

AGENDA

Review of handout on confounding from c678

Key Points in / Commentary on ALR Ch 3.5 —

Worked Examples

- Salaries of PhD's versus MSc's
 - Y measured in \$.. i.e. interval scale
 - same principles whether confounder is binary or interval
 - extreme confounding : example of Simpson's paradox
- Who is more likely to receive death penalty?
 - 2 binary "X" variables
 - extreme confounding: example 2 of Simpson's paradox
- Neonatal Mortality (Brand and Keirse article)
 - good e.g. of adjustment formula (crude -> adjusted logit)
 - age confounds M:F comparison, but not vice versa!
- Down's syndrome in relation to parity & maternal age
 - using age as an interval versus categorical variable
- Autism & MMR
 - modeling age-specific dx-rates as log-linear doesn't reduce the confounding (see exercise)

Demo

- Resting on a knife-edge: collinearity
 - Excel spreadsheet in session 4 of course 678

Readings (* = ?insufficient)

* Chapter 3.5 of Hosmer and Lemeshow

Other Resources

texts

Statistical methods for comparative studies: methods for bias reduction, by Anderson et al. -- 6 authors, whom I abbreviate to 'aahovw'. Extract from the relevant chapter 3 on the c622 website. Diagram for c678 adapted from these authors

articles

- Brand & Keirse: pair of very good expository articles (from Pediatric and Perinatal Epidemiology 1990; 4: 22-38) on logistic regression [c678, Resources/Material for sessions 9-11]
- "Appropriate Uses of Multivariate Analysis" (JH) from the Annual Review of Public Health in 1983.
under Other Resources in the c622 website,
See in particular the sections on two major uses of multiple regression:- to make comparisons
FAIRER (reduce bias)
SHARPER (increase precision).
- Figure in the Blood Pressure and Altitude article under resources for session 5 of course 678

Clear and present dangers

1 In the article CLINICAL PRACTICE GUIDELINE: ENDPOINTS OF RESUSCITATION, Tisherman reviews a number of studies, including Maynard, N 1993 "Assessment of splanchnic oxygenation by gastric tonometry in patients with acute circulatory failure. JAMA 270:1203-10." and tells us that

Although a variety of resuscitation endpoints correlated with surviving critical illness, only pHi at 24h proved an independent predictor of *death by logistic regression*

2 A title of a Table in another article..

ADJUSTED ESTIMATES OF ASSOCIATION OF PATIENT CHARACTERISTICS WITH INITIALLY MISSED DIAGNOSIS, DELAYED TREATMENT, AND DEATH FROM MULTIVARIATE LOGISTIC REGRESSION

Key Points in / Commentary on ALR Ch 3.5 —

• **The multivariable model** (3.5 page 65, para 1)

"One goal of such an analysis is to statistically adjust the estimated effect of *each* variable in the model for differences in the distributions of, and associations among the other independent variables " (italics by JH)

Remark 1: yes, in some analyses (but certainly not in all, and not in cases of diagnostic and prognostic functions) this might be a goal

Remark 2: in many studies, there is not the symmetry of role among the X's that the authors imply: more often, there is *one* X variable ('exposure' for want of a better word) that is of *primary* interest, and the *other* X's are a *nuisance*. The authors' use of the word *each* makes the aim broader than it really is. Indeed, in situations where the focus is on one X, it would be helpful to rewrite the equation as

$$g[\mu \text{ or } | \mathbf{X} \mathbf{Z}] = \mu_0 + \beta_X X + \beta_Z \mathbf{Z}$$

where X is the 'eXposure' variable, and **Z** is the collection of the confounding and other explanatory variables. [H&L do so in group & age e.g on p66)

Mathematically speaking, there is a certain *symmetry* to the regression equations. In particular, in c621 situations, X2 being a confounder of the observed X1=>Y relationship implies that X1 is *also* a confounder of the observed X2=>Y relationship. For example, when we study the effect of duration of occupational exposure to noise on the prevalence hearing loss, we are, because of practical research design constraints, forced to confront, and 'adjust for', the simultaneous effect of age. But we would not add as a study aim "to adjust the estimated effect of age for differences in the distributions of exposure". Here, age is an extraneous variable, of no intrinsic interest, one that -- if we could do the study experimentally, on say animals, we would control directly, rather than by regression techniques. [This is not necessarily so in c681, where rules that apply in the 'regular' or identity scale, do not always hold in the logit scale: for example, in Brand and Keirse's first example (Table 9), age confounds the sex=>mortality relationship (if work with odds) but sex does *not* confound the age=>mortality relationship! This goes back to the rules for 'collapsibility, and confounding: these rules are different for different effect measure scales (i.e. the rules for confounding are different for odds ratios than they are for risk differences. I expect you will cover this again in the more advanced epidemiology classes).

