

<div> <div>Chapter 12: Strategies for Variable Selection (Class 2 of 2)</div> <div>Class 20, 4/22/09 W</div> </div>	<div>Slide 1 Chapter 12: Strategies for Variable Selection (Class 2 of 2)</div> <div>NOTES:</div>
<div> <div>HW 12 due Friday 4/24/09</div> <div>Submit as Myname-HW12.doc (or *.rtf)</div> <ul style="list-style-type: none"> <li>WIMBA sessions: <ul style="list-style-type: none"> <li>Tonight &amp; every Weds 10-11 pm</li> <li>Thursday Noon - 1 pm (log on from anywhere)</li> </ul> </li> <li>New Homework due dates <ul style="list-style-type: none"> <li>HW 12 10:28: El Niño and Hurricanes</li> <li>Due <b>Friday 4/24/09</b> Noon</li> <li>Note: There will 2 WIMBA sessions available on this topic</li> </ul> </li> <li>HW 13 Cammen's ingestion rate data. Note that this was a 2003 final exam problem <ul style="list-style-type: none"> <li>Read Cammen (1980) &amp; evaluate his regression model</li> <li>Due <b>Weds 4/29/09</b> Noon This problem will count double!</li> </ul> </li> <li>Read Chapter 12: Selection of variables</li> <li>Run my overfitting syntax: overfitting.sps</li> <li>Read Campbell &amp; Kenney Chapters 4 &amp; 5 on the regression artefact and gender inequities</li> <li>Run my Campbell &amp; Kenny syntax: RTMCK.sps</li> </ul> </div>	<div>Slide 2 HW 12 due Friday 4/24/09</div> <div>NOTES:</div>
<div> <div>HW12: Cammen model</div> <p>Cammen (1980) compiled data from the literature on the ingestion rates of 22 deposit feeders. Deposit feeders are organisms that live in mud and sand and ingest mud and sand. Deposit feeders use the organic matter in the mud and sand for growth. Table 1 shows the species from the literature, their ingestion rates, the fraction organic matter in sediment, and the body weights of individual deposit feeders. Cammen (1980) used regression to estimate the ingestion rate of deposit feeders (<b>ING</b>) (mg dry weight/day) using the fraction organic matter in the sediment (<b>OM</b>) and body weight of the deposit feeder (<b>WT</b>). He regressed <math>\log_{10}</math> (<b>ING</b>) as the response variable with two explanatory variables <math>\log_{10}</math> (<b>WT</b>) and <math>\log_{10}</math> (<b>OM</b>). He deleted the three bivalves from his analyses because they appeared to be outliers, and based his regressions on the 19 non-bivalve species.</p> </div>	<div>Slide 3 HW12: Cammen model</div> <div>NOTES:</div>

Table 1. Data from **Cammen (1980)**. Loaded on **Fromotheta** as **cammen.csv**, in case you wanted to examine the data (optional). The last 5 highlighted species are bivalve molluscs (indicated under **Taxon**). **WT** is the body weight of the deposit feeder (dry weight of the animal) in milligrams. **ING** is the ingestion rate in mg dry weight/day. Cammen scaled the ingestion rate to account for temperature effects (higher ingestion at higher temperatures). **CPI** is the organic matter content (% weight organic matter / % total sediment dry weight), expressed as %.

Species	Taxon	WT	ING	CPI
1. <i>Nereis virens</i>	Gastropod mollusc	0.2	0.57	16
2. <i>Nereis virens</i>	Gastropod mollusc	0.2	0.66	17
3. <i>Tubificoides</i>	Oligochaete (annelid)	0.27	0.46	29.7
4. <i>Nereis virens</i>	Crustacean	0.32	0.46	50
5. <i>Phoronopsis</i>	Gastropod mollusc	0.46	2.7	14.4
6. <i>Nereis virens</i>	Gastropod mollusc	0.9	0.67	13
7. <i>Nereis virens</i>	Polychaete (annelid)	5.8	20.2	6.8
8. <i>Phoronopsis</i>	Crustacean (annelid)	8.4	1.46	50
9. <i>Chironomus</i>	Crustacean	12.4	4.4	88
10. <i>Alvinicella</i>	Polychaete (annelid)	20.4	24.0	2.2
11. <i>Tricorophus</i>	Polychaete (annelid)	46	25.0	1
12. <i>Ampelisca</i>	Crustacean	53	30.0	4.2
13. <i>Uca</i>	Crustacean	63.9	19.9	51
14. <i>Scapharca</i>	Crustacean	65	50	23.6
15. <i>Phoronopsis</i>	Polychaete (annelid)	80	10.07	0.7
16. <i>Alvinicella</i>	Polychaete (annelid)	280	34.00	1.2
17. <i>Alvinicella</i>	Polychaete (annelid)	380	34.00	0.4
18. <i>Alvinicella</i>	Polychaete (annelid)	880	47.00	0.64
19. <i>Alvinicella</i>	Crustacean	2050	46.00	2.1
20. <i>Macoma</i>	Bivalve mollusc	5.1	4.46	20
21. <i>Macoma</i>	Bivalve mollusc	19.9	3.94	6.8
22. <i>Macoma</i>	Bivalve mollusc	280	4.3	5.4

## Slide 4

NOTES:

## HW12: Cammen model

- Answer each question and address each issue.**
- Was Cammen (1980) justified in dropping the three bivalve molluscs from his regression equation?
  - Consider both the case-wise diagnostic tests (residuals vs. predicted values, Cook's D, studentized residuals, and leverage values), and the results of fitting bivalves as a dummy variable.
  - Discuss the problems in using Cook's D, leverage, and studentized residuals in detecting outliers when more than one datum may be an outlier.
  - There is no strictly right or wrong answer to this question, but you must justify your choice with evidence from the regression analyses.
  - There were 5 groups of animals in Cammen's data. Is there evidence that the ingestion rates as a function of weight and organic matter differ among these 5 groups?
  - Based on your analyses, produce a graph showing the relationship between ingestion rate, body weight and organic matter.
  - Write the regression equation expressing the relationship between ingestion rate, organic matter, and body weight. Pay attention to significant figures, and include an estimate of the standard error of the coefficients.
  - If you found that the animal groups differed in ingestion rate, your final graphs and model should reflect this full model

## Slide 5 HW12: Cammen model

NOTES:

## Homework Presentations

- William Walker for HW 8
- Steven Kichefski for HW 9 and
- Lisa Greber for HW10



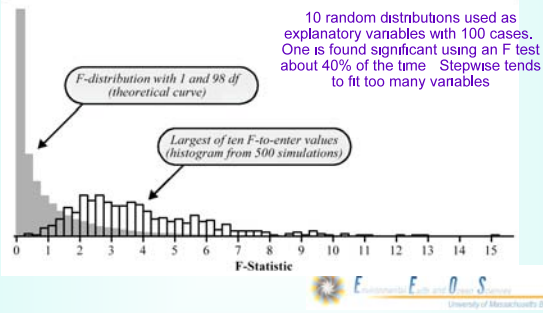
## Slide 6 Homework Presentations

NOTES:

<p><b>Chapter 12: Strategies for variable selection (continued)</b></p>	<p><b>Slide 7 Chapter 12: Strategies for variable selection (continued)</b></p> <p>NOTES:</p>
<p><b>Overfitting: why stepwise procedures should not be used to estimate p values.</b></p>	<p><b>Slide 8 Overfitting: why stepwise procedures should not be used to estimate p values.</b></p> <p>NOTES:</p>
<p><b>Covariates: overfitting &amp; multicollinearity</b></p> <ul style="list-style-type: none"> <li>• Overfitting.sps <ul style="list-style-type: none"> <li>▸ 32 random variables, 100 cases</li> <li>▸ Stepwise, forward &amp; backward regression will usually always find a significant regression</li> <li>▸ One solution: use 40 times as many cases as covariates (&gt;1200 for a 32-variable model!)</li> </ul> </li> <li>• More guns, less crime <ul style="list-style-type: none"> <li>▸ Ayres, I and JJ Donohue (2003) Shooting down the more guns, less crime hypothesis. Stanford Law Review.</li> <li>▸ Including many covariates, many correlated with the key explanatory variable (gun control laws) produces an artifact, showing an effect when none existed</li> </ul> </li> <li>• Peterson's voucher studies <ul style="list-style-type: none"> <li>▸ Kreuger critique: Including pre-test scores as a covariate produces an effect when none existed</li> </ul> </li> </ul>	<p><b>Slide 9 Covariates: overfitting &amp; multicollinearity</b></p> <p>NOTES:</p>

Display 12.7

## Simulated distribution of the largest of ten F-statistics



## Slide 10

## NOTES:

## Gallagher's overfitting.sps

\* Overfitting simulation, inspired by  
 \* Nontechnical Introduction to Overfitting in Regression-Type Models, Babyak (2004).  
 \* Michael A Babyak What You See May Not Be What You Get: A Brief, Nontechnical  
 \* Introduction to Overfitting in Regression-Type Models.  
 \* Psychosom Med 2004 66: 411-421.  
 \* Written by E Gallagher, revised 4/12/05.  
 \* Generate 100 cases, with 32 normally distributed variates.  
 new file.  
 input program.  
 loop # = 1 to 100.  
 COMPUTE V1 = RV.normal (0,1) .  
 COMPUTE V2 = RV.normal (0,1) .  
 ..  
 COMPUTE V32 = RV.normal (0,1) .  
 end case.  
 end loop.  
 end file.  
 end input program.  
 formats V1 to V32 (f4,2).  
 exe.

