

Module 3: Multiple Regression

R Practical

Camille Szmaragd and George Leckie¹
Centre for Multilevel Modelling

Pre-requisites box

- Modules 1-2

Contents

| | |
|---|-----------|
| Introduction to the Scottish Youth Cohort Trends Dataset | 2 |
| P3.1 Regression with a Single Continuous Explanatory Variable | 3 |
| P3.1.1 Examining the data | 3 |
| P3.1.2 A simple linear regression analysis | 8 |
| P3.2 Comparing Groups: Regression with a Single Categorical Explanatory Variable | 18 |
| P3.2.1 Comparing attainment for girls and boys | 18 |
| P3.2.2 Attainment by parental social class | 20 |
| P3.2.3 Fitting a non-linear relationship to attainment and cohort | 23 |
| P3.3 Regression with More than One Explanatory Variable (Multiple Regression) | 25 |
| P3.4 Interaction Effects | 29 |
| P3.4.1 Model with fixed cohort effect for boys and girls | 29 |
| P3.4.2 Fitting separate models for boys and girls | 32 |
| P3.4.3 Allowing for sex-specific trends in a pooled analysis: interaction effects | 34 |
| P3.4.4 Allowing the trend in attainment to depend on social class | 37 |
| P3.5 Checking Model Assumptions in Multiple Regression | 43 |
| P3.5.1 Checking the normality assumption | 44 |
| P3.5.2 Checking the homoskedasticity assumption | 45 |
| P3.6 References | 47 |

¹ This R practical is adapted from the corresponding MLwiN practical: Steele, F. (2008) Module 3: Multiple Regression MLwiN Practical. LEMMA VLE, Centre for Multilevel Modelling. Accessed at <http://www.cmm.bris.ac.uk/lemma/course/view.php?id=13>.

Some of the sections within this module have online quizzes for you to test your understanding. To find the quizzes:

EXAMPLE

From within the LEMMA learning environment

- Go down to the section for **Module 3: Multilevel Modelling**
- Click "[3.1 Regression with a Single Continuous Explanatory Variable](#)" to open Lesson 3.1
- Click **Q1** to open the first question

Pre-requisites

- Understanding of types of variables (continuous vs. categorical variables, dependent and explanatory); covered in Module 1.
- Correlation between variables
- Confidence intervals
- Hypothesis testing, p-values
- Independent samples t-test for comparing the means of two groups

Online resources:

<http://www.sportsci.org/resource/stats/>

<http://www.socialresearchmethods.net/>

<http://www.animatedsoftware.com/statglos/statglos.htm>

<http://davidmlane.com/hyperstat/index.html>

The aim of these exercises is to gain practical experience of the application and interpretation of multiple regression.

Introduction to the Scottish Youth Cohort Trends Dataset

You will be analysing data from the Scottish School Leavers Survey (SSLS), a nationally representative survey of young people. We use data from seven cohorts of young people collected in the first sweep of the study, carried out at the end of the final year of compulsory schooling (aged 16-17) when most sample members had taken Standard grades.² These are subject-based examinations, typically taken in up to eight subjects. Each subject is graded on a scale from 1 (highest) to 7 (lowest). The dependent variable is a total attainment score calculated by assigning 7 points for a '1', 6 for a '2' and so on.

The analysis dataset contains the following five variables:

| Variable name | Description and codes |
|-----------------|--|
| Caseid | Anonymised student identifier |
| Score | Point score calculated from awards in Standard grades. Scores range from 0 to 75, with a higher score indicating a higher attainment |
| cohort90 | The sample includes the following cohorts: 1984, 1986, 1988, 1990, 1996 and 1998. The cohort90 variable is calculated by subtracting 1990 from each value. Thus values range from -6 (corresponding to 1984) to 8 (1998), with 1990 coded as zero |
| Female | Sex of student (1 = female, 0 = male) |
| Sclass | Social class, defined as the higher class of the mother or father (1 = managerial and professional, 2 = intermediate, 3 = working, 4 = unclassified) |

There are 33,988 students in the dataset.

² We are grateful to Linda Croxford (Centre for Educational Sociology, University of Edinburgh) for providing us with these data. The dataset was constructed as part of an ESRC-funded project on Education and Youth Transitions in England, Wales and Scotland 1984-2002. Further analyses of the data can be found in Croxford and Raffe (2006).

P3.1 Regression with a Single Continuous Explanatory Variable

We will begin by looking at the relationship between attainment (**score**) and cohort (**cohort90**). Has attainment changed over time and, if so, is the trend linear?

P3.1.1 Examining the data

Download the R dataset for this lesson:

From within the LEMMA Learning Environment

- Go to **Module 3: Multiple Regression**, and scroll down to **R Datasets and R files**
- Right click “3.1.txt” and select **Save Link As...** to save the dataset to your computer

Read the dataset into R using the `read.table` function and create a dataframe object named **mydata**:³

```
> mydata <- read.table(file = "3.1.txt", sep = ",", header = TRUE)
```

We specify the `sep = ","` argument as the file is a comma separated variable file while we specify the `header = TRUE` argument as the first line of this file contains the names of the variables. If we did not specify this second argument, the variable names would be incorrectly treated as the first row of data in the dataframe.

