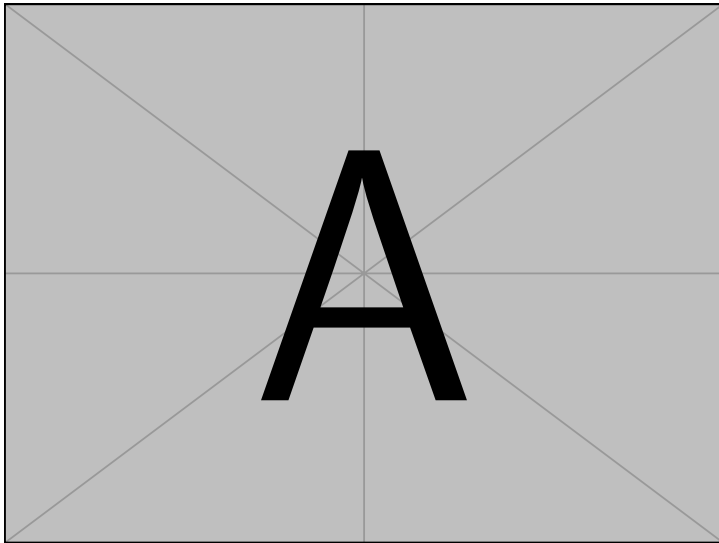


*CEE384—Numerical Methods*

## **Regression**

*Least Square Regression*



ARIZONA STATE UNIVERSITY

*Brian Chevalier*

Updated: November 27, 2019

### **Contents**

<b>1</b>	<b>Linear Regression</b>	<b>1</b>
<b>2</b>	<b>General Polynomial Regression</b>	<b>2</b>
	<b>References</b>	<b>3</b>

## 1 Linear Regression

A general linear equation takes the form:

$$y = a_0 + a_1x \quad (1)$$

where  $y$  is the function output,  $a_0$  is the constant term,  $a_1$  is the slope, and  $x$  is the function input. This equation can be used to interpolate between two points, or create a series of linear spline interpolation equations. However, here we will be looking not at precisely crossing all points in a data set, but instead, minimizing the error associated with creating a linear model for a set of data.

The data are a series of  $x$  and  $y$  coordinates in the form shown in the following equations:

$$\mathbf{x} = [x_0 \quad x_1 \quad \dots \quad x_{n-1}] \quad (2)$$

$$\mathbf{y} = [y_0 \quad y_1 \quad \dots \quad y_{n-1}] \quad (3)$$

where  $n$  is the number of elements in the data set.

The linear regression equation is:

$$y = a_0 + a_1x + e \quad (4)$$

where  $e$  is the error (sometimes called the residual), or the difference between the equation and the actual data set. Now, we want to *minimize* the sum of the squares of the error.

$$S_r = \sum_{k=0}^{k<n} (e_k)^2 = \sum_{k=0}^{k<n} (y_k - a_0 - a_1x_k)^2 \quad (5)$$

where the subscript  $k$  denotes the  $k$ th element of the  $x$  and  $y$  data set. Taking the derivative with respect to each coefficient term yields:

$$\frac{\partial S_r}{\partial a_0} = -2 \sum_{k=0}^{k<n} (y_k - a_0 - a_1x_k) = 0 \quad (6)$$

$$\frac{\partial S_r}{\partial a_1} = -2 \sum_{k=0}^{k<n} x_k (y_k - a_0 - a_1x_k) = 0 \quad (7)$$

Distributing the sum and rearranging:

$$na_0 + a_1 \sum x_k = \sum y_k \quad (8)$$

$$a_0 \sum x_k + a_1 \sum x_k^2 = \sum x_k y_k \quad (9)$$

Note that the bounds on the summations have been dropped for the sake of brevity and is short for  $\sum_{k=0}^{k<n}$ .

Putting the system of equations in matrix form:

$$\begin{bmatrix} n & \sum x_k \\ \sum x_k & \sum x_k^2 \end{bmatrix} \begin{Bmatrix} a_0 \\ a_1 \end{Bmatrix} = \begin{Bmatrix} \sum y_k \\ \sum x_k y_k \end{Bmatrix} \quad (10)$$

Solving the system of equations yields:

$$a_1 = \frac{n \sum x_k y_k - \sum x_k \sum y_k}{n \sum x_k^2 - (\sum x_k)^2} \quad (11)$$

$$a_0 = \bar{y} - a_1 \bar{x} \quad (12)$$

where  $\bar{y} = \frac{\sum y_k}{n}$  and  $\bar{x} = \frac{\sum x_k}{n}$ .

