

# Homework 1

Samuel Rosenberg

Due September 13, 2021 at 11:59pm

## 1 Instructions

**Setup.** Pull the latest version of this assignment from Github and set your working directory to `stat-961-fall-2021/homework/homework-1`. Consult the [getting started guide](#) if you need to brush up on R, LaTeX, or Git.

**Collaboration.** The collaboration policy is as stated on the Syllabus:

“Students are permitted to work together on homework assignments, but solutions must be written up and submitted individually. Students must disclose any sources of assistance they received; furthermore, they are prohibited from verbatim copying from any source and from consulting solutions to problems that may be available online and/or from past iterations of the course.”

In accordance with this policy,

*Please list anyone you discussed this homework with:*

*Please list what external references you consulted (e.g. articles, books, or websites):*

**Writeup.** Use this document as a starting point for your writeup, adding your solutions between `\begin{sol}` and `\end{sol}`. See the [preparing reports guide](#) for guidance on compilation, creation of figures and tables, and presentation quality. Show all the code you wrote to produce your numerical results, and include complete derivations typeset in LaTeX for the mathematical questions.

**Programming.** The `tidyverse` paradigm for data manipulation (`dplyr`) and plotting (`ggplot2`) are strongly encouraged, but points will not be deducted for using base R.

```
library(tidyverse)
```

**Grading.** Each sub-part of each problem will be worth 3 points: 0 points for no solution or completely wrong solution; 1 point for some progress; 2 points for a mostly correct solution; 3 points for a complete and correct solution modulo small flaws. The presentation quality of the solution for each problem (as exemplified by the guidelines in Section 3 of the [preparing reports guide](#)) will be evaluated out of an additional 3 points.

**Submission.** Compile your writeup to PDF and submit to [Gradescope](#).

**Problem 1. Change of basis.** (Adapted from Agresti Ex. 1.17)

Let  $\mathbf{X}$  and  $\mathbf{X}'$  be full-rank  $n \times p$  model matrices.

- Show that  $C(\mathbf{X}) = C(\mathbf{X}')$  if and only if  $\mathbf{X}' = \mathbf{X}\mathbf{A}$  for some nonsingular  $p \times p$  matrix  $\mathbf{A}$ . In plain language, express what the operation  $\mathbf{X} \mapsto \mathbf{X}\mathbf{A}$  does to the columns of  $\mathbf{X}$  (one sentence is sufficient).
- Let  $\hat{\beta}$  and  $\hat{\beta}'$  be the least squares solutions obtained from regressing a response vector  $\mathbf{y}$  on  $\mathbf{X}$  and  $\mathbf{X}' \equiv \mathbf{X}\mathbf{A}$ , respectively, where  $\mathbf{A}$  is a nonsingular  $p \times p$  matrix. What is the relationship between  $\hat{\beta}$  and  $\hat{\beta}'$  (express the latter in terms of the former)? Justify your answer.
- Consider the linear model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon; \quad \epsilon \sim (0, \sigma^2), \quad (1)$$

so that  $\mathbf{X} = [\mathbf{1}, \mathbf{x}_{*1}, \mathbf{x}_{*2}]$  for columns  $\mathbf{x}_{*j} \equiv (x_{1j}, \dots, x_{nj})^T$ ,  $j \in \{1, 2\}$ . Sometimes it is useful to center the predictors by subtracting their means:

$$\mathbf{x}'_{*j} \equiv \mathbf{x}_{*j} - \bar{x}_j \mathbf{1}; \quad \bar{x}_j \equiv \frac{1}{n} \sum_{i=1}^n x_{ij}, \quad j \in \{1, 2\}.$$

Defining  $\mathbf{X}' \equiv [\mathbf{1}, \mathbf{x}'_{*1}, \mathbf{x}'_{*2}]$ , find the matrix  $\mathbf{A}$  such that  $\mathbf{X}' = \mathbf{X}\mathbf{A}$  ( $\mathbf{A}$  may itself be expressed in terms of  $\mathbf{X}$ ). Express the coefficient estimates from the centered regression ( $\hat{\beta}'_0, \hat{\beta}'_1, \hat{\beta}'_2$ ) in terms of those from the original regression ( $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$ ).

- Let  $w \in \{a, b, c\}$  be a categorical variable with three levels. Define  $x_1 \equiv \mathbb{1}(w = b)$  and  $x_2 \equiv \mathbb{1}(w = c)$ , and consider the linear regression (1). This corresponds to regressing  $y$  on the categorical variable  $w$ , with baseline category  $a$ . Sometimes a different baseline category may make more sense, e.g. category  $b$ . In this case, we would define  $x'_1 \equiv \mathbb{1}(w = a)$  and  $x'_2 \equiv \mathbb{1}(w = c)$ . Defining  $\mathbf{X} \equiv [\mathbf{1}, \mathbf{x}_{*1}, \mathbf{x}_{*2}]$  and  $\mathbf{X}' \equiv [\mathbf{1}, \mathbf{x}'_{*1}, \mathbf{x}'_{*2}]$ , find the matrix  $\mathbf{A}$  such that  $\mathbf{X}' = \mathbf{X}\mathbf{A}$ . Express the coefficient estimates from the transformed regression ( $\hat{\beta}'_0, \hat{\beta}'_1, \hat{\beta}'_2$ ) in terms of those from the original regression ( $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$ ). What are the interpretations of the original and transformed coefficients, and why do the relationships between these coefficients derived above make sense in terms of these interpretations?

**Solution 1.**

- Because each column space is full rank (rank  $p$ ), the columns of each model matrix form a basis; we write each of these bases as follows:  $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$  and  $\{\mathbf{x}'_1, \dots, \mathbf{x}'_p\}$  respectively.  
 $\Rightarrow$ : Suppose that  $C(\mathbf{X}) = C(\mathbf{X}')$ . Since  $C(\mathbf{X}) = C(\mathbf{X}')$ ,  $\mathbf{x}'_j \in C(\mathbf{X})$ . This means that  $\mathbf{x}'_j = \sum_{k=1}^p a_{kj} \mathbf{x}_k = \sum_{k=1}^p a_{kj} (x_{1k}, \dots, x_{nk})^T$  for some  $a_{kj}$ . So,  $x'_{ij} = \sum_{k=1}^p a_{kj} x_{ik} = (x_{i1}, \dots, x_{ip})(a_{1j}, \dots, a_{pj})^T$ . This is equivalent to saying that  $\mathbf{X}' = \mathbf{X}\mathbf{A}$ , for the  $p \times p$  matrix  $\mathbf{A}$  with entry  $a_{ij}$  in the  $i$ -th row and  $j$ -th column.  
 Suppose for contradiction that  $\mathbf{A}$  is singular, then  $\text{rank}(\mathbf{A}) < p$ . So, there is some column  $\mathbf{a}_\ell = (a_{1\ell}, \dots, a_{p\ell})^T$  that can be written as a linear combination of the other  $\mathbf{a}_j$ ; i.e.  $\mathbf{a}_\ell = \beta_1 \mathbf{a}_1 + \dots + \beta_{\ell-1} \mathbf{a}_{\ell-1} + \beta_{\ell+1} \mathbf{a}_{\ell+1} + \dots + \beta_p \mathbf{a}_p$  and  $a_{i\ell} = \beta_1 a_{i1} + \dots + \beta_{\ell-1} a_{i\ell-1} + \beta_{\ell+1} a_{i\ell+1} + \dots + \beta_p a_{ip}$ .

