

Multiple regression

Bert van der Veen

Department of Mathematical Sciences, NTNU

Introduction

Going to largely omit code in presentation. But see .Rmd files.

Outline Today

- ▶ Multiple linear regression
- ▶ Model validation
- ▶ Introduction to GLMs

Questions about yesterday?



What if we have >1 explanatory variable?

We often want to look at the impacts of several variables together

- ▶ they may all have some effect
- ▶ we might be doing an experiment where factors interact
- ▶ we might want to model one variable as a polynomial

The model

This is our model for simple regression

$$y_i = \alpha + \beta x_i + \epsilon_i$$

How can we extend it to more than one variable?

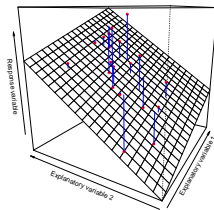
Two explanatory variables



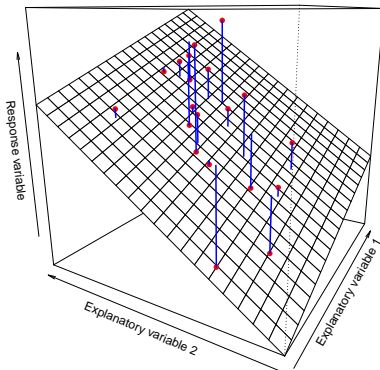
$$y_i = \alpha + \beta_1 x_{1i} + \beta_2 x_{2i} + \epsilon_i$$

This is a plane

With more than two covariates it is a **hyperplane**



Plane



Fitting in R

In R we can just use the same function as we did before.

The only change is in the formula. It was

$y \sim x$

now it is

$y \sim x1 + x2$

and the same for categorical and continuous covariates.

More than two: general linear regression

$$y_i = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \cdots + \beta_k x_{ik} + \epsilon_i$$

$$y_i = \alpha + \sum_{k=1}^p \beta_k x_{ik} + \epsilon_i$$

- ▶ we have p covariates, labelled from $k = 1$ to p
- ▶ the k^{th} covariate values for the i^{th} observation is x_{ik}
- ▶ we also have p covariate effects β_k

Design Matrices

We can write this more compactly. First, we turn the intercept into a covariate by filling a column of 1s for every data point. Then we write all of the covariates in a matrix, \mathbf{X} :

$$\mathbf{X} = \begin{pmatrix} \overline{x_1} & x_2 & x_3 \\ 1 & 2.3 & 3.0 \\ 1 & 4.9 & -5.3 \\ 1 & 1.6 & -0.7 \\ \vdots & \vdots & \vdots \\ 1 & 8.4 & 1.2 \end{pmatrix}$$

This is called the *Design Matrix*: it helps to write down the model

Writing the Model

Using matrix algebra, the regression model becomes

$$\mathbf{Y} = \mathbf{X}\beta + \epsilon$$

where \mathbf{Y} , β and ϵ are now all vectors of length n , where there are n data points. \mathbf{X} is an $n \times p$ matrix.

We will not look at the mathematics in any detail: the point here is that the model for the effect of covariates can be written in the design matrix.

The Solution (just so you can see it)

After a bit of matrix algebra, one can find the MLE:

$$\hat{\mathbf{b}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}(??) \quad (1)$$

where \mathbf{b} is the MLE for β .

In practice:

- ▶ you don't have to calculate this: the computer does it
- ▶ the computer actually doesn't use equation (??)

Writing the Model: continuous covariates

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$$

is

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 2.3 & 3.0 \\ 1 & 4.9 & -5.3 \\ 1 & 1.6 & -0.7 \\ \vdots & \vdots & \vdots \\ 1 & 8.4 & 1.2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

α is the intercept, β_1 is the slope parameter for x_1 , and so on.

Categorical variables

Categorical variables need to be turned into something numerical.

$$\mathbf{x} = \begin{pmatrix} \text{Species} \\ \text{Orchid} \\ \text{Orchid} \\ \text{Dandelion} \\ \vdots \\ \text{Daisy} \end{pmatrix} \Rightarrow \mathbf{X} = \begin{pmatrix} \text{Orchid} & \text{Dandelion} & \text{Daisy} \\ 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{pmatrix}$$

But do we need each column?

Contrasts

There are many ways to construct a design matrix for categorical variables.

`contrasts` and `constr.treatment`

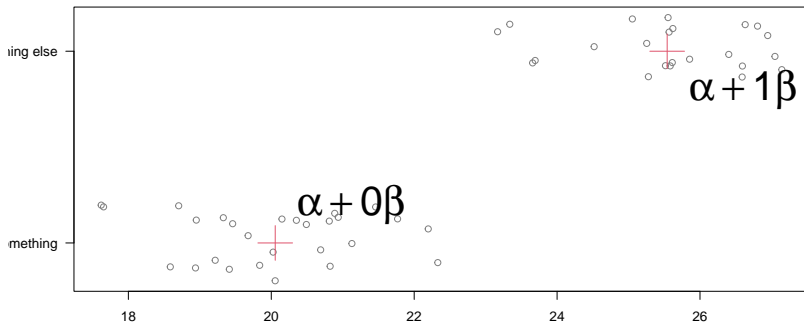
- ▶ Treatment contrast are default in R (“dummy”)
- ▶ Sum-to-zero
- ▶ Polynomial
- ▶ Difference
- ▶ Etc.