There may be *some situations* where *two or more* X's (e.g. *risk factors*) are of *simultaneous interest*, and one is keenly interested in the true *net* effect of each. An example is the article "Parental periconceptional smoking and male:

female ratio of newborn infants (on the 678 website under resources for session 9-11), where the authors assessed whether the smoking habits of (1) the mother and (2) the father around the time of conception affects the likelihood of the offspring being male or female, and one would like to *isolate* the two effects (mother, father) from each other.. i.e. have narrow interval estimates for each beta, i.e.. we seek simultaneous CI's for the two betas such that the two estimates are not highly correlated (often negatively, if the two 'exposures' are positively correlated in their distribution in the sample, and their true effects are in the same direction), or even if there are, are within narrow ranges

• **The multivariable model** (3.5 page 65, para 4)

"If the age (Z) distribution is also the same for the two groups, then a univariate analysis *would suffice* and we would compare the mean weight (Y) in the two groups (X=0 and X=1, " (italics by JH)

Remark 1: YES, in c621, one would get the same estimate whether one omitted age or not. BUT including age would (if age contributes to the variation in Y, and presumably it does or else one would not be considering it) reduce the standard error (i.e. increase the precision) of the estimate.

Remark 2: Paradoxically, this again is one of the places where the message from c621 does not automatically carry over to c681.. Because, at any X, the variance of a binary Y is very large relative to its mean, and because this variance [(1-)] changes with the mean [i.e., with], it is sometimes difficult to say whether including a Z in a logistic model will decrease or increase the SE of the beta for the X of interest. Given that, and given that a sample of n binary Y's does not mean we have n effective degrees of freedom to fit as many parameters as we wish, we have to be more prudent about inclusion of Z's that in c621 "might not do much good, but would not do much harm either."

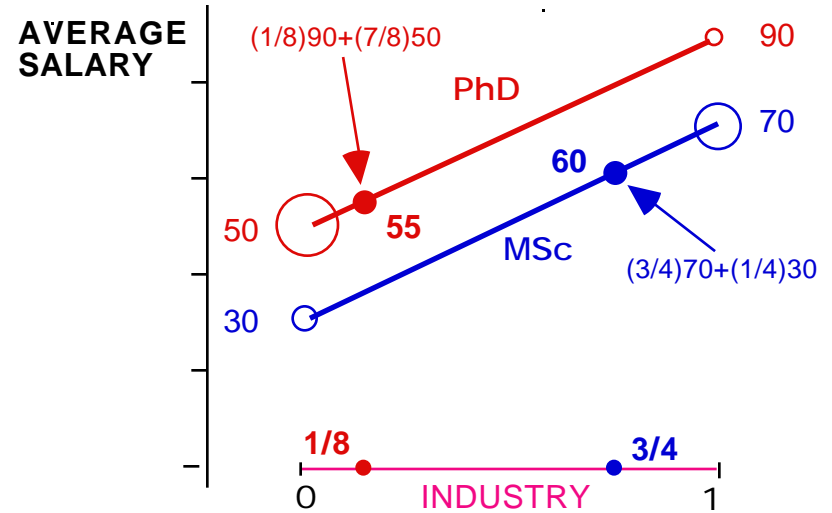
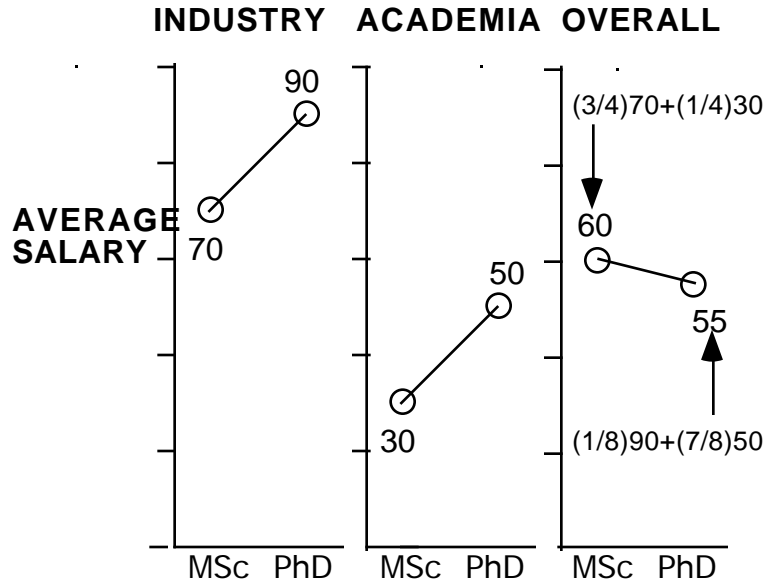
• **The multivariable model (analysis of covariance)** (3.5 page 65, para 5)

