- p. 8-1 Mean structure • idea: data are generated from an underlying system, which is assumed to have the form:  $\underline{y} = \underline{f}(\underline{x}_1, \dots, \underline{x}_m) + \underline{\varepsilon}$ , where  $\underline{f}$  is <u>unknown</u>. • regression *approximates* the mean structure f by a linear combination of (known) <u>base functions</u>  $g_i(x_1, ..., x_m)$ 's, i=1, ..., p, i.e.,  $f \leftarrow \sum_{i=1}^{p} \beta_i \cdot g_i(x_1, \dots, x_m)$  $\triangleright$  when the structure of f is simple and almost linear, it can be approximated by a simple structure with fewer terms, e.g.,  $E(y) = f \cong \beta_0 + \beta_1 \underline{x_1} + \dots + \beta_m x_m$ • Q: nature is simple? • Q: are there sufficient data to support/fit a complex model?  $\blacktriangleright$  when f is complex and non-linear  $\Rightarrow$  need more terms to get a good approximation • more parameters, need more degrees of freedom, i.e., more data • e.g., 2 levels, only linear effects; 3 levels, linear and quadratic effects • Q: what other complex models? • base functions for quantitative and qualitative predictors  $x_i$ 's are defined in different ways NTHU STAT 5410, 2022, Lecture Notes made by S.-W. Cheng (NTHU, Taiwan) p. 8-2 location and/or scale change •  $\underline{x}_i \rightarrow \underline{(x_i+a)/b}$  or  $\underline{y} \rightarrow \underline{(y+a)/b}$ , where  $\underline{a}$  and  $\underline{b}$  are given constants. *a*: change of location, *b*: change of scale • **Q**: why we might want to do this? predictors of similar magnitude are easier to compare  $\triangleright$  rescaling may make  $\hat{\beta}$  easier to read and may aid interpretability > numerical stability is enhanced when all predictors are on a similar scale  $\blacktriangleright$  for experimental data, it's often that we code two levels (say, 20°C, 30°C)  $\rightarrow$ (-1, 1); three levels (say, 20°C, 30°C, 40°C)  $\rightarrow$  (-1, 0, 1)• <u>influence</u> caused by <u>location/scale change</u> on  $\underline{x}_i$  (i.e.,  $\underline{x}_i \rightarrow (x_i + a)/b$ ) (under a model with intercept) overall *F*-test, *t*-test,  $\underline{R^2}$ ,  $\hat{\sigma}$  all unchanged  $\succ \hat{\boldsymbol{\beta}} \underline{\hat{\boldsymbol{\beta}}} \underline{\text{change}} : \hat{\beta}_i \to b \hat{\beta}_i , \quad \hat{\beta}_0 \to \hat{\beta}_0 - a \hat{\beta}_i$ • influence caused by location/scale change on y (i.e.,  $y \rightarrow (y+a)/b$ )  $\blacktriangleright$  (under a model with intercept) overall F-test, t-test,  $R^2$  unchanged
  - $\blacktriangleright \hat{\underline{\beta}}$  and  $\underline{\hat{\sigma}}$  change:  $\underline{\hat{\sigma}} \rightarrow \hat{\sigma}/b$ ,  $\underline{\hat{\beta}_i} \rightarrow \hat{\beta}_i/b$  for each  $i, \underline{\hat{\beta}_0} \rightarrow (\hat{\beta}_0 + a)/b$

✤ Reading: Faraway (2005, 1<sup>st</sup> ed.), 5.2

**♦** Further reading: D&S, 16.2, 16.3



$\blacktriangleright \text{ define } \underline{z_0} = \underline{1}, \ \underline{z_1} = \underline{a_1} + \underline{b_1 x_1}, \ \underline{z_{11}} = \underline{a_2} + \underline{b_2 x_1} + \underline{c_2 x_1^2}, \ \underline{z_{111}} = \underline{a_3} + \underline{b_3 x_1} + \underline{c_3 x_1^2} + \underline{d_3 x_1^3}, \dots$	p. 8-5
Find $\underline{a_i}, \underline{b_i}, \underline{c_i},$ , that make $\underline{z_j} \underline{z_k} = 0$ if $j \neq k$ (and $  \underline{z_i}   = 1$ sometimes)	
• can apply regression to obtain $z_0, z_1, z_2, \dots$ (note: $\hat{y}^T \hat{\varepsilon} = 0$ ), e.g.,	
regress $x_1$ on $z_0$ , then the residuals is proportional to $z_1$ ; regress	
$\overline{x_{I}^{2}}$ on $\overline{z_{0}}, \overline{z_{I}}$ and the residuals is proportional to $\overline{z_{II}}$ . In $\overline{R}$ , built-	
in function is provided to construct orthogonal polynomials.	
