

R workshop #2: multivariate regression analysis and factor interactions

Nicola Romanò

Introduction

Last year we have talked about linear models and their use to perform regression and analysis of variance (ANOVA). We only considered simple situations with one independent variable influencing the output variable (one-way ANOVA) or two factors (two-way ANOVA) that do not interact with each other. In the lectures we have now talked about interactions and how they change our interpretation of linear models. In this workshop we will have a look at how to deal with interactions in R.

Learning objectives

After completing this workshop you will be able to:

- Use linear models to perform multiple regression and analysis of variance with multiple factors
- Correctly interpret the output of a linear model
- Compare two models to choose the one that fits the data better
- Correctly interpret the results of your analysis in the presence of interactions

Section 1 - A refresher on linear models

We start this workshop with a little refresher of linear models. **This also includes some more details about linear models that we only briefly touched upon last year**, so please go through this carefully. A linear model is a statistical model that relates the changes in a dependent variable (Y) with the changes in one or more independent variables (X_1, X_2, \dots, X_n).

The general equation for such model is:

$$Y = X_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

Where:

- Y is our measured/outcome variables
- X_1, \dots, X_n are the factors (or predictors) that influence Y . They are generally the other variables in your dataset, or transformations/combination of them ¹.

¹ For instance, we may have collected the weight of the subjects in our study, but use $\log(\text{weight})$ as a predictor for our model. Or we may have collected two different values and use their ratio as a model parameter

- β_1, \dots, β_n are the regression coefficients, scaling factors for the predictors.
- ϵ is the error, or residual. It represents the difference between what is explained by the model prediction, and what we have observed. It includes the effect of all the factors that we did not measure in our experimental setup, as well as measurement errors. We generally assume that it is normally distributed².

When we use R (or any other software!) to generate the model, what it does is estimating the coefficients β in such a way to minimise the error³.

In this formula each predictor acts independently from the others. In other words, if we have two predictors, X_1 and X_2 , the effect of X_1 on Y will always be the same, independently of the value of X_2 . As we have seen in the lecture this is not always the case.

Simple regression

As a first example let's consider the dataset *pressure-workshop2.csv*. In this study the effect of a drug on reducing blood pressure (measured in mmHg) has been investigated on 150 patients of different age, weight (measured in kg), and sex.

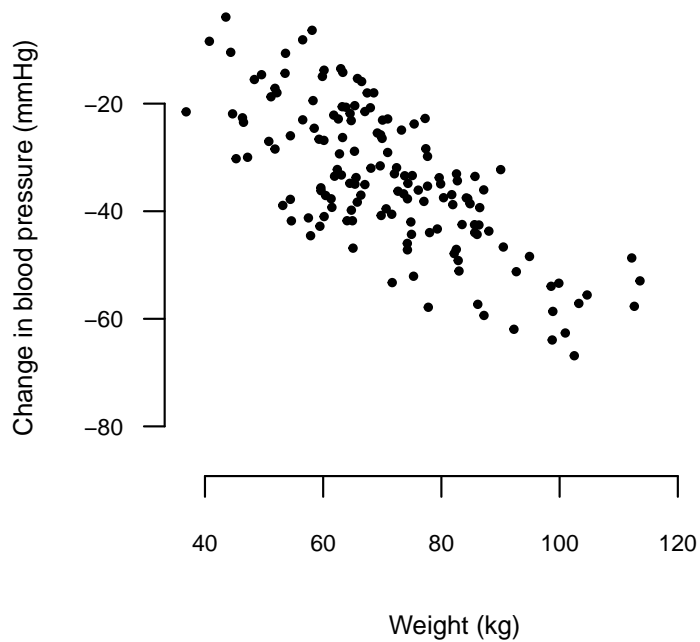
Start by familiarising with the data. How many men and women are there? What range of age and weight? Plot the various variables against each other and see if any particular patterns emerge⁴.

Let's forget for a moment about the other variables and concentrate on the relation between Weight and Response; it looks like the largest effect is seen in heavier patients.

² Note that although having normally distributed residuals makes things easier... you can still have a good, usable model when residuals are not normally distributed, especially if the sample size is big enough. Mostly, it boils down to critical observation of your data and experience.

