

# The general linear model

## Lecture 5 – Multiple predictors

Jan Vanhove

<https://janhove.github.io>

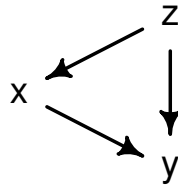
Ghent, 14–16 July 2025

As we’ve seen in the previous lectures, the general linear model can accommodate multiple predictor variables at once, yielding **multiple regression models** as opposed to simple regression models. Mathematically and computationally, nothing really changes when adding more predictor variables to the model: The model estimates the  $\beta$  parameters by minimising the sum of squared residuals and inference is achieved by making assumptions about the distribution of the errors. Conceptually, however, multiple regression models often pose researchers a number of challenges. These mainly relate to the decision whether to include several predictor variables in the model and, relatedly, to how to interpret the estimated regression coefficients. Here, I don’t use the verb ‘interpret’ to refer to subject-matter interpretations, but to more basic statistical interpretations: What do all these numbers literally mean? Clearly, before we can lend subject-matter interpretations to the output of statistical models, we need to understand what they mean literally.

In this lecture, we’ll write simulations to illustrate the consequences of including several predictors in a handful of key scenarios. The following box anticipates the main insight.

The parameter estimates that the general linear model produces first and foremost describe associations in the data. Often, however, researchers want to lend a causal interpretation to these associations. For instance, we’re usually not content with finding out that the experimental group outperforms the control group on average – we want to know whether the experimental group outperforms the control group on average *because* they’re the experimental group.

**Statistical tools cannot prove claims about causality.** This goes even for so-called causal models. But what we can do is make assumptions about the causal connections between our variables and reason about whether and how we can model these statistically. Whether



**Figure 5.1:** In this scenario,  $z$  influences both  $x$  and  $y$ . As a result,  $z$  acts as a confounding variable if we're interested in the causal influence of  $x$  on  $y$ .

these causal assumptions are reasonable is a subject-matter issue, not a statistical one.

A useful tool when discussing assumptions about causal connections are **directed acyclic graphs** (DAGs). For an introduction to DAGs, see Rohrer (2018) and McElreath (2020). We will use DAGs throughout this lecture to represent causal connections between variables ( $x, y, z, \dots$ ). Throughout, we're interested in estimating the causal influence that  $x$  exerts on  $y$ . Our guiding questions are, first, whether this is at all possible, and, second, if so, which variables we should include in the model.

## 1 Taking into account confounding variables

First, we consider the scenario shown in Figure 5.1. We're interested in the causal influence of  $x$  on  $y$ . However, it is possible that a third variable,  $z$ , influences both  $x$  and  $y$ . If you want to, you can replace these abstract variable names by something more specific (e.g.,  $x$  = participation in a content- and language-integrated programme,  $y$  = proficiency in the target language,  $z$  = the parents' socio-economic status). But I think it's ultimately more useful to embrace the abstraction, as this makes it clearer that the lessons drawn from this scenario are valid more generally.

Nevertheless, we'll make this scenario a bit more concrete by fleshing out what the causal connections between the three variables are. To keep this simple, we'll assume that the  $z$  variable was drawn from a normal distribution with mean 0 and standard deviation 1; but this isn't too important:

$$z_i \stackrel{\text{i.i.d.}}{\sim} \text{Normal}(0, 1^2),$$

$$i = 1, \dots, n.$$

Next, we assume that a one-unit increase in the  $z$  variable causes an increase of 1.2 units in the  $x$  variable. There is, however, some variability in the  $x$  variable that is unrelated to  $z$ . We

express this additional variability in an error term  $\tau$ , which we assume is also drawn from a normal distribution with mean 0 and standard deviation 1:

$$\begin{aligned}x_i &= 0 + 1.2 \cdot z_i + \tau_i, \\ \tau_i &\overset{\text{i.i.d.}}{\sim} \text{Normal}(0, 1^2),\end{aligned}\tag{1}$$

$i = 1, \dots, n$ . The numbers 1.2 in the first line and 1 in the second line were chosen arbitrarily – you can play around with these numbers at home. The 0 in the first line merely means that the mean of the  $x$  variable is 0, but little hinges on this.

Additionally, we assume that the  $y$  variable is described by the equation below. The influence of  $x$  on  $y$  is such that a one-unit increase in  $x$  leads to a 0.6-unit increase in  $y$ . The  $z$  variable, however, has a negative causal effect on  $y$ , but we don't want to estimate this effect. Again, the numbers 5.2, 0.6 and  $-1.3$  were chosen arbitrarily.

$$\begin{aligned}y_i &= 5.2 + 0.6 \cdot x_i - 1.3 \cdot z_i + \varepsilon_i, \\ \varepsilon_i &\overset{\text{i.i.d.}}{\sim} \text{Normal}(0, 1^2),\end{aligned}\tag{2}$$

$i = 1, \dots, n$ .

Let's now simulate a dataset with 100 observations of these three variables. If we don't specify any further parameters, the `rnorm(n)` call generates  $n$  observations from a normal distribution with mean 0 and standard deviation 1:

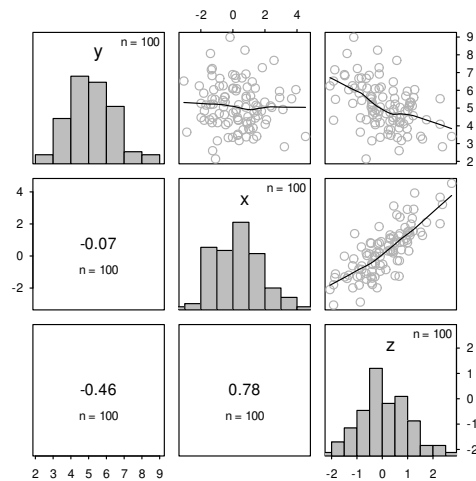
```
library(tidyverse)
library(here)

n <- 100
z <- rnorm(n)
x <- 0 + 1.2*z + rnorm(n)
y <- 5.2 + 0.6*x - 1.3*z + rnorm(n)
d <- tibble(y, x, z)
```

In order to visualise the pairwise associations between several continuous variables at once, we can draw a scatterplot matrix. The file `scatterplot_matrix.R` in the `functions` directory defines a custom-made function for plotting scatterplot matrices. Let's read it in and plot the simulated data; see Figure 5.2.

```
source(here("functions", "scatterplot_matrix.R"))
scatterplot_matrix(d)
```

We will now compare two strategies for analysing these simulated data. If we were analysing



**Figure 5.2:** A scatterplot matrix of the three simulated variables. The numbers in the lower triangle are Pearson correlation coefficients and the number of data pairs they are based on. The trend lines in the scatterplots in the upper triangle are scatterplot smoothers. For more information, see <https://janhove.github.io/posts/2019-11-28-scatterplot-matrix/>.

real data, we wouldn't do this – we'd only run the analysis that makes most sense. But in this lecture, we want to use simulated data in order to find out what the most sensible strategy is.