Then

$$\begin{aligned}
 \mathbf{x}'_\ell &= \sum_{k=1}^p a_{k\ell} \mathbf{x}_k \\
 &= a_{1\ell} \mathbf{x}_1 + \sum_{k=2}^p a_{k\ell} \mathbf{x}_k \\
 &= \left( \sum_{k=2}^p \beta_k a_{1k} \mathbf{x}_1 \right) + \sum_{k=2}^p a_{k\ell} \mathbf{x}_k \\
 &= \sum_{k=2}^p (\beta_k a_{1k} + a_{k\ell}) \mathbf{x}_k.
 \end{aligned}$$

But this means that  $\mathbf{x}'_\ell$  is a non-unique linear combination of the  $\mathbf{x}_k$ , so  $\text{rank}(C(\mathbf{X})) < p$ . This is a contradiction since we assumed  $\mathbf{X}$  was full rank, so  $\mathbf{A}$  must be nonsingular.

$\Leftarrow$ : Now suppose that  $\mathbf{X}' = \mathbf{X}\mathbf{A}$  for some nonsingular  $p \times p$  matrix  $\mathbf{A}$ . Then each entry  $x'_{ij}$  of  $\mathbf{X}'$  has the following form:  $x'_{ij} = (x_{i1}, \dots, x_{ip})(a_{1j}, \dots, a_{pj})^T = \sum_{k=1}^p a_{kj} x_{ik}$ . Consequently, the  $j$ -th column of  $\mathbf{X}'$  is  $\mathbf{x}'_j = \sum_{k=1}^p a_{kj} \mathbf{x}_k$ . Thus,  $\text{span}\{\mathbf{x}'_1, \dots, \mathbf{x}'_p\} = C(\mathbf{X}') \subseteq C(\mathbf{X})$ . But  $\dim(C(\mathbf{X}')) = \dim(C(\mathbf{X}))$ , so we must have that  $C(\mathbf{X}) = C(\mathbf{X}')$ .

Thus,  $C(\mathbf{X}) = C(\mathbf{X}')$  if and only if  $\mathbf{X}' = \mathbf{X}\mathbf{A}$  for some nonsingular  $p \times p$  matrix  $\mathbf{A}$ .

The operation  $\mathbf{X} \mapsto \mathbf{X}\mathbf{A}$  converts the columns of  $\mathbf{X}$  to use the same units as  $\mathbf{X}'$  (i.e. it is a change of bases).

- (b) Recall that the least square estimators for the linear regression of  $\mathbf{y}$  on  $\mathbf{X}$  and  $\mathbf{X}'$  are  $\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$  and  $\hat{\beta}' = ((\mathbf{X}')^T \mathbf{X}')^{-1} (\mathbf{X}')^T \mathbf{y}$  respectively (e.g. by equation 2.3 in Agresti). Taking  $\mathbf{X}' \equiv \mathbf{X}\mathbf{A}$  and using properties of the matrix transpose and matrix inverse (since  $\mathbf{X}$ ,  $\mathbf{X}'$ ,  $\mathbf{A}$  are full rank), we have that

$$\begin{aligned}
 \hat{\beta}' &= ((\mathbf{X}')^T \mathbf{X}')^{-1} (\mathbf{X}')^T \mathbf{y} \\
 &= ((\mathbf{X}\mathbf{A})^T \mathbf{X}\mathbf{A})^{-1} (\mathbf{X}\mathbf{A})^T \mathbf{y} \\
 &= (\mathbf{A}^T \mathbf{X}^T \mathbf{X} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{X}^T \mathbf{y} \\
 &= (\mathbf{X}^T \mathbf{X} \mathbf{A})^{-1} (\mathbf{A}^T)^{-1} \mathbf{A}^T \mathbf{X}^T \mathbf{y} \\
 &= (\mathbf{X}^T \mathbf{X} \mathbf{A})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= \mathbf{A}^{-1} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \\
 &= \mathbf{A}^{-1} \hat{\beta}.
 \end{aligned}$$

- (c) We take

$$\mathbf{A} = \begin{pmatrix} 1 & -\bar{x}_1 & -\bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

Then

$$\begin{aligned}
 \mathbf{X}\mathbf{A} &= \begin{pmatrix} 1 & x_{11} & x_{12} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{pmatrix} \begin{pmatrix} 1 & -\bar{x}_1 & -\bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \\
 &= \begin{pmatrix} 1 & x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 \end{pmatrix} \\
 &= \mathbf{X}'
 \end{aligned}$$

as desired.

Note that

$$\mathbf{A} \begin{pmatrix} 1 & \bar{x}_1 & \bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & -\bar{x}_1 & -\bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & \bar{x}_1 & \bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \mathbf{I}_3,$$

so

$$\mathbf{A}^{-1} = \begin{pmatrix} 1 & \bar{x}_1 & \bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}.$$

As we saw in (b),  $\hat{\beta}' = \mathbf{A}^{-1}\hat{\beta}$ , so

$$\begin{aligned}
 \hat{\beta}' &= \mathbf{A}^{-1}\hat{\beta} \\
 &= \begin{pmatrix} 1 & \bar{x}_1 & \bar{x}_2 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \\
 &= \begin{pmatrix} \hat{\beta}_0 + \bar{x}_1\hat{\beta}_1 + \bar{x}_2\hat{\beta}_2 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix}.
 \end{aligned}$$

(d) We take

$$\mathbf{A} = \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & 0 \\ 0 & -1 & 1 \end{pmatrix}.$$

Note that  $x_2 = x'_2$  implies that  $x_{i2} = x'_{i2}$  for all  $i$ . Also,  $x'_1$  is 1 if  $w = a$  (and  $x_1 = x_2 = 0$ ) and 0 otherwise, so  $x'_1 = 1 - x_1 - x_2$  and  $x'_{i1} = 1 - x_{i1} - x_{i2}$ .