## Slide 11 Gallagher's overfitting.sps

## NOTES:

## Results of Stepwise Selection

## 31 Random predictor variables

Model	Coefficients <sup>a</sup>				Standardized Coefficients <sup>a</sup>				95% Confidence Interval for B			
	B	Std. Error	t	Sig.	B	Std. Error	t	Sig.	Lower Bound	Upper Bound	Lower Bound	Upper Bound
1												
(Constant)	275	894										
V1	301	893	.311	.354								
2												
(Constant)	275	893										
V1	328	893	.365	.354								
V2	154	893	.171	.087								
3												
(Constant)	275	893										
V1	301	894	.311	.354								
V2	154	893	.171	.087								
V3	188	893	.210	.034								
4												
(Constant)	275	893										
V1	411	893	.458	.000								
V2	154	893	.171	.087								
V3	188	893	.210	.034								
V4	227	893	.254	.011								
5												
(Constant)	275	893										
V1	420	893	.468	.000								
V2	154	893	.171	.087								
V3	188	893	.210	.034								
V4	227	893	.254	.011								
V5	188	893	.210	.034								

a. Dependent Variable: Y1

Backward (added V23, V19)	25	(Constant)	268	893								
	V1	154	893	.171	.087							
	V2	154	893	.171	.087							
	V3	188	893	.210	.034							
	V4	227	893	.254	.011							
	V5	188	893	.210	.034							
	V6	188	893	.210	.034							
	V7	188	893	.210	.034							
	V8	188	893	.210	.034							
	V9	188	893	.210	.034							
	V10	188	893	.210	.034							
	V11	188	893	.210	.034							
	V12	188	893	.210	.034							
	V13	188	893	.210	.034							
	V14	188	893	.210	.034							
	V15	188	893	.210	.034							
	V16	188	893	.210	.034							
	V17	188	893	.210	.034							
	V18	188	893	.210	.034							
	V19	188	893	.210	.034							
	V20	188	893	.210	.034							
	V21	188	893	.210	.034							
	V22	188	893	.210	.034							
	V23	188	893	.210	.034							
	V24	188	893	.210	.034							
	V25	188	893	.210	.034							
	V26	188	893	.210	.034							
	V27	188	893	.210	.034							
	V28	188	893	.210	.034							
	V29	188	893	.210	.034							
	V30	188	893	.210	.034							
	V31	188	893	.210	.034							
	V32	188	893	.210	.034							

a. Dependent Variable: Y1

## Slide 12 Results of Stepwise Selection

## NOTES:

**Harrell (2002, p. 56-57) on stepwise****Harrell's conclusion: Don't use stepwise!**

- It yields  $R^2$  values that are biased high
- $F$  and  $\chi^2$  distributions don't have their claimed distributions
- SE of regression coefficients are biased low and CI's and predicted values that are falsely narrow
- P-values too small
- Regression coefficients biased high in absolute value and need shrinkage.
- Rather than solving the problem of collinearity, variable selection is made arbitrary by collinearity
- It allows us not to think about the problem

**Slide 13 Harrell (2002, p. 56-57) on stepwise**

NOTES:

**Overfitting: too many covariates****Harrell (2001, p. 60)**

"When a model is fitted that is too complex, that is it has too many free parameters to estimate for the amount of information in the data, the worth of the model (e.g.,  $R^2$ ) will be exaggerated and future observed values will not agree with predicted values. In this situation **overfitting** is said to be present, and some of the findings of the analysis come from fitting noise or finding spurious associations between X and Y"

**Slide 14 Overfitting: too many covariates**

NOTES:

**Number of cases needed for regression (1 of 2)****Harrell (2001, p. 61)**

- Number of predictors should be less than  $m/10$  or  $m/20$  where  $m$  is the limiting sample size shown below
- Candidate variables must include all variables screened for association with response, including nonlinear terms and interactions

TABLE 4.1: Limiting Sample Sizes for Various Response Variables

Type of Response Variable	Limiting Sample Size $m$
Continuous	$n$ (total sample size)
Binary	$\min(n_1, n_2)^c$
Ordinal ( $k$ categories)	$n - \frac{1}{n^2} \sum_{i=1}^k n_i^3$ <sup>d</sup>
Failure (survival) time	number of failures <sup>e</sup>

**Slide 15 Number of cases needed for regression (1 of 2)**

NOTES:

<div data-bbox="293 163 738 231"> <h3>Number of cases for regression (2 of 2)</h3> </div> <div data-bbox="365 235 646 260"> <p>Tabachnik &amp; Fidell (2001, p 117)</p> </div> <ul style="list-style-type: none"> <li>For multiple regression (from Green 1991) <ul style="list-style-type: none"> <li><math>N &gt; 50 + 8m</math>, where <math>m</math> is the number of explanatory variables, for testing <math>R^2</math>, and</li> <li><math>N \geq 104 + m</math> for individual predictors</li> <li>A higher case to explanatory variable ratio is needed when <ul style="list-style-type: none"> <li>Effect sizes are small</li> <li>Data are skewed</li> <li>Measurement error is expected in explanatory variables</li> </ul> </li> <li>Automated selection procedures (statistical regression) <ul style="list-style-type: none"> <li>Cases <math>&gt; 40 \times</math> explanatory variables</li> </ul> </li> <li>Green's more precise rule <ul style="list-style-type: none"> <li><math>N \geq (8 / f^2) + (m-1)</math>, where <math>f^2 = 0.01, 0.15</math>, and <math>0.35</math> for small, medium and large effect sizes.</li> <li><math>f^2 = R^2 / (1-R^2)</math>, where <math>R^2</math> is the expected squared multiple correlation coefficient</li> </ul> </li> </ul> </li> </ul>	<div data-bbox="815 132 1373 170"> <h3>Slide 16 Number of cases for regression</h3> </div> <div data-bbox="815 195 922 231"> <p>(2 of 2)</p> </div> <div data-bbox="815 317 938 352"> <p>NOTES:</p> </div>
<div data-bbox="318 657 719 695"> <h3>Multicollinearity, collinearity</h3> </div> <ul style="list-style-type: none"> <li>If the explanatory variables are strongly correlated <ul style="list-style-type: none"> <li>The regression coefficient estimates have a huge variance</li> <li>They can change in sign and significance with a slight change in the data, bouncing betas</li> </ul> </li> <li>Assessed with Variance inflation factors (VIF) or tolerance <ul style="list-style-type: none"> <li><math>VIF_i = 1 / (1 - R_i^2)</math>, where <math>R_i^2</math> is the squared multiple correlation coefficient between explanatory variable 'i' and the other explanatory variables</li> <li>Neter et al. (1996): VIF's <math>&gt; 10</math> are cause for concern (but smaller VIF's can also be a problem)</li> <li>Marayuma (1998): VIF <math>&gt; 6</math> or <math>7</math>, as a very rough rule, indicate strong multicollinearity</li> </ul> </li> </ul>	<div data-bbox="815 621 1354 659"> <h3>Slide 17 Multicollinearity, collinearity</h3> </div> <div data-bbox="815 743 938 779"> <p>NOTES:</p> </div>
<div data-bbox="276 1146 764 1182"> <h3>Ways of detecting multicollinearity</h3> </div> <div data-bbox="397 1188 605 1213"> <p>Marayuma (1998, p. 64)</p> </div> <ul style="list-style-type: none"> <li>When the variance (standard errors) of beta weights is large</li> <li>When signs on beta weights are inappropriate [e.g., larger classes <math>\Rightarrow</math> higher test scores]</li> <li>When regression weights and signs change radically upon the addition or removal of single variables</li> <li>When the Variance Inflation Factor is high (VIF <math>&gt; 6</math> or <math>7</math> as a very rough rule)</li> <li>When simple correlations are <math>&gt; 0.8-0.9</math></li> <li>When correlations among predictor variables <math>&gt; R^2</math> for response with all predictor variables</li> </ul>	<div data-bbox="815 1110 1195 1184"> <h3>Slide 18 Ways of detecting multicollinearity</h3> </div> <div data-bbox="815 1268 938 1304"> <p>NOTES:</p> </div>

## Solutions to multicollinearity

- If the goal of the model is to produce predicted values for one analysis, then multicollinearity is **not** a problem. All variables can be included.
  - ▶ However, if the equation is to be used for new data, then the model will be badly overfitted, the predicted values will be biased
  - ▶ Significant coefficients could be spurious or nonsense
- Solutions
  - ▶ Reduce the number of explanatory variables using theory & insight into the field
  - ▶ Cluster analysis of variables: Choose 1 from each cluster
  - ▶ Ridge regression (available using syntax for SPSS - Raynald Lavasque's web site)
  - ▶ Principal components regression
    - Principal component scores are usually orthogonal (uncorrelated)
    - Use principal component scores as explanatory variables
  - ▶ Structural equation modeling

## Slide 19 Solutions to multicollinearity

NOTES:

## Ridge regression

Available as a macro in SPSS, LISREL (not AMOS); increase variance for variables not covariance

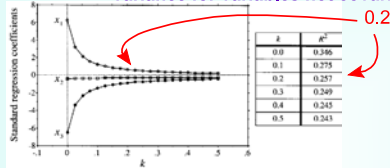


Figure 19.8 Ridge trace diagram showing the estimates of the standardized regression coefficients  $b_j$  as a function of  $k$ . Table: decrease of  $R^2$  as a function of  $k$ .