In the first model, we ignore the z variable:

```
dag1.lm1 <- lm(y ~ x, data = d)
summary(dag1.lm1)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.142553	0.132453	38.82546	2.5995e-61
x	-0.062995	0.086952	-0.72448	4.7050e-01

The resulting parameter estimates are to be interpreted in exactly the same way as explained in Lecture 2:

- If we were to take a large number of observations for which all  $x$  values were 0, then, according to the model, we'd expect the mean of their  $y$  values to be  $5.14 \pm 0.13$ .
- If we were to take a large number of observations for which all  $x$  values were 1, then, according to the model, their mean  $y$  value is expected to be  $0.06 \pm 0.09$  lower than that of the  $x = 0$  group.

*This interpretation is absolutely fine!* Note, however, that this interpretation does not involve any

causal claims.

In the second model, we include  $z$  as a predictor:

```
dag1.lm2 <- lm(y ~ x + z, data = d)
summary(dag1.lm2)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.0321	0.10269	49.0020	3.0861e-70
x	0.6334	0.10731	5.9028	5.2557e-08
z	-1.3475	0.16243	-8.2961	6.2768e-13

These parameter estimates, too, are to be interpreted as outlined in Lecture 2:

- If we were to take a large number of observations for which all  $x$  and  $z$  values were 0, then, according to the model, we'd expect the mean of their  $y$  values to be  $5.03 \pm 0.10$ .
- If we were to take a large number of observations for which all  $x$  values were 1 and all  $z$  values were 0, then, according to the model, their mean  $y$  value is expected to be  $0.63 \pm 0.11$  higher than that of the  $x = 0, z = 0$  group.
- If we were to take a large number of observations for which all  $x$  values were 0 and all  $z$  values were 1, then, according to the model, their mean  $y$  value is expected to be  $1.3 \pm 0.2$  lower than that of the  $x = 0, z = 0$  group.

The second point does *not* contradict the interpretation of the parameter estimates of the first model: Just because parameters with the same name occur in both models ((Intercept),  $x$ ) doesn't mean that these have the same interpretation. (Recall the Greek letter fallacy!) Specifically, the parameter estimates in the second model can only be interpreted correctly when taking into account the other variable in the model,  $z$ . Similarly, to interpret the parameter estimates in the first model, you must not implicitly assume a fixed value for the  $z$  variable that was not included in the model.

Now compare the parameter estimates of the second models with those in Equation 2 on page 3. The estimated parameters are quite close to the true parameters we used to generate the simulated data. In fact, the discrepancies between the estimates and the true values are purely due to chance. We can verify this by simulating lots of datasets using the same parameters as used in Equations 1 and 2 and analysing them in the same way as we did here. Below, we define the function `generate_dag1()`, which by default generates 10,000 such datasets, containing 100 observations each. Each dataset is analysed twice: Once without taking the  $z$  variable into account, and once including this variable in the model. For each simulated dataset and each model, the estimated  $x$  parameter is extracted.

```
generate_dag1 <- function(
  n = 100,          # number of observations per dataset
  sims = 10000,     # number of datasets
  z_x = 1.2,        # influence z -> x
  x_y = 0.6,        # influence x -> y
  z_y = -1.3,       # influence z -> y
  baseline_y = 5.2 # baseline y
) {
  est_lm1 <- vector(length = sims)
  est_lm2 <- vector(length = sims)

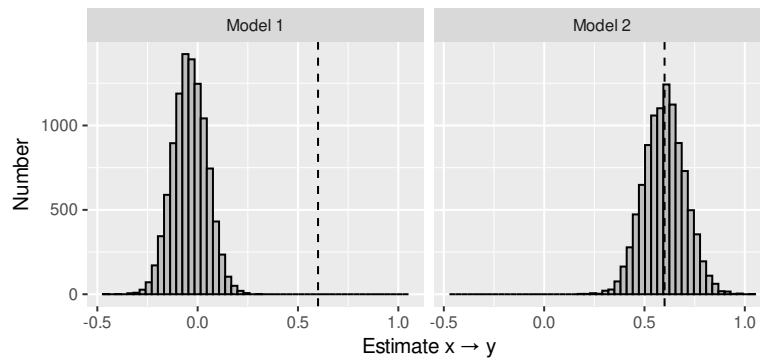
  for (i in 1:sims) {
    z <- rnorm(n)
    x <- z_x*z + rnorm(n)
    y <- baseline_y + x_y*x + z_y*z + rnorm(n)
    mod_lm1 <- lm(y ~ x)
    mod_lm2 <- lm(y ~ x + z)
    est_lm1[[i]] <- coef(mod_lm1)[[2]]
    est_lm2[[i]] <- coef(mod_lm2)[[2]]
  }

  tibble('Model 1' = est_lm1,
        'Model 2' = est_lm2)
}
```

The code below runs this function using the default settings and then plots the estimated parameters obtained by both models; Figure 5.3.

```
est_dag1 <- generate_dag1()

est_dag1 |>
  pivot_longer(cols = everything(),
               names_to = "Model",
               values_to = "Estimate") |>
  ggplot(aes(x = Estimate)) +
  geom_histogram(bins = 50,
                 fill = "grey", colour = "black") +
  geom_vline(xintercept = 0.6, linetype = "dashed") +
  facet_grid(cols = vars(Model)) +
  xlab("Estimate x y") +
```



**Figure 5.3:** Model 1 does not estimate the causal effect of  $x$  on  $y$  in an unbiased way. Depending on the parameter values chosen in Equations 1 and 2, the bias could be an overestimate or, like here, an underestimate. Model 2, by contrast, does provide an unbiased estimate of the effect of  $x$  on  $y$ : On average, the estimated parameter values correspond exactly to the true parameter value.

```
ylab("Number")
```

This simulation confirms that the second model yields an unbiased estimate of the causal effect of  $x$  on  $y$  (0.6), but the first doesn't:

```
apply(est_dag1, 2, mean)

Model 1  Model 2
-0.038681 0.599484
```

The model without the confounding variable ( $z$ ) isn't wrong. If you want to estimate the average difference in the  $y$  variable between groups differing in  $x$ , this is the model you need. But if we assume the causal structure shown in Figure 5.1, we cannot interpret its parameter estimates causally.

