

STATS201.stat201 20 Simple linear regression

Anna Helga Jónsdóttir
Sigrún Helga Lund

December 16, 2012

Copyright This work is licensed under the Creative Commons Attribution-ShareAlike License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-sa/1.0/> or send a letter to Creative Commons, 559 Nathan Abbott Way, Stanford, California 94305, USA.

Contents

1	Simple linear regression	3
1.1	Scatter plot	3
1.2	Scatter plot - continuous variables	3
1.3	The equation of a straight line	3
1.4	The equation of a straight line	4
1.5	Linear relationship	4
1.6	Linear and nonlinear relationship	4
1.7	Sample coefficient of correlation	5
1.8	The size and direction of a linear relationship	5
1.9	The size and direction of a linear relationship	6
1.10	Correlation and causation	6
1.11	The linear regression model	6
1.12	The least squares method	7
1.13	The least squares method	7
1.14	The least squares regression line	8
1.15	Residuals	9
1.16	Residual plot	9
1.17	Interpolation	9
1.18	Extrapolation	10
1.19	Coefficient of determination	10
1.20	Outliers and influential measurements	10
1.21	Outliers and influential measurements	11
1.22	Treatment of outliers and influential measurements	11
1.23	The linear regression model	11
1.24	The random variable ϵ	12
1.25	Confidence interval for β_0	12
1.26	Confidence interval for β_1	13
1.27	Prediction interval	13
1.28	Hypothesis test for the correlation coefficient	14

1 Simple linear regression

1.1 Scatter plot

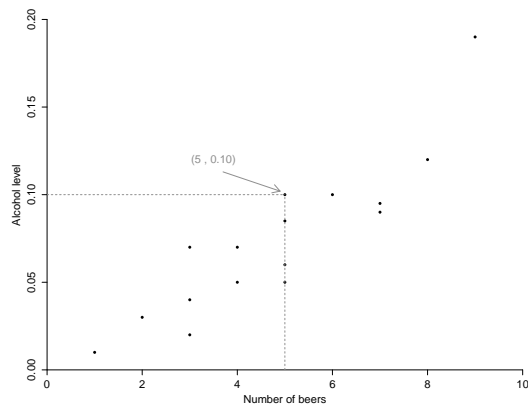
Scatter plot
Scatter plots are used to investigate the relationship between two numerical variables.

The value of one variable is on the y-axis (vertical) and the other on the x-axis (horizontal).

When one of the variable is an explanatory variable and the other one is a response variable, the response variable is always on the y-axis and the explanatory variable on the x-axis.

Response variables and explanatory variables
For every subject, the value of an **explanatory variable** will influence what value the **response variable** receives.

1.2 Scatter plot - continuous variables



1.3 The equation of a straight line

The equation of a straight line
The equation of a straight line describes a linear relationship between two variables, x and y . The equation is written

$$y = \beta_0 + \beta_1 x$$

where β_0 is the **intercept** of the line on the y-axis and β_1 is the **slope** of the line.

1.4 The equation of a straight line

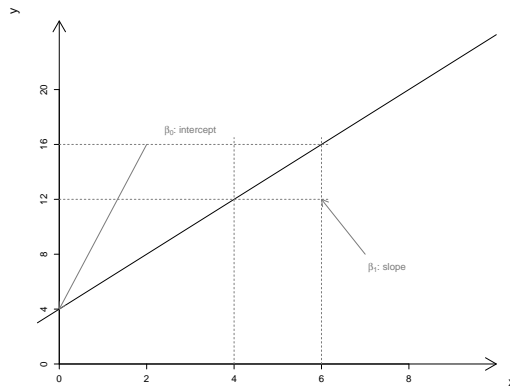


Figure 1: The equation of a straight line.

1.5 Linear relationship

Linear relationship
We say that the relationship between two variables is **linear** if the equation of a straight line can be used to predict which value the response variable will take based on the value of the explanatory variable.

Notice that there can be all sorts of relationship between two variables. For example, the relationship can be described with a parabola, an exponential function and so on. Those relationship are referred to as nonlinear and are not covered in this lecture.

1.6 Linear and nonlinear relationship

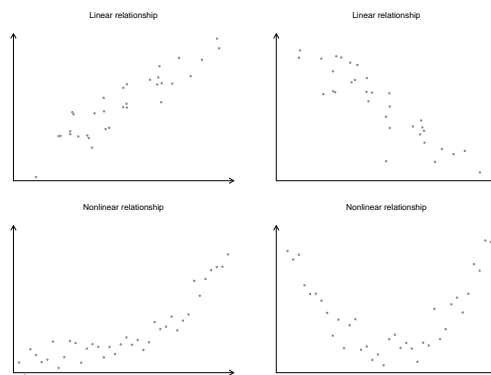


Figure 2: Scatter plot where the relationship between two variables is linear (above) and nonlinear (below).

1.7 Sample coefficient of correlation

Sample coefficient of correlation

Assume that we have n measurements on two variables x and y .

Denote the mean and the standard deviation of the variable x with \bar{x} and s_x and the mean and the standard deviation of the y variable with \bar{y} and s_y .

The sample coefficient of correlation is

$$r = \frac{1}{n-1} \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right);$$

Be careful! We only use correlation to estimate **linear** relationship!

1.8 The size and direction of a linear relationship

The direction of a linear relationship

The sign of the correlation coefficients determines the **direction** of a linear relationship. It is either positive or negative.

- If the correlation coefficient of two variables is positive, we say that their correlation is **positive**.
- If the correlation coefficient of two variables is negative, we say that their correlation is **negative**.

The size of a linear relationship

The absolute value of a correlation coefficient describes the **size** of the linear relationship between the variables.

It tells us how well we can predict the value of the response variable from the value of the explanatory variable.

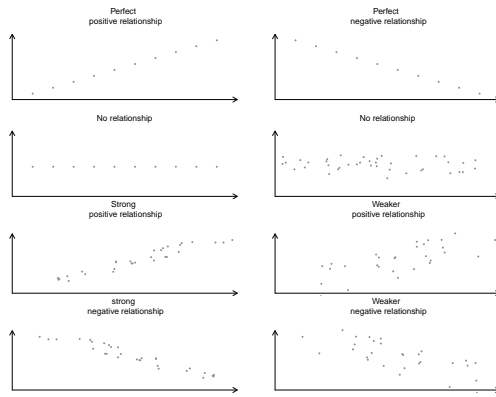


Figure 3: Scatter plot for various values of r .

1.9 The size and direction of a linear relationship

1.10 Correlation and causation

- **Causation** is when changes in one variable **cause** changes in the other variable.
- There is often strong correlation between two variables although there is no causal relationship.
- In many cases, the variables are both influenced by the third variable which is then a **lurking variable**.
- Therefore, high correlation on its own is never enough to claim that there is a causal relationship between two variables.

1.11 The linear regression model

The linear regression model
The simple linear regression model is written

$$Y = \beta_0 + \beta_1 x + \epsilon$$

when β_0 and β_1 are unknown parameters and ϵ is a normally distributed random variable with mean 0.

The aim of the simple linear regression is first and foremost to estimate the parameters β_0 and β_1 with the measurements of the two variables, x and Y .

The method we use is called the least squares method.

1.12 The least squares method

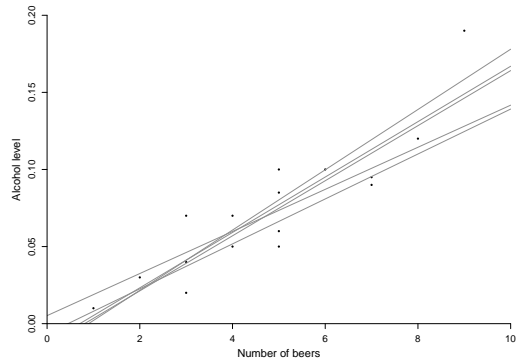


Figure 4: Many lines, but which one is the best?

1.13 The least squares method

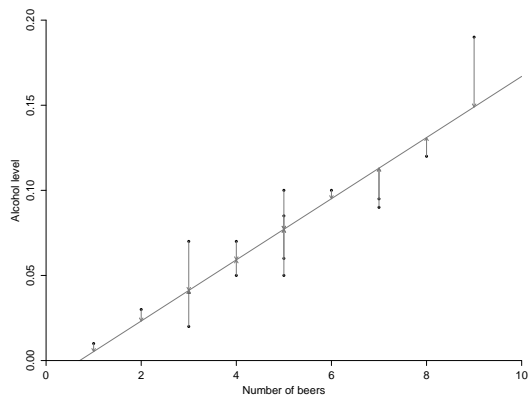


Figure 5: The least squares method.

1.14 The least squares regression line

The least squares regression line
 Denote the mean and standard deviation of the x variable with \bar{x} and s_x and the y variable with \bar{y} and s_y and their correlation coefficient with r .

Let b_0 denote the estimate of β_0 and b_1 denote the estimate of β_1 . Then b_0 and b_1 are given with the equation

$$\hat{\beta}_1 = b_1 = r \frac{s_y}{s_x}$$

and

$$\hat{\beta}_0 = b_0 = \bar{y} - \beta_1 \bar{x}.$$

We use the parameters to **predict** the value of y for a given value of x with the least squares regression line

$$\hat{y} = b_0 + b_1 x$$

Þorgerður and Birna like beer a lot. They decided to make an experiment to investigate the relationship between the alcohol level in blood and the number of consumed beers. 16 students took part in the experiment, the data can be seen below.

2*Nemi	Fjöldi	Alkóhólmagn	2*Nemi	Fjöldi	Alkóhólmagn
	bjóra	í blóði		bjóra	í blóði
1	5	0.100	9	8	0.120
2	2	0.030	10	3	0.040
3	9	0.190	11	5	0.060
4	7	0.095	12	5	0.050
5	3	0.070	13	6	0.100
6	3	0.020	14	7	0.090
7	4	0.070	15	1	0.010
8	5	0.085	16	4	0.050

Use the method of least squares to fit a regression line to the data. From the data we can calculate:

$$\bar{x} = 4.813, \quad s_x = 2.198, \quad \bar{y} = 0.074, \quad s_y = 0.044, \quad r = 0.894.$$

The slope is

$$b_1 = r \frac{s_y}{s_x} = 0.894 \cdot \frac{0.044}{2.198} = 0.018$$

and the intercept is:

$$b_0 = \bar{y} - \beta_1 \bar{x} = 0.074 - (0.018 \cdot 4.813) = -0.013.$$

so the regression line is

$$\hat{y} = -0.013 + 0.018x.$$

1.15 Residuals

Residuals
The vertical distance from our measurements to the regression line are called the **residuals** and are denoted with e . The size of the residuals can be calculated with

$$e_i = y_i - \hat{y}_i$$

Points above the regression line have a positive residue but points below it have a negative.

1.16 Residual plot

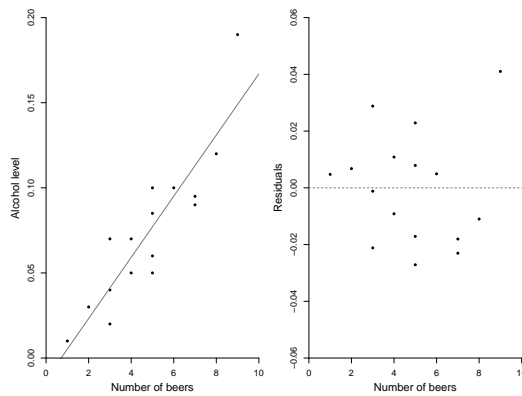


Figure 6: Scatter plot of the data and a residual plot.

1.17 Interpolation

Interpolation
If the regression model is used to predict a value of Y for some value of x which is similar to the x -values that were used to estimate the model is referred to as **interpolating**.

Let us continue with the beer example. Predict the alcohol level in the blood of a person that has drunken 6.5 beers.

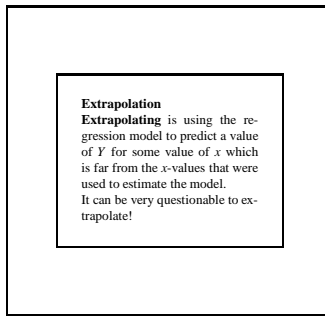
The regression equation is:

$$\hat{y} = -0.013 + 0.018x.$$

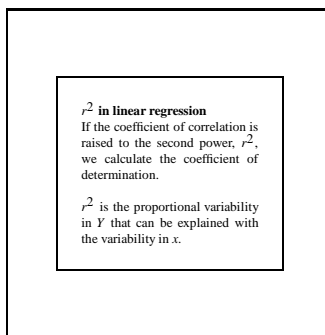
We used data from people drinking from one to nine beers so we are interpolating here. We insert 6.5 in the equation and get:

$$\hat{y} = -0.013 + (0.018 \cdot 6.5) = 0.104.$$

1.18 Extrapolation



1.19 Coefficient of determination



We continue to work with the beer data. How much of the variability in the alcohol level can be explained by the number of consumed beers.

We saw earlier that $r = 0.894$. So we get that $r^2 = 0.894^2 = 0.799$. Around 80% of the variability in alcohol level can be explained by the number of beers consumed.

1.20 Outliers and influential measurements

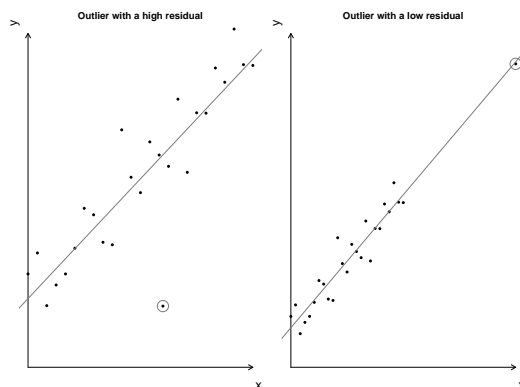


Figure 7: Outliers and their residuals.

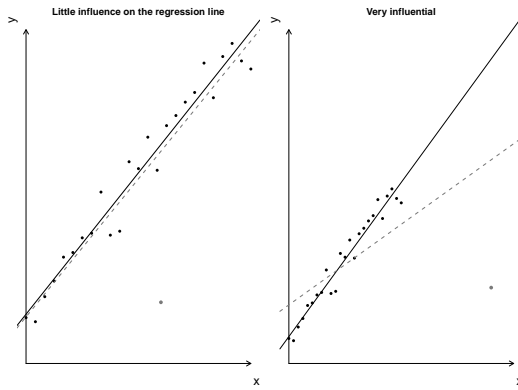


Figure 8: Influential measurements?

1.21 Outliers and influential measurements

1.22 Treatment of outliers and influential measurements

- Outliers and influential measurements shall always be viewed carefully.
- If a mistake has been made, the measurement shall be eliminated.
- If it cannot be shown that a mistake has been made it is often good to show estimates with and without these measurements.
- In some cases it is more appropriate to use the estimates without the outliers/influential measurements.
- In these cases, it shall be pointed out that the model does not fit data outside the range of the measurements used for estimating the model.

1.23 The linear regression model

If we have n paired measurements $(x_1, y_1), \dots, (x_n, y_n)$, the regression model can be written as

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i.$$

- β_0 is the true intercept (population intercept) that we do not know the value of.
- β_1 is the true slope (population slope)
- ε_i are the errors.

β_0 and β_1 are therefore statistics, that we both want to estimate and make inference on.

We do that by applying the least squares method to our data.

1.24 The random variable ε

ε describes the uncertainty in our measurements of Y .

We assume that ε_i are independent and identically distributed random variables that follow a normal distribution with mean 0 and variance σ^2 .

Estimating σ^2 in simple linear regression

The estimate of σ^2 in simple linear regression is denoted with s_e^2 and calculated with

$$s_e^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}.$$

This is the same equation as for the "normal" standard deviation, but now we divide with $n-2$ but not $n-1$.

1.25 Confidence interval for β_0

Confidence interval for β_0

The lower bound of a $1-\alpha$ confidence interval for β_0 is:

$$b_0 - t_{1-\alpha/2, (n-2)} \cdot s_e \sqrt{\frac{1}{n} + \frac{(\bar{x})^2}{s_x^2 \cdot (n-1)}}$$

The upper bound of $1-\alpha$ confidence interval is:

$$b_0 + t_{1-\alpha/2, (n-2)} \cdot s_e \sqrt{\frac{1}{n} + \frac{(\bar{x})^2}{s_x^2 \cdot (n-1)}}$$

where b_0 is calculated the same way as usual, n is the number of paired measurements, \bar{x} is the mean of the explanatory variable, s_x is the standard deviation of the explanatory variable and $t_{1-\alpha/2, (n-2)}$ is in the table for the t-distribution.

1.26 Confidence interval for β_1

Confidence interval for β_1
 The lower bound of $1 - \alpha$ confidence interval for β_1 is:

$$b_1 - t_{1-\alpha/2, (n-2)} \cdot s_e \frac{1}{\sqrt{s_x^2 \cdot (n-1)}}$$

The upper bound of $1 - \alpha$ confidence interval is:

$$b_1 + t_{1-\alpha/2, (n-2)} \cdot s_e \frac{1}{\sqrt{s_x^2 \cdot (n-1)}}$$

where b_1 is calculated the same way as usual, n is the number of paired measurements, s_x is the standard deviation of the explanatory variable and $t_{1-\alpha/2, (n-2)}$ is found in the t-distribution table.

1.27 Prediction interval

Prediction interval
 The lower bound of $1 - \alpha$ prediction interval for Y is:

$$(b_0 + b_1 x_0) - t_{1-\alpha/2, (n-2)} \cdot s_e \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_x^2 (n-1)}}$$

The upper bound of $1 - \alpha$ prediction interval is:

$$(b_0 + b_1 x_0) + t_{1-\alpha/2, (n-2)} \cdot s_e \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{s_x^2 (n-1)}}$$

where b_0 and b_1 are calculated the same way as usual, n is the number of paired measurements, s_x is the standard deviation of the explanatory variable and $t_{1-\alpha/2, (n-2)}$ is found in the t-distribution table.

1.28 Hypothesis test for the correlation coefficient

Hypothesis test for ρ
 The null hypothesis is:

$$H_0 : \rho = 0$$

The test statistic is:

$$T = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}$$

If the null hypothesis is true the test statistic follows the t distribution with $n-2$ degrees of freedom or $T \sim t(n-2)$.

Alternative hypothesis	Reject H_0 if:
$H_1 : \rho < 0$	$T < -t_{1-\alpha}$
$H_1 : \rho > 0$	$T > t_{1-\alpha}$
$H_1 : \rho \neq 0$	$T < -t_{1-\alpha/2}$ or $T > t_{\alpha/2}$

Atli is making an experiment to investigate whether there is a relationship between the icecream sales in a certain shop and the temperature outside. He looks at sales numbers and temperature data on 38 days he chose randomly. He calculated the correlation to be 0.5. Can Atli conclude that the variables temperature and icecream sales are correlated. Use $\alpha = 0.05$.

1. We would like to make a hypothesis test for a correlation.
2. $\alpha = 0.05$.
3. The hypotheses are:

$$H_0 : \rho = 0$$

$$H_1 : \rho \neq 0.$$

4. The test statistic is:

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}}.$$

We have $n = 38$ and $r = 0.5$.

$$t = \frac{0.5\sqrt{38-2}}{\sqrt{1-0.5^2}} = \frac{0.5\sqrt{36}}{\sqrt{1-0.25}}$$

$$= \frac{0.5 \cdot 6}{\sqrt{0.75}} = \frac{3}{\sqrt{0.75}} = 3.46.$$

5. We have $n - 2 = 36$ degrees of freedom. $t_{1-\alpha/2, (n-2)} = t_{0.975, (36)} = 2.028$, so we reject the null hypothesis if $t > 2.028$ or if $t < -2.028$.

We see that $t = 3.46 > 2.028$.

6. We reject the null hypothesis and conclude that temperature and icecream sales are correlated.