"This situation is described graphically in Figure 3.1"

Fig 3.1 is AT THE CORE of the concept of statistical adjustment, and is repeated with colour, and lines that reflect *where the data are*, on the right hand side of page 2 of what was handed out on confounding from c678. JH gives other examples later in the c681 notes that you are reading now, and points here to a textbook and articles where this is explained with greater graphical clarity and impact. The text is Statistical methods for comparative studies: methods for bias reduction, by Anderson et al. -- 6 authors, whom I abbreviate to 'aahovw'. I have put an extract from the relevant chapter 3 on the c622 website. Also, under Other Resources in the c622 website, is JH's article "Appropriate Uses of Multivariate Analysis" from the Annual Review of Public Health in 1983. It has two sections that explain -- graphically and with a few data points-- two big uses of multiple regression.. to make comparisons FAIRER (reduce bias) and SHARPER (increase precision). Finally, to see Fig 3.1 in a real situation, with much more striking graphics, you could examine the Figure in the Blood Pressure and Altitude article under resources for session 5 of course 678.

comments/corrections welcomed .. jh feb 9 2004

Here I adopt, and illustrate **confounding**, & "**adjustment by regression**" for an extreme 'continuous Y' example [607Ch2, Simpson's paradox]



<p>1.</p> <pre>DATA salary; INPUT phd industry salary number; LINES; 0 1 70 750 1 1 90 125 0 0 30 250 1 0 50 875 ;</pre>	<p>2.</p> <pre>PROC MEANS DATA = salary MEAN; CLASS PHD; VAR salary; FREQ number; PHD N Obs Mean_SALARY ----- No 0 1000 60 Yes 1 1000 55</pre>																				
<p>3a</p> <pre>PROC REG DATA = salary ; MODEL salary = Phd ; FREQ number;</pre>	<p>3b. Dep. Var: SALARY</p> <table border="1"> <thead> <tr> <th></th> <th>Estimate</th> <th>SE</th> <th></th> </tr> </thead> <tbody> <tr> <td>INTERCEP</td> <td>60</td> <td>0.49</td> <td></td> </tr> <tr> <td>PHD</td> <td>-5</td> <td>0.69</td> <td>CRUDE Δ</td> </tr> </tbody> </table>		Estimate	SE		INTERCEP	60	0.49		PHD	-5	0.69	CRUDE Δ								
	Estimate	SE																			
INTERCEP	60	0.49																			
PHD	-5	0.69	CRUDE Δ																		
<p>4a.</p> <pre>PROC SORT DATA = salary; by industry; PROC REG DATA = salary; by industry; MODEL salary = Phd; FREQ number;</pre>	<p>4b.</p> <table border="1"> <thead> <tr> <th></th> <th colspan="2">INDUSTRY = 0</th> <th colspan="2">INDUSTRY = 1</th> </tr> <tr> <th></th> <th>Estimate</th> <th>SE</th> <th>Estimate</th> <th>SE</th> </tr> </thead> <tbody> <tr> <td>INTERCEP</td> <td>30</td> <td>0</td> <td>70</td> <td>0</td> </tr> <tr> <td>PHD</td> <td>20</td> <td>0</td> <td>20</td> <td>0</td> </tr> </tbody> </table>		INDUSTRY = 0		INDUSTRY = 1			Estimate	SE	Estimate	SE	INTERCEP	30	0	70	0	PHD	20	0	20	0
	INDUSTRY = 0		INDUSTRY = 1																		
	Estimate	SE	Estimate	SE																	
INTERCEP	30	0	70	0																	
PHD	20	0	20	0																	