It seems clear what conclusion we ought to draw from the considerations above: If confounding variables are at play and we want to make causal claims, we need to include these confounders in the analysis. In practice, however, things aren't so simple.

- We may not know all confounders or we may not have been able to measure all of them.
- Even if we did measure the confounders, we probably didn't measure them perfectly.
- It is possible that the confounders exert some nonlinear influence, whereas we only considered linear effects.

In the exercises below, you explore the first two complicating factors. To anticipate the take-home message:

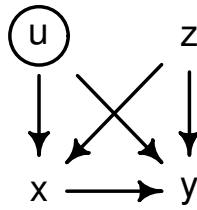


Figure 5.4: The confounder  $z$  was measured; the confounder  $u$  wasn't.

**Tip 5.1. Don't pin your hopes on statistical tools to neutralise the effect of confounding variables.** This would require you to have perfectly measured all confounders and to have specified the functional form of their causal effects correctly. **Statistical tools are no substitute for a solid research design that neutralises confounders.** ◇

**Exercise 5.2** (Unknown confounders). In Figure 5.4, one confounder,  $u$ , was added to our DAG, but for some reason, it wasn't measured. We assume that the following causal relationships are at play:

$$\begin{aligned}
 z_i &\sim \text{Normal}(0, 1^2), \\
 u_i &\sim \text{Normal}(0, 1^2), \\
 x_i &= 0 + 1.2 \cdot z_i + 0.9 \cdot u_i + \tau_i, \\
 y_i &= 5.2 + 0.6 \cdot x_i - 1.3 \cdot z_i + 2.5 \cdot u_i + \varepsilon_i, \\
 \tau_i &\sim \text{Normal}(0, 1^2), \\
 \varepsilon_i &\sim \text{Normal}(0, 1^2),
 \end{aligned}$$

$i = 1, \dots, n$ , with all  $z_i, u_i, \varepsilon_i, \tau_i$  independent. We can generate a dataset governed by these equations that contains 50 observations as follows. Note that while we simulate the  $u$  variable, we don't add it to the dataset since it wasn't measured.

```

n <- 50
z <- rnorm(n)
u <- rnorm(n)
x <- 0 + 1.2*z + 0.9*u + rnorm(n)
y <- 5.2 + 0.6*x - 1.3*z + 2.5*u + rnorm(n)
d <- tibble(y, x, z)

```

While we cannot include  $u$  in the model, we can include  $z$ :

```

exercisel.m <- lm(y ~ x + z, data = d)
summary(exercisel.m)$coefficients

```



	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.2761	0.35514	14.8565	2.2261e-19
x	1.7884	0.26698	6.6989	2.3416e-08
z	-2.1714	0.49826	-4.3580	7.0873e-05

1. Explain what the parameter estimate for  $x$  literally means.
2. Adapt the `generate_dag1()` function and show that the `exercise1.lm` model does not provide an unbiased estimate of the causal effect of  $x$  on  $y$ , even though it includes the measured confounder  $z$ .
3. Would more data help solve this problem? Justify your answer by means of a simulation in which each dataset features 200 instead of 50 observations.  $\diamond$

**Exercise 5.3** (Measurement error on the confounder). Figure 5.5 depicts another scenario where  $z$  confounds the causal relationship between  $x$  and  $y$ . This time, however, we didn't measure  $z$  itself. Instead, we obtained an **indicator**  $z_m$  of  $z$ . This indicator represents the **construct**  $z$  imperfectly. This is quite usual: constructs such as working memory capacity, L2 writing skills, L1 vocabulary knowledge, intelligence, socioeconomic status etc., cannot be observed directly and have to be inferred from test results, questionnaire responses etc. It is therefore crucial to understand the effect of measurement error on statistical control.

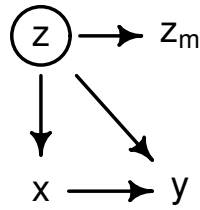
We assume that the same causal links are at play as earlier. The only difference is that we include  $z_m$  instead of  $z$  in the dataset. To construct  $z_m$ , we take the values of  $z$  and add some Gaussian noise to it (from a normal distribution with mean 0 and standard deviation 0.3).

$$\begin{aligned}
 z_i &\sim \text{Normal}(0, 1^2), \\
 z_{m,i} &= z_i + \psi_i, \\
 x_i &= 0 + 1.2 \cdot z_i + \tau_i, \\
 y_i &= 5.2 + 0.6 \cdot x_i - 1.3 \cdot z_i + \varepsilon_i, \\
 \psi_i &\sim \text{Normal}(0, 0.3^2), \\
 \tau_i &\sim \text{Normal}(0, 1^2), \\
 \varepsilon_i &\sim \text{Normal}(0, 1^2),
 \end{aligned} \tag{3}$$

$i = 1, \dots, n$ , with  $z_i, \psi_i, \tau_i, \varepsilon_i$  independent.

