

Module 3: Multiple Regression

MLwiN Practicals

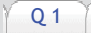
Fiona Steele¹
Centre for Multilevel Modelling

Module 3 (Practice): Multiple Regression

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  to open the first question

Contents

Introduction to the Scottish Youth Cohort Trends Dataset	3
P3.1 Regression with a Single Continuous Explanatory Variable	4
P3.1.1 Examining the data	4
P3.1.2 A simple linear regression analysis	11
P3.2 Comparing Groups: Regression with a Single Categorical Explanatory Variable	23
P3.2.1 Comparing attainment for girls and boys	23
P3.2.2 Attainment by parental social class	24
P3.2.3 Fitting a non-linear relationship to attainment and cohort	28
P3.3 Regression with More than One Explanatory Variable (Multiple Regression)....	31
P3.4 Interaction Effects	35
P3.4.1 Model with fixed cohort effect for boys and girls	35
P3.4.2 Fitting separate models for boys and girls	40
P3.4.3 Allowing for sex-specific trends in a pooled analysis: interaction effects	42
P3.4.4 Allowing the trend in attainment to depend on social class	46
P3.5 Checking Model Assumptions in Multiple Regression	54
P3.5.1 Checking the normality assumption	55
P3.5.2 Checking the homoskedasticity assumption	56

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. The MLwiN software will be used throughout.

¹ With additional material taken from Rasbash, J., Steele, F., Browne, W.J. and Prosser, B. (2005) *A User's Guide to MLwiN version 2.0*. Centre for Multilevel Modelling, University of Bristol. Downloadable from <http://www.cmm.bris.ac.uk/MLwiN/download/manuals.shtml>

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 grades². 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 33988 students in the data file.

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, L. and Raffe, D. (2006) "Education Markets and Social Class Inequality: A Comparison of Trends in England, Scotland and Wales". In R. Teese (Ed.) *Inequality Revisited*. Berlin: Springer.

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?

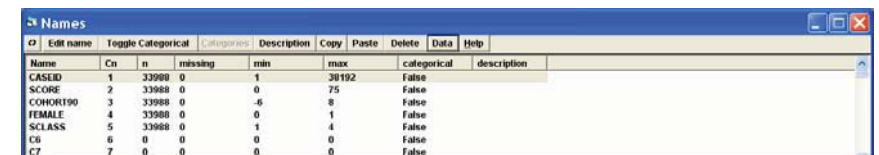
P3.1.1 Examining the data

To access the data files associated with this tutorial, you must have an account with LEMMA. To open the first data file,

From within the LEMMA Learning Environment

- Go to **Module 3: Multiple regression**, and scroll down to **MLwiN Datafiles**
- If you do not already have MLwiN to open the datafile with, click ([get MLwiN](#)).
- Click "[3.1.ws2](#)"

When the worksheet is opened, the filename will appear in the title bar of the main window. The **Names** window will also appear, giving a summary of the data in the worksheet:



Name	Cn	n	missing	min	max	categorical	description
CASEID	1	33988	0	1	38192	False	
SCORE	2	33988	0	0	75	False	
COHORT90	3	33988	0	-6	8	False	
FEMALE	4	33988	0	0	1	False	
SCLASS	5	33988	0	1	4	False	
C6	6	0	0	0	0	False	
C7	7	0	0	0	0	False	

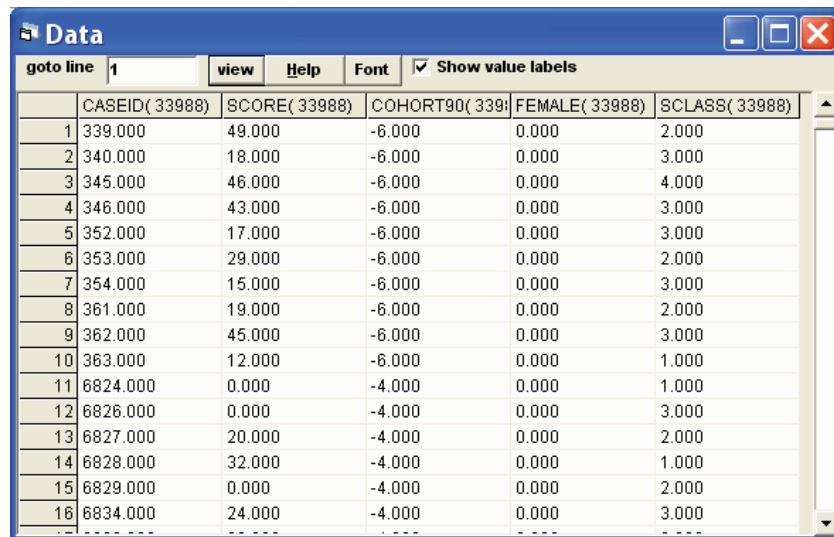
The MLwiN worksheet holds the data and other information in a series of columns, as on a spreadsheet. There are initially named c1, c2, etc. but we recommend that they be given meaningful names to show what their content relates to. This has already been done in the worksheet that you have loaded.