<ul> <li><u>cross-product terms</u> (i.e., <u>interactions</u>) can be <u>defined</u> in a</li> </ul>	
similar manner (e.g., $\underline{z_{12}} = \underline{a+b} \underline{x_1+c} \underline{x_2+d} \underline{x_1x_2}$ , regress	
$\underline{x_1x_2}$ on $\underline{z_0}, \underline{z_1}, \underline{z_2}$ , and the <u>residuals</u> are <u>proportional</u> to $\underline{z_{12}}$ )	
$\succ$ change model based on polynomial terms to model based on z's, e.g.,	
$y = \beta_0 + \beta_1 \underline{x_1} + \beta_2 \underline{x_2} + \beta_{11} \underline{x_1^2} + \beta_{22} \underline{x_2^2} + \beta_{12} \underline{x_1 x_2} + \epsilon$	
$\longrightarrow y = \beta'_0 + \underline{\beta'_1}  \underline{z_1} + \underline{\beta'_2}  \underline{z_2} + \underline{\beta'_{11}}  \underline{z_{11}} + \underline{\beta'_{22}}  \underline{z_{22}} + \underline{\beta'_{12}}  \underline{z_{12}} + \epsilon$	
the two models have same column space $\Omega$ (i.e., same $R^2$ , $\hat{\sigma}$ , overall F), but	
interpretation of $\beta$ 's and $\beta'$ 's are different (i.e., different estimates, <i>t</i> -tests)	
→ orthogonality can save works when selecting model (do not have to refit after	
deleting term), it's more convenient for fitting and testing	
<ul> <li>properties of polynomial model</li> </ul>	
offer more flexible relationship	
$\succ$ remember that it's an <u>approximation</u> , we usually do <u>not</u>	
believe it exactly represents the underlying reality	
NIHU SIAI 5410, 2022, Lecture Notes	
made by S-VV Unend (NTEU Taiwan)	
polynomials have the advantage of smoothness	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the disadvantage that each</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u></li> <li>data point affects the fit <b>globally</b></li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit</u> <b>globally</b></li> <li>For larger values of <i>d</i>, the fitted polynomial curves may become wiggly.</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit</u> <b>globally</b></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become <u>wiggly</u>. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between <u>predictors and response</u>.</li> </ul>	p. 8-6
<ul> <li>polynomials have the advantage of smoothness</li> <li>but, have the disadvantage that each data point affects the fit globally</li> <li>For larger values of d, the fitted polynomial curves may become wiggly. reason: the curve may capture the random variation, rather than the overall shape of the relationship between predictors and response.</li> <li>polynomial model is hard to fit "jump function"</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> data point affects the fit <b>globally</b></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become <u>wiggly</u>. reason: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between <u>predictors and response</u>.</li> <li>polynomial model is <u>hard</u> to <u>fit "jump function</u>"</li> <li><b>* Reading</b>: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For <u>larger values</u> of <u>d</u>, the fitted <u>polynomial curves</u> may become <u>wiggly</u>. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between <u>predictors and response</u>.</li> <li><u>polynomial model</u> is <u>hard</u> to <u>fit "jump function</u>"</li> <li><b>Reading</b>: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2</li> <li><b>broken stick (line) regression (segmented regression)</b></li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between <u>predictors and response</u>.</li> <li>polynomial model is <u>hard</u> to <u>fit "jump function"</u></li> <li><b>Reading</b>: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2</li></ul>	p. 8-6