Let's simulate the data.

```
n <- 50
z <- rnorm(n)
z_m <- z + rnorm(n, sd = 0.3)
```



**Figure 5.5:** The  $z$  variable confounds the causal relationship between  $x$  and  $y$ , but it wasn't measured directly. Instead, we need to make do with a proxy variable  $z_m$  that captures  $z$  imperfectly.

```

x <- 0 + 1.2*z + rnorm(n)
y <- 5.2 + 0.6*x - 1.3*z + rnorm(n)
d <- tibble(y, x, z_m)
  
```

This time, we include  $z_m$  in the analysis:

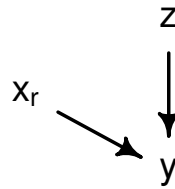
```

exercise2.lm <- lm(y ~ x + z_m, data = d)
summary(exercise2.lm)$coefficients
  
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.39410	0.19468	27.7072	9.1603e-31
x	0.25364	0.19460	1.3033	1.9881e-01
z_m	-0.74224	0.31544	-2.3530	2.2860e-02

1. Explain what the parameter estimate for  $x$  literally means.
2. Adapt the `generate_dag1()` function and show that the `exercise2.lm` model does not provide an unbiased estimate of the causal effect of  $x$  on  $y$ , even though it includes an indicator of the confounder (i.e.,  $z_m$ ).
3. Would more data help solve this problem? Justify your answer by means of a simulation in which each dataset features 200 instead of 50 observations.
4. Would it be better not to take into account the confounder at all? That is, should we just fit the model without the indicator of  $z$ ?
5. What would happen if we had a more reliable indicator for  $z$ ? To answer this question, rerun your simulation but use a standard deviation of 0.1 instead of 0.3 for  $\psi$  in Equation 3. ◇

**Tip 5.4.** If you read about a study that claims that the participants' 'intelligence' or 'socio-economic status' has been controlled for, mentally change this to 'an imperfect indicator of the participants' intelligence/socio-economic status'. ◇



**Figure 5.6:** In this scenario, the values of  $x$  were assigned randomly and independently of  $z$ . This is typical of experiments in which the participants are randomly assigned to the conditions.

## 2 Control variables in randomised experiments

In Section 1, we discussed a scenario in which  $z$  causally affects both  $x$  and  $y$ . This scenario is typical for observational (or correlational) studies and quasi-experiments. In the present section, we turn our attention to randomised experiments, that is, experiments in which the values of  $x$  are manipulated by the researchers based on random assignment. To reflect this, the  $x$  variable is represented as  $x_r$  in Figure 5.6 ( $r$  for *randomised*). A typical example for the more abstract scenario we consider here is an experiment in which participants are randomly assigned to the conditions. In this case,  $x$  would be a categorical variable. But without any loss of generality, we can restrict our discussion to a continuous  $x$ . This makes the simulations a bit easier; it doesn't affect the conclusions we will draw. The  $z$  variable would then be a variable that we're not really interested in but of which we suspect that it correlates with the outcome,  $y$ . The textbook case is a pretest/posttest experiment, where  $x$  represents the experimental condition,  $z$  the pretest performance, and  $y$  the posttest performance.

We assume that the causal links between the variables are described by the following equations:

$$\begin{aligned}x_i &\sim \text{Normal}(0, 1^2), \\z_i &\sim \text{Normal}(0, 1^2), \\y_i &= 5.2 + 0.3 \cdot x_i + 0.9 \cdot z_i + \varepsilon_i, \\ \varepsilon_i &\sim \text{Normal}(0, 1^2),\end{aligned}\tag{4}$$

$i = 1, \dots, n$  with the  $x_i, z_i, \varepsilon_i$  independent. Let's simulate a dataset conforming to these equations:

```
n <- 100
x <- rnorm(n)
z <- rnorm(n)
y <- 5.2 + 0.3*x + 0.9*z + rnorm(n)
d <- tibble(y, x, z)
```

```
# Not shown in script
scatterplot_matrix(d)
```

We again fit one model with and one model without  $z$ . Both models yield similar though different estimates for the  $x$  parameter:  $0.35 \pm 0.14$  and  $0.38 \pm 0.10$ .

```
dag2.lm1 <- lm(y ~ x, data = d)
summary(dag2.lm1)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.40030	0.14198	38.0370	1.7153e-60
x	0.35473	0.14035	2.5275	1.3086e-02

```
dag2.lm2 <- lm(y ~ x + z, data = d)
summary(dag2.lm2)$coefficients
```

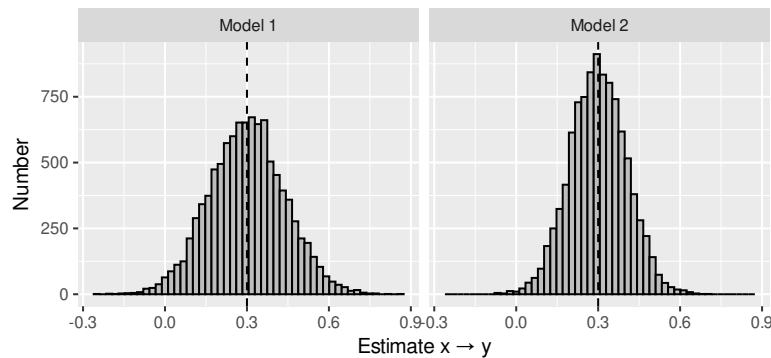
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	5.36358	0.097475	55.0254	5.9391e-75
x	0.38051	0.096327	3.9502	1.4813e-04
z	0.98832	0.093735	10.5438	9.0128e-18

The simulations below will confirm that *both* models yield unbiased estimates of the causal influence of  $x$  on  $y$ . By this standard, neither model is bad. That said, if we have the  $z$  variable at our disposal, the second model is to be preferred. The reason for this will become clearer once we've simulated a couple of thousand datasets. To this end, we could write a new function, `generate_dag2()`. But we can also reuse `generate_dag1()` and set the parameter value of  $z_x$  to 0:

```
est_dag2 <- generate_dag1(x_y = 0.3, z_y = 0.9, z_x = 0)
```

As shown in Figure 5.7, both models yield unbiased estimates of the causal influence of  $x$  on  $y$  in Equation 4. But these estimates vary less from sample to sample in model 2, that is, on average, they are closer to the true parameter value than the estimates that model 1 yields.

```
est_dag2 |>
  pivot_longer(cols = everything(),
               names_to = "Model",
               values_to = "Estimate") |>
  ggplot(aes(x = Estimate)) +
  geom_histogram(bins = 50,
                 fill = "grey", colour = "black") +
  geom_vline(xintercept = 0.3, linetype = "dashed") +
  facet_grid(cols = vars(Model)) +
```