Writing the Model: categorical (ANOVA)

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 0 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Here, α is the intercept for the first category, β_1 the difference of the first and second category, β_2 the difference between the first and third categories.

Examples of linear models: categorical x_i (from yesterday)



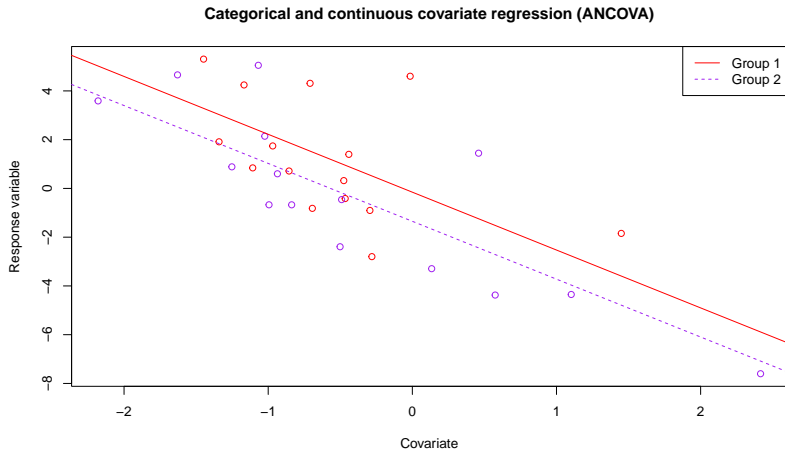
- ▶ α is the mean of the first group
- ▶ β is the deviation from the mean of the first group

Writing the Model: continuous and categorical

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3.0 \\ 1 & 1 & -5.3 \\ 1 & 1 & -0.7 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 1.2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Here, α is the intercept for the first category at $x_3 = 0$, β_1 is the difference for the second category at $x_3 = 0$, and β_2 is the slope parameter for two regression lines.

Writing the Model: continuous and categorical



Example 1: body temperature and heart rate

Data by Mackowiak et al. (1992) via the UsingR package



Figure 1: by Vecteezy

Body temperature: the data

- ▶ 130 individuals
- ▶ Body temperature
- ▶ Resting heart rate

temperature	gender	hr
96.3	1	70
96.7	1	71
96.9	1	74
97.0	1	80
97.1	1	73
97.1	1	75
97.1	1	82

Body temperature: two simple models

$$y_i = \alpha_i + x_i \beta_i + \epsilon_i \quad (2)$$

For this data that means:

$$y_i = \alpha + \textit{female} * \beta + \textit{left.over.stuff}_i$$

and

$$y_i = \alpha + \textit{heart.rate} * \beta + \textit{left.over.stuff}_i$$

Body temperature: fit the two models

```
model1 <- lm(celsius ~ gender, data = normtemp)
coef(model1)
```

```
## (Intercept) genderFemale
## 36.7247863 0.1606838
```

```
model2 <- lm(celsius ~ hr, data = normtemp)
coef(model2)
```

```
## (Intercept) hr
## 35.7259744 0.0146303
```


Body temperature: fit one model and compare

```
model3 <- lm(celsius ~ hr + gender, data = normtemp)
```

##	(Intercept)	hr	genderFemale
## gender	36.7248	NA	0.16
## hr	35.7260	0.0146	NA
## both	35.6949	0.0140	0.15

Interactions

An interaction is when we have the product of two (or more) covariates in the model:

$$y_i = \alpha + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i1} x_{i2} + \epsilon_i$$

$y \sim x_1 + x_2 + x_1 : x_2$ or $y \sim x_1 * x_2$

It means that we expect the effect of two covariates to jointly impact y_i

It does **not** mean we model how x_1 affect x_2 or vice versa!