Then

$$\begin{aligned}
 \mathbf{X}\mathbf{A} &= \begin{pmatrix} 1 & x_{11} & x_{12} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{pmatrix} \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & 0 \\ 0 & -1 & 1 \end{pmatrix} \\
 &= \begin{pmatrix} 1 & 1 - x_{11} - x_{12} & x_{12} \\ \vdots & \vdots & \vdots \\ 1 & 1 - x_{n1} - x_{n2} & x_{n2} \end{pmatrix} \\
 &= \begin{pmatrix} 1 & x'_{11} & x'_{12} \\ \vdots & \vdots & \vdots \\ 1 & x'_{n1} & x'_{n2} \end{pmatrix} \\
 &= \mathbf{X}'
 \end{aligned}$$

as desired.

Note that  $\mathbf{A}^2 = \mathbf{I}_3$ , so  $\mathbf{A} = \mathbf{A}^{-1}$ . As we saw in (b),  $\hat{\beta}' = \mathbf{A}^{-1}\hat{\beta}$ , so

$$\begin{aligned}
 \hat{\beta}' &= \mathbf{A}^{-1}\hat{\beta} \\
 &= \begin{pmatrix} 1 & 1 & 0 \\ 0 & -1 & 0 \\ 0 & -1 & 1 \end{pmatrix} \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \\
 &= \begin{pmatrix} \hat{\beta}_0 + \hat{\beta}_1 \\ -\hat{\beta}_1 \\ -\hat{\beta}_1 + \hat{\beta}_2 \end{pmatrix}.
 \end{aligned}$$

We interpret the coefficients as follows:  $\hat{\beta}_0$  is the mean for the subgroup with  $w = a$ , while  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are the differences in the means of the subgroups where  $w = b$  and  $w = a$  or  $w = c$  and  $w = a$  respectively. Likewise,  $\hat{\beta}'_1$  is the mean for the subgroup with  $w = b$ , while  $\hat{\beta}'_1$  and  $\hat{\beta}'_2$  are the differences in the means of the subgroups where  $w = a$  and  $w = b$  or  $w = c$  and  $w = a$  respectively.

Intuitively, the relationships between the coefficients make sense:  $\hat{\beta}'_0 = \hat{\beta}_0 + \hat{\beta}_1$  is the mean of the subgroup with  $w = a$  plus the difference in the subgroups with  $w = b$  and  $w = a$ ; i.e. it is  $\hat{\beta}'_0$ , the mean of the subgroup with  $w = b$ . Likewise,  $\hat{\beta}'_1 = -\hat{\beta}_1 = \hat{\beta}_0 - (\hat{\beta}_0 + \hat{\beta}_1)$  is the difference between the means of the subgroups where  $w = a$  and  $w = b$ . Finally,  $\hat{\beta}'_2 = -\hat{\beta}_1 + \hat{\beta}_2 = (\hat{\beta}_0 + \hat{\beta}_2) - (\hat{\beta}_0 + \hat{\beta}_1)$  is the difference between the means of the subgroups where  $w = c$  and  $w = b$ .

**Problem 2. Predictor correlation.** (Adapted from Agresti Ex. 2.9)

Consider the linear regression

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon; \quad \epsilon \sim (0, \sigma^2),$$

with observed predictor vectors denoted  $\mathbf{x}_{*1} \equiv (x_{11}, \dots, x_{n1})^T$  and  $\mathbf{x}_{*2} \equiv (x_{12}, \dots, x_{n2})^T$ . (This is the same setup as in Problem 1(c).)

- Suppose  $\mathbf{x}_{*1}$  and  $\mathbf{x}_{*2}$  have sample correlation  $\rho \in (-1, 1)$ . In terms of  $\rho$ , what is the correlation between the estimates  $\hat{\beta}_1$  and  $\hat{\beta}_2$  (which are random variables due to the randomness in  $\epsilon$ )?
- To build intuition for the preceding result, consider the extreme case when  $\mathbf{x}_{*1} = \mathbf{x}_{*2}$ . In this case,  $\rho = 1$  and the regression is not identifiable. For a fixed parameter vector  $(\beta_0^0, \beta_1^0, \beta_2^0)$ , write down the set  $\mathcal{S}$  of parameter vectors  $(\beta_0, \beta_1, \beta_2)$  giving the same value of  $\mathbb{E}[\mathbf{y}]$  as  $(\beta_0, \beta_1, \beta_2) = (\beta_0^0, \beta_1^0, \beta_2^0)$ . In what sense does the result in part (a) reflect the relationship between  $\beta_1$  and  $\beta_2$  for  $(\beta_0, \beta_1, \beta_2) \in \mathcal{S}$ ? (Ignore the fact that the case  $\rho = 1$  is not covered in part (a).)
- Suppose  $z_1, z_2 \stackrel{\text{i.i.d.}}{\sim} N(0, 1)$ , and  $x_1 \equiv z_1 + 0.5z_2$  and  $x_2 \equiv z_1 - 0.5z_2$ . What is the correlation between the random variables  $x_1$  and  $x_2$ ? Suppose the predictors in each row  $\{(x_{i1}, x_{i2})\}_{i=1}^n$  are a sample from this joint distribution. Roughly what do we expect to be the sample correlation between  $\mathbf{x}_{*1}$  and  $\mathbf{x}_{*2}$ ? Fixing  $\mathbf{x}_{*1}$  and  $\mathbf{x}_{*2}$  at their realizations, roughly what do we expect to be the correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$ ?
- To check the conclusions in part (b), run a numerical simulation with  $n = 100$ ,  $\sigma^2 = 1$ ,  $(\beta_0, \beta_1, \beta_2) = (0, 1, 2)$ , and  $\epsilon \sim N(0, \sigma^2)$ . Sample one realization of  $\mathbf{x}_{*1}$  and  $\mathbf{x}_{*2}$ , generate 250 realizations of the response  $\mathbf{y}$ , and for each realization calculate least squares estimates  $\hat{\beta}$ . Summarize the results of your simulation by creating scatter plots of  $\mathbf{x}_{*2}$  versus  $\mathbf{x}_{*1}$  and  $\hat{\beta}_2$  versus  $\hat{\beta}_1$ , with the title of each plot containing the sample correlations of the data it displays. On the scatter plot of  $\hat{\beta}_2$  versus  $\hat{\beta}_1$ , indicate the theoretical expected value of  $(\hat{\beta}_1, \hat{\beta}_2)$  with a red point. Display these two scatter plots side by side using `cowplot::plot_grid`. Do the sample correlations match what you predicted in part (c)?