A ridge regression parameter,  $k$ , is chosen using the ridge trace diagram ( $k=0.2$  in the above example [the base of the horn] from Draper & Smith) that 'shrinks' the regression coefficients, especially those coefficients (Beta's) that are strongly correlated. This offers a partial solution to the problem of collinearity.

## Slide 20 Ridge regression

NOTES:

## Number of cases needed for regression (1 of 2)

Harrell (2001, p. 61)

- Number of predictors should be less than  $m/10$  or  $m/20$  where  $m$  is the limiting sample size shown below
- Candidate variables must include all variables screened for association with response, including nonlinear terms and interactions

TABLE 4.1: Limiting Sample Sizes for Various Response Variables

Type of Response Variable	Limiting Sample Size $m$
Continuous	$n$ (total sample size)
Binary	$\min(n_1, n_2)$ <sup>c</sup>
Ordinal ( $k$ categories)	$n - \frac{1}{k^2} \sum_{i=1}^k n_i^2$ <sup>d</sup>
Failure (survival) time	number of failures <sup>e</sup>

## Slide 21 Number of cases needed for regression (1 of 2)

NOTES:

## Number of cases for regression (2 of 2)

Tabachnik & Fidell (2001, p 117)

- For multiple regression (from Green 1991)
  - $N > 50 + 8m$ , where  $m$  is the number of explanatory variables, for testing  $R^2$ , and
  - $N \geq 104 + m$  for individual predictors
- A higher case to explanatory variable ratio is needed when
  - Effect sizes are small
  - Data are skewed
  - Measurement error is expected in explanatory variables
- Automated selection procedures (statistical regression)
  - Cases  $> 40 \times$  explanatory variables
- Green's more precise rule
  - $N \geq (8 / f^2) + (m-1)$ , where  $f^2 = 0.01, 0.15$ , and  $0.35$  for small, medium and large effect sizes.
  - $f^2 = R^2 / (1-R^2)$ , where  $R^2$  is the expected squared multiple correlation coefficient

## Slide 22 Number of cases for regression

(2 of 2)

NOTES:

## Ayres & Donohue (2003): Too many covariates produces less crime

Lott used 36 demographic covariates, severe collinearity problems

- Lott & Mustard (1997) argue lenient 'will carry' gun law states had less crime
- L&M used 36 demographic variables in their regressions
- The excessive number of covariates produced
  - Multicollinearity effects, changing the sign of the crime terms
  - Note: the sign of a term in a multiple regression is a partial correlation, given the other terms. The sign can change depending on other terms.

1290

STANFORD LAW REVIEW

[Vol. 55:1193]

The question merits investigation: Why do the results of the Lott model (Table 2) support the Lott thesis more than the results of the Zheng model (Table 2)? The reason turns out to be somewhat surprising. As Table 2 documents, the two models include some different explanatory variables that one might think could have important implications. For example, Lott controls for population density and transfer payments to the poor, while Zheng controls for police, poverty, unemployment, and alcohol consumption. But these differences in the substantive controls turn out to be largely unimportant. Interestingly, as we will discuss in the next section, what drives the entire difference between Tables 2 and 3 is that Lott includes a large number of potentially duplicative demographic variables. Indeed, the entry is so extensive as to make multicollinearity a serious issue.<sup>62</sup> Specifically, while Zheng's model controls for percent black and three age groupings, Lott's has thirty-six separate demographic percentages, breaking down each of three difference race categories—black, white, and neither black nor white—and three sexes into six separate age categories from age ten up.<sup>63</sup> The sensitivity of the results to the inclusion or exclusion of an array of nearly collinear demographic variables serves as a cautionary tale to those who conduct or rely upon panel data models of crime. Probably no one examining either Weisberg Zheng's work or that of Lott and Mustard would suspect that conclusions reached from their models would be sensitive to these seemingly second-order demographic controls.

## Slide 23 Ayres & Donohue (2003): Too many covariates produces less crime

NOTES:

## Shooting Down the "More Guns, Less Crime" Hypothesis

Ian Ayres\* & John J. Donohue III

1290

STANFORD LAW REVIEW

[Vol. 55:1193]

The question merits investigation: Why do the results of the Lott model (Table 2) support the Lott thesis more than the results of the Zheng model (Table 2)? The reason turns out to be somewhat surprising. As Table 2 documents, the two models include some different explanatory variables that one might think could have important implications. For example, Lott controls for population density and transfer payments to the poor, while Zheng controls for police, poverty, unemployment, and alcohol consumption. But these differences in the substantive controls turn out to be largely unimportant. Interestingly, as we will discuss in the next section, what drives the entire difference between Tables 2 and 3 is that Lott includes a large number of potentially duplicative demographic variables. Indeed, the entry is so extensive as to make multicollinearity a serious issue.<sup>62</sup> Specifically, while Zheng's model controls for percent black and three age groupings, Lott's has thirty-six separate demographic percentages, breaking down each of three difference race categories—black, white, and neither black nor white—and three sexes into six separate age categories from age ten up.<sup>63</sup> The sensitivity of the results to the inclusion or exclusion of an array of nearly collinear demographic variables serves as a cautionary tale to those who conduct or rely upon panel data models of crime. Probably no one examining either Weisberg Zheng's work or that of Lott and Mustard would suspect that conclusions reached from their models would be sensitive to these seemingly second-order demographic controls.

crime from the Table 3 regressions, one finds utterly bizarre results. For example, an increase of one percentage point in the percentage of black males aged 30-39 would be expected to almost double the violent crime rate, while a similar increase in the percentage of black males aged 40-49 would lead to a drop in violent crime of 66%. Similarly, increasing the percentage of black males aged 50-64 would cause violent crime to jump by 145%, but increasing the percentage of black males over age 65 would lead to a 78% decline in violent crime. These nonsense results prevent us from understanding why the demographic controls can influence the estimates of shall-issue adoption so strongly.

Adding too many covariates can destroy a regression

## Slide 24

NOTES:



## Case 11.2 Gender discrimination

## Slide 25 Case 11.2 Gender discrimination

NOTES:

[illegible]

Is there evidence for sex discrimination AFTER age, education and experience are 'accounted for'?

Note, that Sleuth's approach is subject to 'the regression artifact' (Campbell & Kenny 1999)

## Slide 26

NOTES:

**Display 12.9**

<i>Main Effect Variables</i>	<i>Quadratic Variables</i>	<i>Interaction Variables</i>
s = seniority	t = s <sup>2</sup>	m = s × a    c = a × e
a = age	b = a <sup>2</sup>	n = s × e    k = a × x
e = education	f = e <sup>2</sup>	v = s × x    q = e × x
x = experience	y = x <sup>2</sup>	



## Slide 27

NOTES:

## Slide 28

Bayesian posterior analysis of the difference between male and female log-beginning salaries

Model	p	BIC	posterior probability	Addition of sex indicator		
				coeff	SE	1-sided p-value
saesck	7	-401.40	.7709	-.1196	.0229	6.27E-7
saesvc	7	-398.89	.0625	-.1287	.0236	8.42E-8
saesckq	7	-398.28	.0340	-.1244	.0221	1.18E-7
saesckc	8	-398.08	.0279	-.1173	.0229	9.48E-7
saesvc	6	-397.81	.0213	-.1247	.0238	5.59E-7
saesck	6	-397.51	.0157	-.1135	.0246	6.94E-6
saesckb	8	-396.49	.0057	-.1195	.0229	6.70E-7
saesckc	8	-396.37	.0051	-.1189	.0232	9.10E-7
saesckqb	8	-396.36	.0050	-.1206	.0221	2.41E-7
saesvcn	8	-396.33	.0048	-.1258	.0225	1.37E-7
saesck	6	-396.26	.0045	-.1331	.0221	1.96E-8
saesvcq	6	-396.15	.0040	-.1345	.0201	1.02E-9
saesckf	8	-396.12	.0039	-.1196	.0230	6.93E-7
saesckq	8	-396.05	.0037	-.1208	.0230	5.54E-7
saesvc	5	-395.93	.0032	-.1302	.0211	1.11E-8
saesck	8	-395.91	.0032	-.1257	.0232	2.81E-7
saesvcq	7	-398.89	.0031	-.1328	.0218	1.51E-8
saesckm	8	-395.84	.0030	-.1195	.0231	7.46E-7
saesckv	8	-395.80	.0028	-.1196	.0231	7.31E-7
saesckc	7	-395.20	.0016	-.1230	.0237	6.95E-7

s = seniority  
a = age  
e = education  
x = experience

t = s<sup>2</sup>  
b = s<sup>2</sup>  
f = e<sup>2</sup>  
y = x<sup>2</sup>

m = s \* x \* b  
v = s \* x \* e  
n = s \* x \* x  
q = e \* x \* x

Sleuth (p 343):  
'There is convincing evidence that the median starting salary for females was lower than the median starting salary for males, even after the effects of age, education, previous experience, and time at which the job began are taken into account (1-sided p-value < 0.0001)'