<ul> <li>polynomials have the advantage of smoothness</li> <li>but, have the disadvantage that each data point affects the fit globally</li> <li>For larger values of d, the fitted polynomial curves may become wiggly. reason: the curve may capture the random variation, rather than the overall shape of the relationship between predictors and response.</li> <li>polynomial model is hard to fit "jump function"</li> <li>Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2 </li> <li>Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li>broken stick (line) regression (segmented regression)</li> <li>Recall. polynomial regression: suitable for smooth mean structure, but cannot capture local abrupt change (example?)</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. reason: the <u>curve</u> may capture the random variation, rather than the overall <u>shape</u> of the <u>relationship</u> between <u>predictors and response</u>.</li> <li>polynomial model is <u>hard</u> to <u>fit "jump function"</u></li> <li><b>Reading</b>: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2  Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li><u>broken stick (line) regression (segmented regression)</u></li> <li>Recall. polynomial regression: suitable for smooth mean structure, but <u>cannot</u> capture <u>local abrupt change</u> (example?)</li> <li>Q: how to <u>relax</u> the <u>smoothness</u> restriction?</li> </ul>	p. 8-6
<ul> <li>&gt; polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>&gt; but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>&gt; For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. reason: the <u>curve</u> may <u>capture</u> the random variation, rather than the overall <u>shape</u> of the relationship between predictors and response.</li> <li>&gt; polynomial model is hard to <u>fit "jump function"</u></li> <li>* Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2 * Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li><u>broken stick (line) regression (segmented regression)</u></li> <li>• Recall. polynomial regression: suitable for <u>smooth mean</u> structure, but <u>cannot</u> capture <u>local abrupt change</u> (example?)</li> <li>Q: how to <u>relax</u> the <u>smoothness</u> restriction?</li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the fit <u>globally</u></li> <li>For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the relationship between predictors and response.</li> <li>polynomial model is <u>hard</u> to fit "jump function"</li> <li>Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2  Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li><u>broken stick (line) regression (segmented regression)</u></li> <li>Recall. polynomial regression: suitable for smooth mean structure, but <u>cannot</u> capture <u>local abrupt change</u> (example?)</li> <li><u>Q</u>: how to <u>relax</u> the <u>smoothness</u> restriction?</li> <li><u>one solution</u>: <u>broken line</u> regression.</li> </ul>	p. 8-6
<ul> <li>&gt; polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>&gt; but, have the <u>disadvantage</u> that <u>each</u> data point affects the fit <u>globally</u></li> <li>&gt; For <u>larger values</u> of <u>d</u>, the fitted polynomial curves may become wiggly. reason: the <u>curve</u> may capture the <u>random variation</u>, <u>rather than the</u> overall <u>shape</u> of the relationship between predictors and response.</li> <li>&gt; polynomial model is <u>hard</u> to <u>fit</u> "jump function"</li> <li>&lt; Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2 </li> <li>&gt; Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li>&gt; broken stick (line) regression (segmented regression)</li> <li>• Recall. polynomial regression: suitable for smooth mean structure, but <u>cannot</u> capture <u>local abrupt change</u> (example?)</li> <li>Q: how to <u>relax</u> the <u>smoothness</u> restriction?</li> <li>⇒ <u>one solution</u>: <u>broken line</u> regression.</li> <li>• Q: when to use <u>broken line</u> regression?</li> <li>⇒ believe that <u>different regression models</u> apply in <u>different regions</u></li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For larger values of <u>d</u>, the fitted polynomial curves may become wiggly. reason: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between predictors and response.