

## Index of Lecture 2a: Model building I

Page	Title
1	Practical information
2	Predictor types and forms
3	Continuous predictors
4	Parametrizing categorical predictors (factors)
5	Indicator variables
6	Interaction — introductory remarks (recap)
7	Statistical modelling of interactions
8	Interaction example (VER 14.9)
9	Interactions involving continuous predictors
10	Confounding — introductory remarks (recap)
11	Definition of confounding
12	Multivariable model example (VER 14.12)
13	Additional notes on interpretation and presentation
14	Stata factor notation basics

## PRACTICAL INFORMATION

Today's lecture: follow-up from Lectures 1a–1b, and start of **model building**,

- interpretation of (more complex) models and parameters,
- more focus on the **model equation** and its parameters  $\Rightarrow$  more technical material (but important to understand),
- **interaction** and **confounding**:
  - recap of discussions in previous courses,
  - new features and recommendations (to be practiced as we go on),
- textbook reading: **VER** Sections 14.4–7,
- review of Linear Regression Exercise 1 with software demonstrations, as needed.

**Other updates:**

- **homework for Friday: Exercise 2** in “Linear Regression Exercises”, same setup as for today's exercise:
  - \* use your preferred statistical software for the calculations,
  - \* exercise review with potential Minitab/Stata demo on Friday,
- **scheduling of exam**: can we move to April 16? (currently April 17).

## PREDICTOR TYPES AND FORMS

**Variable types** (review from VHM 801):

- **quantitative** (continuous, discrete),
- **categorical** (dichotomous, ordinal, nominal).

**Predictor forms** (i.e., how predictors enter into models):

- **continuous**  $\sim$  modelled by slope:  $\beta \cdot x$ , (or other functional form),
- **categorical**  $\sim$  modelled by effects for each category:  $\alpha_j$  (for  $j^{\text{th}}$  category).
- relation to variable types:

		Predictor forms	
		continuous	categorical
Variable types			
quantitative		✓	(✓) <sup>a</sup>
dichotomous		✓ <sup>b</sup>	✓
ordinal		(✓) <sup>c</sup>	✓
nominal		no!	✓

<sup>a</sup> often requires grouping  $\sim$  loss of information, when values in a group are not distinguished,

<sup>b</sup> same results as for categorical modelling, when the values are coded as 0 and 1,

<sup>c</sup> implies a non-trivial scale assumption.

## CONTINUOUS PREDICTORS

Ways to (sometimes) **improve interpretation** of parameters:

- **scaling** of  $x$ : change the units for  $x$  to make the “1 unit change” for the **slope** more meaningful, e.g.
  - \* herd sizes are often in the hundreds, so the effect of a 1 unit change can be very small  $\Rightarrow x/10$  or  $x/100$  may be preferable,
  - \* if the range of the predictor is very small, a 1 unit change may be too large  $\Rightarrow 10 \cdot x$  or  $100 \cdot x$  may be preferable,
- **centring** of  $x$ : subtract a value from all values of  $x$  to make the **intercept** (corresponding to predictor value 0) more meaningful, e.g.
  - \* a parity of 0 does not exist  $\Rightarrow$  (parity  $- 1$ ) or (parity  $- 3$ ) may be preferable,
  - \* in dates, a calendar year of 0 is often a wild extrapolation  $\Rightarrow$  subtract the first study or a central (mean/median) study year,

**note:** centring reduces collinearity, but usually not an important consideration.

**Beware:** Scaling or centring does not change the model! (the model’s fit, inference and validity of assumptions are not affected).

## PARAMETRIZING CATEGORICAL PREDICTORS (FACTORS)

**Basic problem:** too many parameters  $\Rightarrow$  cannot meaningfully estimate them all.

**Simplest example:** 1-way ANOVA, say with factor  $A$  and  $a(\geq 2)$  groups:

- most natural parameters: the **group means**  $\mu_j$ ,  $j = 1, \dots, a$ ,
- can also introduce a **mean/intercept**  $\mu$ , so that  $\mu_j = \mu + \alpha_j$ ,  $j = 1, \dots, a$ ,
  - \*  $\alpha_j$  is the deviation from  $\mu$  for group  $j$ ,
  - \* notation with  $\mu$  and  $\alpha_j$ 's extends to multiple factors and slopes,<sup>1</sup>
  - \* results in  $(a+1)$  parameters from  $a$  groups — **one parameter too many!**