## NOTES:

## SPSS output using forward, backward or stepwise

Model Summary<sup>a</sup>

Model	Akaike Information Criterion	Amemiya Prediction Criterion	Selection Criteria		Schwarz Bayesian Criterion
			Prediction	Prediction	
1	-395.813 <sup>a</sup>	.858	38.600		-390.747
2	-407.042 <sup>b</sup>	.761	23.681		-399.444
3	-410.713 <sup>c</sup>	.731	19.134		-400.582
4	-415.957 <sup>d</sup>	.691	13.330		-403.294
5	-419.552 <sup>e</sup>	.665	9.706		-404.356
6	-421.538 <sup>f</sup>	.651	7.716		-408.876
7	-427.248 <sup>g</sup>	.612	2.501		-412.053

a. Predictors: (Constant), f (e<sup>2</sup>)b. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e)c. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e), v (s \* x)d. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e), v (s \* x), k (a \* x)e. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e), v (s \* x), k (a \* x), x (Experience)f. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e), k (a \* x), x (Experience)g. Predictors: (Constant), f (e<sup>2</sup>), n (s \* e), k (a \* x), x (Experience), q (e \* x)

h. Dependent Variable: ln (Salary)

## Slide 29 SPSS output using forward, backward or stepwise

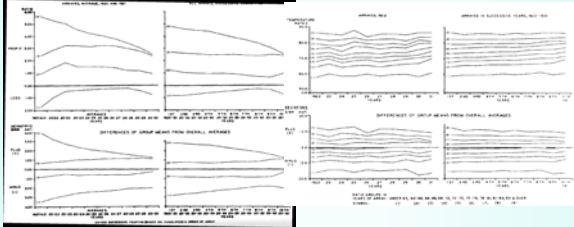

## NOTES:

## Slide 30 Has gender equity really been rejected?

## Has gender equity really been rejected?

Campbell & Kenny: statistical equating often produces gender discrimination when there is none, and racial differences when there are none

## NOTES:

<div data-bbox="228 163 776 212"> <h3>Statistical Equating &amp; RTM</h3> </div> <div data-bbox="228 216 776 562"> <p><b>Campbell &amp; Kenny: The regression artifact</b></p> <ul style="list-style-type: none"> <li>• The sophomore jinx</li> <li>• Spontaneous remission of depression</li> <li>• Misclassification of individuals using standardized tests</li> <li>• Perhaps:             <ul style="list-style-type: none"> <li>▸ Ashland cancer study</li> <li>▸ Washington D.C. vouchers</li> <li>▸ Sander's analysis of African-American failure on the bar exam</li> </ul> </li> <li>• Statistical equating             <ul style="list-style-type: none"> <li>▸ Regression to the mean leads to a bias in estimating gender differences using "equating"</li> <li>▸ Page 84: Ethnic differences in intellectual ability:                 <ul style="list-style-type: none"> <li>■ "We believe that the bias in statistical equating for ethnic differences in achievement and intelligence testing is underadjustment"</li> </ul> </li> </ul> </li> </ul> </div>	<div data-bbox="821 134 1416 182"> <h3>Slide 31 Statistical Equating &amp; RTM</h3> </div> <div data-bbox="821 258 1416 306"> <p>NOTES:</p> </div>
<div data-bbox="228 657 776 705"> <h3>Poor Horace Secrist (1933)</h3> </div> <div data-bbox="228 709 776 772"> <p>Identify companies that had lower than average profits and invest in them; he was aware of RTM Profits (left), temperature (right)</p> </div> <div data-bbox="228 777 776 1035">  </div>	<div data-bbox="821 627 1416 676"> <h3>Slide 32 Poor Horace Secrist (1933)</h3> </div> <div data-bbox="821 751 1416 800"> <p>NOTES:</p> </div>
<div data-bbox="228 1150 776 1199"> <h3>Hotelling's (1933) JASA review</h3> </div> <div data-bbox="228 1220 776 1402"> <ul style="list-style-type: none"> <li>• Business varies, but average temperatures don't vary nearly as much             <ul style="list-style-type: none"> <li>▸ Secrist chose cities spread out throughout the country and looked at interannual variability</li> <li>▸ Small year-to-year variations compared to the big city-to-city variations</li> </ul> </li> <li>• Secrist rebuttal (1934)</li> </ul> </div> <div data-bbox="532 1476 787 1524">  </div>	<div data-bbox="821 1121 1416 1169"> <h3>Slide 33 Hotelling's (1933) JASA review</h3> </div> <div data-bbox="821 1245 1416 1293"> <p>NOTES:</p> </div>

### Hotelling's (1934) rejoinder

Quoted in Stigler's "Statistics on the Table"

"To 'prove' such a mathematical result [regression to the mean in annual reports] by a costly and prolonged numerical study of many kinds of business profit and expense ratios is analogous to proving the multiplication table by arranging elephants in rows and columns, and then doing the same for numerous other kinds of animals. The performance, though perhaps entertaining, and having a certain pedagogical value, is not an important contribution to either zoology or to mathematics."

### Slide 34 Hotelling's (1934) rejoinder

NOTES:

### Statistical Equating

Effects on gender bias & racial differences

**"Including a covariate, like socioeconomic status, can produce a racial or gender bias, when none really exists!"**

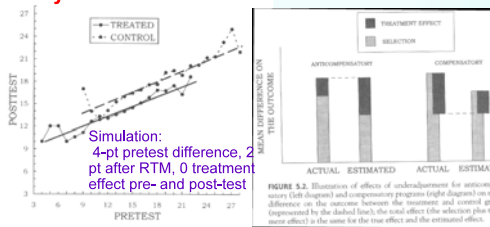


FIGURE 5.1. Parallel regression lines for treatment and control groups adjusted by statistical equating (multiple regression).

### Slide 35 Statistical Equating

NOTES:

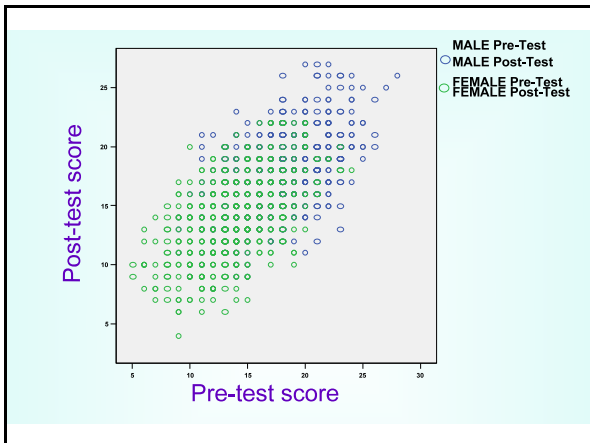
### A hypothetical test of gender effects

Read Campbell & Kenny Chapters 4 & 5

- Are women inferior in mathematics?
- Randomly select 500 women & 500 men for admission to an intense workshop on advanced mathematics.
- Give both groups a pretest of mathematical ability
  - In the simulation (rtm-ck.sps) generate test scores by 4 tosses of a die. Assign males 4 units higher score in both pre & post test
    - Males: sum of 4 dice + 4
    - Females: sum of 4 dice + 0.
- Assume that the workshop does NOTHING to improve ability for either group
- Retest each student, the post-test, which is modeled to have a correlation of 0.5 between pre- & post-test
  - 2 dice the same, 2 new dice throws for each student
- Test whether males did better than females in this advanced workshop, even after controlling for their previous math background

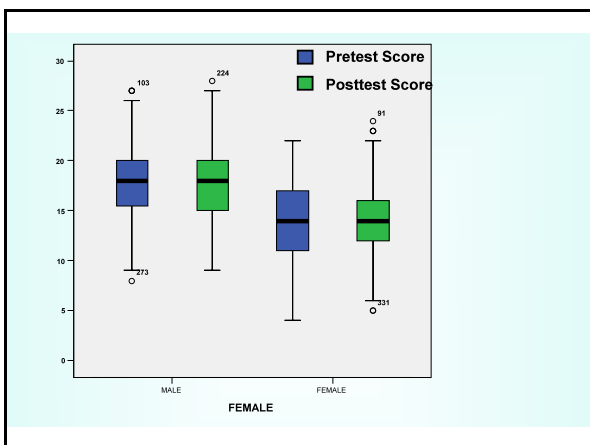
### Slide 36 A hypothetical test of gender effects

NOTES:



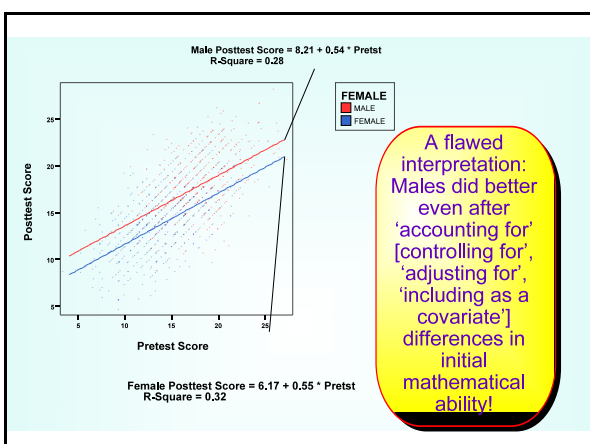
### Slide 37

NOTES:



### Slide 38

NOTES:



### Slide 39

NOTES:

### Flawed interpretation: Females score 2 points less ( $1.9 \pm 0.4$ ) on the post-test, after 'supposedly' controlling for the effect of previous mathematical ability ( $p < 10^{-18}$ )

Model		Unstandardized Coefficients		Standardized Coefficients		t	Sig.	95% Confidence Interval for B	
		B	Std. Error	Beta				Lower Bound	Upper Bound
1	(Constant)	8.146	.486			16.777	7.8E-056	7.194	9.099
	Pretest Score	.545	.026	.542		20.736	1.1E-079	.493	.596
	FEMALE	-1.927	.209	-.240		-9.198	2.1E-019	-2.338	-1.516

a. Dependent Variable: Posttest Score

But: the simulation is set so that the workshop didn't have any effect on either group!