**Solution 2.**

- (a) Note that the model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon, \quad \epsilon \sim (0, \sigma^2)$$

is equivalent to the model

$$y = \beta'_0 + \beta_1(x_1 - \bar{x}_1) + \beta_2(x_2 - \bar{x}_2), \quad \epsilon \sim (0, \sigma^2), \quad \text{with } \beta'_0 = \beta_0 + \beta_1 \bar{x}_1 + \beta_2 \bar{x}_2.$$

So, the coefficients  $\hat{\beta}_1, \hat{\beta}_2$  are unaffected by mean centering. Note that the regression  $y = \beta_1(x_1 - \bar{x}_1) + \beta_2(x_2 - \bar{x}_2)$  has the same coefficients of this second model by orthogonality;  $\bar{x}_1$  is the projection of  $x_1$  onto 1 and  $\bar{x}_2$  is the projection of  $x_2$  onto 1, so  $x_1 - \bar{x}_1, x_2 - \bar{x}_2 \perp 1$ .

Thus, the covariance matrix  $\Sigma$  for the initial model has the following form:

$$\Sigma = \begin{pmatrix} s_{11} & s_{12} & s_{13} \\ s_{21} & \mathbf{\Sigma}' & \\ s_{31} & & \end{pmatrix},$$

where  $\mathbf{\Sigma}'$  is the covariance matrix for  $\hat{\beta}_1$  and  $\hat{\beta}_2$  under the model  $y = \beta_1(x_1 - \bar{x}_1) + \beta_2(x_2 - \bar{x}_2)$ .

Recall that for an identifiable regression with model matrix  $\mathbf{X}$ , the covariance matrix is given by  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ . Note that our model matrix for the third (reduced) regression is

$$\mathbf{X} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 \\ \vdots & \vdots \\ x_{n1} - \bar{x}_1 & x_{n2} - \bar{x}_2 \end{pmatrix}.$$

Then

$$\mathbf{X}^T \mathbf{X} = \begin{pmatrix} \|\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1}\|^2 & (\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}) \\ (\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}) & \|\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}\|^2 \end{pmatrix}$$

and

$$(\mathbf{X}^T \mathbf{X})^{-1} = \frac{1}{\det(\mathbf{X}^T \mathbf{X})} \begin{pmatrix} \|\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}\|^2 & -(\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}) \\ -(\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}) & \|\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1}\|^2 \end{pmatrix}.$$

We have that

$$\text{cov}(\hat{\beta}_1, \hat{\beta}_2) = \sigma^2 \Sigma_{23} = \sigma^2 \Sigma'_{12} = -\frac{\sigma^2}{\det(\mathbf{X}^T \mathbf{X})} (\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}),$$

$$\sigma_{\hat{\beta}_1}^2 = \sigma^2 \Sigma_{22} = \sigma^2 \Sigma'_{11} = \frac{\sigma^2}{\det(\mathbf{X}^T \mathbf{X})} \|\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}\|^2,$$

and

$$\sigma_{\hat{\beta}_2}^2 = \sigma^2 \Sigma_{33} = \sigma^2 \Sigma'_{22} = \frac{\sigma^2}{\det(\mathbf{X}^T \mathbf{X})} \|\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1}\|^2.$$

Define  $C := \frac{\sigma^2}{\det(\mathbf{X}^T \mathbf{X})}$ . Then

$$\begin{aligned} \text{corr}(\hat{\beta}_1, \hat{\beta}_2) &= \frac{\text{cov}(\hat{\beta}_1, \hat{\beta}_2)}{\sigma_{\hat{\beta}_1} \sigma_{\hat{\beta}_2}} \\ &= -\frac{C(\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1})^T (\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1})}{\sqrt{C} \|\mathbf{x}_{*2} - \bar{x}_2 \mathbf{1}\| \sqrt{C} \|\mathbf{x}_{*1} - \bar{x}_1 \mathbf{1}\|} \\ &= -\rho. \end{aligned}$$

That is, the correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$  is  $-\rho$ .

- (b) Suppose that we have that  $(\beta_0^0, \beta_1^0, \beta_2^0)$  is a parameter vector for the case when  $\mathbf{x}_{*1} = \mathbf{x}_{*2}$ . Consider  $\mathcal{S} := \{(\beta_0, \beta_1, \beta_2) \in \mathbb{R}^3 : \beta_0 = \beta_0^0, \beta_1 + \beta_2 = \beta_1^0 + \beta_2^0\}$ . Consider the models where  $y^0$  and  $y$  have the same formula as (1c), with  $(\beta_0^0, \beta_1^0, \beta_2^0)$  and  $(\beta_0, \beta_1, \beta_2) \in \mathcal{S}$  as their respective parameter vectors. Then

$$\mathbb{E}[\mathbf{y}^0] = \mathbb{E}[\beta_0^0 \mathbf{1} + \beta_1^0 \mathbf{x}_{*1} + \beta_2^0 \mathbf{x}_{*2} + \boldsymbol{\epsilon}] = \beta_0^0 \mathbf{1} + \beta_1^0 \mathbf{x}_{*1} + \beta_2^0 \mathbf{x}_{*2} = \beta_0^0 \mathbf{1} + (\beta_1^0 + \beta_2^0) \mathbf{x}_{*1},$$

where the last equality holds since  $\mathbf{x}_{*1} = \mathbf{x}_{*2}$ . But

$$\beta_0^0 \mathbf{1} + (\beta_1^0 + \beta_2^0) \mathbf{x}_{*1} = \beta_0 \mathbf{1} + (\beta_1 + \beta_2) \mathbf{x}_{*1} = \beta_0 \mathbf{1} + \beta_1 \mathbf{x}_{*1} + \beta_2 \mathbf{x}_{*2} = \mathbb{E}[\mathbf{y}],$$

where the first equality holds by construction of  $\mathcal{S}$ . So,  $\mathbb{E}[\mathbf{y}^0] = \mathbb{E}[\mathbf{y}]$ ; i.e. any parameter vector in  $\mathcal{S}$  will yield the same expectation  $\mathbb{E}[\mathbf{y}]$  as the original fixed parameter vector.

Take  $c := \beta_1^0 + \beta_2^0$ , which is one of the constraints that appears when defining  $\mathcal{S}$  (i.e.  $\beta_1 + \beta_2 = c$ ). When we rearrange this expression, we get  $\beta_1 = c - \beta_2$ . In other words, perfect correlation of the predictors  $x_1$  and  $x_2$  leads to a constraint of perfect anticorrelation of the coefficients  $\beta_1$  and  $\beta_2$ , where the slope of the linear relationship between  $\beta_1$  and  $\beta_2$  is determined to be the negative of the slope of the linear relationship between  $x_1$  and  $x_2$ . This corresponds to  $\rho = 1$  implying that  $\text{cor}(\hat{\beta}_1, \hat{\beta}_2) = -1$ , if the result from (a) could be extended to  $\rho = \pm 1$ .

