

# EMAT30007 Applied Statistics

## Lab 6: Linear Models (Simple Linear Regression)

Nikolai Bode

This lab focusses on simple linear regression models covered in lecture 6. These models take the general form:  $Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$ .

We will explore three data sets in this lab. The first two are from real-world studies and the last one is a famous statistical data set (references are given in square brackets):

(1) DRILLROCK.csv

[Penner, R., Watts D. G. (1991). "Mining information". The American Statistician. 45(1): 6]

(2) APPLIANCES.csv

[Candanedo, L. M., Feldheim, V., Deramaix, D. (2017) "Data driven prediction models of energy use of appliances in a low-energy house". Energy and Buildings, Volume 140, 1 April 2017, Pages 81-97]

(3) ANSCOMBE.txt

[Anscombe, F. J. (1973). "Graphs in Statistical Analysis". American Statistician. 27 (1): 17–21]

... you can download all of the data sets on Blackboard.

### (1) Dataset DRILLROCK

This data is from a rock drilling experiment. There are two processes for hydraulic drilling of rock: dry drilling and wet drilling. In a dry hole, compressed air is forced down the drill rods to flush the cuttings and drive the hammer. In a wet hole, water is used instead of air. An experiment was conducted to determine whether the time it takes to dry drill a distance of 5 feet in rock increases with depth. The results for one portion of the experiment are provided in the file.

**1.1. Read in the data, see how many data points there are and plot a scatter plot of the data.**

```
% read data from file:
delimiterIn = ',';
A = importdata('DRILLROCK.csv',delimiterIn);

% how many data points are there?
A.data
size(A.data)

%% plot scatterplots
clf
```

```
x = A.data(:,1); % depth data
y = A.data(:,2); % time data
plot(x,y, '.')
title('Scatterplot') % title for plot
xlabel('depth') % x-axis label
ylabel('time') % y-axis label
```

There are  $n = 17$  data points and from the scatterplot it already looks like there is a relationship between depth and time.

## 1.2. Fit a simple linear model with intercept to this data.

To demonstrate the model fitting process, we will use three different approaches in Matlab to fit the model. Of course, we could simply use ready-made model fitting tools in Matlab (as we will do most of the time), but it is instructive to see how model fitting works. The general structure of simple linear models is given above and in the lecture notes. We will use 'depth' as an explanatory variable and 'time' as our response variable.

We begin with the first method using the *exact equations* for estimating the model parameters,

$\hat{\beta} = (X'X)^{-1}X'Y$  and  $\hat{\sigma}^2 = \frac{1}{n-2}(Y - X\hat{\beta})'(Y - X\hat{\beta})$ . In Matlab  $x'$  is the transpose of a matrix  $x$  and `inv(X)` gives the inverse of a matrix.

```
Y = y;
% define the matrix X. Don't forget about the intercept.
X = [ones(17,1) x];
% ... continue from here using the equations above...
```

You should find  $\hat{\beta}_0 = 4.7896$ ,  $\hat{\beta}_1 = 0.0144$  and  $\hat{\sigma} = 1.4322$ .

In our next approach to fit the model, we perform a *numerical maximisation of the likelihood function*. Recall from lecture 6 that the likelihood function for our model is the product of the normal probability densities for all data points, given the model. In practice, rather than maximising this likelihood directly, we work with the log-likelihood. This avoids *underflow errors* where numbers become too small for the computer to cope with them (this happens a lot if we multiply many small probability densities). In fact, because Matlab has convenient options for finding the minimum of functions, we will find the minimum of the negative log-likelihood.

To streamline our code, we can write a separate Matlab file that computes the negative log-likelihood for our model given parameter values and data (i.e. values for response and explanatory variables). This is provided for you and is called `negloglik.m`. Have a look at this file and make sure you understand what it does.

```
% define values for parameters that we use at the start of the optimisation.
% We have 3 parameters: sigma, beta_0, beta_1
% x0 is the same length as the vector 'paras' in the function 'negloglik'.
x0 = [sqrt(var(y)) 1 0];

% check that computing the negative log-likelihood using 'negloglik.m' works:
negloglik([1 2 0],1,2) % uses sigma=1, beta_0=2 and beta_1=0 for data: Y=1 and x=2.
```

```

negloglik(x0,y,x) % uses starting values of parameters and all data.

% maximum likelihood fit:
fun = @(s)negloglik(s,y,x); % this sets up the function we want to minimise.
s = fminsearch(fun,x0); % here we performe the optimisation
% (i.e. find parameters for which negloglik is minimal)

% The fitted parameters are stored in 's':
s(1) % sigma
s(2) % beta_0
s(3) % beta_1

```

We (should) find the same values for  $\hat{\beta}_0$  and  $\hat{\beta}_1$ . However, the variance estimate is different to the value of  $\hat{\sigma}$  we found before. In fact, it is equal to  $\frac{1}{n} (Y - X\hat{\beta})'(Y - X\hat{\beta})$ . We say that the maximum likelihood estimator for the variance is *biased*. A detailed explanation for this is beyond the scope of the course, but notice that the difference between the estimate we use above and the maximum likelihood estimate will become negligible for large data sets (i.e. large  $n$ ).