</li> <li>polynomial model is <u>hard</u> to <u>fit "jump function</u>"</li> <li>Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2  Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li><u>broken stick (line) regression (segmented regression)</u></li> <li>Recall, polynomial regression: suitable for <u>smooth</u> mean structure, but <u>cannot</u> capture local abrupt change (example?)</li> <li><u>Q</u>: how to <u>relax</u> the <u>smoothness</u> restriction?</li> <li><u>one solution</u>: <u>broken line</u> regression.</li> <li><u>Q</u>: when to use <u>broken line</u> regression?</li> <li><u>believe that different regression models</u> apply in <u>different regions</u> of <u>data</u>, and the <u>fit</u> should be <u>continuous</u> at the <u>broken points</u></li> </ul>	p. 8-6
<ul> <li>polynomials have the <u>advantage</u> of <u>smoothness</u></li> <li>but, have the <u>disadvantage</u> that <u>each</u> <u>data point</u> affects the <u>fit globally</u></li> <li>For larger values of <u>d</u>, the fitted polynomial curves may become wiggly. <u>reason</u>: the <u>curve</u> may <u>capture</u> the <u>random variation</u>, <u>rather than</u> the overall <u>shape</u> of the <u>relationship</u> between predictors and response.</li> <li>polynomial model is hard to <u>fit</u> "jump function"</li> <li><b>Reading</b>: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2 <b>Further reading</b>: D&amp;S, 12.1, 12.3, 22.2 <b>broken stick (line) regression (segmented regression)</b></li> <li>Recall. polynomial regression: suitable for <u>smooth</u> mean structure, but <u>cannot</u> capture local abrupt change (example?)</li> <li>Q: how to <u>relax</u> the <u>smoothness</u> restriction?</li> <li>⇒ one solution: broken line regression.</li> <li>Q: when to use <u>broken line</u> regression?</li> <li>⇒ believe that <u>different regression</u> models apply in <u>different regions</u> of <u>data</u>, and the <u>fit</u> should be <u>continuous</u> at the <u>broken points</u></li> <li>suppose the <u>break</u> occurs at the <u>known</u> value c, define the <u>base function</u> (where c in the points)</li> </ul>	p. 8-6
<ul> <li>polynomials have the advantage of smoothness</li> <li>but, have the disadvantage that each data point affects the fit globally</li> <li>For larger values of d, the fitted polynomial curves may become wiggly. reason: the curve may capture the random variation, rather than the overall shape of the relationship between predictors and response.</li> <li>polynomial model is hard to fit "jump function"</li> <li>Reading: Faraway(2005, 1<sup>st</sup> ed.), 7.2.2  Further reading: D&amp;S, 12.1, 12.3, 22.2</li> <li>broken stick (line) regression (segmented regression)</li> <li>Recall. polynomial regression: suitable for smooth mean structure, but cannot capture local abrupt change (example?)</li> <li>Q: how to relax the smoothness restriction?</li> <li>&gt; one solution: broken line regression.</li> <li>Q: when to use broken line regression?</li> <li>&gt; believe that different regression models apply in different regions of data, and the fit should be continuous at the broken points</li> <li>suppose the break occurs at the known value c, define the base function (where c is called a knot):     <ul> <li>(1, if x &gt; c,</li> </ul> </li> </ul>	p. 8-6
$ \frac{1}{2} 1$	p. 8-6







• model 5: $y = \beta_0 + \beta_1 d + \beta_2 x + \beta_3 (d \cdot x) + \epsilon$	p. 8-13
$C = c_1$ : $\mu_{1,x} = E(y d = 0, x) = \beta_0 + \beta_2 x$	ž – – – – – – – – – – – – – – – – – – –
$\overline{C = c_2}: \qquad \overline{\mu_{2,x}} = E(y d = 1, x) = (\beta_0 + \beta_1) + (\beta_2 + \beta_3)x$	
$\frac{2}{\beta_0} = \frac{\mu_{1,0}}{\mu_{1,0}} \text{ (intercept of category } c_1)}$	р
$\beta_2 = \text{slope of category } c_1$	00 02 04 06 08 10
$\Rightarrow \frac{1}{\beta_1} = \frac{1}{\text{difference in intercepts}}$	2
$\beta_3 = \text{difference in slopes}$	5
$\rightarrow$ alternative coding of dummy variable (better orthogonality)	α- Σ-
$(-1)  \text{if } C = c_1.$	2 0 - 0 -
$d(\underline{C}) = \begin{cases} \frac{1}{2}, & \text{if } \underline{C} = c \end{cases}$	00 02 04 06 08 10 x
<b>O</b> : how to interpret $\beta_i$ 's in models $1 \sim 5$ under this coding?	