**Solution for too many parameters:** put some restrictions on them, e.g. (again for the 1-way ANOVA example),<sup>2</sup>

- $\alpha_1 = 0$  ( $\mu$  becomes value of first category, the **reference category**, Stata/Minitab/R default),
- $\alpha_a = 0$  ( $\mu$  becomes value of last category, the **reference category**, SAS default)
- $\alpha_1 + \dots + \alpha_a = 0$  ( $\mu$  becomes average value, Minitab/R option).

**Beware:** Choice of parametrization does not change the model! (the model's fit, inference and validity of assumptions are not affected).

<sup>1</sup> Corresponding to  $(\mu + \alpha_j + \beta_i)$  for two factors and to  $(\mu + \alpha_j + \beta \cdot x)$  for an added slope.

<sup>2</sup> Section 14.4.3 of VER describes hierarchical restrictions for an ordinal predictor — not part of core curriculum.

## INDICATOR VARIABLES

**Fact:** A categorical predictor, say  $A$ , is represented in the model in its usual multiple regression form, for observation  $i$ ,

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_k x_{ki} + \varepsilon_i,$$

when the corresponding  $x_j$ 's have a special form: with  $\alpha_1 = 0$ ,<sup>3</sup>

- $x_{ji} = \begin{cases} 1 & \text{if } A_i = j \\ 0 & \text{if } A_i \neq j \end{cases}$ , for  $j = 2, \dots, a$
- $x_2, \dots, x_a$  are called **indicator** (or “dummy”) variables,
- a total of  $(a-1)$  variables (none for the reference category).

Therefore, the model output will display (only) the estimates  $\hat{\beta}_1, \dots, \hat{\beta}_a$ .

**Choice of reference category:**

- still does not change the model (fit, inference, assumptions),
- ideally a category for which the comparisons with reference are (most) relevant, because
  - \*  $\beta_j \sim$  difference between  $j^{\text{th}}$  and reference category,
- avoid a reference category with small sample size (if unbalanced design).

---

<sup>3</sup> Minitab notation: “(1,0)” coding, as opposed to “(-1,0,1)” coding for  $\mu \sim$  average value.

## INTERACTION — INTRODUCTORY REMARKS (RECAP)

**Interaction** (synergism/antagonism, covariation):

- the **combined effect of two predictors** (often factors, but not necessarily so) **on the outcome** is not predictable from isolated effects of each of them (on the outcome) when examined separately,
- the effect of one factor **depends on the level** of the other factor,
- **non-parallel lines** in **interaction plot** (i.e., plot of combined means versus one factor), where parallel lines  $\sim$  **additive effects** (and hence no interaction).
- is **not the same** as collinearity: can occur with both independent and dependent predictors,<sup>4</sup>
- **note**: interaction is dependent on scale (of outcome), so affected by transformation.

**Why explore interactions?:**

- biological mechanisms of interest, either in addition to main (separate) effects or as the primary interest (e.g., with time),
- leads to more complex understanding of predictor effects  $\Rightarrow$  of potential use/interest,
- improves model fit, but beware of multiple testing issues.

---

<sup>4</sup> Strong dependence between predictors, e.g. incompleteness of  $A \times B$ , will make the assessment of interaction difficult and/or little meaningful.

## STATISTICAL MODELLING OF INTERACTIONS

Interaction in the **model equation**:

- additional terms to the main effects (**beware** to always include main effects),
- specific form depends on types of predictors involved and the chosen parametrization (discussed in the following),
- statistical significance assessed by a suitable  $F$ -test.<sup>5</sup>

**Interpretation of interaction**:

- guided by the **interaction plot** (strongly recommended),
- quantitative interpretation by estimates and CI — depends on type and form of interaction and often requires additional analysis, e.g. pairwise comparisons.