³ In the case of `lm`, this is called a least-square estimation. In statistics books and publications you may see the estimated parameters indicated as $\hat{\beta}$ (read as "beta hat"). This is to indicate that this is the result of an estimation, that is an approximation of the true value of β in the population, which remains unknown. We can obtain confidence intervals for these estimates by using `confint(model)`

⁴ If you do not remember how to do that, see Workshop 1.



We can use a linear model to test whether such a relation exists.

As always, we start by stating our null hypothesis _____

Do you remember how to perform a linear regression in R? Try it,
if you don't remember see the following page!

```
model <- lm(Response ~ Weight, data = pressure)
```

This generates the model

$$\text{Response} = \beta_0 + \beta_1 * \text{Weight} + \epsilon$$

What are the assumption of this model? Do you remember how to verify that they are satisfied? ⁵

This is one of the simplest linear models we can generate, where the value of the outcome depends on a single parameter. This is called *simple regression*.

Let's look at the output of the model

```
summary(model)

##
## Call:
## lm(formula = Response ~ Weight, data = pressure)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -19.190  -6.354   0.874   6.568  19.554
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  11.78143    3.32703   3.541 0.000533 ***
## Weight       -0.64852    0.04521 -14.346 < 2e-16 ***
## ---
## Signif. codes:
##  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9.133 on 148 degrees of freedom
## Multiple R-squared:  0.5817, Adjusted R-squared:  0.5789
## F-statistic: 205.8 on 1 and 148 DF, p-value: < 2.2e-16
```

The summary gives us a lot of information.

First of all, it tells us the parameters β (coefficients) that have been estimated by the model.

$$\hat{\beta}_0 = 11.78 \text{ and } \hat{\beta}_1 = -0.65$$

Therefore

$$\text{Response} = 11.78 - 0.65 * \text{Weight} + \epsilon$$

This means that for any increase of 1 Kg in weight there is a decrease of 0.65 mmHg in blood pressure following the intake of the drug. The effect of weight on the response to the drug is statistically significant ($F_{1,148} = 205.8, p = 2 * 10^{-16}$)⁶.

⁵ Let's discuss this in the forum! I would say that for this case the assumptions are generally satisfied, what do you think?

⁶ R also reports a p-value for the intercept; this is the result of a one sample t-test comparing the intercept to 0. In other words, in this case the intercept is statistically different from 0. The intercept is the value corresponding to a change in blood pressure where all of the factors (in this case weight) are equal to zero. Since a weight of 0 is not biologically meaningful we can ignore this value in this instance

Another important value is the coefficient of determination (R^2 , sometimes called deviance). This is a measure of how good the model is, or how much of the variation in the data it explains. R^2 is not a great way to compare two different models, since it depends on the number of parameters; that is, if we add an extra descriptor to our model R^2 will always increase. For this reason, R reports also an “adjusted” version of it. In this case $adj.R^2 = 0.5789$; this means that our model describes/explains ~57.9% of the variability in our data, which is OK but not great. It means that there are other factors that we have not considered accounting for >40% of the variability! ⁷ So, what are these other factors?

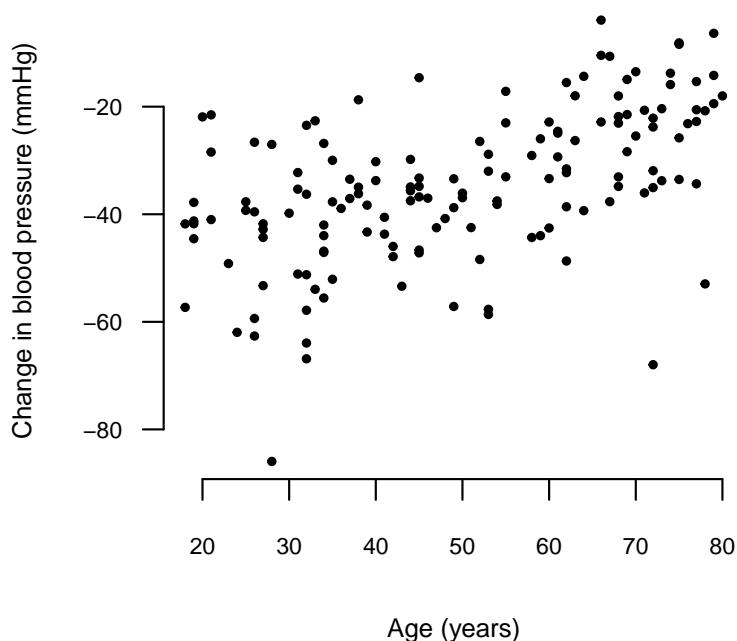
⁷ What do you think is the maximum value of R^2 ? Why?