• model 1: $y = \beta_0 + \beta_1 \underline{d} + \epsilon$	
$C = c_1:$ $\mu_1 = E(y d = -1) = \beta_0 - \beta_1$ $\beta_0 = (\mu_1 - 1) = \beta_0 - \beta_1$	$+ \mu_2)/\underline{2}$
$\overline{C = c_2}: \qquad \overline{\mu_2} = E(y \underline{d = 1}) = \beta_0 + \beta_1 \implies \overline{\beta_1} = (\overline{\mu_2} - \overline{\beta_1}) = \beta_0 + \beta_1 = \beta_0 + \beta_0 + \beta_0 + \beta_0 = \beta_0 = \beta_0 + \beta_0 = \beta_0 + \beta_0 = \beta_$	$(-\mu_1)/2$
> analysis strategy: start from the full model (model 5) if there are en	nough
degrees of freedom, and then test if some terms can be eliminated	
identical methodology applies for more than 2	
categories and more quantitative predictors	
$\succ$ <b>Q</b> : what if <u>data</u> in the <u>two categories</u> have <u>different variance</u> ?	
NTHU STAT 5410, 2022, Lecture Notes	
made by S _VV_(Chend (NTHLT Laiwan)	
made by SW. Cheng (NTHU, Taiwan)	p. 8-14
Made by SW. Cheng (NTHU, Taiwan) $ \underline{ANalysis} \text{ of } \underline{COVAriance}: \text{ testing model 3} (\Omega) \text{ against model 2} (\Omega) \text{ against model 2} (\Omega) \text{ against model 2} (\Omega) \text{ than 2 categories and more quantitative predictors is } $	p. 8-14
$ \xrightarrow{\text{Malysis}} of \underline{COVAriance}: \text{ testing model 3} (\Omega) \text{ against model 2} ( (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is$	( <u>(</u> )
$ \xrightarrow{\text{Malysis}} of \underline{COVAriance}: \text{ testing model 3} (\Omega) \text{ against model 2} (\Omega) \text{ against model 2}$	( <u>(</u> ))
$ \xrightarrow{\text{Made by SW. Cheng (NHO, Tawan)}} $ $ \xrightarrow{\text{ANalysis} of COVAriance:} testing model 3 (\Omega) against model 2 (O) against mo$	( <u>(</u> ))
<ul> <li><u>ANalysis</u> of <u>COVAriance</u>: testing model 3 (Ω) against model 2 ( (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called <u>covariate</u> and is expected to have the <u>same effect</u> in <u>all categories</u>. The <u>difference</u> between <u>categories</u> is assumed to be an <u>additive effect</u>.</li> <li><u>one polytomous predictor</u>: more than two <u>categories</u></li> </ul>	( <u>ω</u> ) <u><u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u> <u></u></u>
<ul> <li>ANalysis of <u>COVAriance</u>: testing model 3 (Ω) against model 2 ( (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called <u>covariate</u> and is expected to have the <u>same effect</u> in all categories. The <u>difference</u> between <u>categories</u> is assumed to be an <u>additive effect</u>.</li> <li>one polytomous predictor: more than two categories &gt; for <u>k categories</u>, <u>k-1</u> dummy variables are needed to depict the <u>difference</u></li> </ul>	p. 8-14 ( <u>(</u> )) <u>(</u> ) <u>(</u> ) <u></u>
<ul> <li><u>ANalysis</u> of <u>COVAriance</u>: testing model 3 (Ω) against model 2 ( (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called <u>covariate</u> and is expected to have the <u>same effect</u> in all categories. The <u>difference</u> between <u>categories</u> is assumed to be an <u>additive effect</u>.</li> <li><u>one polytomous</u> predictor: more than two <u>categories</u></li> <li>for <u>k categories</u>, <u>k-1</u> dummy variables are needed to depict the <u>d</u> between <u>categories</u> (one parameter is used to represent <u>constant</u>)</li> </ul>	$(\underline{\omega})$ $($