<p>5a.</p> <pre>PROC REG DATA = salary ; MODEL salary = Phd industry; FREQ number;</pre>	<table border="1"> <thead> <tr> <th></th> <th>Estimates</th> <th>SE</th> </tr> </thead> <tbody> <tr> <td>INTERCEP</td> <td>30</td> <td>0</td> </tr> <tr> <td>PHD</td> <td>20</td> <td>0</td> </tr> <tr> <td>INDUSTRY</td> <td>40</td> <td>0</td> </tr> </tbody> </table>		Estimates	SE	INTERCEP	30	0	PHD	20	0	INDUSTRY	40	0
	Estimates	SE											
INTERCEP	30	0											
PHD	20	0											
INDUSTRY	40	0											

Average Salary = 30 + 20 if PhD + 40 if Industry

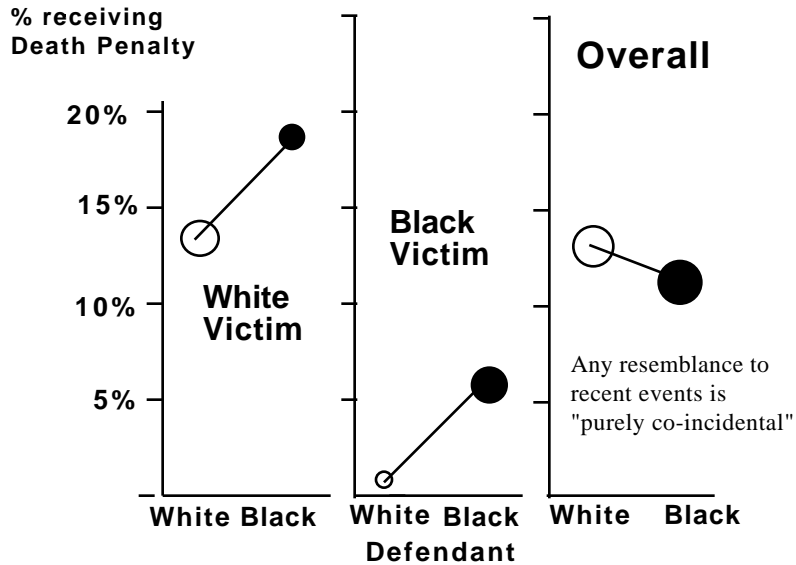
<p>6a.</p> <pre>PROC MEANS DATA = salary MEAN; CLASS PHD; VAR industry; FREQ number;</pre>	<p>6b. Distrn. of CONFOUNDER</p> <table border="1"> <thead> <tr> <th>PHD</th> <th>Mean of (Indicator)</th> <th>variable INDUSTRY</th> </tr> </thead> <tbody> <tr> <td>No 0</td> <td>0.750</td> <td>(3/4)</td> </tr> <tr> <td>Yes 1</td> <td>0.125</td> <td>(1/8)</td> </tr> </tbody> </table>	PHD	Mean of (Indicator)	variable INDUSTRY	No 0	0.750	(3/4)	Yes 1	0.125	(1/8)
PHD	Mean of (Indicator)	variable INDUSTRY								
No 0	0.750	(3/4)								
Yes 1	0.125	(1/8)								

From 2 or from 3b ("CRUDE") \$value of PhD: -5
 From 4b or from 5b ("NET") \$value of PhD: +20

