>

implicitly suggests:

$$2\text{-sided } \begin{cases} H_0: \hat{RR} = 1 & (\text{no effect}) \\ H_1: \hat{RR} \neq 1 \end{cases}$$

yet in epid, there may be no biologic or soc. basis (plausibility) for hypothesizing E as protective (or causal) ∴ can increase power of test by employing 1-sided test

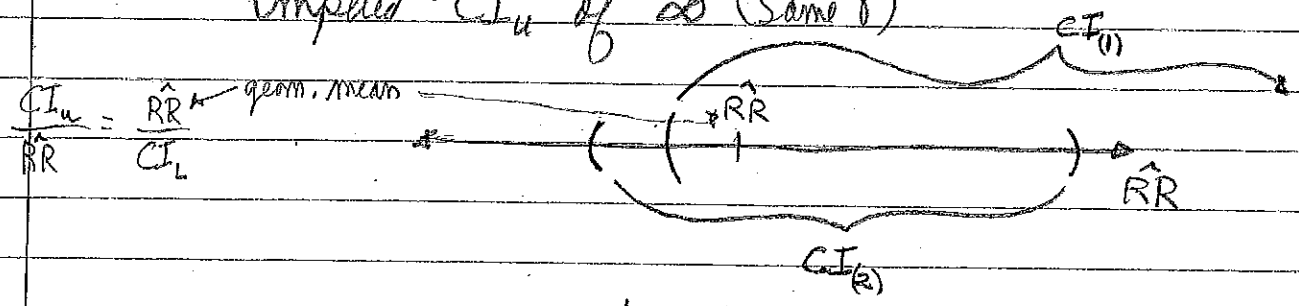
$$H_0: \hat{RR} = 1$$

$$H_1: \hat{RR} > 1 \text{ or } \hat{RR} < 1$$

∴ p-value will be cut in half (given same  $\alpha$ )

$$Z_{1-\alpha/2} \rightarrow Z_{1-\alpha}$$

OM recommends, that even a CI reflect the one-sidedness of most epid H's, by computing a tighter 1-sided CI<sub>L</sub> w/ an implied CI<sub>u</sub> of  $\infty$  (same  $\alpha$ )



∴ increasing ability of rejecting the H<sub>0</sub>

$$CI_{\alpha}[\hat{RR}] = [CI_L, \infty] = \left[ RR \left( 1 - \frac{Z_{\alpha}}{X} \right) \right]$$

H. Example - CR study

	D	$\bar{D}$	
E	40	23	63
$\bar{E}$	60	77	137
	100	100	200

III.

H. let:  $\delta = .95$

$\therefore \alpha = .05$  (1-sided test)

$Z_{.95} = 1.645$

$$\hat{I}_{OR} = OR = \frac{ad}{bc} = \frac{40(77)}{23(60)} = 2.232$$

$$\chi_{MH} = \frac{40 - \frac{63(100)}{200}}{\sqrt{\frac{63(137)(100)(100)}{200^2(199)}}} = 2.581$$

(use Z table  
p = .005  
(1-sided))

$$CI_L [RR] = 2.232 \left(1 - \frac{1.645}{2.581}\right) = 1.338$$

( $CI_u = \infty$ )

if we had used traditional method for calculating CI by est. SE,  $CI_L$  would have been somewhat less (closer to 1)

I. advantages of test-based method - p. 48

J. Summary of  $\chi$  calculations for different study designs - pp. 28-32, #10

next: Stratified analysis

Addendum

Hypergeometric

	D	$\bar{D}$	
$E_1$	a	b	$n_1$
$E_0$	c	d	$n_0$
	$m_1$	$m_0$	$n$

$a \sim \text{Hyp}$

$N = n$   
 $n_1 = n$   
 $D = m_1 \left( \frac{m_1}{n} = p \right)$

$E[a] = \frac{nD}{N} = \frac{n_1 m_1}{n}$

correcting factor

$np(1-p) \frac{(N+n)}{(N-1)} = V[a] = \frac{nD(N-D)(N-n)}{N^2(N-1)}$

$V[a] = \frac{n_1 m_1 m_0 n_0}{n^2(n-1)}$

$m_1 \sim \text{Bin}$  or  $a_1 \sim \text{Bin}$

$p = m_1/n$

$m_1 \sim \text{Nor}$  if  $n$  is large

$\frac{m_1 - np}{\sqrt{np(1-p)}} = \frac{\frac{m_1}{n} - p}{\sqrt{\frac{p(1-p)}{n}}} = \frac{\hat{p} - p}{\sqrt{V[\hat{p}]}}$

Can  $p$  be  $I_d = \frac{m_1}{T}$ ?

Binomial

	D	$p-n$
$E_1$	a	$T_1$
$E_0$	c	$T_0$
	$m_1$	$T$

$n = m_1$   
 $p = T_1/T$   
 $a \sim \text{Bin}$

$E[a] = np = \frac{m_1 T_1}{T}$

$V[a] = \bar{p} p (1-p) =$

$\frac{m_1 T_1 T_0}{T^2}$

$\frac{a}{T} = \frac{m_1}{T} \quad I_{d_1} - I_{d_0}$

$\sqrt{\frac{m_1 T_0}{T_1 T^2}}$

$$\hat{I}_d = \frac{10 \text{ cases}}{1000 \text{ p-yr.}} = .01$$

$$CI[\hat{I}_d] = (.01) \pm 1.96 \cdot SE[\ln(.01)]$$

$$SE[\ln(.01)] = \sqrt{V[\ln(.01)]} = \sqrt{\frac{1}{10} + \dots}$$

*typer*

$$CI = .01 e^{\pm \frac{1.96}{\sqrt{10}}} = [.005 - .019]$$

$$CI = .01 \pm 1.96 \sqrt{\frac{.01(.99)}{1000}} = [.004 - .016]$$

	D	p-yr.
E <sub>1</sub>	15	500
E <sub>0</sub>	10	500
	25	1000

.03 - .02

$$\hat{I}_R = 1.5$$

$$CI[\hat{I}_R] = 1.5 \left(1 \pm \frac{1.96}{\sqrt{7}}\right) = 1.5 \pm 1.96$$

$$\lambda = \frac{15 - \frac{25(500)}{1000}}{\sqrt{\frac{25(500)(500)}{1000^2}}} = 1$$

[.68] - 2.32

$$\hat{I}_{DD} = .01$$

$$CI[\hat{I}_{DD}] = .01(1 \pm 1.96) = [.0096 - .0296]$$

$$CI[\hat{I}_{DD}] = .01 \pm 1.96 \sqrt{\frac{.03(.97)}{500} + \frac{.02(.98)}{500}} = [.0093 - .0293]$$

2/11/76 Outline + 2/18

I. Introd.

A. extend principles of simple analysis to stratified analysis

II. Hypothesis testing - MH test

A. consider g strata of C (pot. conf. factors) - j=1...g

$$\chi_{MH_j} = \frac{a_j - E[a_j]}{\sqrt{V[a_j]}} = \frac{\hat{M}H_j}{SE[\hat{M}H_j]} \quad \begin{matrix} \hat{M}H_j \sim N \\ \epsilon = 0 \end{matrix}$$

$$\sum_{j=1}^g \chi_{j^2} = \sum \frac{\hat{M}H_j^2}{SE^2[\hat{M}H_j]} = \text{"total chi"} \quad (\text{or } \sum \chi_{j(g)}^2)$$

H<sub>0</sub>: no assoc. in any strata i.e. no good because 1 strata may distort results

B. what is needed is an overall test statistic:

X<sub>0</sub>, which is a function of an overall meas of assoc. where A<sub>0</sub> combines strata to form a pooled est.