**Slide 40 Flawed interpretation: Females score 2 points less ( $1.9 \pm 0.4$ ) on the post-test, after 'supposedly' controlling for the effect of previous mathematical ability ( $p < 10^{-18}$ )**

NOTES:

### Classic Analysis of covariance

Huge Male-female difference in post-workshop scores, after 'controlling' for pre-test ability

- \* Classic analysis of covariance (ANCOVA)
  - \* to test for treatment effect
  - \* with pretest as the covariate.
- ANOVA postst BY treat(0,1) with pretst  
/STATISTICS=ALL.

		ANOVA <sup>a, b</sup>					
		Sum of Squares	df	Mean Square	F	Sig.	B
Posttest Score	Covariates	3630.112	1	3630.112	429.995	1.08E-079	.545
	Main Effects	714.243	1	714.243	84.604	2.08E-019	
	Model	7654.148	2	3827.074	453.326	9.44E-141	
	Residual	8416.888	997	8.442			
	Total	16071.036	999	16.087			

a. Posttest Score by FEMALE with Pretest Score

b. All effects entered simultaneously

**Slide 41 Classic Analysis of covariance**

NOTES:

### Repeated measures designs (Chapter 16) produce the correct solution: No effect of gender on post-test

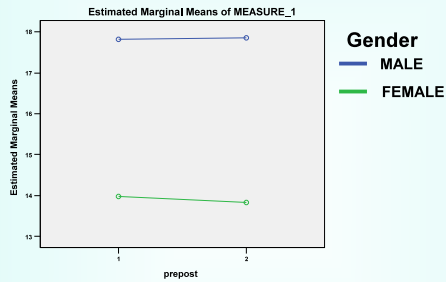
There is no pre-test to post-test x gender interaction

		Tests of Within-Subjects Effects				
Measure: MEASURE_1		Type III Sum of Squares	df	Mean Square	F	Sig.
Source	Sphericity Assumed	1.458	1	1.458	.266	.606
	Greenhouse-Geisser	1.458	1.000	1.458	.266	.606
	Huynh-Feldt	1.458	1.000	1.458	.266	.606
	Lower-bound	1.458	1.000	1.458	.266	.606
prepost * treat	Sphericity Assumed	4.232	1	4.232	.771	.380
	Greenhouse-Geisser	4.232	1.000	4.232	.771	.380
	Huynh-Feldt	4.232	1.000	4.232	.771	.380
	Lower-bound	4.232	1.000	4.232	.771	.380
Error(prepost)	Sphericity Assumed	5476.310	998	5.487		
	Greenhouse-Geisser	5476.310	998.000	5.487		
	Huynh-Feldt	5476.310	998.000	5.487		
	Lower-bound	5476.310	998.000	5.487		

**Slide 42 Repeated measures designs (Chapter 16) produce the correct solution: No effect of gender on post-test**

NOTES:

### Profiles from Repeated Measures ANOVA



### Slide 43 Profiles from Repeated Measures ANOVA

NOTES:

### Change score: Do paired t tests on males & females separately

Paired Samples Test									
Paired Differences									
	Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference		t	df	Sig. (2-tailed)	
				Lower	Upper				
Pair 1 Pretest Score - Posttest Score	-.038	3.365	.150	-.334	.258	-.253	499		.801

Paired Samples Test									
Paired Differences									
	Mean	Std. Deviation	Std. Error Mean	95% Confidence Interval of the Difference		t	df	Sig. (2-tailed)	
				Lower	Upper				
Pair 1 Pretest Score - Posttest Score	.146	3.260	.146	-.140	.432	1.002	499		.317



### Slide 44 Change score: Do paired t tests on males & females separately

NOTES:

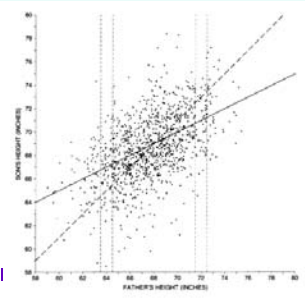
### Why didn't regression & ANCOVA work?

See Cambell & Kenny (Ch 4-5) for full analysis

- Whenever there is less than perfect correlation between the covariate and the response, the effect of the covariate on the response is **not** removed by regression (=Analysis of covariance)
- This is due to regression to the mean
- Since the correlation between pre-test and post-test was set at  $r=0.5$ , only 50% of the pre-test effect can be 'explained' or accounted for by multiple regression
- Whenever the covariate is less than perfectly correlated with the response, multiple regression does not fully 'control for' or 'account for' or 'adjust for' the effects of the covariate.
  - ▶ Note that if the pre-test score had a correlation with the post-test score of 0.25, then only 1/4 of the pre-test difference would be accounted for by including pre-test as a covariate. There would be a 3-point advantage for males after including pre-test as a covariate

### Slide 45 Why didn't regression & ANCOVA work?

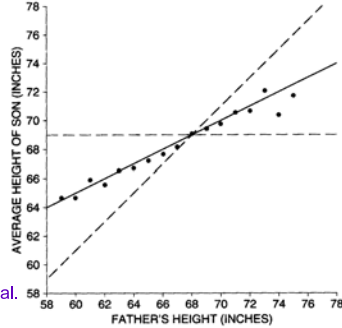
NOTES:

**Galton's regression to the mean**Son's height 1" taller than father's,  $r=0.5$ ,  $SD=2.5$ "Figure from  
Freedman et al**Slide 46 Galton's regression to the mean**

NOTES:

**RTM effect  $\propto 1/r$** 

From Freedman et al

Figure from  
Freedman et al.**Slide 47 RTM effect  $\propto 1/r$** 

NOTES:

**Galton squeeze**

If you naively use pretest as a covariate, you'll introduce an artifact in the analysis.

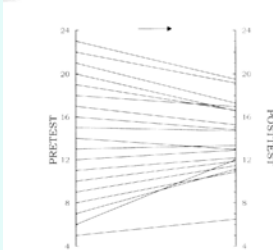


FIGURE 1.8. Galton-squeeze diagram for the data set with 500 cases using pretest to predict posttest.

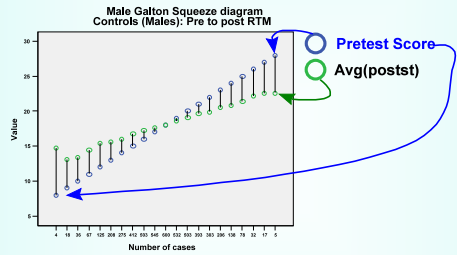
Using pre-test to predict post-test will be subject to 'regression to the mean.' If  $r$  between pre- and post is 0.5, only half of the pre-test [gender] effect will be accounted for.**Slide 48 Galton squeeze**

NOTES:



### Galton squeeze

Only about  $\frac{1}{2}$  the pretest effect is removed if the correlation is 0.5 between covariate and response. The other half appears as the male-female difference in the post-test scores

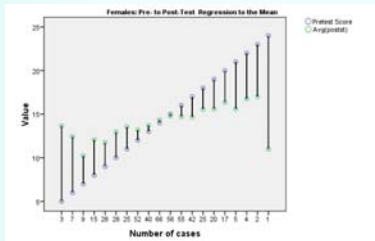


### Slide 49 Galton squeeze

NOTES:

### Galton squeeze, if $r=0.25$

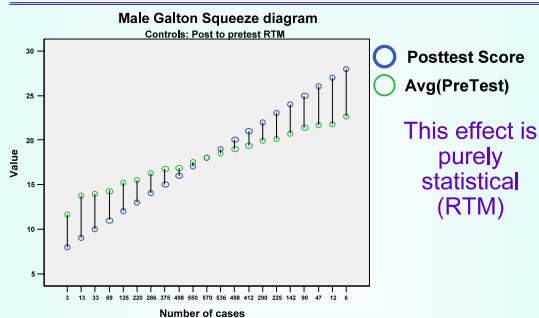
Only about  $\frac{1}{4}$  of the pretest effect is removed if the correlation is 0.25 between covariate and response. The other  $\frac{3}{4}$  appears as the male-female difference in the post-test scores



### Slide 50 Galton squeeze, if $r=0.25$

NOTES:

### Regression to the mean applies forward & backward



### Slide 51 Regression to the mean applies forward & backward

NOTES:

<div data-bbox="342 170 677 207" data-label="Section-Header"> <h3>The Regression Fallacy</h3> </div> <div data-bbox="282 214 747 239" data-label="Text"> <p>Stigler (1999) Chapter 9 Regression toward the mean</p> </div> <div data-bbox="230 239 773 495" data-label="List-Group"> <ul style="list-style-type: none"> <li>• "I suspect that <b>the regression fallacy</b> is the most common fallacy in the statistical analysis of economic data." Milton Friedman (1992) [emphasis added]</li> <li>• "The recurrence of <b>regression fallacies</b> is testimony to its subtlety, deceptive simplicity, and I speculate, to the wide use of the word regression to describe least squares fitting of curves, lines, and surfaces. Researchers may err because they believe they know about regression, yet in truth have never fully appreciated how Galton's concept works. History suggests that this will not change soon. Galton's achievement remains one of the most attractive triumphs in the history of statistics, but it is one that each generation must learn to appreciate anew, on that seemingly never loses its power to surprise."</li> </ul> </div> <div data-bbox="531 499 786 543" data-label="Image"> </div>	<div data-bbox="815 132 1269 172" data-label="Section-Header"> <h3>Slide 52 The Regression Fallacy</h3> </div> <div data-bbox="815 256 940 291" data-label="Text"> <p>NOTES:</p> </div>
<div data-bbox="298 655 738 695" data-label="Section-Header"> <h3>Statistical matching &amp; equating</h3> </div> <div data-bbox="318 699 711 726" data-label="Text"> <p>Creates 'bias' in assessing treatment effects</p> </div> <div data-bbox="230 730 773 976" data-label="List-Group"> <ul style="list-style-type: none"> <li>• <b>Matching:</b> If a covariate (e.g., pretest scores) is used to select groups, and there is less than perfect correlation between pre-and post-test assessments, then there will be regression to the mean.       <ul style="list-style-type: none"> <li>▸ Each group will regress to its own mean</li> <li>▸ <b>The regression to the mean effect will produce a treatment difference due to the treatment when none may have existed.</b></li> </ul> </li> <li>▸ Scaling College math performance vs. Gender based on categorical variables like (high school algebra I, Algebra I &amp; II, Algebra I, II &amp; Calculus) is still prone to the regression artifact</li> <li>• <b>Equating:</b> If the covariate is weakly correlated with the presumed factor that it is controlling for (SES), &amp; the covariate is positively associated with the response, then differences among groups can be magnified by the addition of the covariate.</li> </ul> </div> <div data-bbox="531 984 786 1029" data-label="Image"> </div>	<div data-bbox="815 621 1388 659" data-label="Section-Header"> <h3>Slide 53 Statistical matching &amp; equating</h3> </div> <div data-bbox="815 741 940 777" data-label="Text"> <p>NOTES:</p> </div>
<div data-bbox="282 1144 750 1182" data-label="Section-Header"> <h3>Structural modeling vs. ANCOVA</h3> </div> <div data-bbox="279 1186 763 1213" data-label="Text"> <p>Cook &amp; Campbell 1979. Primer on Regression artifacts</p> </div> <div data-bbox="230 1213 763 1488" data-label="List-Group"> <ul style="list-style-type: none"> <li>• "The usefulness of analysis of covariance is closely coupled to the assumption that each covariate be measured without error"       <ul style="list-style-type: none"> <li>▸ Other assumptions too</li> <li>▸ Violation of this assumption could be disastrous</li> </ul> </li> <li>• Using unreliable covariates can produce treatment effects that do not exist and can mask strong treatment effects.       <ul style="list-style-type: none"> <li>▸ Gender discrimination</li> <li>▸ Racial differences on standardized tests</li> </ul> </li> <li>• Really unreliable covariates can change the sign of a treatment effect</li> </ul> </div> <div data-bbox="531 1472 786 1516" data-label="Image"> </div>	<div data-bbox="815 1108 1273 1182" data-label="Section-Header"> <h3>Slide 54 Structural modeling vs. ANCOVA</h3> </div> <div data-bbox="815 1266 940 1302" data-label="Text"> <p>NOTES:</p> </div>

### Solutions to Equating & matching problems

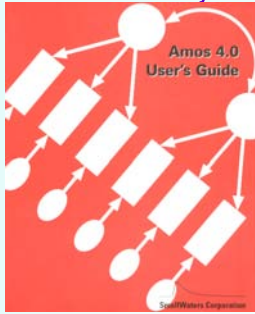
- Need a procedure that can adjust for the effect of the covariate, to correct for the 'bias' due to the regression to the mean phenomenon
- Equating & ANCOVA, may be ok when
  - Randomized assignment of subjects to cases
    - Equating not needed at all for reliability, but only for Increasing 'power'
  - **If there is little correlation between the treatment groups and the covariate.**
- Alternatives to multiple regression: Structural equation modeling, change-score analysis (Campbell & Kenny 1999), Hierarchical linear models, James-Stein (empirical Bayes) estimators

### Slide 55 Solutions to Equating & matching problems

NOTES:

### Structural equation modeling

AMOS: Analysis of moment structures



Covered in EEOS612:  
No time in EEOS611

### Slide 56 Structural equation modeling

NOTES:

### Path analysis and regression

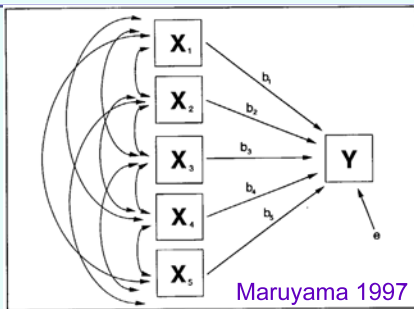


Figure 2.1. Regression Model With Five Predictor Variables

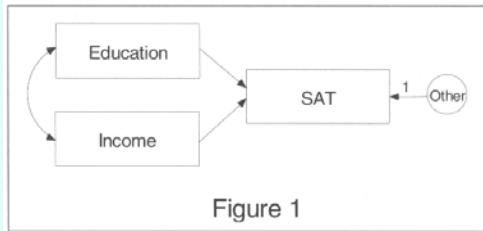
Maruyama 1997

### Slide 57 Path analysis and regression

NOTES:

### Regression: a subset of structural equation modeling

The path diagram in **Figure 1** shows a model for these data:



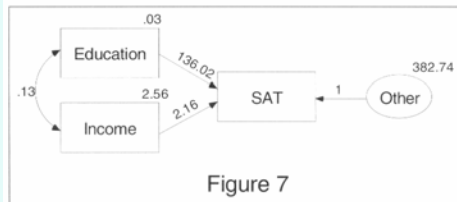
### Slide 58 Regression: a subset of structural equation modeling

NOTES:

### AMOS graphical solutions

#### Path coefficients (unstandardized or standardized)

Now to see the unstandardized estimates, simply click on **Unstandardized estimates** in the open dialog box next to your drawing area. Your path diagram should now look like **Figure 7**:



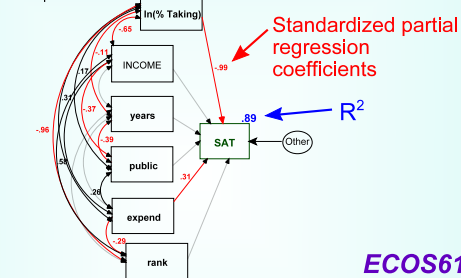
### Slide 59 AMOS graphical solutions

NOTES:

### Predicting SAT scores from states

Ramsey & Schafer (2001) "Statistical Sleuth" Ch. 12  
Modeled with AMOS

Example 12.01



### Slide 60 Predicting SAT scores from states

NOTES:

### Results from a standard OLS regression

	Unstandardized Coefficients		Standardized Coefficients		t	Sig.	95% Confidence Interval for B		Lower Bound	Upper Bound
	B	Std. Error	Beta	Partial Beta						
(Constant)	1008.6	16.7			60.500	.000	974.9	1042.3		
EXPERIENCE	.46	.08	.03		6.000	.001	.31	.61		
LEADERSHIP	-.652	.34	-.10		-1.910	.079	-1.34	.04		



### Slide 61 Results from a standard OLS regression

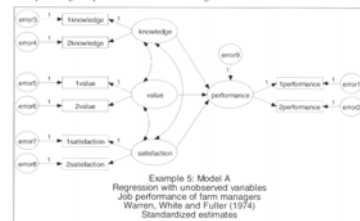
NOTES:

### From Path to Factor analysis

Latent variables (=unmeasured variables, Factors)

Model A

This path diagram presents a model for the eight subtests:



Four ellipses in the figure are labeled knowledge, value, satisfaction and performance. They represent the unobserved variables that are indirectly measured by the eight split-half tests.

### Slide 62 From Path to Factor analysis

NOTES:

### Measurement & Structural submodels

#### Measurement model

The set of connections between the observed and unobserved variables is often called the measurement model. The current problem has four distinct measurement submodels:



#### Structural model

The model component connecting the observed variables is called the structural model:



### Slide 63 Measurement & Structural submodels

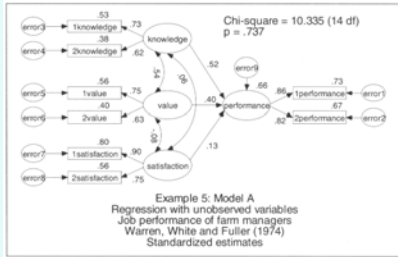
NOTES:

## AMOS Results

Chi-square, under  $H_0$ ,  $\approx$  d.f.