Multiple regression

Our dataset contains two other descriptors: Age and Sex. It is very biologically plausible that these would affect blood pressure, so we should add them to our model⁸. To keep things simple, we will start with Age, and consider gender later.

It is useful, at this point, to also plot the change in blood pressure against age.

⁸ Note that, although for the sake of simplicity we are adding these descriptors one at a time, in practice we would probably start from a complete model, including all of the descriptors that we measured. That is why we measured, isn't it?



We see a possible in the response to the drug depending on age. Let's incorporate age in our model.

```
model.2 <- lm(Response ~ Weight + Age, data = pressure)
```

This will generate a model that considers the effect of weight and

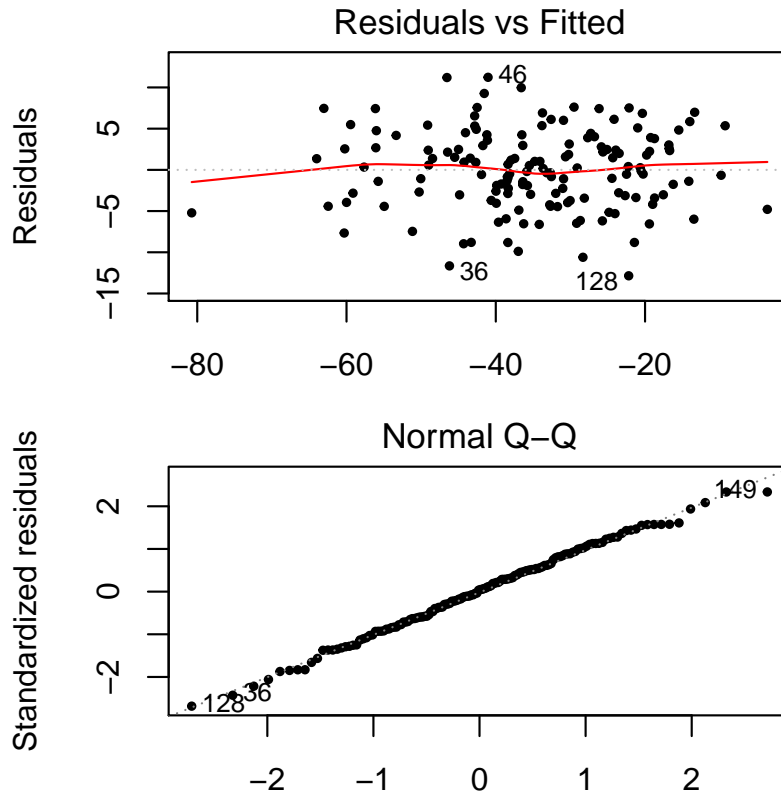
the effect of age, independently of each other ⁹

What are the null hypotheses¹⁰ that this model is testing?

Again, we want to check the assumptions of the model by using diagnostic plots.

⁹ This means that the model will look at the effect of the weight of the individual on his/her response to the drug, independently of his/her age, and vice versa for age.

¹⁰ There is more than one!!!



Since the diagnostic plots look good (equal variance throughout and normally distributed residuals), we can continue with this model.

```
summary(model.2)
```

```
##
## Call:
## lm(formula = Response ~ Weight + Age, data = pressure)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.8545  -3.3021   0.1972   3.1070  11.2320
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -8.60214    2.04592  -4.205 4.53e-05 ***
## Weight       -0.65779    0.02392 -27.501 < 2e-16 ***
## Age           0.42390    0.02169  19.541 < 2e-16 ***
## ---
## Signif. codes:
```

```
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.832 on 147 degrees of freedom
## Multiple R-squared: 0.8837, Adjusted R-squared: 0.8821
## F-statistic: 558.6 on 2 and 147 DF, p-value: < 2.2e-16
```

You can interpret the result of this model just like you did for the previous one. It tells us that there is a statistically significant effect of weight ($p < 2 * 10^{-16}$) and of age ($p < 2 * 10^{-16}$) on the response to the drug ($F_{2,147} = 558.6$ ¹¹) Note that now the model explains 88% of the variability!