1. 2 ways of getting A<sub>0</sub>

- a. ML est. esp. when #'s are small (pencil)
- b. weighted average of A<sub>j</sub> - requires large #'s in each cell

2. (1b.) →  $A_0 = \frac{\sum w_j A_j}{\sum w_j}$  where:  $w_j = \frac{1}{V[A_j]} = \frac{1}{SE^2[A_j]}$

ie, wts. maximize efficiency of the estimate (unlik. SH)

$$\therefore SE[A_j] = \frac{1}{\sqrt{w_j}}$$

since A<sub>0</sub> is a wted. sum of A<sub>j</sub>,  $SE[A_0] = \frac{1}{\sqrt{\sum w_j}}$

II, B. 3. using  $A_j = \hat{M}H_j$ , w/o proof (in paper)

$$\chi_{0(1)}^2 = \frac{a - E[a]}{\sqrt{V[a]}} = \frac{\sum a_j - \sum E[a_j]}{\sqrt{\sum V[a_j]}}$$

[NB: #'s in cells need not be large]

∴ same form as in simple analysis

where:  $\sum E[a_j] = \sum \frac{m_{Dj} n_{Ej}}{n_j}$  in  $I_c$  or CR

$$= \sum \frac{m_{Dj} T_{Ej}}{T_j}$$

in  $I_D$

$$\sqrt{\sum V[a_j]} = \sqrt{\sum \frac{n_{Ej} n_{\bar{E}j} m_{Dj} m_{\bar{D}j}}{n_j^2 (n_j - 1)}} \quad \text{in } I_c \text{ or CR (hypergeom. V)}$$

$$= \sqrt{\sum \frac{m_{Dj} T_{Ej} T_{\bar{E}j}}{T_j^2}} \quad I_D \text{ (bin V.)}$$

### III. Test for uniformity

A.  $\chi_0^2$  (or  $\chi_0^2$ ) is ~~not~~ applicable to any study design, where C is conf. —  
 but, if C is also an effect mod.,  $A_0$  (or  $MH_0$ ) is a deceptive measure because different estimates are pooled —  
 so, you don't get  $A_0$  unless there is no sign. EM — ie, all strata-sp. est. need not be ~~exactly~~ exactly equal, but the data must be consistent with an assumption of uniformity — (ie, not sign. diff. from uniformity)

III. B. the test:

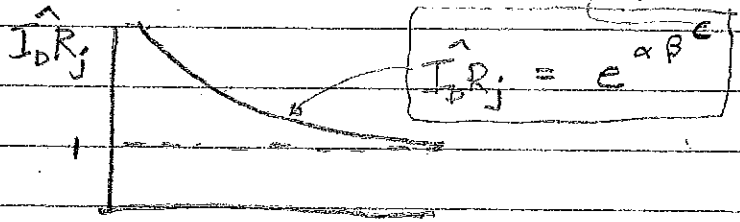
$$\chi^2_{u(q-1)} = \sum \chi^2_{j(q)} - \chi^2_{0(q)}$$

computed as described

if p is small  $\Rightarrow$  reject  $H_0$  of uniformity; don't pool strata to get  $A_0$  - tho can still use  $\chi^2_0$  - can get SMR, SRR  
if p is large  $\Rightarrow$  no EM  $\therefore$  use  $A_0$  as pt. est.

\* if  $A_0 = \hat{M}H_0$  (ie,  $\chi_j = \chi_{MH}$ ), then  $\chi^2_u$  is only good in a C-R study - I guess because  $\hat{M}H$  isn't a stable est. in a CH study (like OR)

C. recall that  $\chi^2_u$  isn't a very sensitive test because it only requires nominal (unordered) data we might fit the data to a model instead - eg, (OM et al, 1976) - cig. sm.  $\rightarrow$  acute MI (non-fatal)



C. (Composite Score)

- testing:
- ①  $\alpha > 0 \Rightarrow$  assoc between cig sm + AMI
  - ②  $\beta < 1 \Rightarrow$   $\hat{R}R_j$  approaches unity w/ ↑ C (ref. risk)
  - ③ fit of model

IV. Pt. est. of overall effect ( $A_0$ ) - don't usually get  $\hat{M}H_0$  because it has little meaning (not RD, not RR) must have large #'s in each strata, or use ML est.

who Kleinbaum computes it in 287

IV. A.  $I_c$  study -

$$\hat{RD}_0 = A_0 = \hat{I}_c \hat{D}_0 = \frac{\sum w_j \hat{I}_c \hat{D}_j}{\sum w_j} \quad \hat{RD}_j \sim N$$

$$w_j = \frac{n_{Ej} n_{Ej} n_j}{m_{Dj} m_{Ej}} = \frac{1}{V}$$

$$\hat{RR}_0 = A_0 = \hat{I}_c \hat{R}_0 \quad \hat{RR}_j \sim \ln N$$

$$\therefore \ln \hat{I}_c \hat{R}_0 = \frac{\sum w_j (\ln \hat{I}_c \hat{R}_j)}{\sum w_j}$$

$$\hat{I}_c \hat{R}_0 = \exp \left[ \frac{\sum w_j (\ln \hat{I}_c \hat{R}_j)}{\sum w_j} \right]$$

$$w_j = \frac{m_{Dj} n_{Ej} n_{Ej}}{m_{Ej} n_j} = \frac{1}{V}$$

B.  $I_D$  study

$$\hat{RD}_0 = A_0 = \hat{I}_D \hat{D}_0 = \frac{\sum w_j \hat{I}_D \hat{D}_j}{\sum w_j}$$

$$w_j = \frac{T_{Ej} T_{Ej}}{m_{Dj}}$$

$$\hat{RR}_0 = A_0 = \hat{I}_D \hat{R}_0 = \exp \left[ \frac{\sum w_j (\ln \hat{I}_D \hat{R}_j)}{\sum w_j} \right]$$

$$w_j = \frac{m_D T_{Ej} T_{Ej}}{T_j^2}$$



### IV. C. C-R Study

$\hat{RR}_0$  — could use same technique as above,  
 w/  $OR_j \equiv I_{D_j} R_j$   
 but there is a simpler way:  
M-H Procedure for est. RR in a C-R study

$$\hat{RR}_0 = I_{D_j} \hat{R}_0 = \frac{\sum \frac{a_j d_j}{n_j}}{\sum \frac{b_j c_j}{n_j}} = \frac{\sum w_j \cdot OR_j}{\sum w_j} \quad w_j = \frac{b_j c_j}{n_j}$$

can still use  $\chi_0$  ( $\chi_{MH}$ )

D. Table 1, #10, pp. 28-32 — Summary of  $\chi_{MH} + A_0$  for different study designs + analysis, including Matched analysis, to be covered next week

begin  
 ↓  
 For  $\chi_{MH}$  formula  
 p. 3  
 $\hat{A}_0 = \frac{\sum w_j \hat{A}_j}{\sum w_j}$   
 $\hat{RR}_0 =$  C-R Study

### V. Interval Estimation — CI

A. for overall meas. of assoc. —  $A_0$   
 use test-based intervals exactly as in Simple anal.

eg,  $CI_{95}[\hat{RR}_0] = \hat{RR}_0 (1 \pm \frac{z_{\alpha/2}}{x})$

where:  $\hat{RR}_0$  is one of the previous weighted formulas or a ML est.

CI [EP] — p. 5 Simple analysis

this method (CI) is only good when the A has maximized effic. — which occurs when there is uniformity of effect over strata (of C)

Similar  
to SMR

CH  $CI_y[SRD] = \hat{SRD} \pm Z_\alpha \sqrt{\hat{V}[SRD]}$

where:  $\hat{V}[SRD] \cong \frac{\sum_j n_{ij} I_{c_{ij}} (1 - I_{c_{ij}})}{\left(\sum_j n_{ij} I_{c_{oj}}\right)^2} + \frac{\sum_j n_{oj} I_{c_{oj}} (1 - I_{c_{oj}})}{\left(\sum_j n_{oj} I_{c_{oj}}\right)^2}$

external  
CI only

$$CI_y [SRD] = SRD \pm Z_\alpha \sqrt{\hat{V} [SRD]} \quad \left( \frac{S_j}{S} = \frac{n_{0j}}{n_j} \right) \quad (6)$$

$$\hat{V} [SRD] = \sum_j \left( \frac{n_{0j}}{n_j} \right)^2 \left[ \frac{I_{01j} (1 - I_{01j})}{n_{1j}} + \frac{I_{02j} (1 - I_{02j})}{n_{2j}} \right]$$

V

B. for standardized measures - SMR, SRR, which haven't max. effic. - must estimate SE

assume:  $I_{0j}$  are  
Bin & indep. for  
each category

the SMR method

eg,  $CI_y [SMR] = SMR \pm Z_\alpha \cdot SE [SMR]$

$$\hat{V} [SMR] \approx \frac{\sum n_{Ej} \hat{R}_{Ej} (1 - \hat{R}_{Ej})}{\left( \sum n_{Ej} \hat{R}_{Ej} \right)^2}$$

? SRR  
1 by dropping

but can't compute SE (or  $\hat{R}$ ) in C-R stud  
so I don't know if SE has been estimated  
for C-R study

VI. Example - C-R study

	C <sub>1</sub>			C <sub>2</sub>		
	D	$\bar{D}$		D	$\bar{D}$	
E	35	55	90	25	35	60
$\bar{E}$	20	90	110	10	30	40
	55	145	200	35	65	100

A. Is there conf.?

$$CRR = \frac{60(120)}{90(30)} = 2.667$$

$$SMR = \frac{60}{\frac{55(20)}{90} + \frac{35(10)}{30}} = 2.512$$

RR\* = 1.062      RCE = .041      (Some pos. conf.)