Simplest example: **two categorical predictors** (factors, say  $A$  and  $B$ ):

- **model equation**:  $\alpha_j + \beta_l + \gamma_{jl}$ , where  $j = 1, \dots, a \sim A$  and  $l = 1, \dots, b \sim B$ ,
- **restrictions** on  $(\gamma_{jl})$  needed — with 1 as the reference category for both  $A$  and  $B$ :

$$\gamma_{11}, \gamma_{21}, \dots, \gamma_{a1} = 0 \quad \text{and} \quad \gamma_{11}, \gamma_{12}, \dots, \gamma_{1b} = 0.$$

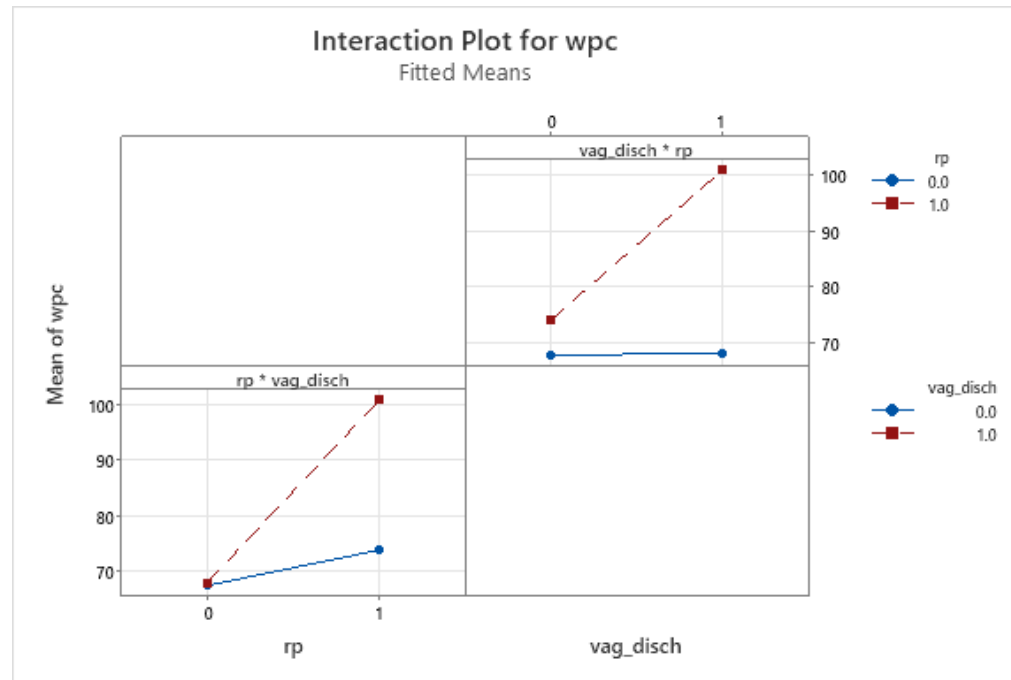
---

<sup>5</sup> The numerator degrees of freedom should equal the product of the two main effects, i.e.  $DF_{A*B} = DF_A \cdot DF_B$ ; caution advised if not the case!

## INTERACTION EXAMPLE (VER 14.9)

- Model summary:**
- **outcome:** wpc (time from wait period to conception),
  - **model terms:** rp, vag\_disch, and their interaction ( $P=0.039$ ).

**Interaction plot:**



**Parameter estimates:**

- “constant” (67.67): wpc for  $rp=0$  and  $vag\_disc=0$ ,
- “rp 1” (6.34): rp difference  $(1-0)$  in wpc for  $vag\_disc=0$ ,
- “vag\_disch 1” (0.54): vag\_disch difference  $(1-0)$  in wpc for  $rp=0$ ,
- “(vag\_disch,rp)(1,1)” (26.35): added rp difference  $(1-0)$  in wpc at  $vag\_disc=1$ .<sup>6</sup>

<sup>6</sup>Also: added vag\_disch difference  $(1-0)$  in wpc at  $rp=1$ .