**Figure 5.7:** Both models yield unbiased estimates of the causal influence of  $x$  on  $y$  (0.3). But these estimates vary less from sample to sample in the second model than in the first model. In other words, on average, the second model yields more precise estimates.

```
xlab("Estimate x → y") +
ylab("Number")
```

We can check this numerically. The means of the estimates in both models correspond to the true value (making allowance for simulation error). But compared to the estimates by the first model, the standard deviation of the estimates by the second model is about 25% lower.

```
apply(est_dag2, 2, mean)

Model 1 Model 2
0.30060 0.30018

apply(est_dag2, 2, sd)

Model 1 Model 2
0.13736 0.10291
```

What's happened here is that by including the  $z$  variable in the model – even though it doesn't act as a confounder for the causal link between  $x$  and  $y$  – we have reduced the error variance and have hence increased the precision of our estimate. So even if you don't actually care about  $z$ , it can pay dividends to still collect it and include it in your analysis.

**Tip 5.5** (Choosing control variables). When designing an experiment, one or two additional variables that you suspect to be strongly correlated to the outcome may be useful, particularly if they're not too strongly related to each other. (Else they would both essentially be doing the same job.) Such variables typically are *so* obviously related to the outcome that they are completely uninteresting in and of themselves (e.g., pretest performance). But don't collect umpteen additional variables on the off-chance that they might be related to the outcome and

help you reduce the error variance. Moreover, the decision whether or not to include some additional variable in your analysis should be taken *before* you analyse your data – don't run multiple models and then report the one that works 'best'. You completely invalidate your inferential results this way. ◇

**Exercise 5.6** (A pretest/posttest experiment (Part I)). The data for this exercise stem from a study by Hicks (2021). She investigated how well 260 children could learn German–English cognates. There were three waves of data collection: T1, T2, and T3. After the first wave, an intervention was undertaken with 120 of the children, the goal of which was to impart to them a greater awareness of cognate correspondences. The other 140 children served as the control group. Here we're interested in the question if the intervention bore fruit, that is, that children who took part in the intervention were better able to learn German–English cognates.

For the time being, we'll pretend that the assignment of children to experimental condition was done at random and on an individual basis. We're interested in the T3 data (T3cog); the pretest scores from T1 serve as the control variable (T1cog); we ignore the T2 data.

Read in the data, retaining just the columns that we actually need:

```
hicks <- read_csv(here("data", "hicks2021.csv")) |>
  select(ID, Class, Group, T1cog, T3cog)
```

One option to visualise these data is to draw a scatterplot with the pretest scores along the  $x$ -axis and the posttest scores along the  $y$ -axis and with different colours depending on the experimental condition. Further, regression lines for both conditions can be added. These are derived from separate simple regression models for the two conditions:

```
ggplot(hicks,
       aes(x = T1cog, y = T3cog,
           colour = Group)) +
  geom_point(shape = 1) +
  geom_smooth(se = FALSE, method = "lm") +
  xlab("Pretest score") +
  ylab("Posttest score")
```

Based on this plot, how would you answer the research question? What aspect of the visualisation did you base your answer on?

Now fit a linear model of the form  $\text{outcome} \sim \text{condition} + \text{control}$ . Don't forget to create a dummy variable for the condition variable. Interpret the model in terms of the research question. ◇

**Exercise 5.7** (A pretest/posttest experiment (Part II)). The children in the study by Hicks

(2021) were pupils in classes. They were assigned to the experiment's conditions in whole classes rather than on an individual basis. This induces a dependency between different data points, threatening the validity of the inferential results obtained in the previous part of the exercise.

There are a couple of possible solutions to this problem (see Vanhove, 2020, for a comparison). The easiest – and possibly the best – is to compute the mean pretest and mean posttest score for each cluster (class) and run the analysis using these averages instead. That is, assuming you named the dummy variable `n.Group`:

```
hicks_byclass <- hicks |>
  group_by(n.Group, Class) |>
  summarise(mean_T1 = mean(T1cog),
            mean_T3 = mean(T3cog))
byclass.lm <- lm(mean_T3 ~ n.Group + mean_T1, hicks_byclass)
```

Interpret the output of this model in terms of the research question. How would you report the finding in an article? Focus only on what's important. ◇

**Exercise 5.8** (A pretest/posttest experiment (Part III)). In Part II, we still assumed that random assignment was used – but on the level of the classes rather than on the level of the individual pupils. Look up in Hicks' article how the children were actually assigned to the different conditions. Briefly discuss plausible consequences. ◇

### 3 Posttreatment variables

We now turn our attention to scenarios where  $z$  is a posttreatment variable. That is, when we draw a DAG, we can arrive at  $z$  starting in  $x$  by following arrows. In the context of a randomised experiment with  $x$  as the predictor of interest, this means that  $z$  was collected *after* the randomisation. As a consequence,  $z$  may be influenced by  $x$ .

**Mediators.** First consider the DAG in Figure 5.8. Notice that there are two causal paths from  $x$  to  $y$ : one is direct ( $x \rightarrow y$ ), one is **mediated** by  $z$  ( $x \rightarrow z \rightarrow y$ ). A dataset conforming to this DAG can be simulated as follows:

```
n <- 200
x <- rnorm(n)
z <- 0.4*x + rnorm(n)
y <- 0.6*x + 1.2*z + rnorm(n)
d <- tibble(x, y, z)
```

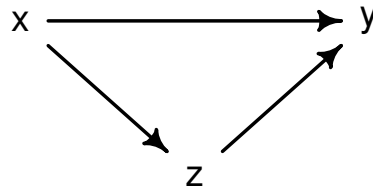


Figure 5.8

What's the causal influence of  $x$  on  $y$ ? That is, if you increase  $x$  by one unit, what change in  $y$  does this cause?