B. EM?       $OR_1 = 2.86$  } is this diff  
                  $OR_2 = 2.14$  } stat. sign?

(1) II c. But first, is there an overall assoc. -

$$\chi_{(0)}^2 = \frac{a - E[a]}{\sqrt{V[a]}} = \frac{(35+25) - \left(\frac{55(90)}{200} + \frac{35(60)}{100}\right)}{\sqrt{\frac{90(110)(55)(145)}{200^2(199)} + \frac{60(40)(35)(65)}{100^2(99)}}$$

$$= 3.627 \quad p = \underline{.00014} \text{ (1-sided)}$$

$$\chi_0^2 = 13.157$$

D. check to see whether EM occurred by chance  
(to determine ~~what~~ the type of pt. est. to use)

$$\chi_{u(q-1)}^2 = \sum \chi_{j(q)}^2 - \chi_{0(c)}^2$$

$$\sum \chi_j^2 = \chi_1^2 + \chi_2^2 = \frac{\left(35 - \frac{55(90)}{200}\right)^2}{\frac{90(110)(55)(145)}{200^2(199)}} + \frac{\left(25 - \frac{35(60)}{100}\right)^2}{\frac{60(40)(35)(65)}{100^2(99)}}$$

$$= 13.494$$

$$\therefore \chi_u^2 = 13.494 - 13.157 = .337$$

$$\chi_u = \sqrt{.337} = .587 \quad p = \underline{.561} \text{ (2-sided)}$$

\(\therefore\) no EM

E.  $\Rightarrow$  can get  $\hat{RR}_0$  using MH procedure

$$\hat{RR}_0 = \frac{\sum \frac{a_j d_j}{n_j}}{\sum \frac{b_j c_j}{n_j}} = \frac{\frac{35(90)}{200} + \frac{25(30)}{100}}{\frac{55(20)}{200} + \frac{35(10)}{100}} = \underline{2.583}$$

VI, F. CI<sub>y</sub> - γ = .95

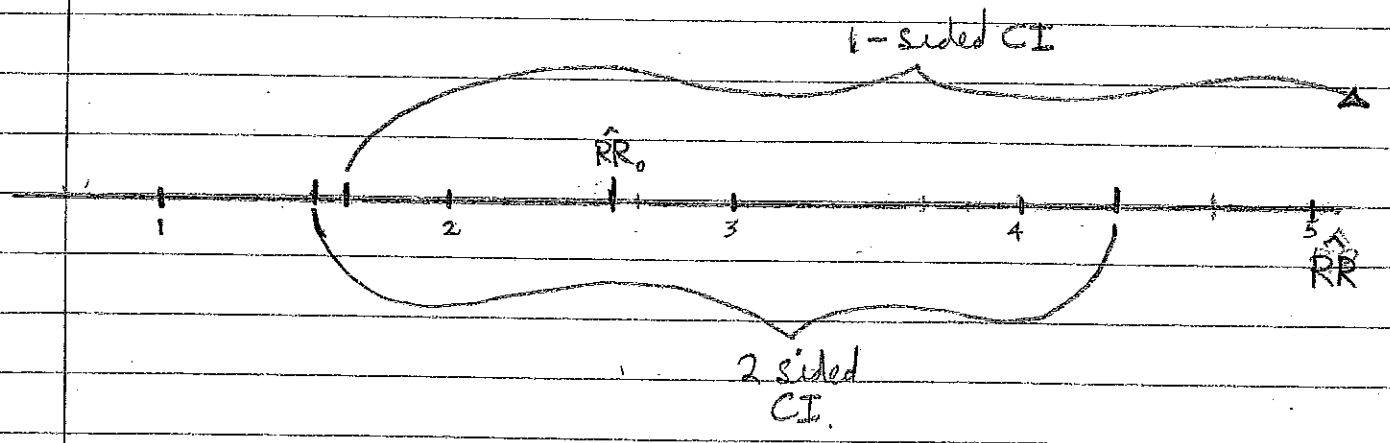
α = .05 (1-sided) Z<sub>γ</sub> = 1.645

95% CI<sub>L</sub> = RR̂ (1 - Z<sub>γ</sub>/k) = 2.583 (1 - 1.645/3.627) =

= 1.68 ∴ 95% confid. that "true RR" is > 1.68

2-sided 95% CI = 2.583 (1 ± 1.96/3.627)

= [1.55 ; 4.31]



# 2/18/76 Outline

## I. Intro. → cl. 8 (OM) -

A.  $CI_0[\hat{E}F]$   
 $CI_0[\hat{I}_0] + CI[\hat{I}_c]$  } other notes

B. today:

1. finish stratified analysis
2. Matching + M. analysis (#12)

## II. Defn's

A. design: predictor + outcome (response) var's

B. predictor var → test + referent groups

C. M (in selection of subjects in obs. study) ≡ restricting the selection of referent series by making it similar to the test group (selected indep.) w/ re: to C (pat. conf.)

M - one way to control for conf. factors in selection of subjects

1. indiv. M: each test subject is grouped with one or more referent subjects according to their value in C. → making ref. series comparable to test series w/ re: to C - each set of 1 test + 1 or more ref.'s ≡ stratum

a. fixed-ratio (r) M - all sets of M. subjects contain same ratio of ref.: test subj.

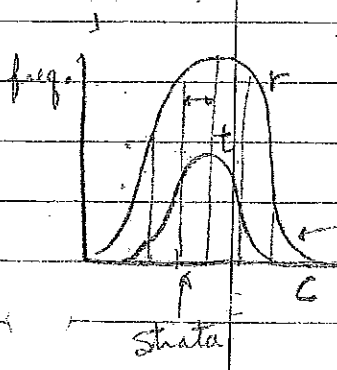
i. pairwise M - r=1

ii. non-pairwise - r > 1

b. variable-ratio M - r is not uniform across strata

2. freq. M - forces comparability of distrib's of test + ref. series w/ re: to C -

OM: freq M is like other side of continuum from pairwise M - stratum, defined over C, as in non-M  
 ∴ must have lg. #'s w/in each category of C - analyzed like gen. stratified



D. procedure: all 'M' factors are divided into a few (2-4) classes (k<sub>i</sub>) combining factors produces M. categories = Π k<sub>i</sub>  
 ex. sex + race → 4 M. cates.

II. E. freq. M has less power than Ind. M. because  
~~the~~ freq. M. has fewer df (= # strata - 1)  
 where strata are larger - than Ind. M (df = # strata - 1)  
 but strata are smaller (as little as 2)  
 ∴ in planning a study there is little reason to use freq. M unless Ind. M. is too cumbersome  
 say because of too many M categories (> 12) -  
 on the other hand, we will see that such extensive M. may not be called for  
 - ∴ only consider Ind. M.

III. Analysis

A. Cohort - pairwise

- n subjects - n/2 pairs - each E is Mcd w/ 1  $\bar{E}$   
 ∴ n/2 strata of 4 types

	D	$\bar{D}$		D	$\bar{D}$		D	$\bar{D}$		D	$\bar{D}$
E	1	0	1	1	0	1	0	1	1	0	1
$\bar{E}$	1	0	1	0	1	1	1	0	1	0	1
	2		0	2			1	1	2	0	2
	redundant						redundant				

Since all pairs of similar type give same info. (Validity + precision) eg w/re:  $\hat{RR}$  - they can be combined, which amounts to stratification in analysis -  
 Summarized by 2x2 table

$$\hat{ML} \hat{RR} = \frac{P+Q}{P+R}$$

		D	$\bar{D}$	
E	P	P	Q	"
$\bar{E}$	$\bar{D}$	R	S	
		n/2 (pairs)		

- if stratification had been ignored

$$\hat{RR} = \frac{\frac{P+Q}{n/2}}{\frac{P+R}{n/2}} = \frac{P+Q}{P+R} = \hat{ML} \text{ w/stratification}$$

	D	$\bar{D}$	
E	P+Q	R+S	n/2
$\bar{E}$	P+R	Q+R	n/2
	n (subj.)		

so failure to stratify in M.c.H. study doesn't result in bias

III. A. 3.  $H_0$  test -  $\chi_{MH}$

$$\chi_{MH} = \frac{a - E[a]}{\sqrt{V[a]}} = \chi_0$$

where:  $a = \sum a_j = P+Q$

$$E[a] = P + \frac{Q+R}{2}$$

$$V[a_j] = \frac{n_{Ej} n_{\bar{E}j} m_{Dj} m_{\bar{D}j}}{n_j^2 (n_j - 1)} \quad \text{Hypergeom. V}$$

$$= \frac{0}{2^2(1)} = 0 \quad \text{for concordant pairs}$$

$$= \frac{1(1)(1)(1)}{2^2(1)} = \frac{1}{4} \quad \text{for disc. pairs}$$

$$\therefore V[a] = \sum V[a_j] = \frac{Q+R}{4}$$

$$\chi_{MH} = \chi_0 = \frac{(P+Q) - \left(P + \frac{Q+R}{2}\right)}{\sqrt{\frac{Q+R}{4}}} = \frac{Q+R}{\sqrt{Q+R}} \quad \text{"McNemar test"}$$

can use this  $\chi$  to compute CI [ $\hat{RR}$ ] as described

4. effect of fixed  $K-M$  - no assoc. between  $E+C$ ,  
but can still study assoc. between  $C+D$  (can risk ind.)  
after having H.