Amos Graphics output

The path diagram with standardized parameter estimates inserted is:



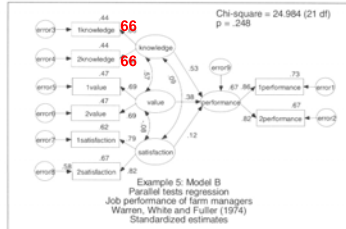
## Slide 64 AMOS Results

NOTES:

## Full vs Reduced models

Test equal slopes model, fewer parameters

Now, here are the corresponding standardized estimates and squared multiple correlations in Amos Graphics output:



Testing Model B against Model A

## Slide 65 Full vs Reduced models

NOTES:

## How to handle covariates in RTM

213 11-year olds, pre- & post-test with training

### The data

Olsson (1973) administered a battery of eight tests to 213 11-year-old students on two occasions. We will employ two of the eight tests, *Synonyms* and *Opposites*, in this example. Between the two administrations of the test battery, 108 of the students (the experimental group) received training that was intended to improve performance on the tests. The other 105 students (the control group) did not receive any special training. As a result of taking two tests on two occasions, each of the 213 students obtained four scores:

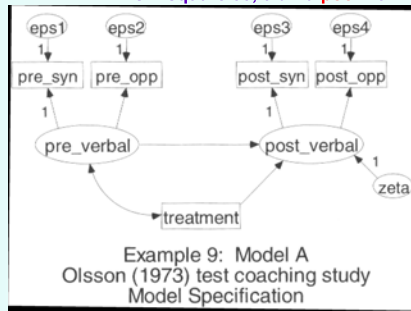
scores	explanation
pre_syn	Pretest scores on the Synonyms test
pre_opp	Pretest scores on the Opposites test
post_syn	Post-test scores on the Opposites test
post_opp	Post-test scores on the Synonyms test
treatment	A dichotomous variable taking on the value 1 for students who received the special training, and 0 for those who did not. This variable was created especially for the analyses in this example.

## Slide 66 How to handle covariates in RTM

NOTES:

### A reduced model

Chi square 33, 3 df: a **poor fit**

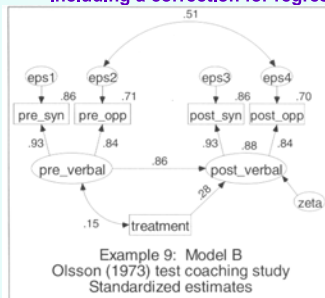


### Slide 67 A reduced model

NOTES:

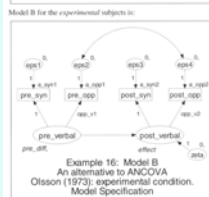
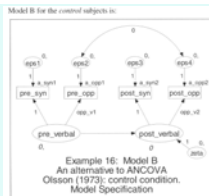
### SEM, testing between groups

Including a correction for regression to the mean

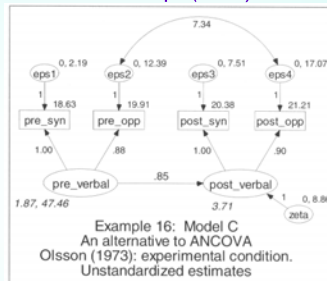


### Slide 68 SEM, testing between groups

NOTES:



Testing treatment vs.  
Control with regression to  
mean; can estimate  
intercept (3.71)



### Slide 69

NOTES:

<div data-bbox="261 315 774 352" data-label="Section-Header"> <h3>MCAS Analyses and the thrip fallacy</h3> </div>	<div data-bbox="815 132 1364 210" data-label="Section-Header"> <h4>Slide 70 MCAS Analyses and the thrip fallacy</h4> </div> <div data-bbox="815 294 940 325" data-label="Text"> <p>NOTES:</p> </div>
<div data-bbox="306 693 716 730" data-label="Section-Header"> <h3>Applications to SAT &amp; MCAS</h3> </div> <ul style="list-style-type: none"> <li>• SAT scores: can be analyzed using SEM <ul style="list-style-type: none"> <li>▸ % Taking exams and expenditure per students are the most important variabls</li> </ul> </li> <li>• How should socioeconomic factors be included in evaluating schools with MCAS <ul style="list-style-type: none"> <li>▸ Strong collinearity among socio-economic variables</li> <li>▸ Gaudet &amp; UMASS Donahue Instiute <ul style="list-style-type: none"> <li>▪ Socioeconomic variables are strongly correlated</li> <li>▪ Used principal component regression (didn't need to)</li> <li>▪ Could have used ridge regression</li> </ul> </li> <li>▸ Tuerck, Beacon Hill Institute <ul style="list-style-type: none"> <li>▪ Class size increases MCAS scores: probably an artifact, but need original data.</li> </ul> </li> <li>▸ Chen &amp; Ferguson (2002) simultaneous spatial autoregressive model (SAR)</li> </ul> </li> </ul>	<div data-bbox="815 657 1367 697" data-label="Section-Header"> <h4>Slide 71 Applications to SAT &amp; MCAS</h4> </div> <div data-bbox="815 781 940 814" data-label="Text"> <p>NOTES:</p> </div>
<div data-bbox="284 1182 748 1220" data-label="Section-Header"> <h3>Gaudet's Ranking of MA Schools</h3> </div> <p>1998 UMASS/Amherst Ph.D. and Donahue Institute Annual reports</p> <ul style="list-style-type: none"> <li>• Gaudet's method for evaluating school quality <ul style="list-style-type: none"> <li>▸ Socioeconomic variables from the 1990 census database, per student expenditure from MA DOE, MEAP results</li> <li>▸ 6 variables used in a "Major Axis" or principal components regression <ul style="list-style-type: none"> <li>▪ average education level, average income, poverty rate, single-parent status, language spoken, and percentage of school-age population enrolled in private schools.</li> </ul> </li> <li>▸ 86% of the variation in 1998 MCAS score is due to socioeconomic background of the students</li> <li>▸ Reduced to 85%, 83%, 81% and 81%MA</li> </ul> </li> <li>• Rerank 240 communities after controlling for 6 socioeconomic factors.</li> </ul>	<div data-bbox="815 1146 1412 1184" data-label="Section-Header"> <h4>Slide 72 Gaudet's Ranking of MA Schools</h4> </div> <div data-bbox="815 1268 940 1302" data-label="Text"> <p>NOTES:</p> </div>



### The best 10th grade classes

Gaudet's ranking for President Bulger's office

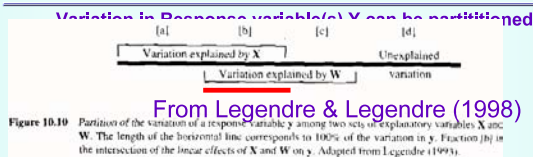
District	ELA 10 Score	Overscore	District	Math 10 Score	Overscore
Berlin	255	10	Harvard	254	10
Boylston	251	8	Lenox	250	9
Lenox	250	8	Newburyport	251	8
Stoneham	250	8	Westborough	253	8
Northampton	248	8	Amesbury	246	8
Harvard	254	8	Northampton	245	7
Nauset	250	8	Gardner	240	7
Braintree	250	8	Nauset	247	7
Clinton	245	7	Shrewsbury	249	7
Wareham	244	7	Berlin	250	7
Shrewsbury	251	7	Boylston	250	7
Pentucket Rte	250	6	Braintree	247	6
Norwood	248	5	Nashoba	250	6
Westborough	251	5	Tyngsborough	245	6

Similar to Case Study 12.1, the residual after fitting covariates (Socio-economic factors) is used to assess teaching Quality

### Slide 73 The best 10<sup>th</sup> grade classes

NOTES:

### The thrip/regression fallacy



From Legendre & Legendre (1998)

Andrewartha & Birch (1954) on 'weather' vs. Biological interactions controlling thrip abundance and Smith's critique

### Slide 74 The thrip/regression fallacy

NOTES:

### Chen & Ferguson (2002)

Evaluating school quality

$$Y_i = \beta_0 + \sum_{j=1}^4 \beta_j X_{ij} + \varepsilon_i \quad (A5.1)$$

where,  $Y_i, i = 1, 2, \dots, 226$  is the grand average of MCAS scores for years 1998, 1999, and 2000 for district  $i$ , and  $X_{ij}, j = 1, 2, 3, 4$  are the covariates of economic and demographic factors. They are AFRICAN-AMERICAN, PERCAP, TWOPHLD, and TAFDCPER. (LIM.ENG, which might quite reasonably be deemed a non-school related variable, is not used in this equation, since in combination with these variables alone it is not significant.) Once again, however, a Moran test indicates that the residuals of (A5.1) are spatially autocorrelated.



### Slide 75 Chen & Ferguson (2002)

NOTES:

## Slide 76 Chen &amp; Ferguson (2002)

Just as in the earlier equation we employ spatial models. Here the model is:

$$Y_i = \beta_0 + \sum_{j=1}^k \beta_j X_{ij} + \delta_i + \varepsilon_i \quad (A5.2)$$

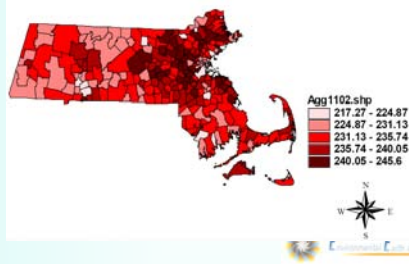
Again, as in Appendix 3, we estimate both a Conditional Spatial Autoregression (CAR) model using S-Plus and a Bayesian spatial approach estimated with WinBUGS. The estimated coefficients and p-values are listed in Table A5.3.