Use the `dim` function to display the number of rows of data and the number of columns of variables in the dataframe:

```
> dim(mydata)
[1] 33988 5
```

and use the `str` function to produce a basic description of the dataframe, which includes some general information on the number of variables, along with a list of the variables and their first few values:

```
> str(mydata)
'data.frame':  33988 obs. of  5 variables:
 $ caseid  : int  339 340 345 346 352 353 354 361 362 363 ...
 $ score   : int  49 18 46 43 17 29 15 19 45 12 ...
 $ cohort90: int  -6 -6 -6 -6 -6 -6 -6 -6 -6 -6 ...
 $ female  : int  0 0 0 0 0 0 0 0 0 0 ...
 $ sclass  : int  2 3 4 3 3 2 3 2 3 1 ...
```

The dataframe contains 33,988 observations on 5 variables.

³ At the beginning of your R session, you will need to set R's working directory to the file location where you saved the dataset. This can be done using the command line and the `setwd` command:

```
> setwd("C:\\userdirectory\\")
```

Or through selecting Change Dir... on the File menu.

Module 3 (R Practical): Multiple Regression

Dataframes may be displayed in matrix form. You can view subsets of the dataframe using standard matrix indexing conventions. Here we specify rows `1:20` to display only the first 20 rows of observations in the data. Note, we have not specified which columns we wish to display and so the values of all the variables are displayed:

```
> mydata[1:20, ]
  caseid score cohort90 female sclass
1     339   49      -6     0       2
2     340   18      -6     0       3
3     345   46      -6     0       4
4     346   43      -6     0       3
5     352   17      -6     0       3
6     353   29      -6     0       2
7     354   15      -6     0       3
8     361   19      -6     0       2
9     362   45      -6     0       3
10    363   12      -6     0       1
11   6824    0      -4     0       1
12   6826    0      -4     0       3
13   6827   20      -4     0       2
14   6828   32      -4     0       1
15   6829    0      -4     0       2
16   6834   24      -4     0       3
17   6836   23      -4     0       2
18  13206    7      -2     0       3
19  13209   38      -2     0       3
20  13215   46      -2     0       1
```

For example, the 10th student in the data belongs to the 1984 cohort and scored 12 out of 75. This student is a boy from a managerial social class background.

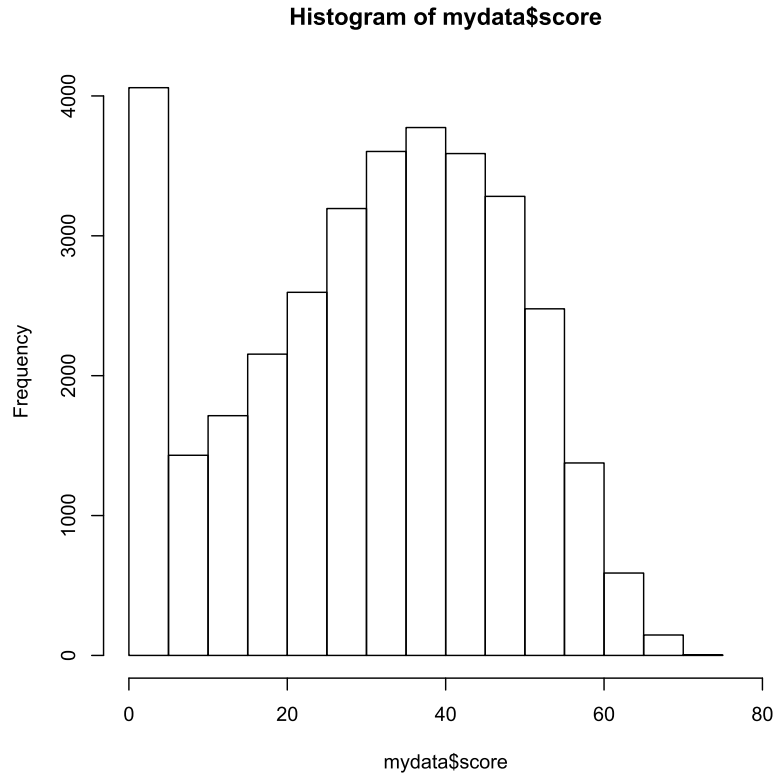
Having viewed the data we will examine **score** and **cohort90**, the variables to be considered in our first regression analysis.

Distribution of score

We will begin by obtaining a histogram and descriptive statistics for the dependent variable, **score**. We obtain a histogram with the `hist` command:

```
> hist(mydata$score, xlim = c(0,80))
```

where we have referred to the **score** variable within the **mydata** dataframe as `mydata$score`



The histogram should look like the above figure. The `xlim` argument is used to scale the x-axis from zero to 80. Apart from a peak at around zero, the distribution looks approximately normal. Remember that in a linear regression model it is the residuals that are assumed to be normal; we will check this assumption at the end of the exercise.

The `summary` command can be used to calculate and display a variety of univariate summary statistics for the variables in the dataset. To obtain summary statistics only for `score`:

```
> summary(mydata$score)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
  0.00  19.00   33.00   31.09  45.00   75.00
```

The standard deviation of `score` can also be obtained with the `sd` command

```
> sd(mydata$score)
[1] 17.31437
```

We see that `score` has a mean of 31.09, a standard deviation of 17.31 and can range between a minimum and maximum value of 0 and 75.

Distribution of cohort90

Because `cohort90` contains only six distinct values, we will look at its distribution in a frequency table rather than graphically. The `table` command produces one-way

This document is only the first few pages of the full version.

To see the complete document please go to learning materials and register:

<http://www.cmm.bris.ac.uk/lemma>

The course is completely free. We ask for a few details about yourself for our research purposes only. We will not give any details to any other organisation unless it is with your express permission.