(c) Note that  $\mathbb{E}[x_1] = \mathbb{E}[z_1 + 0.5z_2] = \mathbb{E}[z_1] + 0.5\mathbb{E}[z_2] = 0$ ; likewise,  $\mathbb{E}[x_2] = 0$ . Then  $\text{cov}(x_1, x_2) = \mathbb{E}[(x_1 - \mathbb{E}[x_1])(x_2 - \mathbb{E}[x_2])] = \mathbb{E}[x_1x_2] = \mathbb{E}[(z_1 + 0.5z_2)(z_1 - 0.5z_2)] = \mathbb{E}[z_1^2] - 0.25\mathbb{E}[z_2^2]$ . But we have that  $z_1^2, z_2^2 \sim \chi_1^2$ , so  $\mathbb{E}[z_1] = \mathbb{E}[z_2] = 1$  and  $\text{cov}(x_1, x_2) = 3/4$ . Also,  $\sigma_{x_1}^2 = \sigma_{x_2}^2 = \text{var}(z_1) + 0.25\text{var}(z_2) = 5/4$ . So, the correlation between  $x_1$  and  $x_2$  as random variables is  $\text{corr}(x_1, x_2) = \frac{3/4}{\sqrt{5/4}} = 3/5$ .

We expect the sample correlation between  $\mathbf{x}_{*1}$  and  $\mathbf{x}_{*2}$  to be roughly equal to the correlation between  $x_1$  and  $x_2$  as random variables; that is,  $3/5$ . This is because the sample correlation is an asymptotically unbiased estimator of the correlation between random variables. Using our result from (a), we thus expect the correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$  to be approximately  $-3/5$ .

```
(d) library(cowplot)
library(latex2exp)

### Set simulation parameters
# Set values of model parameters
beta_0 <- 0
beta_1 <- 1
beta_2 <- 2

n <- 100
sigma <- 1

num_sim <- 250

### Fix one realization of  $x_{*1}$ ,  $x_{*2}$  based on the distribution from (c)
# Generate n standard normal IID values for  $z_1$ ,  $z_2$ 
z_1 <- rnorm(n, mean = 0, sd = 1)
z_2 <- rnorm(n, mean = 0, sd = 1)

# Compute  $x_{*1}$ ,  $x_{*2}$ 
x_1 <- z_1 + 0.5 * z_2
x_2 <- z_1 - 0.5 * z_2

### Generate num_sim realizations of the response y and calculate  $\widehat{\beta}$ 
# Vectors to store  $\widehat{\beta}$  estimates
beta_1_hat <- c()
beta_2_hat <- c()

# Run loop num_sim times
```



```

for(i in 1:num_sim){
  # Get \epsilon and calculate y
  epsilon <- rnorm(n, mean = 0, sd = sigma)

  y <- beta_0 + beta_1 * x_1 + beta_2 * x_2 + epsilon

  # Combine data into data frame and run regression
  sim_data <- data.frame(y = y, x_1 = x_1, x_2 = x_2)
  sim_lm <- lm(y ~ x_1 + x_2, data = sim_data)

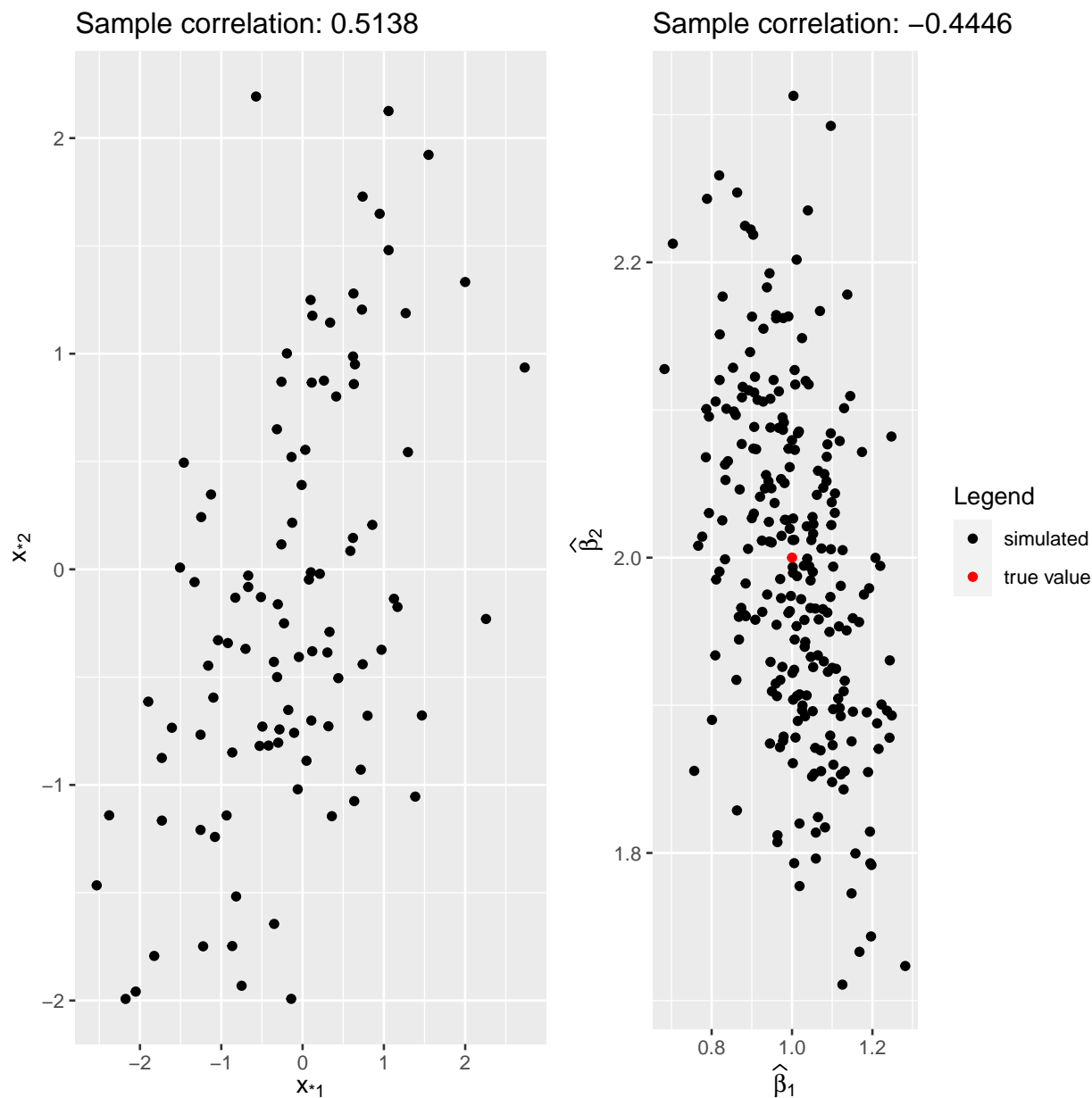
  # Extract and store \widehat{\beta} estimates
  beta_hat <- sim_lm$coefficients
  beta_1_hat <- c(beta_1_hat, beta_hat[2])
  beta_2_hat <- c(beta_2_hat, beta_hat[3])
}

x_plt <-
  ggplot(data = data.frame(x_1 = x_1, x_2 = x_2),
    mapping = aes(x = x_1, y = x_2)) +
    geom_point() +
    xlab(TeX("$x_{*1}$")) +
    ylab(TeX("$x_{*2}$")) +
    ggtitle(paste("Sample correlation:",
      round(cor(x = x_1, y = x_2), digits = 4)))

beta_plt <-
  ggplot(data = data.frame(beta_1_hat = c(beta_1_hat, beta_1),
    beta_2_hat = c(beta_2_hat, beta_2),
    labels = c(rep("simulated", times = num_sim),
      "true value")),
    mapping = aes(x = beta_1_hat,
      y = beta_2_hat,
      color = labels)) +
    geom_point() +
    xlab(TeX("$\\widehat{\\beta}_1$")) +
    ylab(TeX("$\\widehat{\\beta}_2$")) +
    ggtitle(paste("Sample correlation:",
      round(cor(x = beta_1_hat, y = beta_2_hat), digits = 4))) +
    scale_color_manual(name = "Legend", values = c("black", "red"))

plot_grid(plotlist = list(x_plt, beta_plt), nrow = 1, ncol = 2)