Qualitative predictors and dummy variables

Let's now consider gender and add it to our model. Plot the data, do you think gender affects the response to the drug?

In this case, we are dealing with a discrete qualitative variable, with two levels, F and M. All that we said so far still applies, and `lm` is able to deal with this type of variables with no issue. However, the way we deal with these type of variables is slightly different.

```
model.3 <- lm(Response ~ Weight + Age + Sex, data = pressure)
```

This is modelling the following:

$$\text{Response} = \beta_0 + \beta_1 * \text{Weight} + \beta_2 * \text{Height} + \beta_3 * D$$

We introduce a new variable D , called a “dummy variable”, that codes Sex in this way:

$$D = \begin{cases} 1, & \text{if Sex} = M \\ 0, & \text{otherwise} \end{cases}$$

By default R assigns 0 to the first level of the variable (the *reference level*, in this case F), and 1 to the second¹²

Therefore, for observations at the reference level (so from female subjects), the third term $\beta_3 * D$ will be 0; for male subjects that will be $\beta_3 * 1 = \beta_3$. Thus, β_3 represents the **difference between the response of a male and a female**, keeping all of the other factors constant.

Let's have a look at the summary of the model to better clarify this.

¹¹ Note that the F statistic reported by `summary` refers to the whole model. If you wanted to know the F statistic for specific components of the model you could run `anova(model.2)`, which would give you F and DF for the various descriptors of your model.)

¹² Levels are ordered alphabetically; see Workshop 1 for how to change level ordering.

```
summary(model.3)

##
## Call:
## lm(formula = Response ~ Weight + Age + Sex, data = pressure)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12.6250  -3.2145   0.1914   3.2024  11.0588
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -8.36394     2.11449  -3.956 0.000119 ***
## Weight       -0.66298     0.02645 -25.062 < 2e-16 ***
## Age          0.42224     0.02204  19.154 < 2e-16 ***
## SexM         0.41103     0.88467   0.465 0.642898
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.845 on 146 degrees of freedom
## Multiple R-squared:  0.8839, Adjusted R-squared:  0.8815
## F-statistic: 370.5 on 3 and 146 DF, p-value: < 2.2e-16
```

The output is not much different from what we had before. Does Sex have a statistically significant effect on the response? What percent of variance is explained by this model?¹³

Consider the estimates

$$\beta_1 = -0.66; \beta_2 = 0.42; \beta_3 = 0.41$$

These mean that:

- For every increase in 1 kg of weight, the response decreases of 0.66 mmHg (keeping age and sex the same)
- For every increase in 1 year of age, the response increases of 0.42 mmHg (keeping weight and sex the same)
- If the patient is male the response increases of 0.41 mmHg

So, what would you predict the response of a 50 year old male weighing 82 kg will be? Write your response on the forum¹⁴.

¹³ Note how, although minimally, R^2 has increased, since we have added an extra parameter, however adj. R^2 has decreased!

¹⁴ Remember to include the intercept as well!

Dummy variables for multiple levels

You may have realised at this point, that dummy variables are what R uses to code for groups or other discrete factors when doing an ANOVA!

In some cases, however, you will have more than two levels; the reasoning is the same, however multiple dummy variables will be used to define the different levels.

For example, suppose you measured the levels of LH in three different species of fish: mackerel, salmon, and trout.