Each line in the body of the **Names** window summarises a column of data. In the present case only the first five of the 400 columns of the worksheet contain data. Each column contains 33988 values, one for each student represented in the data set. There are no missing values, and the minimum and maximum value in each column are shown. It is possible to define a variable as categorical (we shall do this later) and to add variable descriptions.

You can view individual values in the data using the **Data** window as follows:

- From the **Data Manipulation** menu, select **View or edit data**

The following window appears:



	CASEID(33988)	SCORE(33988)	COHORT90(339)	FEMALE(33988)	SCLASS(33988)
1	339.000	49.000	-6.000	0.000	2.000
2	340.000	18.000	-6.000	0.000	3.000
3	345.000	46.000	-6.000	0.000	4.000
4	346.000	43.000	-6.000	0.000	3.000
5	352.000	17.000	-6.000	0.000	3.000
6	353.000	29.000	-6.000	0.000	2.000
7	354.000	15.000	-6.000	0.000	3.000
8	361.000	19.000	-6.000	0.000	2.000
9	362.000	45.000	-6.000	0.000	3.000
10	363.000	12.000	-6.000	0.000	1.000
11	6824.000	0.000	-4.000	0.000	1.000
12	6826.000	0.000	-4.000	0.000	3.000
13	6827.000	20.000	-4.000	0.000	2.000
14	6828.000	32.000	-4.000	0.000	1.000
15	6829.000	0.000	-4.000	0.000	2.000
16	6834.000	24.000	-4.000	0.000	3.000

Because there are only five variables in the data file, all columns can be seen. When there are more variables, you can view any selection of columns, spreadsheet fashion, as follows:

- Click the **View** button
- Select columns to view
- Click **OK**

You can select a block of adjacent columns either by pointing and dragging or by selecting the column at one end of the block and holding down 'Shift' while you select the column at the other end. You can add to an existing selection by holding down 'Ctrl' while you select new columns or blocks. Use the scroll bars of the **Data** window to move horizontally and vertically through the data, and move or resize the window if you wish. You can go straight to line 1000, for example, by typing 1000 in the **goto line** box, and you can highlight a particular cell by pointing and clicking. This provides a means to edit data.

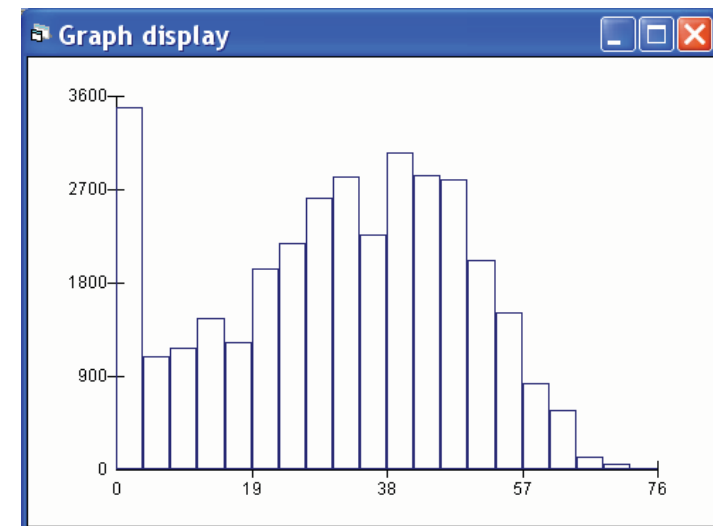
Having viewed the data we will examine SCORE and COHORT90, the variables to be considered in our first regression analysis.

Distribution of SCORE

We will begin by obtaining a histogram and descriptive statistics for the dependent variable, SCORE.