Fortunately, Matlab does all the hard work in fitting simple linear models for us. All we need to know are a few simple commands and how to use them:

```

% create a table for the data (x,y, defined above):
data = table(x,y,'VariableNames',{'depth','time'});
% show first few rows of table
data(1:5,:)

% fit simple linear model
m1 = fitlm(data, 'time~depth');

```

Note how the structure of the model is defined using the variable names of the table `data` we defined. The standard notation in Matlab is that '~' separates the response variable from the predictor. An intercept is included into the model by default.

The object `m1` stores a lot of useful information about the model fit (see <https://uk.mathworks.com/help/stats/linear-model-class.html>) and we can obtain a fairly comprehensive report on the model simply by calling `m1`:

```
m1
```

An explanation for this output is provided in the lecture notes.

Before interpreting the parameter estimates and the hypothesis tests, we need to ensure our model is appropriate for the data (model checking).

### **1.3. Compute the residuals for the model fit and use residual plots to check that no model assumptions are violated.**

In lecture 6, we discussed the assumptions we make when using simple linear models. One way to check these assumptions are not violated is to look at residual plots. Once we have fitted our model to the data and obtained parameter estimates, it is easy to compute residuals  $\hat{\varepsilon} = Y - X\hat{\beta}$ . Compute

residuals in this way. Plot their distribution (using `hist()`), to see if they roughly follow a normal distribution. In addition, plot the fitted values,  $X\hat{\beta}$ , against the residuals, to ensure homoscedasticity.

```
% use parameter estimates and data to find residuals.  
% you can use the matrices and parameters from question 1.2 here.  
  
% plot a histogram of the residuals, use hist()  
  
% plot the residuals versus the fitted values
```

In Matlab we can use a few simple commands to obtain the key residual plots. To show the residual plots discussed in lectures, you can use the following code:

```
% plot distribution of residuals (for outliers etc.)  
clf  
subplot(2,2,1)  
plotResiduals(m1)  
% plot to check normality  
subplot(2,2,2)  
plotResiduals(m1,'probability')  
% residuals versus fitted values (check for homoscedasticity)  
subplot(2,2,3)  
plotResiduals(m1,'fitted')  
% auto-correlation (via lagged residuals)  
subplot(2,2,4)  
plotResiduals(m1,'lagged') % want no trend in this!
```

Looking at these residual plots, can you see any obvious reasons for why the model assumptions may not hold?

Answer: overall, the model assumptions seem appropriate. The data form a fairly straight line in the normal probability plot and the histogram of the residuals is somewhat mound-shaped. There is no indication that the data are not normally distributed. The plot of the residuals against the fitted values does not show evidence for heteroscedasticity (e.g. funnel shape or a football shape). There is therefore no indication that the assumption of constant variance is violated. There is also no evidence that residuals are auto-correlated.

#### 1.4. Check the model fit visually

We have ensured the model assumptions are appropriate. Now we need to check that the model fit is sensible. For simple linear regression models, it is always worth it to plot the model we have fitted alongside the data.

```
% check model fit:  
clf  
% plot the raw data:  
plot(x,y,'.')  
hold on  
% plot model fit line using the estimates parameters:  
xx = [0 400];  
yy = [m1.Coefficients{1,1}+0*m1.Coefficients{2,1} m1.Coefficients{1,1}+400*m1.Coefficients{2,1}];  
plot(xx,yy,'k','LineWidth',2)  
xlabel('depth') % x-axis label  
ylabel('time') % y-axis label  
title('Model fit line')
```

This looks sensible, as the straight line appears to capture the trend in the data reasonably well.

We have discussed alternative measures for model fit in the lecture (e.g.  $R^2$ ) and Matlab computes some of these automatically. We won't focus on these for now, as we will look at such measures in more detail in lecture 7 & lab 7.

### 1.5. Does the time it takes to dry drill a distance of 5 feet in rock increase with depth?

We have now fitted our model and checked it thoroughly. Now we can investigate the parameter estimates and interpret what they mean. To answer the question above, we look at the parameter estimate  $\hat{\beta}_1$ , or the slope of our model. It is positive, suggesting that as the depth increases, the time to drill 5 feet increases. We should formally test this assertion by computing a p-value for the null hypothesis  $H_0 : \beta_1 = 0$ . While the result of this test is already provided as part of the Matlab output above, it is instructive to repeat this computation manually.

Use the formula for the T-statistic for the model parameter  $\hat{\beta}_1$  given in the lecture notes and the cumulative distribution function of the Student's T-distribution in Matlab (`tcdf()`) to find the p-value for the hypothesis above.

```
% re-calculate the error estimate
% you can use the matrices and parameter estimates from question 1.2.
sigmahat = sqrt((Y-X*betahat)'*(Y-X*betahat)/(17-2));
% calculate the quantity SS_xx
ssxx = sum((x-mean(x)).*(x-mean(x)));
% ... continue from here...
```

### 1.6. Compute the prediction and estimation intervals for the model fit.

We briefly covered this in lectures. Once we have fit a model to data, we can then use it to estimate the mean value of the response, or to predict the response, for a given value of the explanatory variable. The confidence intervals for prediction and estimation can easily be obtained in Matlab. The default are 95% confidence intervals.

```
% prediction and estimation intervals:
xnew = 0:0.1:400;
% confidence bounds for fitted mean values (estimation)
[ypred,yci] = predict(m1,xnew','Prediction','curve');
% confidence bounds for new observations (prediction)
[ypred2,yci2] = predict(m1,xnew','Prediction','observation');

% plot these intervals:
xx = [0 400];
yy = [m1.Coefficients{1,1}+0*m1.Coefficients{2,1}...
      m1.Coefficients{1,1}+400*m1.Coefficients{2,1}];
clf
plot(xx,yy,'k','LineWidth',2)
% estimation and prediction limits
hold on
plot(xnew,yci(:,1),'k',xnew,yci2(:,1),'k--')
hold on
plot(xnew,yci(:,2),'k',xnew,yci2(:,2),'k--')
legend('fit','95% estimation limits','95% prediction limits')
xlabel('depth') % x-axis label
```

```
ylabel('time') % y-axis label
```

Above, we have calculated estimation and prediction intervals for many values of the depth. Now compute these intervals for a value of the depth of 200 and 800. What is problematic about computing estimates or making predictions using the latter value of the depth?

```
% give estimation and prediction values for Y for two values: depth=200 and depth=800.
```

## (2) Dataset APPLIANCES

This data was collected to investigate the energy use of appliances in low-energy buildings. The data set is from one house. Outside temperature and appliance energy use readings were taken every 10 minutes for about 4.5 months. Thus, there are a lot more data points than in the first data set and this quantity of data is more representative of what statistical analysis currently looks like.

Use the techniques we have covered so far to answer the following exercises. You can use the standard Matlab commands for this analysis.

### ***2.1. Read in the data, see how many data points there are and plot a scatter plot of the data.***

Looking at this data, do you think a simple linear model will describe the data well?

```
% ... you can use similar code to question 1.
```

### ***2.2. Fit a simple linear model with intercept to the data.***

Do not start to interpret the model estimates or outcomes of the hypothesis tests yet. We need to check the model assumptions hold first!

```
% ... you can use similar code to question 1.
```

### ***2.3. Use residual plots to check that no model assumptions are violated.***

Referring to the assumptions simple linear models make, explain whether the residual plots show evidence of problems with the model fit. If there are problems with the residual plots, speculate, using your knowledge of the data, what might cause these problems.

```
% ... you can use similar code to question 1.
```

*Discussion of residual plots:*

The distribution of residuals does not look like a normal distribution. The peak around zero could be normally distributed, but there is another peak, making the distribution *bimodal*.

This non-normal distribution is also clearly indicated in the normal probability plot. We could speculate that outside temperature is not the only factor that explains energy use in the house. For example, time of day or day of the week could also be relevant.

Plotting fitted values versus residuals confirms our observation from the scatterplot of the data that there are many high-energy data points our model cannot explain.

Plotting residuals versus lagged residuals shows some evidence of autocorrelation in the residuals (i.e. the plot has some structure along the line indicating positive correlation). In some ways that is not surprising. The data set form a time series (readings every 10 minutes). Consecutive readings are therefore unlikely to be entirely independent (think about appliances being switched on for >10 mins).

Overall, we have to be very careful with interpreting this model fit to the data. Some model assumptions clearly do not hold. This means that the theory we use may give misleading results.

## 2.4. Check the model fit visually

Having seen the residual plots we should re-think our modelling approach. However, it's good to practice. Once you've had a look at the model fit visually, you can also look at the Matlab summary of the model fit. Keeping all the limitations about the model fit in mind, do you think temperature is a factor in determining the energy usage in the house?

```
% ... you can use similar code to question 1.
```

Thinking more broadly, if there is an effect, can we be sure that lower temperatures cause higher energy usage? Or could it be that temperatures are correlated with other factors (e.g. light levels) that are more important in determining the usage and energy consumption of appliances? Remember: *correlation does not imply causation!*

## (3) Dataset ANSCOMBE

This file actually contains four data sets. They were constructed in 1973 by the statistician Francis Anscombe to demonstrate both the importance of graphing data before analyzing it and the effect of outliers on statistical properties.

The file `ANSCOMBE.txt` has eight columns. The first two columns are the explanatory and response variables for the first data set, respectively. The third and fourth columns make up the second data set, and so forth.

**3.1. Read in the data sets and fit a simple linear model to each of them.**

**3.2. Compare the model fits.**

**3.3. Plot scatterplots of the data sets with the fit lines of the corresponding models.**

Without looking at residual plots, do you think it's appropriate to use simple linear models to investigate all of the data sets?