$$\phi_{CD} = \frac{PS - QR}{\sqrt{(P+Q)(R+S)(P+R)(Q+S)}} \quad 0 \leq \phi \leq 1$$



III

B. C-R - pairwise

1. n Subjects - n/2 pairs or strata - each D is  
Match w/ 1  $\bar{D}$  -

Stratified analysis is represented as:

		$\bar{D}$	
		E	$\bar{E}$
$\hat{I}D = \frac{X}{Y}$	D	E	E
		W	X
		$\bar{E}$	Z
			n/2 (pairs)

2. Suppose: no stratif.

		D	$\bar{D}$
	E	W+X	W+Y
	$\bar{E}$	Y+Z	X+Z
		n/2	n/2
			n (subj)

$$\hat{I}D = OR = \frac{(W+X)(X+Z)}{(W+Y)(Y+Z)} \neq \frac{X}{Y}$$

$\therefore$  failure to stratify in analysis can bias RR.  
often done

3.  $H_0$  test - similar derivation of  $\chi_{MH}$

$$\chi_{MH} = \frac{(X-Y)^2}{\sqrt{X+Y}} \quad \text{"McNemar test"}$$

can use to compute test-based interval (CI)

4. effect fixed: R. M. - no assoc. between C+D so  
cant study C as risk indic of D after matching  
but can get assoc. between C+E

$$\phi_{EC} = \frac{WZ - XY}{\sqrt{(W+X)(Y+Z)(W+Y)(X+Z)}}$$

MH est. of  $\hat{RR}$  in fixed-r. c-c study

		$\bar{D}$ # exposed (g)		
		2	1	0
$D$	# exp. (f)	1	0	
	0		$Z_{fg}$	

$Z$  (sets of  $r$  controls + 1 case)

$$f = 0, 1$$

$$g = 0, 1, \dots, r-1$$

$$(\text{not ML}) \quad I_{DR} = \sum_f \sum_g Z_{fg} (r-g) f$$

$$\hat{I}_{DR} = \frac{\sum_f \sum_g Z_{fg} (r-g) f}{\sum_f \sum_g Z_{fg} (1-f) g}$$

Last time:

defined M. - restricting selection of ref. Series  
analysis of pairwise M. in CH + C-R - RR,  $\chi_0$ , CI  
stratify in analysis - CH  
C-R

III

### C. Fixed-ratio M - ( $r > 1$ )

$$\chi_{MH} = \frac{a - E[a]}{\sqrt{V[a]}}$$

1. can still extend basic MH test to  $r > 1$   
for either C-R or CH design - the calculations are a little tricky (see: my paper)

2. NB: stratif. analysis involves  $> 4$  cells  
(eg) M-CR study  $r=2$

		2E	E/E	2E
D	E	U	V	W
	E	X	Y	Z

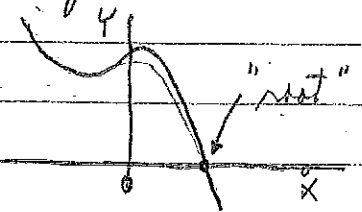
how many exp.?  
 $a+b = a+V+W + 2(U+X) + Y = 3U+2V+2X+U+Y$

so each stratum has 3 subjects (strata)  $1D + 2D$   
3. RR - not simple algebraic soln. but requires iterative soln. to get  $\hat{M}$   $r > 2$

Suppose we set up an equation to solve for an unknown ( $X$  in RR), in which there is no alg. soln. - eg,

Mantel-Haenszel do have simple algebraic soln. but not  $\hat{M}$ .  
 $\hat{RR} = \frac{2W+V}{2X+Y}$

$$X^3 + X - 1 = 0$$



good approach: get successive approx's that narrow in on the exact root - how?

w/ recursive function of form:  $X_{n+1} = f(X_n)$   
(keep replacing  $X_n$  w/  $X_{n+1}$ )

$$X(X+1) = 1$$

$$X = \frac{1}{X^2+1} \Rightarrow X_{n+1} = \frac{1}{X_n^2+1}$$

eg, start  $X_0 = 0 \therefore X_1 = 1$   
 $X_1 = 1 \therefore X_2 = .5$   
 $X_2 = .5 \therefore X_3 = .8 \dots$

III

C. 3. in the case of RR, OM, has shown how to get MI for Med C-R study - in order to derive the proper recursive function, need diff. calculus - use: N-R technique

D. Var-r M. - still use basic  $\chi^2_{MH}$  also can use iterative soln. to get RR, treating each type of strata separately -  $r = 1, 2, \dots$

when needed:  
eg ① M. Subs.  
② limited # of ref. Subjects in certain M. categories

still use CI

E. optimum r: (similar Q. in unmatched designs)

- w/  $r > 4$  or 5, little gain in size-eff. & thus informativeness - NB: usually limited in # cases
- mostly a function of cost-eff. - ie, the relative costs of getting data on test and ref. subjects

$r_0 \approx \sqrt{\frac{\#_t}{\#_r}}$  (one rule of thumb)

$\therefore$  use  $r < 4$  only when:  $\#_t > \sqrt{\frac{\#_t}{\#_r}}$   
 $4 > \sqrt{16}$

$\therefore$  use  $r < 4$  when ref. subj. costs  $> \frac{1}{16}$  of test  
ie. when ref. subj. are not readily available  
~~if using secondary data, use  $r < 4$  when there are fewer than 16~~

IV. When & on what var. should one M?

write on board

A. M. effects accuracy or inform. of study, according to

- Assoc's in source pop
  - $C \sim E$
  - $C \sim D$
  - among C's
- Methodology (most important)
  - type of study design - CH - CR
  - " " M. - indiv (r) vs freq.
  - selection of subj. - add bias? (do assoc. from source pop.  $\Delta$ )
  - type of analysis - stratified?

IV.

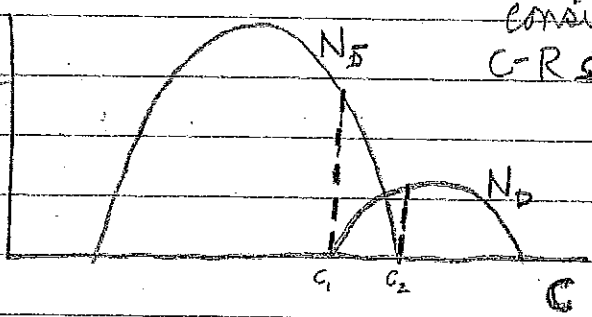
B. 2 criteria for evaluating effects of M: 2 components of in  
1. Validity (internal)  $\rightarrow$  systematic errors  $\rightarrow$  <sup>inference</sup> <sup>in</sup> <sup>accuracy</sup> <sup>of</sup> <sup>est.</sup>  $\rightarrow$  lack of bias  
2. Size-eff. - lack of random errors, reflected in range of CI for a given n.

not definitive statements of M, but based on crucial assumption: no additional bias introduced in selection of study pop. from the source pop. described (as assoc's) in table 1 - eg, Suppose age is risk factor of D in source pop, but restrict subj. to all 50-55 yr olds at entry; age could not be predictive of dis. in study, nor conf.

C. See: Table 1, pp 30 (#12) for summary - explain  
1. NB: (1):  $\downarrow$  indicates "overmatching"  $\uparrow$  M was more effective than not M.  
(2): only  $E_s$  (exact val.) is affected by this is not to say that M. is not done to control for conf. factors - only that bias may also be eliminated thru methods in analysis (eg, standardizing.)  
explain how M. can improve  $E_s$

i. Table 1 is a guide, designed to illustrate principles of M, not a system for making decisions  
 $\therefore$  table 1 represents expected effect of M, given the type of design, assoc. in source pop. + type

# persons (source pop)



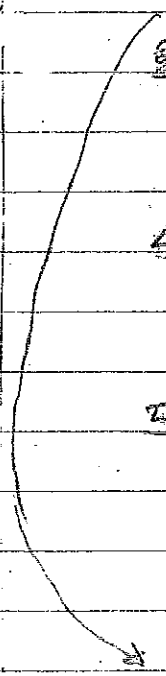
considers C-R study  
 $\downarrow$   
NB: C is a assoc. w/ D

in C-R study (or CH), info. is achieved by comparing D w/  $\bar{D}$ ; stratific. works by holding C constant + making comparisons w/in each stratum - now suppose, we took a simple random sample of D +  $\bar{D}$  in our C-R study: most noncases lie to left of  $c_1$  and most cases to right of  $c_2$   $\therefore$  thus no info. would be gained from these subjects since they cannot be compared in their respect strata - the effect of M. is to ensure that all subjects lie w/in "overlap" region thus achieving

### IV. C: Reasons for M

- \* 1. gain eff. (Es) - to prevent empty cells in stratified analysis - esp, important when C is meas. on nominal scale w/ many categories (can't combine) eg, "neighborhood"
- 3. cost of M. can be great, but cost-eff. can be improved if there is an inexpensive supply of ref. Subjects - assuming benefit in Es exists
- 4. in primary data - M obviates need to collect info on C => ↓ compute processing - can't study modifying effect of C
- 5. M. simplifies analysis conceptually, making it intuitively appealing as a means of controlling for conf - easy for non-stat. to see how M. removes bias (i.e. multivar. analysis)
- 2. In hybrid or C-R study (with new cases), there is a practical advantage in matching since the comparison series (noncases) should ideally be comparable to cases w/ respect to time of diagnosis. If cases are thus identified over time, the most feasible way to control for time of diagnosis is to individual match. Since if all noncases are selected at one time (say end of case detection), the time of diagnosis can't be controlled for thus stratification.

relates to (1)



IV

c. 1. (Ent) mbr info./ subject - i.e. improving  $E_s$   
\*  $\therefore$  M. reduces the loss in  $E_s$  that must result from stratif.