## INTERACTIONS INVOLVING CONTINUOUS PREDICTORS

**Essential change** (compared to factors only):

interactions involve differences in slopes (instead of differences in means),

- interpretation is often more difficult<sup>7</sup> and technical,
- biologically, non-parallel lines ( $\sim$  interaction) is often **more natural** than parallel lines ( $\sim$  no interaction)!
- interaction plots are possible, but not always as easily accessible in software.<sup>8</sup>

Example: **one categorical predictor** (say  $x$ ) and **one factor** (say  $A$ ):

- **model equation**:  $\alpha_j + \beta \cdot x + \gamma_j \cdot x$ , where  $j = 1, \dots, a \sim A$ ,
- **restrictions** on  $(\alpha_j)$  and  $(\gamma_j)$  needed — with 1 as the reference category for  $A$ :  
 $\alpha_1 = 0$  and  $\gamma_1 = 0$ ,
- **interpretation**: regression line for  $A = j$ : (“constant” +  $\alpha_j$ ) +  $(\beta + \gamma_j) \cdot x$ ,
  - \* slopes:  $\beta$  for  $A = 1$  and  $(\beta + \gamma_j)$  for  $A = j$ ,
  - \* intercepts: “constant” for  $A = 1$  and (“constant” +  $\alpha_j$ ) for  $A = j$ .
- VER Example 14.10 with  $x = \text{milk120}$  and  $A = \text{dyst}$ ; also Example 28.12 in PSLS.

---

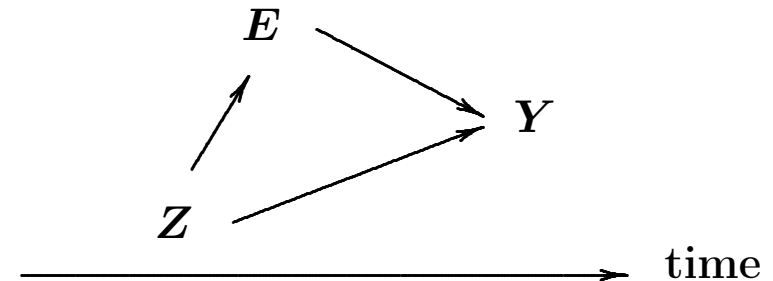
<sup>7</sup> Interaction between two continuous predictors is difficult to visualize and interpret (see VER Example 14.11), and therefore often not attractive in practice.

<sup>8</sup> The most general approach is to construct the desired plot(s) from suitably fitted values from the model.

## CONFOUNDING — INTRODUCTORY REMARKS (RECAP)

Basic notation/terminology and causal diagram (VER 13.5):

- outcome  $Y$  (of any type, but continuous for now),
- exposure  $E$  (of any type),
- extraneous factor of interest  $Z$  (measured or unmeasured; of any type),
- confounder (or lurking variable)  $Z$ : extraneous factor that exerts confounding of the relation  $E \rightarrow Y$ .<sup>9</sup>



3 necessary conditions for  $Z$  to confound the relation  $E \rightarrow Y$ :

- 1)  $Z$  must be a risk factor for  $Y$ ; more precisely:
  - \* at the reference level of  $E$ , i.e. within “exposure-negative subjects” (because the risk must not be caused by a link with  $E$ ),
- 2)  $Z$  must be associated with  $E$  in the source population,<sup>10</sup>
- 3)  $Z$  must not be affected by  $E$  (which would make  $Z$  an **intermediate (mediating) variable**<sup>11</sup> between  $E$  and  $Y$ ), and  $Z$  must not be an effect of  $Y$ .

<sup>9</sup> Note that a confounder is always tied to **both** outcome and exposure.

<sup>10</sup> See VER 13.5 for specific conditions in cohort and case-control studies.

<sup>11</sup> The VER textbook uses the (non-standard) term **intervening variable**.

## DEFINITION OF CONFOUNDING

### Definition of confounding:

- mathematically not easy; best attempt uses counterfactual arguments, but often infeasible in practice to realistically assume that all its conditions are met,
- literature agrees on **necessary** (but not sufficient) conditions for confounding.

**Pragmatic solution:** define confounding by  $z$  for the relation  $x \rightarrow y$  as present when both of the conditions (i)–(ii) below are met:

- (i)  $z$  meets the 3 necessary conditions, 1)–3) on previous page, to be a confounder,
- (ii) the difference between a crude (“total”) measure of association/effect<sup>12</sup> and a confounding-adjusted measure of association/effect<sup>13</sup> is “substantial”, i.e.<sup>14</sup>
  - \* a bias above 20–30% (arbitrary cut-off set in VER) measured relative to the crude estimate.