Interactions: continuous-continuous

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 2.3 & 3.0 & 2.3 * 3.0 \\ 1 & 4.9 & -5.3 & 4.9 * -5.3 \\ 1 & 1.6 & -0.7 & 1.6 * -0.7 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 8.4 & 1.2 & 8.4 * 1.2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Interactions: continuous-continuous

$$y_i = \alpha + x_{i1}\beta_1 + x_{i2}\beta_2 + x_{i1}x_{i2}\beta_3 + \epsilon_i$$

- ▶ β_1 is the slope for x_1
- ▶ β_2 is the slope of x_2
- ▶ β_3 is their joint parameter

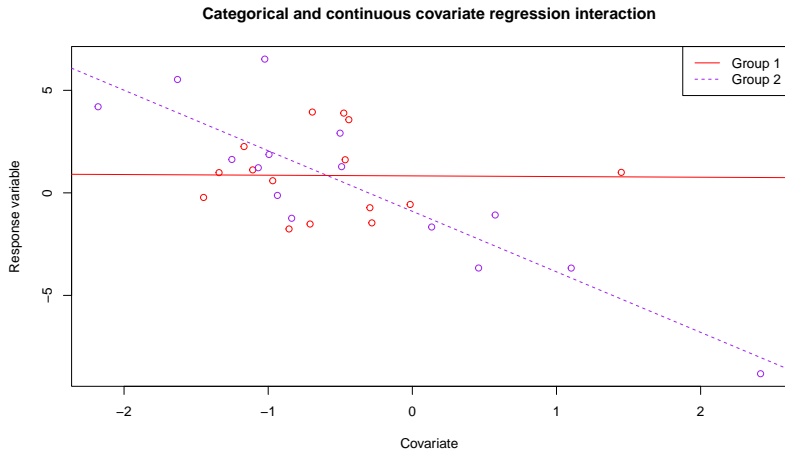
β_3 represents how the effect of x_1 or x_2 changes with the other covariate. E.g., body size and heart rate on body temperature

Interactions: categorical-continuous

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & 3.0 & 0 \\ 1 & 1 & -5.3 & -5.3 \\ 1 & 1 & -0.7 & -0.7 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1.2 & 0 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta_1 \\ \beta_2 \\ \beta_3 \\ \beta_4 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

A separate regression line for each category. Here, α and β_2 are the slope and intercept for the regression line of the first category. $\alpha + \beta_1$ is the intercept and $\beta_3 + \beta_4$ is the slope of the regression line for the second category.

Interactions: categorical-continuous

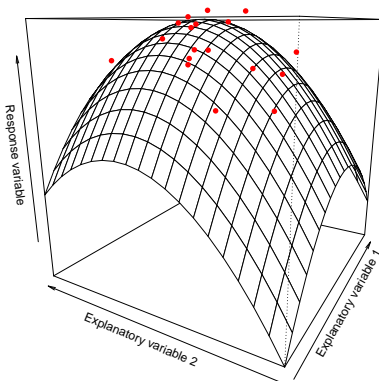


Other functions of explanatory variables

As long as the model is linear in the parameters, we can also have functions:

- ▶ Quadratic: $y_i = x_i\beta + x_i^2\beta_2 + \epsilon_i$
- ▶ Centering: $y_i = (x_i - \bar{\mathbf{x}})\beta + \epsilon_i$
- ▶ Exponential: $y_i = \exp(x_i)\beta + \epsilon_i$
 - ▶ or logarithmic: $y_i = \log(x_i)\beta + \epsilon_i$

Surface: quadratic effects



Wiggly things

$$\mathbf{Y} = s(\mathbf{X}) + \epsilon$$



See [GAM workshop](#)

Finding a “good” model

We do not usually explicitly specify regressions in terms of their imposed hypersurface.

More on how to find a model that fits the data well tomorrow.

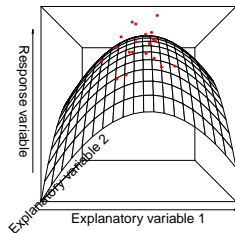
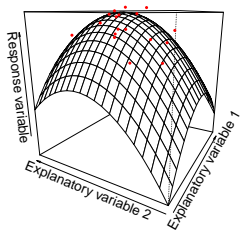
The predict function

In R we can calculate $\hat{y}_i = \hat{\alpha} + x_i\hat{\beta}_1$ with the `predict` function:

```
predict(model, newdata = newX)
```

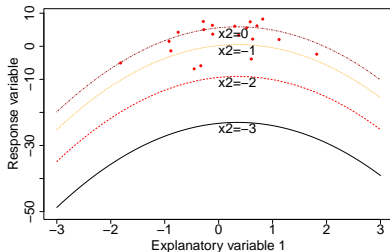
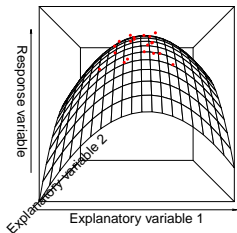
Here, `newX` are the values of the covariate that we want to calculate \hat{y}_i for. For the observed values we leave it empty.

Visualizing a multiple regression



- ▶ We want to look at the regression in 2D anyway
- ▶ So we need to choose what point to do that from

Visualizing a multiple regression



- ▶ We want to look at the regression in 2D anyway
- ▶ So we need to choose what point to do that from

Summary

- ▶ Multiple regression and the design matrix
- ▶ Fortunately we have the `lm()` function in R!

Example code for practical

```
dataset <- read.csv("some_place_on_my_computer/awesomedata.csv")  
lm(y ~ x1+x2, data = dataset)
```

Questions