\* NB. even if distrib. are overlapping in their range, still same benefit accrues

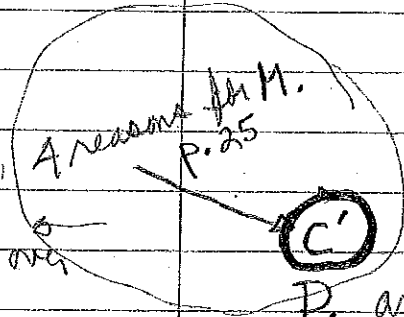
2. this gain in  $E_s$  is greatest when outcome var. is dichotomized

Following empirical (not theoretical) results

3. dividing contin. var. into > 3 or 4 categories produces little gain in  $E_s$

4. analogy: M in obs. studies is like blocking (B) in exp.

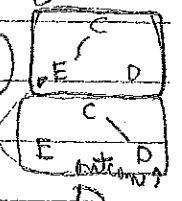
in exp. - control for conf. factors by randomiz[ation]  
Blocking <sup>(akin to R.)</sup> makes randomiz. more eff. by creating homogeneous groups w/ lots of  $\pm$  between-B. variance  $\therefore$  that the effect being measured is not a part of experimental error



D. another p[ro]b. from table is that ~~the~~ effects of M. differ for C-R & CH studies

M factor  
outcome

(I) CR: M  $\rightarrow$   $\downarrow$  Size-eff. (tho not always)  
CH: M  $\rightarrow$   $\uparrow$  " " (usually)



\* why this paradox? - in CH <sup>(as in exp)</sup> subj. develop. D over follow-up period & so D (outcome) is not regulated by M. procedure - so if E is perfectly correlated w/ D, E var. has a better chance of showing its effect - but, in C-R's E is outcome var. ~~so that~~ but is predetermined at onset of study  $\therefore$  E variability is influenced by M. - of perfect cond. exists between C and E, C M. only limits variability of E which then has less of chance of showing its effect on D.

Similar misconceptions:  
M - max. comparability  
not M - representativeness of controls

\* selection of controls in C-R study

\* w/ CR:  $\downarrow$  in  $E_s$  when  $E \sim C$  is strong points out a common fallacy of both M + w/M. C-R studies - D + B should not be alike on all variables except E.  $\Rightarrow$   $\downarrow$   $E_s$  because of less var. in E

IV

D. (3) effects of stratification in analysis - (in M. ~~and~~ design)

CH: (a) no loss in validity - i.e. selection bias can't be introduced by M.

(b) generally only gain in  $E_s$  thru stratification

C-R: (a) can lose ~~sub~~ val. by not stratification in analysis as we already saw occurs when  $\phi_{EC}$  is large enough

(b)  $\uparrow$  or  $\downarrow$  in  $E_s$

E. Table doesn't include X-S studies (of 1 pop) - M. is most costly \* is seldom warranted unless we have a fixed N - otherwise its cheaper to increase N and not match

F. if outcome var. is meas. on interval scale: M. is generally less eff. than no M. & the use of log analysis

\* G. 2 rules of thumb for selecting C's (con which to M.)

1. C's should be strong + predictors of dis. (exp. to avoid overfit in C-R)

2. if ~~there~~ several pot. C's (as they usually are), only M. on those that are ~~not~~ not highly assoc. w/ each other. - even tho the cost of M on ~~a~~ way isn't much more than 2, their intercorrel. will ~~not~~ result in  $\downarrow E_s$  because of M. on redundant factors. -

note example, p. 10 same IdR

3. can use "selective M" - but stratify (by categories of C's) on all variables (including M var.) in the analysis

V

After M. - stratify or not

A. Can't tell if conf. after having M'ed, but can determine one assoc. -  $\phi_{EC}$  in C-R;  $\phi_{CO}$  in CH -

B. see Table 2;

C. function of both val. & eff.   
 make do with unique var. of M. & E.



II. Illustrative Ex. - C-R; hypoth. source pop. - p.36 (#12)

A. do study in 4 ways (each way w/ n=140)

1. <sup>pairwise</sup> matched - strat. (by pairs)

2. " - unstrat. + strat. (by categories of C)

3. unmatched - stratif (since C is conf.)

4. matching - strat. by categories of C results p. 41

B. Note: IRR for M, strat  $\approx$  SMR unmatched strat any diff. is due to rounding errors, but - they would have been different if ~~if there~~ C had also been an effect modifier - in that case, fixing the ratio of test to ref. subj. (by M.) implicitly alters the distrib. of C in the study pop.  $\therefore$  creating a different weighting of ~~stratum~~ category. Specified effects (of course applies to unM. analysis as well)

\* Thus. M influences the IRR in a ~~matched~~ C-R in which C is also an effect modifier  $\therefore$  we might want to make separate analyses of the E  $\rightarrow$  D assoc. for diff. categories of C, but this requires enough subjects per category of C  $\rightarrow$  seriously reducing the power of each test  $\therefore$  M in C-R not only prevents us from studying C as a risk indicator but also makes it more difficult to study EM.

NB. M-strat analysis is more precise

\* C. M. - then forgetting pairs, but stratifying in the analysis by categories of C - get same RR as strat. by pairs, & even slightly better p-value  $\rightarrow$  Perhaps because of rounding errors

2/25 Outline 3/3/76

I. Intro.

- A. today: new method for considering several pot. conf. OM calls: "Confounder Summation"
- B. the problem: we would like to determine assoc. between E + D controlling for several pot. conf.
- C. Solutions (already discussed)

1. sequential stratif - no good because we saw that C's must be considered simultaneously
2. genl. stratified analysis - but when ~~the~~ several C's, we wind up w/ lots of strata (cells) some w/ very few or no subjects - also strata in which there is either mostly test or ref. subjects  $\therefore$  loss of  $E_s$  - these cells w/ very few subjects or w/ mostly only one group provide little or no info.

3. model fitting (multivariate) - eg, covariance - problem
  - a. hard to do, esp. for non-statisticians
  - b. more difficult to interpret - since specific trends get buried in coef's etc. - eg, might want to know how interaction occurs, not just if is sign. - to ~~draw~~ draw inferences
  - c. emphasis is usually on sign. levels, not estimation of effect. - tho, certain ~~an~~ stat. models allow for easy computation of certain measures of Assoc. - CH etc.

$\Rightarrow$  prediction technique

d. step-wise model showing sign. contrib. of variables may be reduced by additional covar. indicators in the data (which would be partialled out first) just like stratification problem of too many strata

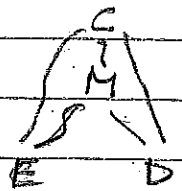
if also conf.

note: (eg) stepwise regression looks for parsimonious model - ie, model w/ least terms & explaining most variance - ie, eliminate redundancy which may have causal significance

c. added assumptions may be needed - eg in linear relationships in partial correl. or normality = variance, or no interaction (in covariance)

- i. linear -  $R = \alpha + \beta X$ ;  $RD = \beta(X_E - X_{\bar{E}})$
- ii. ln linear -  $R = \exp[\alpha + \beta X]$ ;  $RR = \exp[\beta(X_E - X_{\bar{E}})]$
- iii. logistic -  $R = 1 / (1 + \exp[-(\alpha + \beta X)])$ ;  $OR = \exp[\beta(X_E - X_{\bar{E}})]$   
(true even w/  $X_i$  if no interaction terms)

I. c. 3.\* (A) determining interaction creates a problem -  
 if modifier (eg, age) is assoc. w/ C, it is kept "fixed" w/ C in the analysis i.e. rendered meaningless as EM



→ NB: could do partial correl. or regression (interpretation is a statistical problem)

II. New Method - I'll call "composite stratification"

(one alternative - might also use binary variable multiple reg. Feldstein)

A. idea is to summarize all conf.'s into a single multivar. function + use this F to rank all subj. into a small (5) # of strata - then use genl stratif. analysis as if we had 1 conf.