The answer to this question is *not* 0.6 units. While the direct causal effect of  $x$  on  $y$  is indeed such that a one-unit increase results in a 0.6-unit increase in  $y$ ,  $x$  also influences  $y$  via  $z$ . A one-unit increase in  $x$  results in a 0.4-unit increase in  $z$ ; and a one-unit increase in  $z$  results in a 1.2-unit increase in  $y$ . So in addition to the 0.6-unit increase that a one-unit increase in  $x$  causes directly in  $y$ , it also results in a  $0.4 \cdot 1.2 = 0.48$ -unit increase via  $z$ , for a total causal effect of  $0.6 + 0.48 = 1.08$  units increase in  $y$  per one-unit increase in  $x$ !

If we want to estimate the total causal influence of  $x$  on  $y$ , we shouldn't close any causal paths going from  $x$  to  $y$  by controlling for  $z$ , that is, we should fit a simple regression model that does *not* include  $z$ . If, however, we want to estimate the causal effect of  $x$  on  $y$  that is not mediated by  $z$ , then we *do* need to control for  $z$ . Which model you want to run depends entirely on what you want to estimate.

```
total.lm <- lm(y ~ x, data = d)
direct.lm <- lm(y ~ x + z, data = d)
summary(total.lm)$coefficients[, 1:2]
```

	Estimate	Std. Error
(Intercept)	0.21203	0.11484
x	1.18115	0.12344

```
summary(direct.lm)$coefficients[, 1:2]
```

	Estimate	Std. Error
(Intercept)	0.023547	0.071269
x	0.706786	0.080178
z	1.250551	0.069016

**Colliders.** Now consider the DAG in Figure 5.9. Here, both  $x$  and  $y$  causally affect  $z$  ( $x \rightarrow z \leftarrow y$ ). A variable in which two or more causal variables clash together is known as a **collider**. In a sense, colliders are the opposite of confounders: As long as colliders are *not* controlled for,



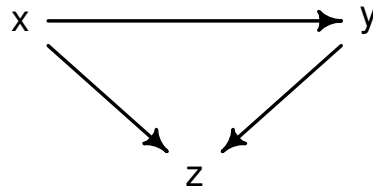


Figure 5.9

their presence does not bias the causal estimates of interest. But once they are controlled for, they may bias these causal estimates.

By way of example, let's assume the data are generated following these equations:

$$\begin{aligned}x_i &\sim \text{Normal}(0, 1^2), \\y_i &= 0.4x_i + \varepsilon_i, \\z_i &= 0.5x_i + 0.8y_i + \tau_i, \\ \varepsilon_i &\sim \text{Normal}(0, 1^2), \\ \tau_i &\sim \text{Normal}(0, 1^2),\end{aligned}$$

$i = 1, \dots, n$  with the  $x_i, \varepsilon_i, \tau_i$  independent.

The function `generate_posttreatment()` is a slight adaptation of `generate_dag1()`.

```
generate_posttreatment <- function(  
  n = 100,      # number of observations per dataset  
  sims = 10000, # number of datasets  
  x_y = 0.4,    # influence x -> y  
  x_z = 0.5,    # influence x -> z  
  y_z = 0.8     # influence y -> z  
) {  
  est_lm1 <- vector(length = sims)  
  est_lm2 <- vector(length = sims)  
  
  for (i in 1:sims) {  
    x <- rnorm(n)  
    y <- x_y*x + rnorm(n)  
    z <- x_z*x + y_z*y + rnorm(n)  
    mod_lm1 <- lm(y ~ x)  
    mod_lm2 <- lm(y ~ x + z)  
    est_lm1[[i]] <- coef(mod_lm1)[[2]]
```

```
est_lm2[[i]] <- coef(mod_lm2)[[2]]  
}  
  
tibble('Model 1' = est_lm1,  
       'Model 2' = est_lm2)  
}
```

As this simulation shows, the model without the collider ( $z$ ) is able to estimate the causal effect of  $x$  on  $y$  (0.4) in an unbiased way. The estimates for the model with the collider, by contrast, are biased.

```
est_collider <- generate_posttreatment()  
apply(est_collider, 2, mean)  
  
Model 1    Model 2  
0.39975704 0.00059706
```

Note again, though, that the model with the collider is well-suited if you want to estimate the mean difference in the  $y$  variable between two groups that differ in the  $x$  variable but whose  $z$  values are all the same. As always, whether the model is justified depends on what you want to find out and on your assumptions.

It also bears pointing out that colliders are sometimes inadvertently controlled for during the design stage. See Rohrer (2018).

**Miscellaneous.** In the DAG in Figure 5.10,  $z$  is also a posttreatment variable. This time, however, it is only indirectly influenced by  $x$ . In order to check whether it would be best to include  $z$  as a variable in the model, let's assume the following equations describe the causal links between the variables:

$$\begin{aligned}x_i &\sim \text{Normal}(0, 1^2), \\y_i &= 0.4x_i + \varepsilon_i, \\z_i &= y_i + \tau_i, \\\varepsilon_i &\sim \text{Normal}(0, 1^2), \\\tau_i &\sim \text{Normal}(0, 1^2),\end{aligned}$$

$i = 1, \dots, n$  with the  $x_i, \varepsilon_i, \tau_i$  independent.

We can reuse the `generate_posttreatment()` function and just specify the new parameter values. Note that, again, the model without the additional variable provides an unbiased estimate of the causal effect of  $x$  on  $y$ , whereas the model with this additional variable doesn't.



Figure 5.10



Figure 5.11

```

est_proxy_y <- generate_posttreatment(x_y = 0.4, x_z = 0, y_z = 1)
apply(est_proxy_y, 2, mean)

Model 1 Model 2
0.39863 0.19938
  
```

In the scenario depicted in Figure 5.10, including the additional variable will bias the estimate of interest towards zero. By spelling out what the literal meaning is of the estimated parameter for  $x$  in the second model, you should see why this is the case.

The situation is only slightly different in the scenario depicted in Figure 5.11. Here,  $z$  is directly affected by  $x$ , but not by  $y$ . We'll simulate datasets conforming to the following equations:

$$\begin{aligned}
 x_i &\sim \text{Normal}(0, 1^2), \\
 y_i &= 0.4x_i + \varepsilon_i, \\
 z_i &= x_i + \tau_i, \\
 \varepsilon_i &\sim \text{Normal}(0, 1^2), \\
 \tau_i &\sim \text{Normal}(0, 1^2),
 \end{aligned}$$

$i = 1, \dots, n$  with the  $x_i, \varepsilon_i, \tau_i$  independent.