Net (adjusted) value of PhD =

$$\begin{aligned}
 & \text{"CRUDE" value} \\
 & \text{MINUS} \\
 & \text{NET value of Industry} \times ([\text{Industry}|\text{PhD}] - [\text{Industry}|\text{non-PhD}]) \\
 & \quad -5 \\
 & \text{MINUS} \\
 & 40 \times (0.125 - 0.750) = -5 -40(-0.625) = +20
 \end{aligned}$$

On the "participation in industry" scale, PhDs are at 0.125, and non_PhD's at 0.750, a difference of 0.625; the salary gap from 'not in industry' (industry = 0, ie academia) to in industry (=) is 40K



<pre>data a; INPUT de_black vi_black n_death n_spared; n_cases = n_death + n_spared; LINES; 0 0 19 132 0 1 0 9 1 0 11 52 1 1 6 97 ;</pre>	<pre>data b; set a; death=1; nmbr = n_death ; output; death=0; nmbr = n_spared; output; PROC PRINT; VAR de_black vi_black death nmbr; DE_BLACK VI_BLACK DEATH NMBR 0 0 1 19 0 0 0 132 0 1 1 0 0 1 0 9 1 0 1 11 1 0 0 52 1 1 1 6 1 1 0 97</pre>
---	---

<pre>PROC FREQ data=b; TABLES death * de_black / cmh norow nopercnt; WEIGHT nmbr; "CRUDE"</pre>	<table border="1"> <thead> <tr> <th colspan="2"></th> <th colspan="2">DE_BLACK</th> <th rowspan="2">Total</th> </tr> <tr> <th>DEATH</th> <th></th> <th>0</th> <th>1</th> </tr> </thead> <tbody> <tr> <td>0</td> <td></td> <td>141</td> <td>149</td> <td>290</td> </tr> <tr> <td></td> <td></td> <td>88.1%</td> <td>89.8%</td> <td></td> </tr> <tr> <td>1</td> <td></td> <td>19</td> <td>17</td> <td>36</td> </tr> <tr> <td></td> <td></td> <td>11.9%</td> <td>10.2%</td> <td>(11%)</td> </tr> <tr> <td colspan="2">Total</td> <td>160</td> <td>166</td> <td>326</td> </tr> </tbody> </table>			DE_BLACK		Total	DEATH		0	1	0		141	149	290			88.1%	89.8%		1		19	17	36			11.9%	10.2%	(11%)	Total		160	166	326
		DE_BLACK		Total																															
DEATH		0	1																																
0		141	149	290																															
		88.1%	89.8%																																
1		19	17	36																															
		11.9%	10.2%	(11%)																															
Total		160	166	326																															
<pre>Type of Study Method OR CI C-C M-H 0.847 0.42 1.69 Logit 0.847 0.42 1.69</pre>																																			

<pre>PROC LOGISTIC data=a; MODEL n_death/n_cases = de_black ; "CRUDE"</pre>	<pre>-2 LOG L 226.513(null) 226.291(model) Estimate SE ChiSq P-value OR INTERCPT -2.00 0.24 67.3 0.0001 . DE_BLACK -0.16 0.35 0.2 0.6382 0.847</pre>
---	--

<pre>PROC FREQ data=b; TABLES vi_black*death*de_black / measures cmh; WEIGHT nmbr;</pre>	<table border="1"> <thead> <tr> <th colspan="2"></th> <th colspan="2">DE_BLACK</th> <th rowspan="2">VI_BLACK=0 Tot</th> </tr> <tr> <th>DEATH</th> <th></th> <th>0</th> <th>1</th> </tr> </thead> <tbody> <tr> <td>0</td> <td></td> <td>132</td> <td>52</td> <td>184</td> </tr> <tr> <td></td> <td></td> <td>87.42</td> <td>82.54</td> <td></td> </tr> <tr> <td>1</td> <td></td> <td>19</td> <td>11</td> <td>30</td> </tr> <tr> <td></td> <td></td> <td>12.58</td> <td>17.46</td> <td></td> </tr> <tr> <td colspan="2">Tot</td> <td>151</td> <td>63</td> <td>214</td> </tr> </tbody> </table> <p style="text-align: right;">OR 1.47 (0.65, 3.3)</p>			DE_BLACK		VI_BLACK=0 Tot	DEATH		0	1	0		132	52	184			87.42	82.54		1		19	11	30			12.58	17.46		Tot		151	63	214
		DE_BLACK		VI_BLACK=0 Tot																															
DEATH		0	1																																
0		132	52	184																															
		87.42	82.54																																
1		19	11	30																															
		12.58	17.46																																
Tot		151	63	214																															