B. Multivar. Statistics (model fitting) is used but it not interpreted as a predictive model, but only to control for conf. in one's data

III. Summarizing all conf. into 1 score<sup>(S)</sup> amounts to finding a way to pool the many strata that would result from genl stratif. - ie, the problem of too many strata will be dealt w/ by pooling them according to some criterion

A. this criterion: when can 2 strata be combined w/o introducing conf. in a study

1.  $P_{ED, \bar{E}}$  for any type of study is requirement for no conf.
2. so can combine any strata in which the expected rate (or %) of D is identical for condition of  $\bar{E}$  - eg, if rate  $(P, I, I_d)$  of D is same for old fem. + young males, cond. on  $\bar{E}$ , can combine these 2 strata

B. So, the Conf. Summary Score (S) will be based on some multivar. function (F) conditional on  $\bar{E}$  - ie, S will be assumed  $\bar{E}$  in order to rank subj.

III. E. It's not necessary that F actually estimate null risk ( $Pr[D|CE]$ ), but only need allow ranking according to the null risk ~~found~~ in the study pop. — in fact this Pr is not even a risk in most C-R studies, only  $Pr[\text{being a case}|C, E]$

IV. Obtaining the Conf. - Summarization Scores (S)

A. Use stat. model to fit data - F

1. type of model (dept. var)

a. Dis dichot. → <sup>linear</sup> discrimin. function

b. Dis interval → <sup>multiple</sup> regression

2. indep. var's

a. interval - use the value, eg, age

b. categorical - create k-1 dummy var's - for k categories of X - eg,

race: W, B, R — i. race = factor

$$\left. \begin{matrix} X_1 = (B=1 \quad \bar{B}=0) \\ X_2 = (R=1 \quad \bar{R}=0) \end{matrix} \right\} = \text{variables}$$

c. let  $E_i =$  opposite ~~to~~ <sup>may be</sup>  $> 1$  var. (E if dichotomous) <sup>if categorical</sup>

d. let  $C_j =$  post-conf. variables - including all dummy var - or some transform of var. — eg,  $\ln \text{age}$

how to select initial  $c_j$  - from RR\* one at a time - the only <sup>TR</sup> a crude guide

B. create F (assume 2 categ.)

$$\hat{D} = a + bE + \sum c_j C_j = F \quad (a, b, c_j \text{ are coef. from model})$$

1. NB: since object will only be to rank Subj, assumption (eg,  $\sim N$ ) need not be met

2. Size of  $c_j$  (coef.) doesn't indicate how conf.  $C_j$  is —

3. Since object is only to rank subjects, a is not needed

IV.

c. Since any stratif. analysis reduces  $E_s$ , we would like as few  $C_j$  in model + still control for conf. to minimize loss in  $E_s$  - i.e., reduction of model F

- Shouldn't be done by step-wise procedure or any technique based on sign. testing - because object is to control for conf., which isn't a matter of making inferences to a source pop, but refers only to study pop.
- to do: criterion - take out each  $C_j$  (one at a time) and obtain new F w/ new coef. if  $b$  (for  $E$ ) does not  $\Delta$  much, ~~at~~ that  $C_j$  was not very conf., so remove from model

Is there any order in removing  $C_j$  - 2 ways:  
 → ① try model without each of the  $C_j$  & drop any that don't  $\Delta$   $b$  much  
 ② Sequential: drop one  $C_j$ , then get new  $b$  for which we would evaluate next  $C_j$  ...

$$F' = a' + b'E + \sum_{j=1}^{r-1} c_j C_j \quad (\text{if } b \approx b', \text{ remove } C_j \text{ from model})$$

D. then must convert F into <sup>indiv</sup> Score (S) on which we will stratify - done by setting  $E$  to null value ( $E=0$ ) -

Reduced (denote w/ symbol)  $F'$ ?

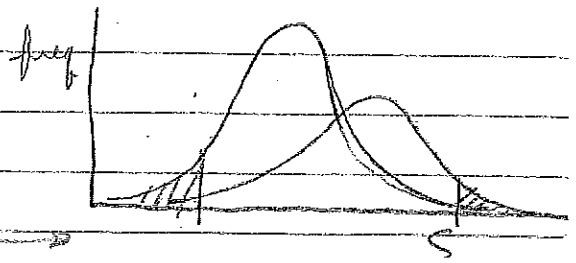
$$S = F' - E = a + \sum c_j C_j \quad (\text{same coef. - reduced } F')$$

(i.e. each subj. is given score  $S$  summarizing all  $C_j$  |  $E$ )

V. Stratify subjects into  $\approx 5$  categories along  $S$  - trickiest part of this whole procedure method

A. prelim. step: eliminate subjects from which we get little or no info. - i.e., delete fringe areas of  $S$  where either  $D$  or  $\bar{D}$  predominate

this is equivalent to subj. restriction in selection so bounded range of  $S$

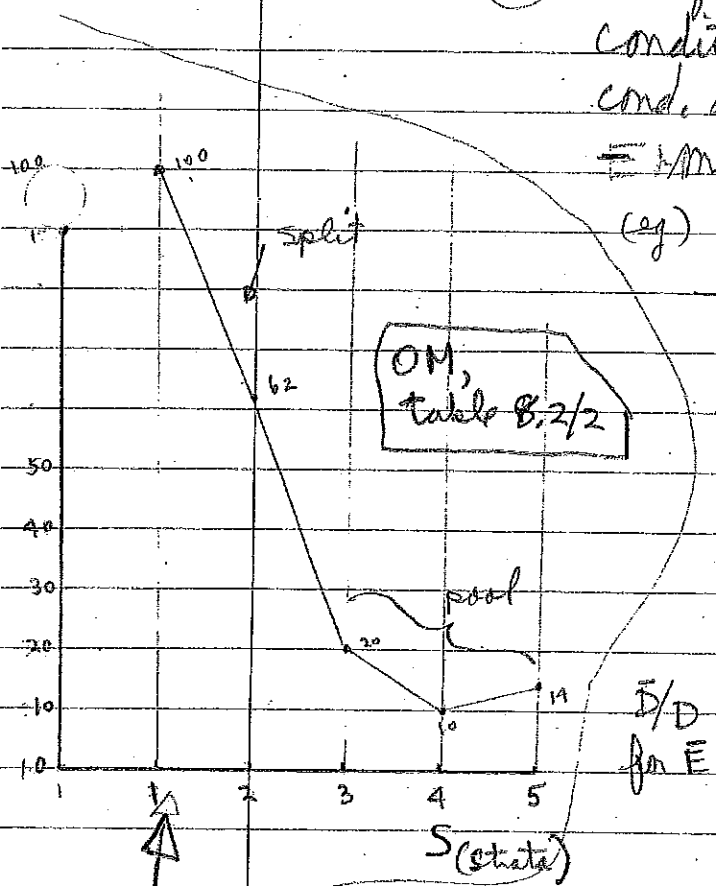


IV

B. make initial strata at stratif. by simply dividing all D (or  $\bar{D}$  if they are fewer) into quintiles along S & let the  $\bar{D}$  fall as they may

C. we would like to create reasonably homogeneous strata <sup>within</sup> w/in range of S ~~just~~ <sup>just</sup> defined, because we would like ~~inter-stratum~~ <sup>inter-stratum</sup> conf. to occur between strata, not w/in them - 2 ways of maximizing the conf. between strata + thereby Min. variability of S w/in strata

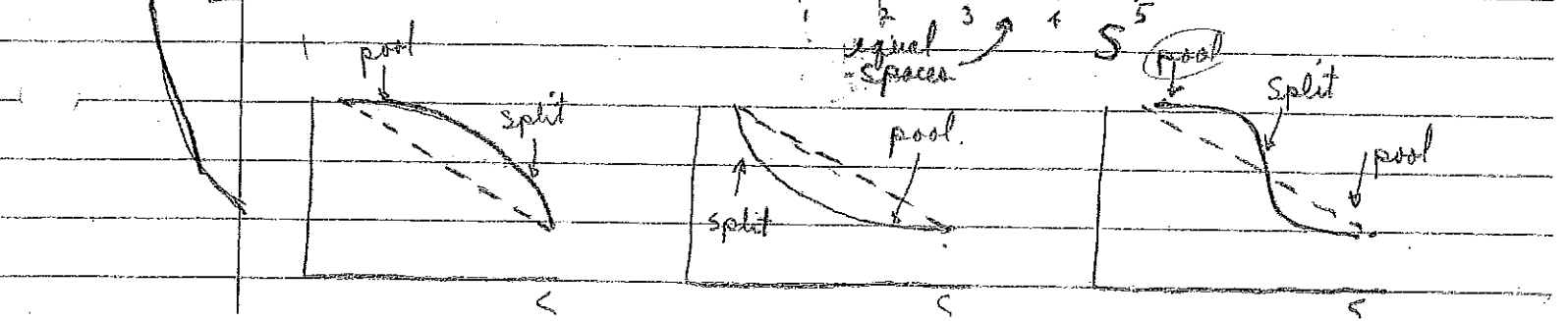
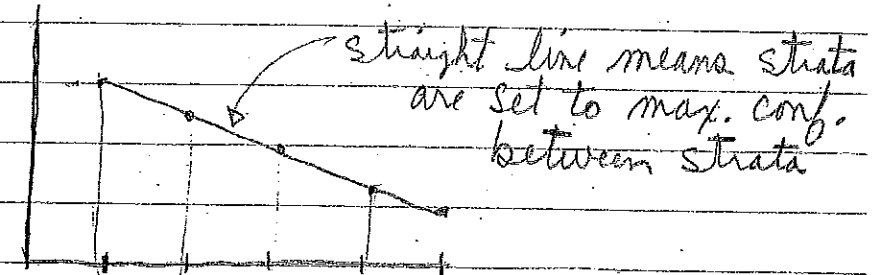
- (topol) 1. uniformity of  $e_{\bar{D}}$  over strata, or  
 2. uniformity of  $(\bar{D}/D)$  ratio over strata, conditional on  $(\bar{E} \text{ or } E)$  - NB. must be cond. on E, so we don't confuse conf. w/E  
 = My procedure: ratio  $(\bar{D}/E) | \bar{E} = \bar{D}/D$