As the simulation results show, both models yield unbiased estimates of the causal effect of  $x$  on  $y$ . But notice that the estimates vary more from sample to sample for the second model

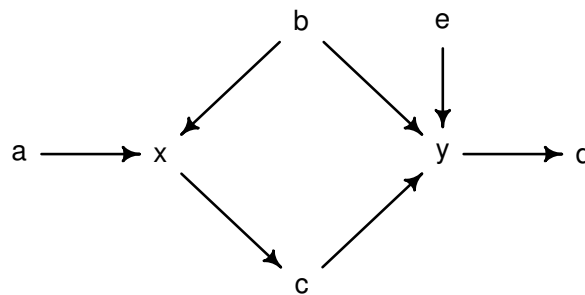


Figure 5.12

than for the first. So a given estimate resulting from the model without the additional variable is more likely to be closer the actual parameter value than a given estimate resulting from the model with the additional variable. It's hard to imagine a research question where the second model would be preferred in this scenario.

```

est_proxy_x <- generate_posttreatment(x_y = 0.4, x_z = 1, y_z = 0)
apply(est_proxy_x, 2, mean)

Model 1 Model 2
0.39975 0.40048

apply(est_proxy_x, 2, sd)

Model 1 Model 2
0.10225 0.14665

```

In sum, including posttreatment variables as predictors in the analysis is usually a bad idea. In randomised experiments, there luckily exists a simple trick for preventing that the predictors you include in the model are posttreatment variables: Just collect the predictors before carrying out the intervention!

**Exercise 5.9.** Consider Figure 5.12. Assume that you want to obtain an unbiased and maximally precise estimate the total causal effect of  $x$  on  $y$  and that all relationships shown are linear and additive. Which variables would you include in a general linear model as predictors? Justify your answer. ◇

**Exercise 5.10.** Same question for Figure 5.13. ◇

**Exercise 5.11.** Let's say that instead of merely observing and measuring all the variables in Figure 5.12, we devise a randomised experiment in which we randomly assign the participants to values of  $x$ .

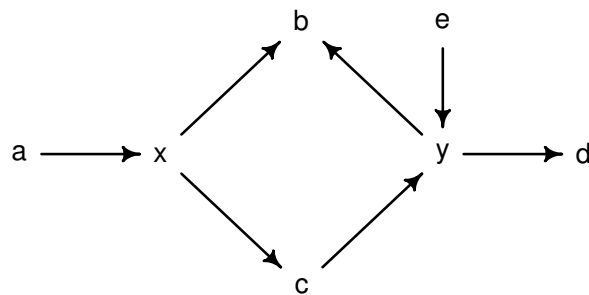


Figure 5.13

1. Draw the updated DAG.
2. Which predictors would you in the general linear model now? Justify your answer. ◇

**Exercise 5.12** (Interpreting model estimates (1)). Vanhove et al. (2019) had 1,000 short French texts written by children rated for their lexical diversity on a 9-point scale by between 2 and 18 raters each. For the purposes of this exercise, we're interested in modelling these human ratings in terms of the length of the text (the logarithmically transformed number of tokens (using base 2), `log2_ntokens`) and the type-token ratio (TTR), an easily computed metric of a text's lexical diversity.

```

d <- read_csv(here("data", "text_ratings.csv"))
lexdiv.lm <- lm(mean_rating ~ log2_ntokens + TTR, data = d)
summary(lexdiv.lm)$coefficients

```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-4.0928	0.458028	-8.9357	1.9186e-18
log2_ntokens	1.2152	0.047348	25.6655	6.1198e-112
TTR	3.8364	0.330558	11.6059	2.6707e-29

1. (Paper and pencil.) For each claim, decide whether it is correct. Justify your answers.
  - The model output shows that human raters are sensitive to differences in the type-token ratio when ratings texts for their lexical diversity.
  - The model output shows that there is a positive linear relationship between the TTR values and the mean ratings: Texts with higher TTR values tend to receive higher ratings than texts with lower TTR values.
2. (With R.) Draw a scatterplot matrix of these variables and revise your answer to the previous question if needed.

3. (Paper and pencil.) Explain what each of the three estimated model parameters literally means.
4. (Paper and pencil.) According to this model, what mean rating would you expect for a text consisting of 64 tokens with a TTR of 0.7? ◇

**Exercise 5.13** (Interpreting model estimates (2)). We're given a data set `d` with the outcome `CorrectSpoken` and four predictors:

- `n.Sex`: 0.5 for men,  $-0.5$  for women (sum coding).
- `NrLang`: Number of languages spoken for each participant (varies from 1 to 5).
- `DS.Span`: Score on the backward digit span, a working memory test (varies from 2 to 8, more is better).
- `Raven.Right`: Score on a test of fluid intelligence (varies from 0 to 35, more is better).

We fit the following model:

```
mod.lm <- lm(CorrectSpoken ~ NrLang + DS.Span + n.Sex*Raven.Right,
             data = d)
summary(mod.lm)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	8.554	1.532	5.58	1.0e-07
NrLang	1.115	0.353	3.16	1.9e-03
DS.Span	-0.097	0.330	-0.29	7.7e-01
n.Sex	-5.017	1.716	-2.92	4.0e-03
Raven.Right	0.286	0.048	6.01	1.3e-08
n.Sex:Raven.Right	0.201	0.088	2.28	2.4e-02

1. What does the estimate of  $-5.02$  for the `n.Sex` parameter refer to? That is, what does it mean literally.
2. Let's say we wanted to fit a similar model but one that has an intercept with a more useful interpretation. Explain how we could achieve this. Also explain how we can interpret the intercept in this new model.
3. Claim: We can glean from the model that, at least in this data set, participants with a high working memory capacity do not outperform participants with a low working memory capacity in terms of the `CorrectSpoken` variable.  
Is this claim correct? (Yes or no.) Justify your answer.
4. What gets computed here?