```



Note that we have that the sample correlation between  $x_{*1}$  and  $x_{*2}$  is roughly  $3/5 = 0.60$ , the correlation between  $x_1$  and  $x_2$  calculated in (c). Likewise, the sample correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$  is roughly  $-3/5 = -0.60$ , which was predicted as the correlation between  $\hat{\beta}_1$  and  $\hat{\beta}_2$  in (c).

### Problem 3. Data analysis: Anorexia treatment. (Adapted from Agresti Ex. 1.24)

For 72 young girls suffering from anorexia, the `Anorexia.dat` file under `stat-961-fall-2021/data` shows their weights before and after an experimental period:

```
setwd("C:/Users/Sam/Documents/School/2021-2022/First Semester/Stat 961/stat-961-fall-2021/home")
anorexia_data = read_tsv("../data/Anorexia.dat", col_types = "ifdd")
print(anorexia_data, n = 5)

## # A tibble: 72 x 4
##   subj therapy before after
##   <int> <fct>    <dbl> <dbl>
## 1     1 b      80.5  82.2
## 2     2 b      84.9  85.6
## 3     3 b      81.5  81.4
## 4     4 b      82.6  81.9
## 5     5 b      79.9  76.4
## # ... with 67 more rows
```

The girls were randomly assigned to receive one of three therapies during this period. A control group received the standard therapy, which was compared to family therapy and cognitive behavioral therapy. The goal of the study is to compare the effectiveness of the therapies in increasing the girls' weights.

- Prepare the data by (1) removing the `subj` variable, (2) re-coding the factor levels of `therapy` as `behavioral`, `family`, and `control`, (3) renaming `before` and `after` to `weight_before` and `weight_after`, respectively, and (4) adding a variable called `weight_gain` defined as the difference of `weight_after` and `weight_before`. Print the resulting tibble.
- Explore the data by (1) making box plots of `weight_gain` as a function of `therapy`, (2) making a scatter plot of `weight_gain` against `weight_before`, coloring points based on `therapy` and (3) creating a table displaying, for each `therapy` group, the mean weight gain, maximum weight gain, and fraction of girls who gained weight (i.e. `weight_gain > 0`). Based on these summaries: What therapy appears overall the most successful and why? How effective does the standard therapy appear to be? What is the greatest weight gain observed in this study? Which girls tended to gain most weight (in the absolute sense), based on their weight before therapy? Why might this be the case?
- Run a linear regression of `weight_gain` on `therapy` and print the regression summary (print in R, without using `kable`). Identify the base category chosen by R and discuss the interpretations of the fitted coefficients. It makes more sense to choose `control` as the base category. Recode the factor levels so that `control` is the first (and therefore will be chosen as the base category), rerun the linear regression, and print the summary again. Do the relationships among the fitted coefficients in these two regressions match what was found in Problem 1d?
- Directly compute the between-groups, within-groups, and corrected total sums of squares (without appealing to the `aov` function or equivalent) and verify that the first two add up to the third. What is the ratio of the between-groups sum of squares and the corrected total sum of squares? What is the interpretation of this quantity, and what quantity in the regression summaries printed in part (c) is it equivalent to?