<ul> <li>ANalysis of COVAriance: testing model 3 (Ω) against model 2 ( (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called <i>covariate</i> and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect.</li> <li>one polytomous predictor: more than two categories</li> <li>for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant to various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exa treatment coding</li> </ul>	$(\underline{\omega})$ $($
$ \begin{array}{c} \hline \label{eq:sW. Cheng (NHU, Tawan)} \\ \hline eq:sW. Cheng (Quantication of the set o$	p. 8-14 ( $\underline{\omega}$ ) <u>ce</u> <u>d</u> <u>d</u> <u>d</u> <u>d</u> <u>d</u> <u>d</u> <u>d</u> <u>d</u>
<ul> <li>ANalysis of <u>COVAriance</u>: testing model 3 (Ω) against model 2 (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called <u>covariate</u> and is expected to have the same effect in all categories. The difference between <u>categories</u> is assumed to be an <u>additive effect</u>.</li> <li>one polytomous predictor: more than two categories</li> <li>for <u>k</u> categories, <u>k-1</u> dummy variables are needed to depict the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding dummy variables: 4 categories c<sub>1</sub>, c<sub>2</sub>, c<sub>3</sub>, c<sub>4</sub> exacted to depice the debetween <u>categories (one parameter is used to represent constant for k various coding dummy variables: 4 categories c<sub>1</sub>, c<sub></sub></u></u></u></u></u></u></u></u></li></ul>	$(\underline{\omega})$
$ \widehat{ANalysis} of \underline{COVAriance}: testing model 3 (\Omega) against model 2 ((more than 2 categories and more quantitative predictors isallowed). The quantitative predictor is called covariate and isexpected to have the same effect in all categories. The differencebetween categories is assumed to be an additive effect.• one polytomous predictor: more than two categories> for k categories, k-1 dummy variables are needed to depict the cbetween categories (one parameter is used to represent constant forvarious coding of dummy variables: 4 categories c1, c2, c3, c4 exatreatment coding Helmert coding sum codi \boxed{\frac{d_1 d_2 d_3}{c_1 - 1 - 1}} \qquad \boxed{\frac{d_1 d_2}{c_2 + 1 - 1}}$	$(\underline{\omega})$
$ \frac{ANalysis}{2} of COVAriance: testing model 3 (Ω) against model 2 (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant for k categories) of dummy variables: 4 categories c1, c2, c3, c4 exacted to difference sum coding reatment coding reatm$	$(\underline{\omega})$ $($
$ \frac{ANalysis}{c_1 + 0} of COVAriance: testing model 3 (Ω) against model 2 (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant for k categories (one parameter is used to represent constant for various coding of dummy variables: 4 categories c_1, c_2, c_3, c_4 exameter categories (add dated da$	$(\underline{\omega})$
$\frac{ANalysis}{(more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c_1, c_2, c_3, c_4 exa treatment coding\frac{d_1 d_2 d_3}{c_1 0 0 0} \frac{d_1 d_2 d_3}{c_1 0 0 0 1} \frac{d_1 d_2 d_3}{c_1 0 0 0 1} \frac{d_1 d_2 d_3}{c_1 0 0 0 1} \frac{d_1 d_2 d_3}{c_1 0 0 0 3} \frac{d_1 d_2 d_3}{c_1 0 0 0 3}$	$(\underline{\omega})$