```
8.554 + 1.115*3 - 0.097*4 + 0.286*25
```

```
[1] 19
```

5. How large, according to the model, is the expected average difference in the `CorrectSpoken` variable between women with a `Raven.Right` score of 20 and women with a `Raven.Right` score of 30, keeping `NrLang` and `DS.Span` constant. It suffices to provide the computation without evaluating it. ◇

## 4 Putting it all together

The following exercise isn't specific to any one of the lectures covered up to this point. Rather, you can draw on several of the lectures to tackle it

**Exercise 5.14** (Analysing an experiment). From the abstract of Vanhove (2016):

“This article investigates whether learners are able to quickly discover simple, systematic graphemic correspondence rules between their L1 and an unknown but closely related language in a setting of receptive multilingualism.

(...)

Eighty L1 German speakers participated in a translation task with written Dutch words, most of which had a German cognate. In the first part of the translation task, participants were shown 48 Dutch words, among which either 10 cognates containing the digraph `oe` (always corresponding to a German word with `u`) or 10 cognates with the digraph `ij` (corresponding to German `ei`). During this part, participants were given feedback in the form of the correct translation. In the second (feedback-free) part of the task, participants were shown another 150 Dutch words, among which 21 cognates with `oe` and 21 cognates with `ij`.”

What I wanted to know was whether the exposure and feedback to ten words containing a specific interlingual orthographic correspondence (i.e., `oe-u` or `ei-ij`) was sufficient for the participants to pick up on this correspondence and use their knowledge in translating new words containing the correspondence.

First, we read in the data and compute, for each participant, the proportion of words in the second (feedback-free) part of the task which they translated correctly in each category (cognates containing `oe`, cognates containing `ij`, other cognates, non-cognates).

```
d <- read_csv(here("data", "correspondencerules.csv"))
d_perParticipant <- d |>
  filter(Block != "training") |>
```

```
group_by(Subject, LearningCondition, Category, WSTRight) |>
summarise(ProportionCorrect = mean(Correct == "yes"),
          .groups = "drop")
```

I've also retained a measure of the participants' L1 (German) vocabulary knowledge, namely the `WSTRight` variable. Use `View(d_perParticipant)` to inspect the structure of the cleaned-up and restructured dataset.

Analyse this dataset to answer the research question.

Hints:

- There are several acceptable ways of analysing these data.
- All of them involve plotting the data :)
- Note that `LearningCondition` varies *between* participants whereas `Category` varies *within* participants.
- First try to work out on a piece of paper what it is you want to do. You can always ask someone to help you implement your idea in R.
- Ask for help and feedback! :) ◇

## Summary

- A multiple regression model is not just multiple simple regressions run at once. The meaning of the parameter estimates changes if you add or remove predictors to or from the model. Also see Morrissey & Ruxton (2018) and Vanhove (2021).
- The decision which predictors to include in a model depends on what it is exactly you want to estimate and how you think different variables may be causally related. In my experience, people's difficulties with regression models aren't so much statistical in nature as due to their not having worked out what they actually want to find out.
- Know what the literal meaning of the estimated model parameters is before you interpret them in terms of the subject matter.
- As shown in the exercises, plots of the actual data may help prevent both you and your readership from interpreting the model output incorrectly.



## Further reading

On the limits of statistical control in observational studies, see Christenfeld et al. (2004), Huitema (2011, Part VII) and Westfall & Yarkoni (2016). On the utility of statistical control in randomised experiments, see Vanhove (2015) and references therein. For an accessible introduction to DAGs, confounders and colliders, see Rohrer (2018).

## References

- Christenfeld, Nicholas J. S., Richard P. Sloan, Douglas Carroll & Sander Greenland. 2004. Risk factors, confounding, and the illusion of statistical control. *Psychosomatic Medicine* 66. 868–875. doi:10.1097/01.psy.0000140008.70959.41.
- Hicks, Nina Selina. 2021. Exploring systematic orthographic crosslinguistic similarities to enhance foreign language vocabulary learning. *Language Teaching Research* doi:10.1177/13621688211047353.
- Huitema, Bradley E. 2011. *The analysis of covariance and alternatives: Statistical methods for experiments, quasi-experiments, and single-case studies*. Hoboken, NJ: Wiley.
- McElreath, Richard. 2020. *Statistical rethinking: A Bayesian course with examples in R and Stan*. Boca Raton, FL: CRC Press 2nd edn.
- Morrissey, Michael B. & Greame D. Ruxton. 2018. Multiple regression is not multiple regressions: The meaning of multiple regression and the non-problem of collinearity. *Philosophy, Theory, and Practice in Biology* 10(3). doi:10.3998/ptpbio.16039257.0010.003.
- Rohrer, Julia M. 2018. Thinking clearly about correlations and causation: Graphical causal models for observational data. *Advances in Methods and Practices in Psychological Science* 1(1). 27–42. doi:10.1177/2515245917745629.
- Vanhove, Jan. 2015. Analyzing randomized controlled interventions: Three notes for applied linguists. *Studies in Second Language Learning and Teaching* 5. 135–152. doi:10.14746/ssllt.2015.5.1.7.
- Vanhove, Jan. 2016. The early learning of interlingual correspondence rules in receptive multilingualism. *International Journal of Bilingualism* 20(5). 580–593. doi:10.1177/1367006915573338.
- Vanhove, Jan. 2020. Capitalising on covariates in cluster-randomised experiments. *PsyArXiv Preprints* doi:10.31234/osf.io/ef4zc.
- Vanhove, Jan. 2021. Collinearity isn't a disease that needs curing. *Meta-Psychology* 5. doi:10.15626/MP.2021.2548.

- Vanhove, Jan, Audrey Bonvin, Amelia Lambelet & Raphael Berthele. 2019. Predicting perceptions of the lexical richness of short French, German, and Portuguese texts using text-based indices. *Journal of Writing Research* 10(3). 499–525. doi:10.17239/jowr-2019.10.03.04.
- Westfall, Jacob & Tal Yarkoni. 2016. Statistically controlling for confounding constructs is harder than you think. *PLOS ONE* 11(3). e0152719. doi:10.1371/journal.pone.0152719.