(eg) ( $\bar{E}$  only)

S	D	$\bar{D}$	$\bar{D}/D$	<del>ratio <math>\bar{D}/D</math></del>
1	10	30	3.00	3.00
2	10	25	2.50	2.25
3	10	20	2.00	2.50
4	10	15	1.50	2.25
5	10	10	1.00	3.00

quintile  
 assume evenly spaced out



IV. D. but ~~this~~ while this process may conf. between strata, it does not insure that there is only a minimal amount of conf. w/in strata, (absolutely) because E itself may not have been adequate - so we ~~obscure~~ check homogen. w/in strata by comparing of D and  $\bar{D}$  as to their distrib. by each of the indiv.  $C_j$ , ~~exp~~ for only  $\bar{E}$  subjects (so as not to confuse EM w/ conf.) - we do this by comparing the %'s of D and  $\bar{D}$  that are ~~split~~ at "high risk" of each indiv.  $C_j$ , within each stratum, in  $\bar{E}$  subjects.

(% of subj. in high risk / stratum)

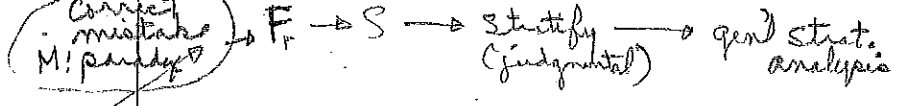
$C_j$	S=1		2		3		4		5	
	D	$\bar{D}$	D	$\bar{D}$	D	$\bar{D}$	D	$\bar{D}$	D	$\bar{D}$
$C_1 = \text{eq. age (40-60)}$	20	20	30	30	40	40	50	50	60	60
$C_2 = \text{sex = M.}$	10	10	15	30	comb.					
$\vdots$			% of exposed cases in $S_1$ who are M							

see: table 8.2/1 (OM) look at: MI - oh age - not, but may be EM

E. Note that getting strata involves a lot of judgment &  $\therefore$  takes practice

VI. Analysis

- A. Proceed as in Gen'l. strat. - as if w/ one  $C = S$
1.  $H_0$  testing -  $\chi_0$
  2. Est. of Effect
    - a. pt. est. - eq,  $\hat{RR}_0$  or  $\hat{SMR}$ ,  $\hat{SRR}$ ,  $\hat{EF}$
    - b. CI



begin:

VI

B. F (or S) does not necessarily represent a risk function - ie, if a  $E_j$  (eg, age) is not included in the reduced F, it doesn't mean that age is not a risk factor of D - S is only a classification score used to rank subjects in the study

1. CH study: 2 ways to see  $S \neq$  risk
  - a. we might select subj. or restrict their selection in such a way as to reduce the variability of a  $E_j$  (eg, age) in the study pop, so that age cannot predict the D in the particular data.
  - b. also, age might be an important predictor of D, but not be assoc. w/ E in the study  $\therefore$  age would not be conf. - ie it would be dropped from F ( $\Delta b$ )

2. C-R study:
  - a. never really considering risk ( $Pr[\text{getting } D | E]$ ) in most C-R studies, but  $Pr[\text{being } D | E]$
  - b. we do not even get this Pr if the ref. series is unrep. of the source pop. of the cases - might occur in 3 ways
    - i. Matching - ie (eg) if age-M, age won't appear in F because it can't discrim. D from  $\bar{D}$
    - ii (even w/o M), selection might occur so that source of ref. subj. is unrep. of the source of cases - eg, dilute role of age ~~can~~ in F by selecting D +  $\bar{D}$  of similar ages <sup>distrib<sup>g</sup></sup> eg from same hospital (Sort of like partial Matching or freq. M.)



VI

C. Suppose: E has > 2 categories i. requiring 2 or + dummy variables, eg, race

E<sub>1</sub> = Blackness - 1 (if B); 0 (if not B)

E<sub>2</sub> = Redness - 1 (if R); 0 (if not R)

W = ref.

F = a + b<sub>1</sub>E<sub>1</sub> + b<sub>2</sub>E<sub>2</sub> + Σ c<sub>j</sub>C<sub>j</sub> (unreduced)

1. one way: 2 separate analyses; one comparing B w/ W and R w/ W

F<sub>1</sub>' - depends on Δb<sub>1</sub> (treating E<sub>2</sub> as conf.)

F<sub>2</sub>' - " " Δb<sub>2</sub> (" E<sub>1</sub> " " )

S<sub>1</sub> = F<sub>1</sub>' - b<sub>1</sub>E<sub>1</sub>

S<sub>2</sub> = F<sub>2</sub>' - b<sub>2</sub>E<sub>2</sub>

problem: if interested in gradient of effect (over E<sub>i</sub>), standardizing requires comparable set of weight for each comparison - i.e., same standard -> SRR - but multiple analyses changes the strata & thus the weights, so...

Correctly recommended

2. F' - remains unreduced & live w/ some inefficiency - or perhaps use conf. of previous analysis (assuming E is dich.)

S = F' - b<sub>1</sub>E<sub>1</sub> - b<sub>2</sub>E<sub>2</sub> (one analysis)

then, standardize externally ("mutually stand.") to get comparable estimates of effect over categories of E<sub>i</sub> - O.M. used in S<sub>m</sub> -> M.I.I

3. can also assume E is ordinal and assign integers (eg, 0, 1, 2...) but this adds extra assumption - not recommended unless E is nearly interval

(recall that main problem w/ partial correl., stepwise regression, or any model fitting is inability of looking for effect mod. of var. that are assoc. w/ conf.)

VI

D. Looking for EM. (in broad sense): uniformity of eff

1. S as EM - not to draw specific inferences about associations but to determine whether strata of S may be pooled to create  $RR_0$  (than  $\hat{M}$  or weighting by inverses of variances of  $\hat{RR}_j$ )

just calculate  $\hat{RR}_j$  - test for uniformity:

a.  $\chi^2_u = \sum X_j^2 - \chi^2_0$   $n_j$

b. a more specific model - eg, using weighted regression

(?) 2. a single component factor (of S) as the effect modifier<sup>(M)</sup> - more concerned w/ specific assoc's - 2 approaches - which is best?

~~a. separate analyses for each category of M - w/ diff. E, S, strata, etc. eg,  $S_m \rightarrow CHD$  do analysis for:~~

- ~~i. all subjects in study~~
- ~~ii. hypertensive~~
- ~~iii. borderline hypert.~~
- ~~iv. normotensive~~

~~ie., we control for confounding separately within each category of M, calculate SMR's and compare~~

recommended

b. single analysis (one S w/ all subj) - calculate SMR's w/in each category of M - may have been what OM used in  $S_m \rightarrow MI$  -

(one) it seems to me that this approach may only be justified on the basis of convenience

VI

E. How to consider interaction of ~~several~~ factors

1. Interaction of  $C_j$ 's - in order to account for more conf. in the data - since the importance of any single conf. is reduced when there are several, interaction of  $C_j$  may not be particularly useful (unlike a prediction model), but under certain conditions, it might be desirable - eg, w/ drugs which are known to interact pharmacologically - eg,  $c_1, c_2, c_3 \rightarrow$  treat as another conf. or age

2. Interaction of E and  $C_j$  - not interested in studying interaction per se, but improving model

(F) - ie, getting planning more variance in D

(eg) E = Smoking  
 $C_1$  = age

$$F = a + bE + b_1 E \cdot C_1 + c_1 C_1 + \sum_{j=2} c_j C_j$$

$$S = F - bE - b_1 E \cdot C_1 \quad (\text{when } E=0)$$

NB: can't reduce F in normal way because both E +  $C_1 \cdot E$  become 0 when  $E=0$ , so would have to observe A in 2 b's -

so, might get  $C_j$  first w/o interactive term then do ~~some~~ ranking w/ it - tho might have to sacrifice some effie.