### Solution 3.

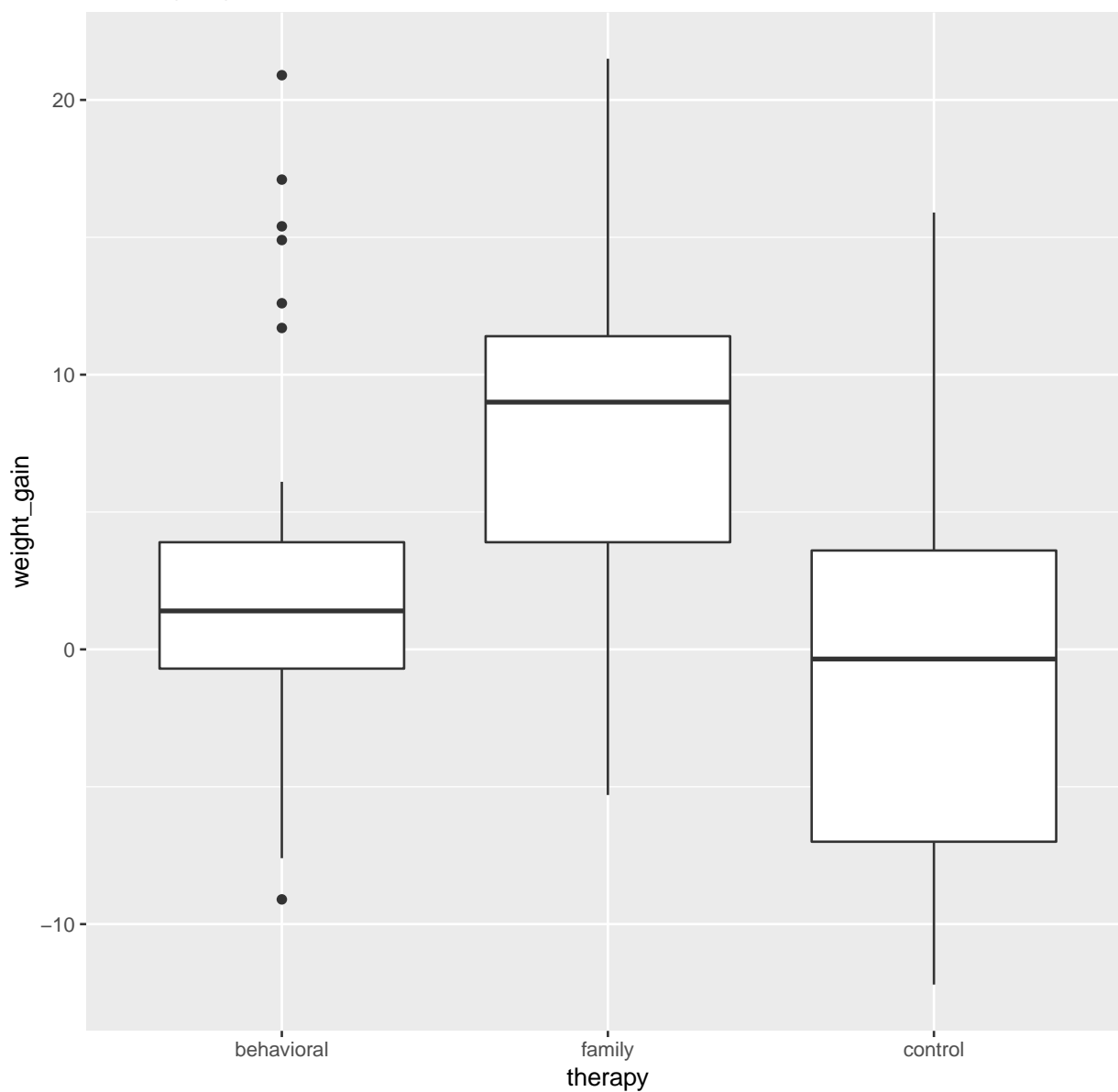
```
(a) anorexia_data_processed <-
  anorexia_data %>%
    select(-subj) %>%
    mutate(therapy =
      recode(therapy,
        b = "behavioral",
        f = "family",
        c = "control")
    ) %>%
    rename(weight_before = "before", weight_after = "after") %>%
    mutate(weight_gain = weight_after - weight_before)

print(anorexia_data_processed, n=5)
```

```
## # A tibble: 72 x 4
##   therapy    weight_before weight_after weight_gain
##   <fct>          <dbl>         <dbl>         <dbl>
## 1 behavioral      80.5           82.2           1.70
## 2 behavioral      84.9           85.6           0.700
## 3 behavioral      81.5           81.4          -0.100
## 4 behavioral      82.6           81.9          -0.700
## 5 behavioral      79.9           76.4          -3.5
## # ... with 67 more rows
```

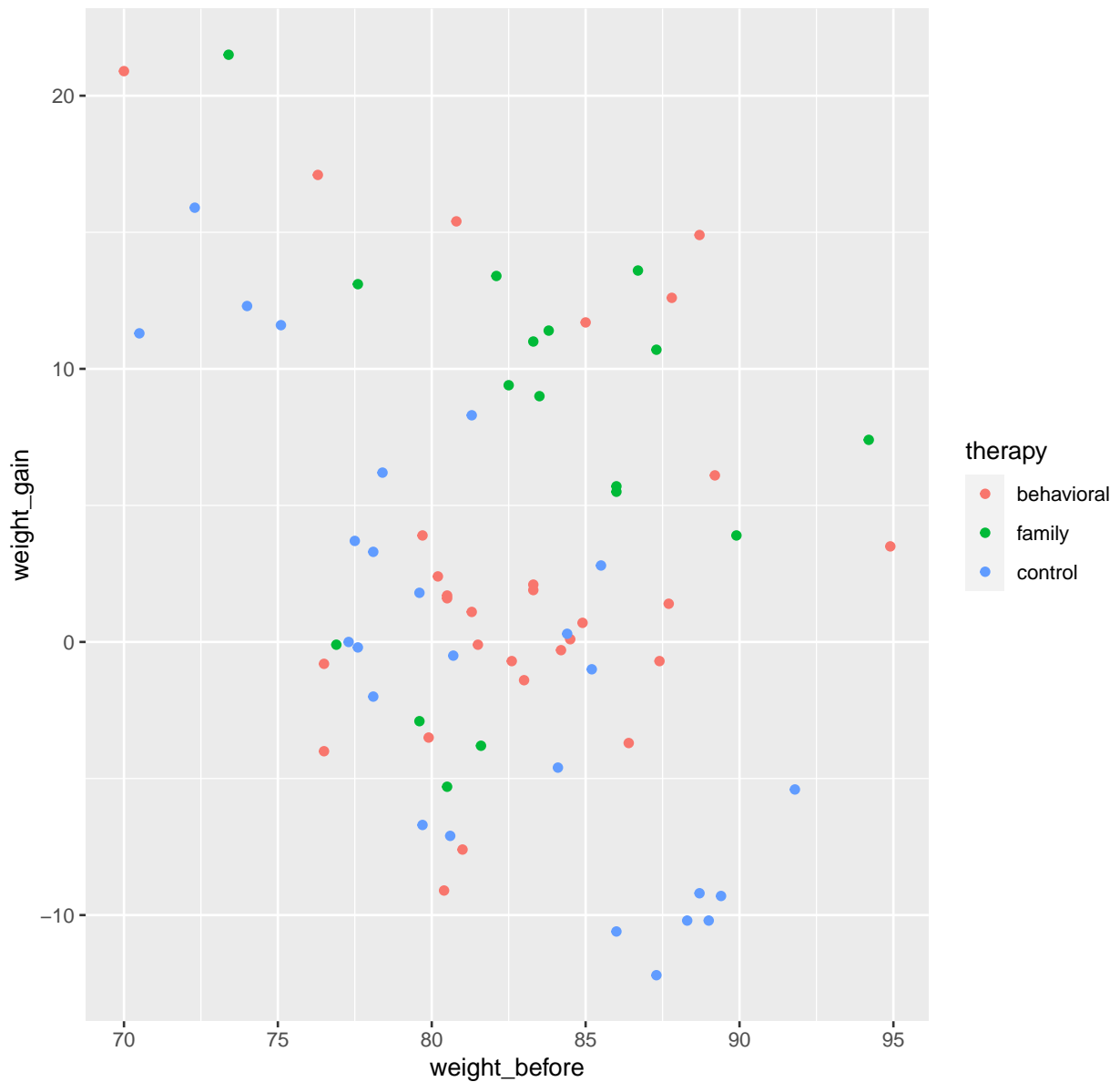
```
(b) anorexia_data_processed %>%
  ggplot() +
  geom_boxplot(mapping = aes(x = therapy, y = weight_gain)) +
  ggtitle("1. Weight gain as a function of therapy")
```

