

Indicator Variables (continued yet again)

prostate.lm3

$$\text{Model: } Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 \cdot X_2 + \varepsilon;$$

where $Y = \text{lcazol} \leftarrow \text{continuous response}$

$X_1 = \text{lpsa} \leftarrow \text{continuous covariate}$

$$X_2 = \text{sui} = \begin{cases} 1 & \text{if sui} \\ 0 & \text{otherwise} \end{cases}$$

If $\text{sui} = 0, X_2 = 0$

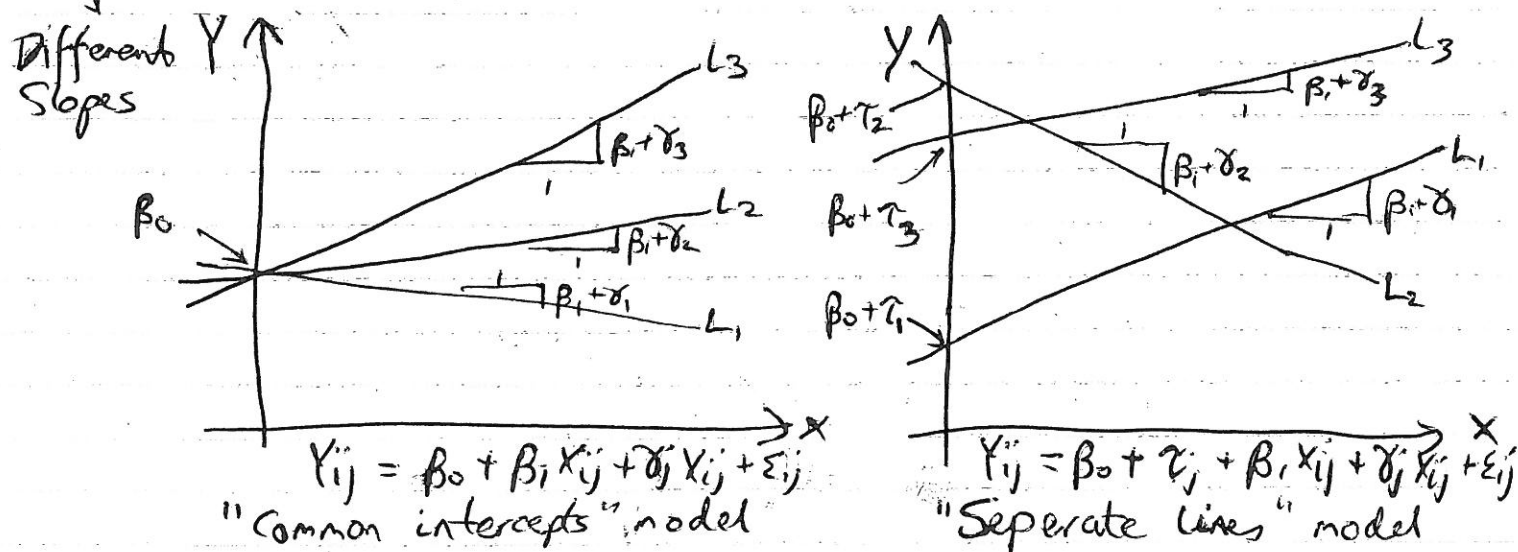
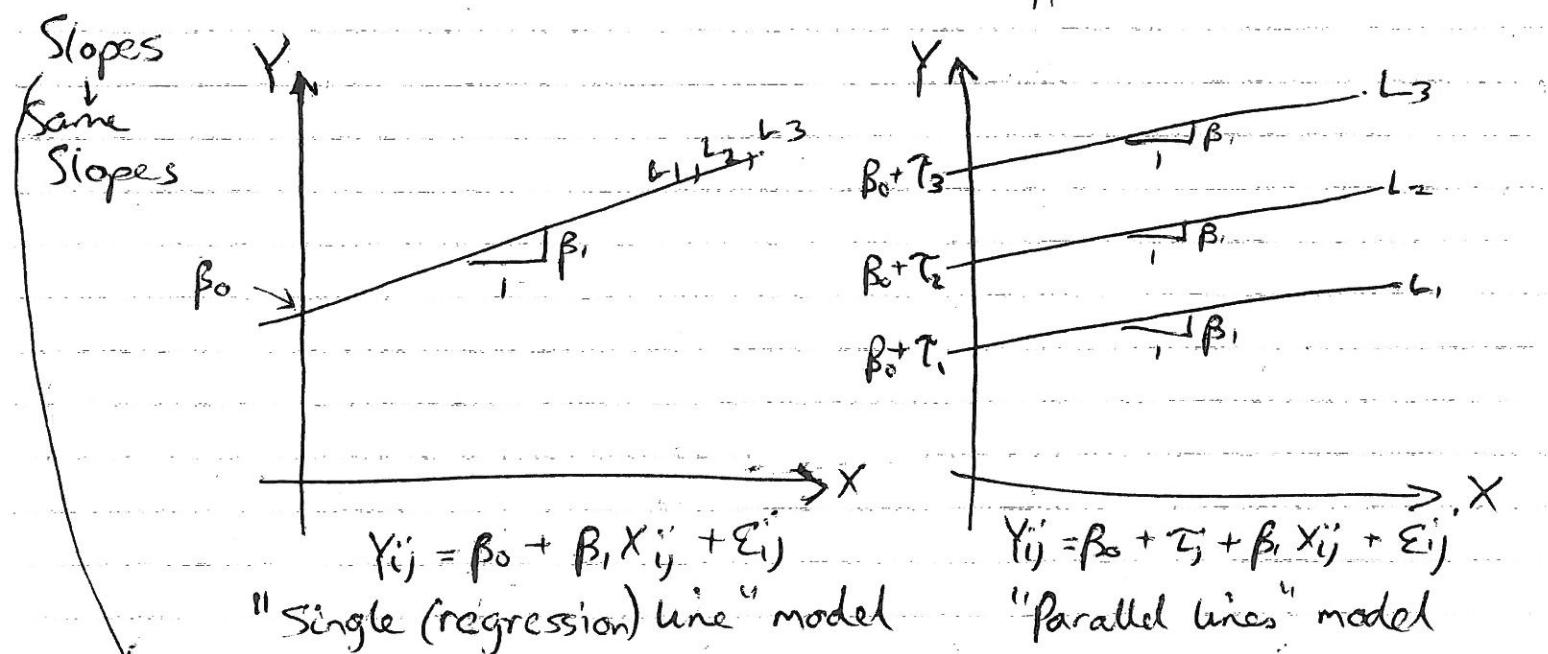
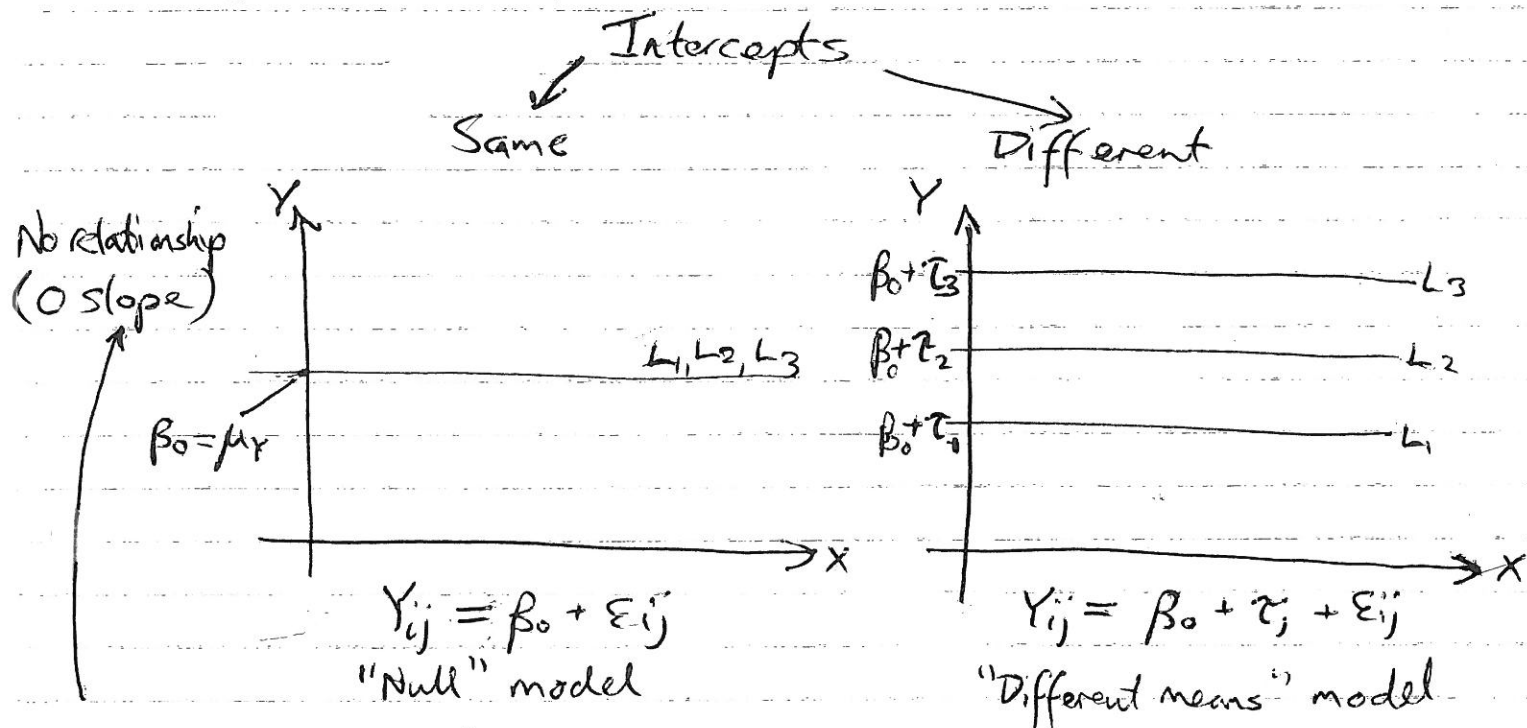
$$\begin{aligned} \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + 0 + 0 \\ &= \hat{\beta}_0 + \hat{\beta}_1 X_1 \quad \text{"base" model} \end{aligned}$$