## 2 General Polynomial Regression

To establish the general n-order polynomial equations, we will derive the quadratic equations and find the pattern between these and the linear equations. The quadratic regression formulas can be derived through the same method as their linear counterparts.

$$y = a_0 + a_1x + a_2x^2 + e \quad (13)$$

$$S_r = \sum_{k=0}^{k < n} (y_k - a_0 - a_1x_k - a_2x_k^2)^2 \quad (14)$$

Taking the derivative with respect to each coefficient:

$$\frac{\partial S_r}{\partial a_0} = -2 \sum (y_k - a_0 - a_1x_k - a_2x_k^2) \quad (15)$$

$$\frac{\partial S_r}{\partial a_1} = -2 \sum x_k (y_k - a_0 - a_1x_k - a_2x_k^2) \quad (16)$$

$$\frac{\partial S_r}{\partial a_2} = -2x_k^2 \sum x_k (y_k - a_0 - a_1x_k - a_2x_k^2) \quad (17)$$

Expanding the equations:

$$a_0n + a_1 \sum x_k + a_2 \sum x_k^2 = \sum y_k \quad (18)$$

$$a_0 \sum x_k + a_1 \sum x_k^2 + a_2 \sum x_k^3 = \sum x_k y_k \quad (19)$$

$$a_0 \sum x_k^2 + a_1 \sum x_k^3 + a_2 \sum x_k^4 = \sum x_k^2 y_k \quad (20)$$

$$\underbrace{\begin{matrix} & 0 & 1 & 2 \\ 0 & \left[ \begin{matrix} n \\ \sum x_k \\ \sum x_k^2 \end{matrix} \right] \\ 1 & \left[ \begin{matrix} \sum x_k \\ \sum x_k^2 \\ \sum x_k^3 \end{matrix} \right] \\ 2 & \left[ \begin{matrix} \sum x_k^2 \\ \sum x_k^3 \\ \sum x_k^4 \end{matrix} \right] \end{matrix}}_A \underbrace{\begin{Bmatrix} a_0 \\ a_1 \\ a_2 \end{Bmatrix}}_x = \underbrace{\begin{Bmatrix} \sum y_k \\ \sum x_k y_k \\ \sum x_k^2 y_k \end{Bmatrix}}_b \quad (21)$$

**Matching the pattern.** There is a very clear pattern for the elements other than  $A_{0,0}$  and  $b_0$ . However, replacing element  $A_{0,0}$  with  $\sum x_k^0$ , and element  $b_0$  with  $\sum x_k^0 y_k$ , creates a much clearer pattern to follow. The system of equations becomes:

$$\underbrace{\begin{matrix} & 0 & 1 & 2 \\ 0 & \left[ \begin{matrix} \sum x_k^0 \\ \sum x_k^1 \\ \sum x_k^2 \end{matrix} \right] \\ 1 & \left[ \begin{matrix} \sum x_k^1 \\ \sum x_k^2 \\ \sum x_k^3 \end{matrix} \right] \\ 2 & \left[ \begin{matrix} \sum x_k^2 \\ \sum x_k^3 \\ \sum x_k^4 \end{matrix} \right] \end{matrix}}_A \underbrace{\begin{Bmatrix} a_0 \\ a_1 \\ a_2 \end{Bmatrix}}_x = \underbrace{\begin{Bmatrix} \sum x_k^0 y_k \\ \sum x_k^1 y_k \\ \sum x_k^2 y_k \end{Bmatrix}}_b \quad (22)$$

Building the equations in a matrix equation can be summarized with the following two equations:

$$A_{i,j} = \sum_{k=0}^{k < n} x_k^{i+j} \quad (0 \leq i < n) \quad (0 \leq j < n) \quad (23)$$

where  $i$  is the row index, and  $j$  is the column index.

$$b_i = \sum_{k=0}^{k < n} x_k^i y_k \quad (0 \leq i < n) \quad (24)$$

These equations can now be programmed, and solved using linear algebra techniques (i.e. gauss elimination).

## References

- [1] Chapra, Steven C. and Canale, Raymond P., *Numerical Methods for Engineers*, 7th ed. McGraw Hill Education, 2015.
- [2] A. K. Kaw, E. E. Kalu, and D. Nguyen. Numerical methods with applications. [Online]. Available: [http://nm.mathforcollege.com/topics/textbook\\_index.html](http://nm.mathforcollege.com/topics/textbook_index.html)