$\frac{ANalysis}{c_1 \text{ of } COVAriance}: \text{ testing model 3} (\Omega) \text{ against model 2} (\Omega) (more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the debetween categories (one parameter is used to represent constant to be various coding of dummy variables: 4 categories c_1, c_2, c_3, c_4 examples the treatment coding \frac{1}{c_1 \text{ of } 0 \text{ of } 0}{\frac{c_1 \text{ of } 0 \text{ of } 0}{\frac{c_2 \text{ 1 o } 0}{\frac{c_3 \text{ 0 1 0 } 0}{\frac{c_4 \text{ 0 0 } 1}{\frac{c_4 \text{ 0 0 } 0}{\frac{c_4 \text{ 0 } 0}$	$(\underline{\omega})$ $($
$\frac{ANalysis}{(more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant) > various coding of dummy variables: 4 categories c1, c2, c3, c4 exa treatment coding Helmert coding sum codi > various coding of 100 additive effect. > consider the model: y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 \underline{d_3} + \epsilon• properties of treatment coding:C = c_1: \mu_1 = E(y d_1 = 0, d_2 = 0, d_3 = 0) = \beta_0$	$(\underline{\omega})$ $($
$\frac{ANalysis}{(more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant to be various coding of dummy variables: 4 categories c_1, c_2, c_3, c_4 examples treatment coding reatment coding rea$	$(\underline{\omega})$ $($
$\frac{ANalysis of COVAriance}{(more than 2 categories and more quantitative predictors is allowed). The quantitative predictor is called covariate and is expected to have the same effect in all categories. The difference between categories is assumed to be an additive effect. • one polytomous predictor: more than two categories > for k categories, k-1 dummy variables are needed to depict the d between categories (one parameter is used to represent constant for k various coding of dummy variables: 4 categories c_1, c_2, c_3, c_4 examples c_1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 +$	$(\underline{\omega})$ $($

p. 8-15  $\square$  treats  $\underline{c}_{1}$  as a reference □ it is convenient if a "standard" categories exists  $\square$   $d_1, d_2$ , and  $d_3$  are mutually orthogonal, but not orthogonal to constant term • properties of <u>Helmert coding</u>:  $y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 d_3 + \epsilon$  $C = c_1: \qquad \mu_1 = E(y|\underline{d_1} = -1, \underline{d_2} = -1, \underline{d_3} = -1) = \beta_0 - \beta_1 - \beta_2 - \beta_3$  $\underline{C = c_2}: \qquad \underline{\mu_2} = E(y | \underline{d_1 = 1}, \underline{d_2 = -1}, \underline{d_3 = -1}) = \beta_0 + \beta_1 - \beta_2 - \beta_3$  $\underline{C = c_3}: \qquad \underline{\mu_3} = E(y|\underline{d_1 = 0}, \underline{d_2 = 2}, \underline{d_3 = -1}) = \underline{\beta_0 + 2\beta_2 - \beta_3}$  $C = c_4$ :  $\mu_4 = E(y|d_1 = 0, d_2 = 0, d_3 = 3) = \beta_0 + 3\beta_3$  $\underline{\beta_0} = \underline{\mu_1 + \mu_2 + \mu_3 + \mu_4}_{\underline{\Lambda}} \equiv \underline{\mu}$  $\Rightarrow \qquad \begin{array}{rcl} \underline{\beta_1} & = & \frac{\mu_2 - \underline{\mu_1}}{2} \\ \Rightarrow & \\ \underline{\beta_2} & = & \frac{\mu_3 - (\underline{(\mu_1 + \mu_2)/2})}{3} \end{array} \end{array}$  $\underline{\beta_3} = -\frac{\mu_4 - ((\mu_1 + \mu_2 + \mu_3)/3)}{4}$  $\Box$  constant term,  $\underline{d_1}$ ,  $\underline{d_2}$ , and  $\underline{d_3}$  are orthogonal when there are equal # of observations in each categories hard to interpret parameters may suitable for ordinal qualitative predictor • properties of sum coding:  $y = \beta_0 + \beta_1 \underline{d_1} + \beta_2 \underline{d_2} + \beta_3 \underline{d_3} + \epsilon$ p. 8-16  $\underline{C = c_1}: \qquad \underline{\mu_1} = E(y|\underline{d_1 = -1}, \underline{d_2 = -1}, \underline{d_3 = -1}) = \underline{\beta_0 - \beta_1 - \beta_2 - \beta_3}$  $\underline{C=c_2}: \qquad \underline{\mu_2}=E(y|\underline{d_1=1},\underline{d_2=0},\underline{d_3=0})=\underline{\beta_0+\beta_1}$ <u> $C = c_3$ </u>: <u> $\mu_3 = E(y|\underline{d_1} = 0, \underline{d_2} = 1, \underline{d_3} = 0) = \underline{\beta_0 + \beta_2}$ </u>  $C = c_4:$  $\underline{\mu_4} = E(y|\underline{d_1} = 0, \underline{d_2} = 0, d_3 = 1) = \beta_0 + \beta_3$  $\underline{\beta_0} \quad = \quad \frac{\mu_1 + \mu_2 + \mu_3 + \mu_4}{{}^{\scriptscriptstyle A}} \equiv \underline{\mu}$  $\Rightarrow \quad \frac{\beta_1}{\beta_2} = \frac{\mu_2 - \bar{\mu}}{\mu_3 - \bar{\mu}}$  $\frac{\beta_2}{\beta_3} = \frac{\mu_4 - \bar{\mu}}{\mu_4 - \bar{\mu}}$  $\square \underline{\beta}_0$  represent overall mean compare each category with the overall mean lesser orthogonal  $\blacktriangleright$  Note: the choice of coding does not affect the  $\underline{R^2}$ ,  $\hat{\sigma}$  and overall F-test (to test  $H_0: \underline{\beta}_1 = \underline{\beta}_2 = \underline{\beta}_3 = 0$ , the three codings have same  $\underline{\omega}$  and  $\underline{\Omega}$ ) the overall F-test is one-way ANOVA (ANalysis Of VAriance)  $\triangleright$  Q: how to work with quantitative predictors?  $\Rightarrow$  identical methodology as in 2 categories case. Q: how to interpret parameters in the case?







$\succ$ a simpler method	. 8-23
• approximate $\underline{x}_i^{\lambda}$ by $\underline{x}_i + (\lambda - l) \underline{x}_i \log(x_i)$ (i.e., first 2 terms in Taylor's expansi	on
of $\underline{x^{\lambda}}$ w.r.t. $\underline{\lambda}$ ) to determine the best $\underline{\lambda} \Rightarrow add$ the terms $\underline{x_i} log(\underline{x_i})$ to this model	
• suppose $\underline{x_i \log(x_i)}$ has regression coefficient $\underline{\eta} \implies \underline{\text{test}} H_0: \underline{\eta} = 0$ .	
If accept, no transformation; if rejected, do transformation	
• $\beta_i^* \underline{x}_i^{\lambda} \approx \beta_i^* [\underline{x}_i + (\lambda - 1) \underline{x}_i \log(\underline{x}_i)] \implies \hat{\eta} = \underline{\hat{\beta}}_i (\lambda - 1) \implies \hat{\lambda} = \underline{(\hat{\eta} / \hat{\beta}_i) + 1}$	
• Some <u>issues</u> in <u>transformation</u>	
transformation can be used to	
<ul> <li>stabilize variance</li> <li>improve fitting</li> </ul>	
<ul> <li>make errors nearly normally distributed</li> </ul>	
<ul> <li>a transformation of scale may also allow use of a simpler model</li> </ul>	
these four goals for transformation will not always be met by	
the same transformation, and compromises may be required	
$\blacktriangleright$ transformation of <u>Y</u> can alter the error structure, e.g.,	
<u>additive</u> $\leftrightarrow$ <u>multiplicative</u> in <u>exp/log</u> . In practice, <u>try different transformation</u>	
and check if the residuals satisfy the conditions required for linear regression	
$\rightarrow$ prediction in <u>Y-space</u> $\Rightarrow$ back-transforming, same for <u>C.I.</u> for the prediction of <u>Y</u>	<u>,</u>
→ It may be difficult to relate the parameters of the untransformed model to the	
parameters of transformed model. After transforming, regression	
coefficients will need to interpreted w.r.t. the transformed scale.	
<ul> <li>Reading: Faraway (1<sup>st</sup> ed.), 7.1</li> <li>Further reading: D&amp;S, chapters 13</li> <li>NTHU STAT 5410, 2022, Lecture Notes</li> </ul>	