<pre>Estimates of Common OR M-H ** 1.574 0.701 3.533 Logit* 1.454 0.667 3.173 ** test-based. * correction of 0.5 in every cell of tables that contain a zero.</pre>	<table border="1"> <thead> <tr> <th colspan="2"></th> <th colspan="2">DE_BLACK</th> <th rowspan="2">VI_BLACK=1 Total</th> </tr> <tr> <th>DEATH</th> <th></th> <th>0</th> <th>1</th> </tr> </thead> <tbody> <tr> <td>0</td> <td></td> <td>9</td> <td>97</td> <td>106</td> </tr> <tr> <td></td> <td></td> <td>100.00</td> <td>94.17</td> <td></td> </tr> <tr> <td>1</td> <td></td> <td>0</td> <td>6</td> <td>6</td> </tr> <tr> <td></td> <td></td> <td>0.00</td> <td>5.83</td> <td></td> </tr> <tr> <td colspan="2">Tot</td> <td>9</td> <td>103</td> <td>112</td> </tr> </tbody> </table> <p style="text-align: right;">OR not computed - 0 cell</p>			DE_BLACK		VI_BLACK=1 Total	DEATH		0	1	0		9	97	106			100.00	94.17		1		0	6	6			0.00	5.83		Tot		9	103	112
		DE_BLACK		VI_BLACK=1 Total																															
DEATH		0	1																																
0		9	97	106																															
		100.00	94.17																																
1		0	6	6																															
		0.00	5.83																																
Tot		9	103	112																															

<pre>PROC LOGISTIC data=a; MODEL n_death/n_cases = de_black vi_black ; OUTPUT OUT=fitted PREDICTED = fitted_p ; *cf 0.266 OR (crude) victim white 1 (ref) victim black 0.347</pre>	<pre>-2 LOG L 226.5(B0) 219.1(B0,B1,B2) ChiSq: 7.4 with 2 DF (p=0.0243) Var. Est. SE ChiSq P OR INTERCPT -1.95 0.25 63.9 0.00 . DE_BLACK 0.44 0.40 1.2 0.27 1.553 VI_BLACK -1.32 0.52 6.5 0.01 0.266*</pre>
--	---

<pre>PROC PRINT data=fitted; VAR de_black vi_black fitted_p; fitted_p;</pre>	<table border="1"> <thead> <tr> <th>DE_BL</th> <th>VI_BL</th> <th>FITTED_P</th> <th>OBS_P</th> <th>#OBS</th> <th>#FITTED</th> </tr> </thead> <tbody> <tr> <td>0</td> <td>0</td> <td>0.1237</td> <td>0.1258</td> <td>19</td> <td>18.7</td> </tr> <tr> <td>0</td> <td>1</td> <td>0.0362</td> <td>0.0000</td> <td>0</td> <td>0.3</td> </tr> <tr> <td>1</td> <td>0</td> <td>0.1798</td> <td>0.1746</td> <td>11</td> <td>11.3</td> </tr> <tr> <td>1</td> <td>1</td> <td>0.0551</td> <td>0.0583</td> <td>6</td> <td>5.7</td> </tr> </tbody> </table> <p>NB: 4 proportions modeled by 3 parameters: close to 'saturated' model.</p>	DE_BL	VI_BL	FITTED_P	OBS_P	#OBS	#FITTED	0	0	0.1237	0.1258	19	18.7	0	1	0.0362	0.0000	0	0.3	1	0	0.1798	0.1746	11	11.3	1	1	0.0551	0.0583	6	5.7
DE_BL	VI_BL	FITTED_P	OBS_P	#OBS	#FITTED																										
0	0	0.1237	0.1258	19	18.7																										
0	1	0.0362	0.0000	0	0.3																										
1	0	0.1798	0.1746	11	11.3																										
1	1	0.0551	0.0583	6	5.7																										

Other eg's: neonatal mortality (Brand/Keirse); Down's, parity, age.