3. Interaction of 2 E's - ie, when we are interested in studying interaction per se - eg,  $E_1; E_2$

$$E_1 = 1 (\text{Sm}) ; 0 (\text{non-Sm})$$

$$E_2 = 1 (\text{dr}) ; 0 (\text{no dr})$$

$$F = a + b_1 E_1 + b_2 E_2 + b_3 E_1 E_2 + \sum c_j C_j$$

$$S = F - b_1 E_1 - b_2 E_2 - b_3 E_1 E_2$$

VI E. 3. (cont) - again, F cannot be reduced in normal way, so might have to sacrifice some E<sub>s</sub> - then analyze as if S were 1 var., which we would like to control in looking for interaction

SERR = SRR - 1 (Cep. stand.)

	$S_m$	$S_{\bar{m}}$	
$D_1$			ie, control for conf. w/in each comparison ref. need large n
$D_2$		0	

VII. "Conf.-Summarization" vs. alternative methods:

A. Confusion exists: some say

1. it is new method
2. it is same thing that is done by researchers who stratify by a MLF (eg. Evans City)
3. (me) actually, the 2 procedures are related, but do not have same utility

B. Differences between conf-sum. and MLF

1. major diff.: MLF is created as a predictor of D (ie, risk function) - while S is a classification score for ranking subj. in order to control for conf
2. often  $\beta$ 's of MLF are applied to another study pop. - not done w/ S; requires different ranking for each pop. - different selection factors in each study
3. MLF is chosen as the model because it creates a function that remains between 0 & 1 (like a Pr) - but not necessary w/ S; they might want to take log transformation - eg. ln Age
4. MLF only ~~more~~ reflects assoc of G<sub>j</sub> w/ D, not G<sub>j</sub> w/ E, which is also a prereq. for conf. ie MLF as a stratification device is liable to be very inefficient

VII

- B. 5. ~~an~~ attempt is made to reduce or limit # terms ( $G_j$ ) in MLF by statistical procedure, eg, step-wise procedure, not done in S.F.
- \* 6. MLF is not generally converted into a score conditional on  $\bar{E}$  - so strong predictive MLF may not be very conf. - ie,  $P_{CD}$  is not criterion for conf.

So, using MLF (as is) to control for conf. may be

1. very inefficient (~~not~~)
2. possibly invalid, in the sense that <sup>max.</sup> conf. is not being controlled, ~~between~~

(VIII)

Selective Matching - eg. M. on age, race, sex  
 Sm  $\rightarrow$  CHD  
 also would like to control for obesity, HT, chol.

So: use M. only as restrictive device  
 i.e. matched 'sets' in the analysis but stratify on 'all' b vari's - using Composite Stratification

3/17 Outline

I. Introd.

A. O.M on 3/24 - will meet w/him

B. today:

1. finish composite stratification
2. class exercise (Syme et al., 1964)

II. Criticisms of Syme et al. (1964)

A. Study design and parameters

1. type of design - had all the ingredients of an ambidex. study but treated as C-R
2. could have computed  $I_d, I_{dE}, I_{dE}, I_{dR}, I_{dD}, EF, I_{c_i}$  -  $\hat{I}_d = \frac{228}{20,000+1}$   $+ CI$  but did not actually compute any of these

B. Matching

1. analyzed data without stratifying for age (the M. var.) - in C-R analysis this could cause both invalidity & loss of sig. effic.
2. what should have been done as suggested by authors

note:  
 ① dynamic pop  
 ② should be age-spec.  $I_d$

		$2E$	$E/E^D$	$2E$
$D$	$E$	$U$	$V$	$W$
	$E$	$X$	$Y$	$Z$

$\hat{M}_L$  of  $I_{dR}$  requires iterative soln. because  $r=2$  - tho,

$$\hat{M}-H(\text{not } \hat{M}_L) = \frac{2W+V}{2X+Y} \quad \text{short-cut}$$

%E ratio  $\Rightarrow a = (U+V+W)$

$$E[a] = \sum \frac{a_j + b_j}{r+1} = U\left(\frac{2}{3}\right) + (X+V)\left(\frac{2}{3}\right) + (Y+W)$$

$$= \frac{3U+2X+2V+Y+W}{3}$$

$$a/E[a] = \frac{3U+3V+3W}{3U+2X+2V+Y+W}$$

II. B. 3. What authors did.

	D	$\bar{D}$	
E	u+v+w	2u+2x+v+y	609 (Subj.)
$\bar{E}$	x+y+z	v+y+2w+z	
	203	406	

O/E ratio  $\rightarrow a = u+v+w$  (same)

$$E[a] = \left( \frac{2u+2x+v+y}{406} \right) 203 =$$

$$= u+x + \frac{v+y}{2} = \frac{2u+2x+v+y}{2}$$

$$a/E[a] = \frac{2u+2v+2w}{2u+2x+v+y} \quad (\neq \text{ as above})$$

4. Clearly the authors got a different R than if they had stratified in the analysis -  
 But suppose they had simply computed the OR - instead of these R's (ignoring age & sex)

	D	$\bar{D}$
E	a	b
$\bar{E}$	c	d

OR =  $\frac{ad}{bc}$  but authors;

$$R' = \frac{R_E}{R_{\bar{E}}} = \frac{a/E[a]}{c/E[c]} =$$

authors never do this formally  
 to give it a name

$$= \frac{\frac{a}{b} \text{ (a+c)}}{(b+d)} = \frac{\frac{ad}{bc}}{\frac{d}{b} \text{ (a+c)}} = \frac{ad}{bc} = OR$$

So,  $R' = R_E/R_{\bar{E}} = OR$

( ) II

B, 5. So forgetting the basic mistake of not stratifying by age in analysis, the authors actually got the  $\hat{I}dR$ , which they "correctly" interpreted - see: p. 3 - but this estimate actually involved 3 mistakes w/ re: to what they said they were doing.

a. They included the missing data in the totals (for D +  $\bar{D}$ ) - thus, they maintained a constant ratio of cases to controls

b. the expected values were computed incorrectly - for a C-R study

$$E[a] = \frac{bc}{d} \neq \frac{b}{(b+d)}(a+c)$$

NB:  $SMR = \frac{a}{\sum \frac{b_j c_j}{d_j}} = \frac{ad}{bc} = OR$  (for only 1 stratum)

if  $R$  had then been computed correctly it would have been an appropriate meas. of assoc. itself + might have even been used to standardize for potential confounders - but were forced to compare

c. 2  $R$ 's

e. each  $R_i$  treats the other  $R$  category as the referent - so,

- $R_1 \rightarrow R_2$  is referent
- $R_2 \rightarrow R_1$  " " "

it wasn't until authors compared  $R$ 's ( $R'$ ) that they arrived at a meas. of assoc.



II, B, 5. (c) cont.

this might result from a <sup>lack of the</sup> fundamental understanding that analytical epid. requires a comparison of groups - the authors treat this idea as secondary & seem to believe that a single category (of a var.) ~~is~~ may be observed to be assoc. with the disease

C. Confounding -

1. never actually determined conf., just whether there was still an assoc. w/in categories of pot. conf. (c) - i.e., as if they were studying Eff. Modif. and searching for "secondary assoc." in fact what they did resembled control table analysis but instead of using I<sub>2</sub> ratios they used "R<sub>i</sub>"
2. more important - did not consider all pot. conf. simultaneously - eg, table 11, p.5 - conf. effect of BP should be seen conditionally on all other conf.
3. also table 11,  $\chi^2_{(5)}$  is meaningless - both SC mob. + Hypert. are treated as one factor in getting a p-value → i.e. they looked for the combined effect of both variables, not for the SC effect, controlling for BP - they could have used  $\chi^2_0$
4. it would be important to see whether certain confounders might represent <sup>intervening</sup> mechanisms - eg, SC → diet → CHD  
this might show diet as math. conf. but does not do that from <sup>standard</sup> mechanism of SC

II. C. 4. (cont.)

might have an indication of this intervening mech

RR (of hi SC mob. | low fat diet) ≈ 1

but RR (of hi SC mob) > 1

5. when comparing > 2 categories of E (eg, 3 classes of SC mob.) would like to identify "dose-response" gradient - table 11, p. 5 requires SRR, not SMR to have comparable standards for each comparison

SC<sub>1</sub> = ref.

CRR<sub>3</sub> = 2.28

SRR<sub>3</sub> = 2.64 ∴ neg. conf.

CRR<sub>2</sub> = 1.34

SRR<sub>2</sub> = 1.33 ∴ no conf.

	stable	SC mod. mobile	highly mob!	
SRR	1.00	1.33	2.64	⇒ dose-response

D. Effect Modif. + Interaction

could have used SMR

1. authors seem to confuse conf. + EM.
2. looked explicitly for st. interaction (table 15, p. 6) but failed to actually determine whether effects were nearly additive and failed to control for other potential conf. - could have used SRR to analyze the interaction between any 2 indep. var. (eg, SC + BP)

	Norm	Hyp
(1)	act.	
SC (2)		
(3)		

SERR = SRR - 1