**Illustration:** confounding by herds ( $z$ ) for “effect” of vag\_disch ( $x$ ) on wpc ( $y$ ):

- in a causal diagram, herds are “before” exposure and outcome — condition 3),
- the associations  $z \rightarrow x$  and  $z \rightarrow y$  are significant — conditions 1)–2),
- the relative change in the vag\_disch estimate is:  $|12.0 - 17.8| / |12.0| = 0.48 \sim 48\%$ .

<sup>12</sup> In regression models: a model with  $x$  as predictor for  $y$ , but  $z$  not included.

<sup>13</sup> In regression models: a model with both  $x$  and  $z$  as predictors for  $y$ .

<sup>14</sup> Rothman *al* (2008), p. 262, note that usually 50% would be considered substantial, and 5% would not...

## MULTIVARIABLE MODEL EXAMPLE (VER 14.12)

- **outcome** of interest: time from wait period to calving (wpc) in daisy2red dataset,
- **exposures** of interest (from study objective focusing on diseases/conditions related to calving): rp, vag\_disch, dyst,
- findings of causal diagram (VER Figure 14.4):
  - \* potential **confounders**: herd\_size, parity, autumn calving (aut\_calv), twin,
  - \* **intermediate** variables: milk120, calving to first service (cf),
- special effects: quadratic term for herd\_size, interaction rp\*vag\_disch.

### Interpretation of model's “effects” ...

herd_size	$-36.06 \cdot x + 11.14 \cdot x^2$ (predictor: herd_size/100)
parity	$1.14 \cdot x$ predictor: (parity-1)
aut_calv	-8.26 for autumn calving
twin	20.68 for twin calving
dyst	11.70 for dystocia
rp	5.99 for retained placenta
vag_disch	1.23 for vaginal discharge
rp*vag_disch	22.86 for retained placenta and vaginal discharge
_cons	84.66

## ADDITIONAL NOTES ON INTERPRETATION AND PRESENTATION

### Categorical predictors:

- assess significance by **overall  $F$ -test**, not individual coefficients,
- focus interpretation on **all pairwise comparisons** (possibly adjusting for multiple comparisons), not only those with the reference category.

### Interactions:

- in presence of an important interaction, the main effects are often of limited interest,
- when a reference category parametrization is used, the coefficients for each predictor in an interaction are **not** the main effects.

### Confounding:

- implies **collinearity** between  $x$  and  $z$ , and is often helpful in understanding how to deal with collinearity,
- common practice to **include potential confounders** in a model, even if insignificant,
- if  $z$  is confounder for  $x$ , then  $x$  is typically intermediate for  $z$ ,
- confounding and interaction are **different things**;  
if they exist together, confounding may be explored for each factor level,
- to say  $z$  is a confounder without mentioning the exposure, is quite meaningless.

## STATA FACTOR NOTATION BASICS

For predictors  $x, z$  and an outcome  $y$ , Stata uses the notation:<sup>15</sup>

- $i.x \sim$  categorical effect of  $x$  ( $x$  must be integer),
- $c.x \sim$  continuous (slope) effect of  $x$  ( $x$  must be numerical),
- the **default** depends on the command:
  - \* `regress y x`  $\sim$  `regress y c.x` (i.e., default is  $c.x$ ),
  - \* `anova y x`  $\sim$  `anova y i.x` (i.e., default is  $i.x$ ).

**Combined effects:**

- $x\#z \sim$  interaction  $x \times z$ ; in all commands, the default is  $x\#z \sim i.x\#i.z$ ,
- $c.x\#c.z \sim$  multiplication  $x \cdot z$ ,
- always:  $x\#\#z \sim x\ z\ x\#z$
- **quadratic regression:**  $c.x\#\#c.x \sim$  continuous terms  $x$  and  $x^2$ .

Factor terms can (must) be **used in tests**, e.g. the Stata commands,

- `testparm i.x`
- `testparm c.x\#\#c.x`

---

<sup>15</sup> In the Stata help files, the information is listed under “fvvarlist”.