You can code the species variable with two dummy variables (so, number of levels - 1) D_1 and D_2 such as:

$$D_1 = \begin{cases} 1, & \text{if Species} = \text{"Salmon"} \\ 0, & \text{otherwise} \end{cases}$$

$$D_2 = \begin{cases} 1, & \text{if Species} = \text{"Tuna"} \\ 0, & \text{otherwise} \end{cases}$$

Therefore:

Species	D_1	D_2
Mackerel	0	0
Salmon	1	0
Tuna	0	1

Our model may be something like:

$$\text{LH} = \beta_0 + \beta_1 * D_1 + \beta_2 * D_2 + \epsilon$$

where β_1 represents the difference between LH levels in salmons and mackerels, and β_2 the difference between LH levels in tuna and mackerel.

Section 2 - Choosing a model

Going back to our initial example, we have fitted three models:

1. Response = $\beta_0 + \beta_1 * \text{Weight} + \epsilon$
2. Response = $\beta_0 + \beta_1 * \text{Weight} + \beta_2 * \text{Age} + \epsilon$
3. Response = $\beta_0 + \beta_1 * \text{Weight} + \beta_2 * \text{Age} + \beta_3 * D_{\text{male}} + \epsilon$

We could argue that #2 is better than #1, as it explains a much larger percentage of the variance (88% vs 58%), but what about #3?

Is it correct to say that since Sex does not have a statistically significant effect on the response, and since the value of adjusted R^2 is lower (albeit by a very small amount) we should drop Sex from the model, and only consider Age and Weight as predictors? One way of deciding this is to use the anova function to compare the two models. This tests the null hypothesis that the first model does not fit the data better than the second one.

```
anova(model, model.2)

## Analysis of Variance Table
##
## Model 1: Response ~ Weight
## Model 2: Response ~ Weight + Age
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     148 12346.1
## 2     147  3431.7   1    8914.4 381.85 < 2.2e-16 ***
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

As expected, the p-value is very low, indicating that the second model fits the data better than the first one, and should be preferred to it. See also how much the residual sum of squares (RSS) has decreased, indicating that the model is much closer to the real data (hence the residuals (and their squared sum) are smaller).

Conversely

```
anova(model.2, model.3)

## Analysis of Variance Table
##
## Model 1: Response ~ Weight + Age
## Model 2: Response ~ Weight + Age + Sex
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1     147  3431.7
## 2     146 3426.7   1    5.0665 0.2159 0.6429
```

The p-value is 0.64, indicating that we should refute the null hypothesis, meaning that our third model (with Age, Weight, and Sex) is not better than the simpler model.

R also provides a convenient function, called `drop1`, that removes one predictor at a time from a larger model. You can see that this confirms what we have seen above.

```
drop1(model.3, test = "F")

## Single term deletions
##
## Model:
## Response ~ Weight + Age + Sex
##      Df Sum of Sq    RSS    AIC  F value Pr(>F)
## <none>          3426.7 477.31
## Weight  1   14741.9 18168.5 725.52 628.1081 <2e-16 ***
## Age     1    8610.5 12037.2 663.77 366.8706 <2e-16 ***
## Sex     1         5.1  3431.7 475.53   0.2159 0.6429
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

This confirms what we saw before, a ## Exercise To consolidate what explained so far, consider the dataset xxx.csv.

TODO: DATASET TO ANALYSE

Section 3 - Interactions between factors

In the lectures you have learnt about interactions amongst factor in multiple regression. We will now see how to analyse interactions in R.

```
model.4 <- lm(Response ~ Weight * Age, data = pressure)
summary(model.4)
```

```
##
## Call:
## lm(formula = Response ~ Weight * Age, data = pressure)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
	-12.8931	-3.3649	0.1206	3.2447	11.2083

```
##
## Coefficients:
```

	Estimate	Std. Error	t value
(Intercept)	-1.052e+01	4.699e+00	-2.238
Weight	-6.301e-01	6.571e-02	-9.588
Age	4.643e-01	9.178e-02	5.058
Weight:Age	-5.823e-04	1.286e-03	-0.453

```
##
## Pr(>|t|)
```

	Pr(> t)
(Intercept)	0.0267 *
Weight	< 2e-16 ***
Age	1.25e-06 ***
Weight:Age	0.6514

```
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.845 on 146 degrees of freedom
## Multiple R-squared: 0.8839, Adjusted R-squared: 0.8815
## F-statistic: 370.5 on 3 and 146 DF, p-value: < 2.2e-16
```