## 1. Weight gain as a function of therapy



```
anorexia_data_processed %>%  
  ggplot() +  
  geom_point(mapping = aes(x = weight_before, y = weight_gain, color = therapy)) +  
  ggtitle("2. Weight gain as a function of weight before")
```

## 2. Weight gain as a function of weight before



```
anorexia_data_summary <-
  anorexia_data_processed %>%
    group_by(therapy) %>%
    summarise(mean_weight_gain = mean(weight_gain),
              max_weight_gain = max(weight_gain),
              frac_weight_gain = sum(weight_gain > 0)/length(weight_gain))

print(anorexia_data_summary)
```

```
## # A tibble: 3 x 4
##   therapy    mean_weight_gain max_weight_gain frac_weight_gain
##   <fct>          <dbl>          <dbl>          <dbl>
## 1 behavioral      3.01            20.9            0.621
```

## 2 family	7.26	21.5	0.765
## 3 control	-0.450	15.9	0.423

Note that the family therapy seems the most successful as it has the highest average weight gain, highest max weight gain, and highest proportion of individuals who gained weight as compared to the other two treatments. The standard (control) therapy seems quite ineffective, as less than half of the recipients of this therapy gained weight and the average weight gain for this group is negative (i.e. these participants had a tendency to *lose* weight). The greatest weight gain observed in the study was 21.5 pounds, for an individual in the family therapy group. Girls with lower weights prior to the study tended to have higher weight gains following the study than those who had higher initial weights. This is negative correlation is possibly a sort of regression to the mean, where it is easier for individuals who are far from their “equilibrium” body weight to move toward that weight than away (i.e. an increase in weight for those below and a decrease/smaller increase for those above).

```
(c) lm1 <- lm(weight_gain ~ therapy, data = anorexia_data_processed)
print(summary(lm1))

##
## Call:
## lm(formula = weight_gain ~ therapy, data = anorexia_data_processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.565  -4.543  -1.007   3.846  17.893
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      3.007      1.398   2.151  0.0350 *
## therapyfamily     4.258      2.300   1.852  0.0684 .
## therapycontrol   -3.457      2.033  -1.700  0.0936 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.528 on 69 degrees of freedom
## Multiple R-squared:  0.1358, Adjusted R-squared:  0.1108
## F-statistic: 5.422 on 2 and 69 DF,  p-value: 0.006499
```

The base category chosen by R was behavioral therapy. The intercept of 3.007 means that we expect the average weight gain of the behavioral therapy to be about 3.007 pounds, in agreement with the 3.01 number from (2). The coefficient of 4.258 for the family therapy indicator means that the mean weight gain of the family therapy group should be 4.258 pounds higher than that of the behavioral therapy group (7.265 pounds, again agreeing with the figure from (2)). Finally, the coefficient of -3.457 pounds for the control therapy indicator means that the weight gain of the control therapy group should be -3.457 pounds lower than that of the behavioral therapy group; i.e. the mean weight gain for the group is -0.450 pounds, in agreement with the figure found in (2).

```
anorexia_data_processed <-
  anorexia_data_processed %>%
```

```

mutate(therapy = relevel(therapy, ref = "control"))

lm2 <- lm(weight_gain ~ therapy, data = anorexia_data_processed)
print(summary(lm2))

##
## Call:
## lm(formula = weight_gain ~ therapy, data = anorexia_data_processed)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.565  -4.543  -1.007   3.846  17.893
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    -0.450      1.476  -0.305   0.7614
## therapybehavioral  3.457      2.033   1.700   0.0936 .
## therapyfamily     7.715      2.348   3.285   0.0016 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7.528 on 69 degrees of freedom
## Multiple R-squared:  0.1358, Adjusted R-squared:  0.1108
## F-statistic: 5.422 on 2 and 69 DF,  p-value: 0.006499

```

Taking  $a := \text{behavioral}$ ,  $b := \text{control}$ , and  $c := \text{family}$ , we see that we are in a specific case of (1d) where we have changed the baseline value of a three level categorical variable.

Note that  $\hat{\beta}'_0 = -0.450 = 3.007 + (-3.457) = \hat{\beta}_0 + \hat{\beta}_1$ ,  $\hat{\beta}'_1 = 3.457 = -(-3.457) = -\hat{\beta}_1$ , and  $\hat{\beta}'_2 = 7.715 = -(-3.457) + 4.258 = -\hat{\beta}_1 + \hat{\beta}_2$ , so the coefficients have the same relationship as predicted by (1d).

```

(d) y <- anorexia_data_processed$weight_gain
y_bar <- mean(y)
y_bar_by_gp <- list(
  behavioral = lm1$coefficients[1],
  family = lm1$coefficients[1] + lm1$coefficients[2],
  control = lm1$coefficients[1] + lm1$coefficients[3])
y_bar_ci <-
  anorexia_data_processed %>%
  select(therapy) %>%
  summarise(gp_mean = lm1$coefficients[1] +
    ifelse(therapy == "behavioral",
      0,
      ifelse(therapy == "family",
        lm1$coefficients[2],
        lm1$coefficients[3])))

```



```

# SST = \sum_{i=1}^n (y_i - \overline{y})^2
sst <- sum((y - y_bar)^2)

# SSB = \sum_{i=1}^n (\overline{y}_{C(i)} - \overline{y})^2
ssb <- sum((y_bar_ci - y_bar)^2)

# SSW = \sum_{i=1}^n (y_i - \overline{y}_{C(i)})^2
ssw <- sum((y - y_bar_ci)^2)

print(paste("SST:", sst))

## [1] "SST: 4525.38611111111"

print(paste("SSB:", ssb))

## [1] "SSB: 614.64366892044"

print(paste("SSW:", ssw))

## [1] "SSW: 3910.74244421907"

```

Note that we indeed have that  $SST = SSB + SSW$ . The ratio of SSB to SST is  $614.6437/4525.386 = 0.1358$ . This corresponds to  $R^2$  in both regression summaries and can be interpreted as saying that about 13.6% of the variation in the data can be captured just by accounting for the type of therapy each individual received.