

## FACTORS AND INTERACTIONS

**Example:** Testing a new teaching method.

We use a treatment and a control group. But some students will do better because of prior learning or because of aptitude. So we give a pretest to measure prior learning, a post-test to measure learning at the end of instruction, and an aptitude test. We ask: Do students do better under the experimental method, other things being equal?

*Variables:*

Response:  $y = (\text{post-test score}) - (\text{pretest score})$

Predictors:

$x_1 = \text{score on aptitude test}$  (a *covariate*)

$$x_2 = \begin{cases} 0 & \text{control group} \\ 1 & \text{experimental group} \end{cases}$$

*Possible models:*

Model I.  $E(y \mid x_1, x_2) = \eta_0 + \eta_1 x_1 + \eta_2 x_2$

This says:

For the control group ( $x_2 = \underline{\quad}$ ),  $E(y \mid x_1) = \underline{\quad}$

For the treatment group ( $x_2 = \underline{\quad}$ ),  $E(y \mid x_1) = \underline{\quad}$

Picture:

If the model is correct, then

$\eta_2 > 0$  says:

$\eta_2 = 0$  says:

$\eta_2 < 0$  says:

Will this model fit if, for example, the new method helps poorer students more than better students?

Model II. Adding the *interaction term*  $x_1x_2$  to Model I gives:

$$E(y \mid x_1, x_2) = \eta_0 + \eta_1x_1 + \eta_2x_2 + \eta_3x_1x_2$$

This says:

For the control group ( $x_2 = \text{---}$ ),  $E(y \mid x_1) = \text{---}$

For the treatment group ( $x_2 = \text{---}$ ),  $E(y \mid x_1) = \text{---}$

If  $\eta_2 > 0$  and  $\eta_3 < 0$ , we have the picture:

This says:

If  $\eta_2 < 0$  and  $\eta_3 > 0$ , we have the picture:

This says:

## CATEGORICAL VARIABLES WITH MORE THAN TWO CATEGORIES

Examples:

- Two treatments plus control (3 categories)
- Socioeconomic class (high, medium, low)
- Amount of water received by plants might be set at high, medium, low

### Terminology:

*Level* = one category of a categorical variable.

*Factor* = a set of indicator variables used to describe a single categorical predictor. (*Note*: There are other uses of the term “factor” – e.g., sometimes it is used to refer to the categorical variable, and sometime to any variable contributing to causality.)

*Example*: If a categorical variable has three levels, we could define three indicator variables:

$$v_1 = \begin{cases} 1 & \text{level 1} \\ 0 & \text{level 2} \\ 0 & \text{level 3} \end{cases} \quad v_2 = \begin{cases} 0 & \text{level 1} \\ 1 & \text{level 2} \\ 0 & \text{level 3} \end{cases} \quad v_3 = \begin{cases} 0 & \text{level 1} \\ 0 & \text{level 2} \\ 1 & \text{level 3} \end{cases}$$

However,

$$v_1 + v_2 + v_3 = 1,$$

so using all three

- is not necessary
- introduces (strict) multicollinearity.

Which two we use will depend on which level we prefer to act as "baseline". If we choose  $v_2$  and  $v_3$ , then level 1 is our baseline: It is the "default" case when the coefficients of the terms involving the indicator variables are all zero. In the case of treatment and control groups, typically the control is chosen as baseline.