	S-PLUS	WinBuGS
INTERCEPT	221.54(.00)	224.20
AFRICAN	-0.160(.00)	-0.162
PERCAP	0.594(.00)	0.602
TWOPHLD	0.122(.00)	0.125
TAFDCPER	-2.124(.00)	-2.213

NOTES:

## Spatially correlated residuals

MCAS Three Year Grand Average Scores 1998-2000



## Slide 77 Spatially correlated residuals

NOTES:

## Slide 78 Chen &amp; Ferguson (2002)

(See Text – Details of Economic/Demographic Equation Below)

RANK	SCHOOL	GRADV890	RES890	BAYRE890	3 YR TOT STUDS
1	AMHERST PELHAM	237.35	6.73	7.29	8709
2	LENOX	239.78	6.03	5.37	1701
3	HARVARD	245.60	5.57	6.10	2451
4	WESTBOROUGH	242.76	5.01	5.16	6251
5	BELMONT	243.95	4.61	4.02	7083
6	NAUSET	238.97	4.46	3.13	6837
7	NORTH READING	241.52	4.42	4.32	4468
8	NORTHAMPTON	234.23	4.20	4.21	6101
9	ACTON BOXBOROUGH	245.32	4.20	4.49	10307
10	HAMILTON WENHAM	241.37	3.89	3.62	5026
11	SANDWICH	240.02	3.79	3.04	7946
12	ARLINGTON	239.12	3.45	3.26	8153
13	NEWTON	243.60	3.34	2.59	22600
14	HADLEY	238.29	3.33	4.00	1274

NOTES:

### Slide 79 Chen & Ferguson (2002)

208	MARBLEHEAD	238.98	-2.94	-2.82	5612
209	BELLINGHAM	232.54	-3.00	-2.65	5359
210	SOUTH HADLEY	231.13	-3.01	-2.47	4797
211	SAUGUS	231.55	-3.04	-3.39	6638
212	WINCHENDON	227.46	-3.09	-3.43	3995
213	TAUNTON	226.04	-3.33	-3.20	15658
214	EASTHAMPTON	228.51	-3.49	-2.80	3731
215	MARLBOROUGH	231.44	-3.55	-3.57	8028
216	CAMBRIDGE	226.07	-3.67	-3.25	13788
217	LAWRENCE	217.50	-3.68	-2.76	21674
218	HAVERTHILL	226.92	-4.24	-3.86	16712
219	MAYNARD	231.77	-4.45	-4.00	2656
220	AVON	227.79	-4.45	-4.12	1660
221	LOWELL	222.05	-4.60	-4.56	29854
222	WESTPORT COMMUNITY	229.16	-5.02	-4.19	3871
223	NARRAGANSETT	227.87	-5.17	-5.27	2941
224	SOUTHERN BERKSHIRE	228.34	-6.08	-5.65	2157
225	DOVER SHERBORN	244.12	-6.11	-5.96	3760
226	WESTON	245.15	-6.40	-6.83	4087

NOTES:

---

### What factors affect test scores?

Figure 13.27 Partition of the variation of a response matrix Y between environmental causes X<sub>e</sub> and spatial (causal) W explanatory variables. The length of the horizontal bar corresponds to 100% of the variation in Y. Compare to Fig. 13.13. (Adapted from Rosen et al., 1992, and Ferguson, 1993).

NOTES:

---

### Beacon Hill Institute Study

Goal to rank schools & to evaluate educational policy

- Use 2000 MCAS scores as response variables
- Variables in Multiple regression:
  - Policy: % change in per pupil spending, percentage change in student-teacher ratios, number of students per computer
  - Socioeconomic: crime rates, % of workers that are professionals, % households headed by single females, Urban or non-urban
  - Choice variables: % students in charter schools, % students in METCO
  - Previous performance: 1994 MEAP scores

NOTES:

## Beacon Hill Results

### Increase class sizes for "good schools"

- SES
  - School performance rises with % professionals or managers
  - School performance drops as the crime rate increases
  - School performance drops with higher % single parent households
  - Urbanized school districts have poorer performance
- Choice
  - Charter schools 'spur schools to do better'
  - METCO has no effect
  - % of students attending public schools positively associated with scores
- Policy implications
  - Spending doesn't improve performance
  - Increased class size for "good districts" Improves performance
  - "Win-win situation" Increase class size in good districts by decreasing their funding and shift to poorer districts



## Slide 82 Beacon Hill Results

NOTES:

## The 15 best schools?

### The 15 Best-Performing Massachusetts School Districts

DISTRICT (number of ratings for which district fell in the top 10)	Achieving Good Performance (G Rating)			Reducing Poor Performance (P Rating)		
	4 <sup>th</sup>	8 <sup>th</sup>	10 <sup>th</sup>	4 <sup>th</sup>	8 <sup>th</sup>	10 <sup>th</sup>
Hadley (5)	X	X	X		X	X
Clinton (3)	X	X		X		
Methuen (3)	X			X	X	
Stoneham (3)		X	X			X
Tyngsborough (3)	X		X			X
Nantucket (2)		X			X	
Chelsea (2)				X		X
Dighton-Rehoboth (2)		X			X	
Eastham (2)	X			X		
Everett (2)	X			X		
Hanover (2)		X			X	
Oxford (2)	X			X		
Provincetown (2)			X			X
Shrewsbury (2)			X			X
Sutton (2)	X			X		

## Slide 83 The 15 best schools?

NOTES:

## The 12 worst schools?

Beacon Hill Inst: Weighted average of 4th, 8th & 10th grades

### The 12 Worst-Performing Massachusetts School Districts

DISTRICT (number of ratings for which district fell in the bottom 10)	Achieving Good Performance (G Rating)			Reducing Poor Performance (P Rating)		
	4 <sup>th</sup>	8 <sup>th</sup>	10 <sup>th</sup>	4 <sup>th</sup>	8 <sup>th</sup>	10 <sup>th</sup>
Narragansett (4)	X		X	X		X
Gateway (3)		X	X			X
Somerset (3)			X		X	X
Chesterfield-Goshen (2)	X			X		
Adams Cheshire (2)	X			X		
Hudson (2)		X			X	
Leicester (2)		X			X	
Millis (2)	X			X		X
Mount Greylock (2)		X			X	
Randolph (2)			X			X
Swampscott (2)			X			X
Watertown (2)		X	X			

## Slide 84 The 12 worst schools?

NOTES:

### The Worst 10th grade schools

#### Beacon Hill Institute

Foxborough	86	Taunton	210
Weston	22	Winchendon	192
Quabbin	128	Wareham	186
North Attleborough	171	Melrose	113
Berkshire Hills	133	Carver	187
Uxbridge	170	Leicester	142
Quabog Regional	168	Winthrop	188
Harvard	17	Westford	63
Peabody	193	Lunenburg	104
Longmeadow	46	Randolph	200
Southwick Tolland	199	Littleton	67
North Middlesex	88	Lincoln-Sudbury	36
Sutton	152	Watertown	132
Hopedale	135	Bellingham	174
Mount Greylock	60	Somerset	196
Douglas	172	Narragansett	191
Saugus	197	Swampscott	141
Taunton	210	Gateway	207

### Slide 85 The Worst 10<sup>th</sup> grade schools

NOTES:

### The Beacon Hill Institute Report

#### Would increasing class size improve performance?

- Beacon Hill study
  - No attempt was made to assess colinearity among the many strongly correlated explanatory variables
  - Multicollinearity would invalidate many of their interpretations of betas, especially class size
    - The authors should have calculated VIF's
    - Solutions
      - Do ridge regression or principal components regression
      - Create a structural equation model for the hypotheses
  - A major conclusion from the study that increased class size improves MCAS performance runs counter to controlled experiments
- Experiments or quasi-experiments performed on class size indicate a negative correlation between class size and performance
  - STAR
  - SAGE



### Slide 86 The Beacon Hill Institute Report

NOTES:

### Class size and test scores


#### Inference: reduced class size causes improved performance

- The Tennessee Star Study
  - A controlled experiment
  - Students randomly assigned to class sizes of 15 or 24
  - Long-lasting effects
- The Wisconsin SAGE study
  - Students randomly assigned to small and large classes.
- Analysis of covariance (i.e., multiple regression) IS NOT a valid alternative to a randomized experiment



### Slide 87 Class size and test scores

NOTES:

Conclusions	Slide 88 Conclusions
<ul style="list-style-type: none"> <li>Regression to the mean will be present whenever an explanatory variable (covariate) exhibits less than perfect correlation with the response variable. The higher the variability in the covariate, the more the regression to the mean effect</li> <li>For pre-test vs. Post-test analyses, regressing with pretest score as an explanatory variable DOES NOT remove the effects of pre-test differences.</li> <li>Better approaches: Repeated measures designs, hierarchical linear longitudinal models, or subtract pretest from posttest (called change score analysis)</li> </ul> 	
	NOTES: