

Linear Regression

Fabian Farelo, Autar Kaw, Jamie Trahan

University of South Florida

United States of America

kaw@eng.usf.edu

Introduction

This worksheet demonstrates the use of *Mathematica* to illustrate the procedure to regress a given data set to a straight line.

Linear Regression is the most popular regression model. In this model we wish to predict response points to n data points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ by a regression model given by:

$$y = a_0 + a_1 x \quad (1.1)$$

where a_0 and a_1 are the constants of the regression model. A measure of goodness of fit, that is, how $a_0 + a_1 x$ predicts the response variable y is the magnitude of the residual, ε_i , at each of the n data points

$$\varepsilon_i = (\text{observed value at } x_i - \text{predicted value at } x_i) = y_i - (a_0 + a_1 * x) \quad (1.2)$$

Ideally, if all the residuals ε_i are zero, one will find an equation in which all the points lie on the model. Thus, minimization of the residual is an objective of obtaining regression coefficients. The most popular method to minimize the residual is the least squares method, where the estimates of the constants of the models are chosen such that the sum of the squared residuals, S_r is minimized, that is minimize

$$S_r = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - (a_0 + a_1 x))^2 \quad (1.3)$$

Let us use the least squares criterion where we minimize the sum of the squared residual, S_r :

$$\frac{\partial S_r}{\partial a_0} = -2 \sum_{i=1}^n (y_i - a_0 - a_1 x_i) (-1) = 0 \quad (1.4)$$

$$\frac{\partial S_r}{\partial a_1} = -2 \sum_{i=1}^n (y_i - a_0 - a_1 x_i) (-x_i) = 0 \quad (1.5)$$

Once S_r is minimized with respect to the regression coefficient, a_0 and a_1 (see Linear Regression notes for complete derivation), the coefficients can be solved for:

$$a_1 = \frac{S_{xy}}{S_{xx}} \quad (1.6)$$

$$a_0 = y_{ave} - a_1 * x_{ave} \quad (1.7)$$

where S_{xy} and S_{xx} can be defined as:

$$S_{xy} = \sum_{i=1}^n (x_i y_i) - n * x_{ave} * y_{ave} \quad (1.8)$$

$$S_{xx} = \sum_{i=1}^n x_i^2 - n * (x_{ave})^2 \quad (1.9)$$

and the average values can be defined as:

$$x_{ave} = \frac{\sum_{i=1}^n x_i}{n} \quad (1.10)$$

$$y_{ave} = \frac{\sum_{i=1}^n y_i}{n} \quad (1.11)$$

Section 1: Input

The following are the input parameters to begin the simulation. The user can change those values that are highlighted and the worksheet will calculate an approximate solution to the system of equations

- Number of Data Points, n

$$n = 5$$

5

- Data Points

$$\mathbf{x} = \{1, 2, 3, 4, 5\}$$

{1, 2, 3, 4, 5}

$$\mathbf{y} = \{1, 4, 9, 16, 25\}$$

{1, 4, 9, 16, 25}

Section 2: Linear Regression Method

- Calculating average values using Equations (1.10) and (1.11):

```
Off [General::spell1]
```

$$\bar{X} = \frac{\sum_{i=1}^n X[[i]]}{n}$$

```
3
```

$$\bar{Y} = \frac{\sum_{i=1}^n Y[[i]]}{n}$$

```
11
```

- Calculating S_{xy} and S_{xx} using Equations (1.8) and (1.9):

```
Sxy = 0;
Sxx = 0;
For[i = 1, i ≤ n, i++,
  Sxy = Sxy + X[[i]] * Y[[i]] - Xave * Yave;
  Sxx = Sxx + (X[[i]] ^ 2) - Xave ^ 2]
```

```
Sxy
```

```
60
```

```
Sxx
```

```
10
```

- S_{xx} , S_{xy} , \bar{x} , and \bar{y} can now be used to calculate the regression coefficients, a_0 and a_1 using Equations (1.6) and (1.7).

```
a1 = N[Sxy / Sxx]
```

```
6.
```

```
a0 = N[Yave - a1 * Xave]
```

```
-7.
```

Section 3: Results

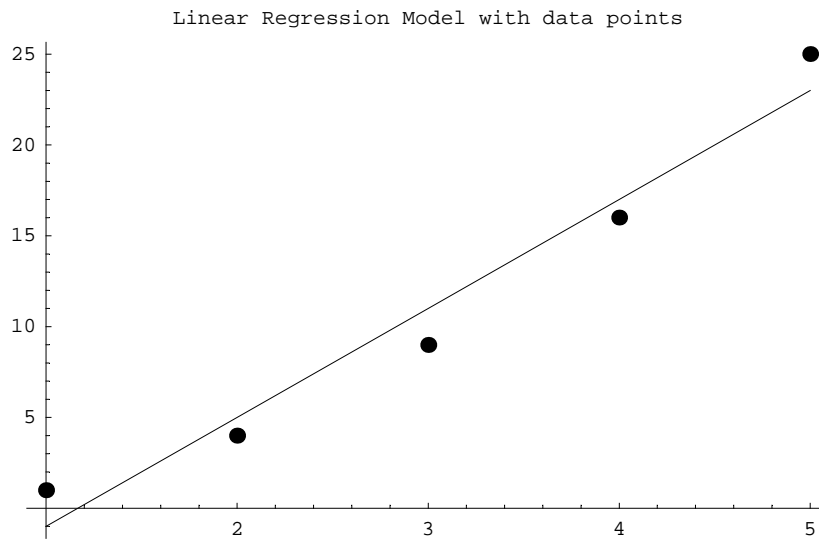
The linear model is described as:

```
y = a0 + a1 * X;
Print["y=", a0, "+", a1, "x"]
```

```
y=-7.+6.x
```

The following figure demonstrates the data points as well as the least squares regression line:

```
Datatable = Table[{X[[i]], Y[[i]]}, {i, 1, n}];  
points = ListPlot[Datatable, PlotStyle -> PointSize[0.02], DisplayFunction -> Identity];  
lin = Plot[a0 + a1 * x, {x, Min[X], Max[X]}, PlotLabel ->  
  {"Linear Regression Model with data points"}, DisplayFunction -> Identity];  
  
Show[points, lin, PlotLabel -> "Linear Regression Model with data points",  
  DisplayFunction -> $DisplayFunction];
```



Section 4: Coefficient of determination

One of the major indicators of how well least squares characterizes or predicts the whole data is a quantity called the **coefficient of determination**, r^2

$$r^2 = \frac{S_t - S_r}{S_t} \quad (4.1)$$

where

S_r = the sum of the squares of the residuals (a value that quantifies the spread around the regression line)

and

S_t = the sum of the squares of deviation from the mean (a value that measures the spread between the data and its mean).

The objective of least squares method is to obtain a compact equation that best describes all data points. The mean can also be used to describe only data points. The magnitude of the sum of squares of deviation from the mean or from the least squares line is therefore a good indicator of how well the mean or least squares characterizes the whole data.

The difference between the two parameters ($S_t - S_r$) measures the error due to describing or characterizing the data in one form instead of the other. A relative comparison of this difference with the sum of squares deviation associated with the mean, (i.e. r^2), describes the proportion of variation in the response data that is explained by the regression model. When all the points in a data set lie on the regression model, the largest possible value of $r^2 = 1$ is obtained, while a minimum possible value of $r^2 = 0$ is obtained when there is only one data point or if the straight line model is a constant line. (Note that $0 \leq r^2 \leq 1$)

Calculation of the coefficient of determination:

- Sum of the difference between observed values and average values, S_t :

$$st = \sum_{i=1}^n ((Y[[i]] - Yave)^2)$$

374

- Sum of the square of the residuals, S_r :

$$sr = \sum_{i=1}^n (Y[[i]] - \hat{Y}[[i]])^2$$

14.

- Coefficient of determination, r^2 :

$$r2 = (st - sr) / st$$

0.962567

Conclusion

Mathematica helped us regress a given data set to a straight line and determine how accurate the linear regression model fits the data.

Question 1: In the table below is given the instantaneous thermal expansion coefficient as a function of temperature. Find the linear regression model that relates the Instantaneous Thermal Expansion as a function of temperature. What is the coefficient of determination of the model?

Table 1: Instantaneous thermal expansion coefficient as a function of temperature.

<i>Temperature</i>	<i>Instantaneous Expansion</i>
Farenheit	E - 06 in / (in F)
80	6.47
60	6.36
40	6.24
20	6.12
0	6.00
-20	5.86
-40	5.72
-60	5.58
-80	5.43
-100	5.28
-120	5.09
-140	4.91
-160	4.72
-180	4.52
-200	4.30
-220	4.08
-240	3.83
-260	3.58
-280	3.33
-300	3.07
-320	2.76
-340	2.45

Question 2: In the table below is given the stress-strain data for a tensile test of a unidirectional composite material.

Strain (%)	Stress (MPa)
0	0
0.183	306
0.36	612
0.5324	917
0.702	1223
0.867	1529
1.0244	1835
1.1774	2140
1.329	2446
1.479	2752
1.5	2767
1.56	2896

Find the straight line regression model that finds the relationship between the stress and strain. Note that the straight line has no intercept. The slope of the straight line is the longitudinal Young's modulus of the composite material.

References

[1] Autar Kaw, *Holistic Numerical Methods Institute*, <http://numericalmethods.eng.usf.edu/nbm>, See Linear Regression