If $\text{sui} = 1, X_2 = 1$

$$\begin{aligned} \hat{Y} &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 \cdot 1 + \hat{\beta}_3 X_1 \cdot 1 \\ &= \underbrace{(\hat{\beta}_0 + \hat{\beta}_2)}_{\text{new intercept}} + \underbrace{(\hat{\beta}_1 + \hat{\beta}_3)}_{\text{new slope}} X_1 \end{aligned}$$

ANALYSIS OF COVARIANCE (ANCOVA) MODELS

(eg one 3 level factor variable (L_1, L_2, L_3) & 1 covariate, X)



Model Selection

A "good" model is one which we can use to address the research question, which may:

- involve certain variables, which we must include in the model, so we can observe and/or "control for" the effects of these variables
- other variables (included in the data) may also be included in the model, if they help to explain some of variation (ie they turn out to be "significant")
- ultimately the research question may require some predictions; preferably predictions that hold general validity

Note if we have already chosen some scale for the variables in the model and a particular form for the model, we can then experiment with models that include other X variables in the data as predictors, as well as derived variables (X^2 , $\log X$, interaction terms involving X , ...)

If we have k possible predictors, then the number of candidate models is $O(2^k)$ as a minimum (as we can also allow for different orders of the predictors)
ie $k=1$, 2 possible models; $k=10 \Rightarrow 1,024$ models
 $k=20$, $2^{20} = 1,048,576$

For observational covariates (optional X 's) we use:

Principle of Parsimony (Occam's Razor)

Of two similar models, we will tend to prefer the simpler one (esp. if there is no significant difference between them)