To obtain a histogram:

- From the **Graphs** menu, select **Customised Graph(s)**
- Next to **y**, select **SCORE** from the drop-down list
- Next to **plot type**, select **histogram**
- Click **Apply**



The histogram should look like the above figure. 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.

To obtain descriptive statistics for SCORE:

- From the **Basic Statistics** menu, select **Averages and Correlations**
- Under **Operation**, retain the default of **Averages**
- Highlight **SCORE** in the variable list
- Click **Calculate**

An **Output** window opens, showing the number of cases, number of missing values, mean and standard deviation of SCORE. The mean is 31.095 and the standard deviation is 17.314.

Distribution of COHORT90

Because COHORT90 contains only six distinct values, we will look at its distribution in a frequency table rather than graphically.

- From the **Basic Statistics** menu, select **Tabulate**
- Under **Output Mode**, retain the default of **Counts**
- Under **Display**, check **Percentages of row totals**
- From the drop-down list next to **Columns**, select **COHORT90**
- Click **Tabulate**

You should see the following table in the **Output** window.

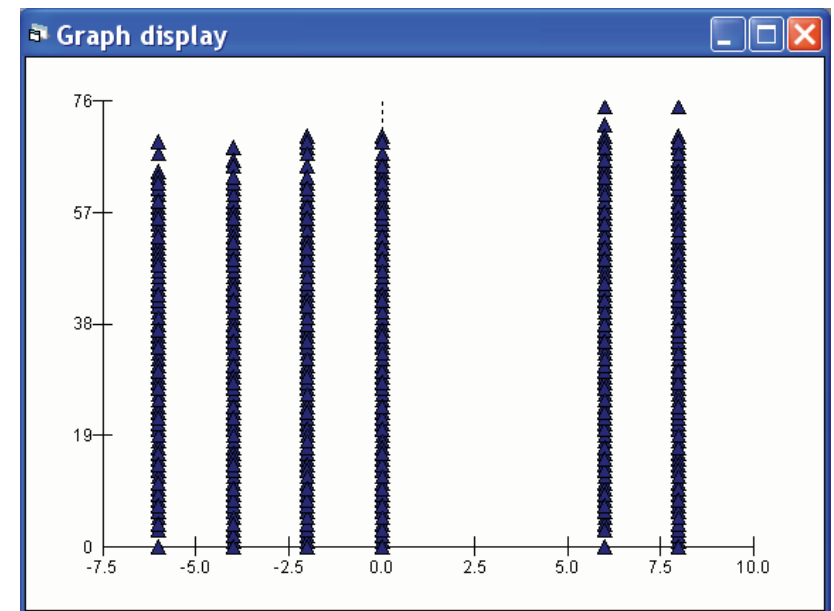
	-6	-5	-4	-3	-2	
N	6478	0	6325	0	5245	
%	19.1	0.0	18.6	0.0	15.4	
	-1	0	1	2	3	
N	0	4371	0	0	0	
%	0.0	12.9	0.0	0.0	0.0	
	4	5	6	7	8	TOTALS
N	0	0	4244	0	7325	33988
%	0.0	0.0	12.5	0.0	21.6	100.0

The number of observations in each category from -6 (year 1984) to 8 (year 1998) are shown. Some are empty because the cohorts in the sample are not from consecutive years. Shown below the number of students (N) is the percentage in each category. The largest proportion of students are from the 1998 cohort, with somewhat smaller proportions from 1990 and 1996.

Relationship between SCORE and COHORT90

Before fitting a linear regression model with attainment and cohort, we will examine the nature of their relationship using a scatter plot.

- From the **Graphs** menu, select **Customized Graph(s)**
- From the drop-down list labelled **plot type**, select **point**
- From the drop-down list labelled **y**, select **SCORE**
- From the drop-down list labelled **x**, select **COHORT90**
- Click **Apply**



Clicking anywhere in the **Graph display** window will bring up the **Graph options** window from which you can, for example, add titles and change the axes scales.

Although there is some suggestion of a positive linear trend, it is difficult to see the relationship because of the small number of distinct values of COHORT90. We will therefore supplement the scatterplot with a table of the mean attainment score for each value of COHORT90.

Before requesting this table, we will compute a recoded version of COHORT90 with consecutive values 1, 2, ..., 6. This will avoid empty cells in the table where there are gaps between cohorts.

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.