

LinkedIn: [linkedin.com/in/stanley-sayianka-8a6450170](https://www.linkedin.com/in/stanley-sayianka-8a6450170)

GitHub: [github.com/stanleyrazor](https://github.com/stanleyrazor)

# **ADDING NOISE TO LINEAR REGRESSION PREDICTIONS USING THE NEAREST NEIGHBOUR ALGORITHM.**

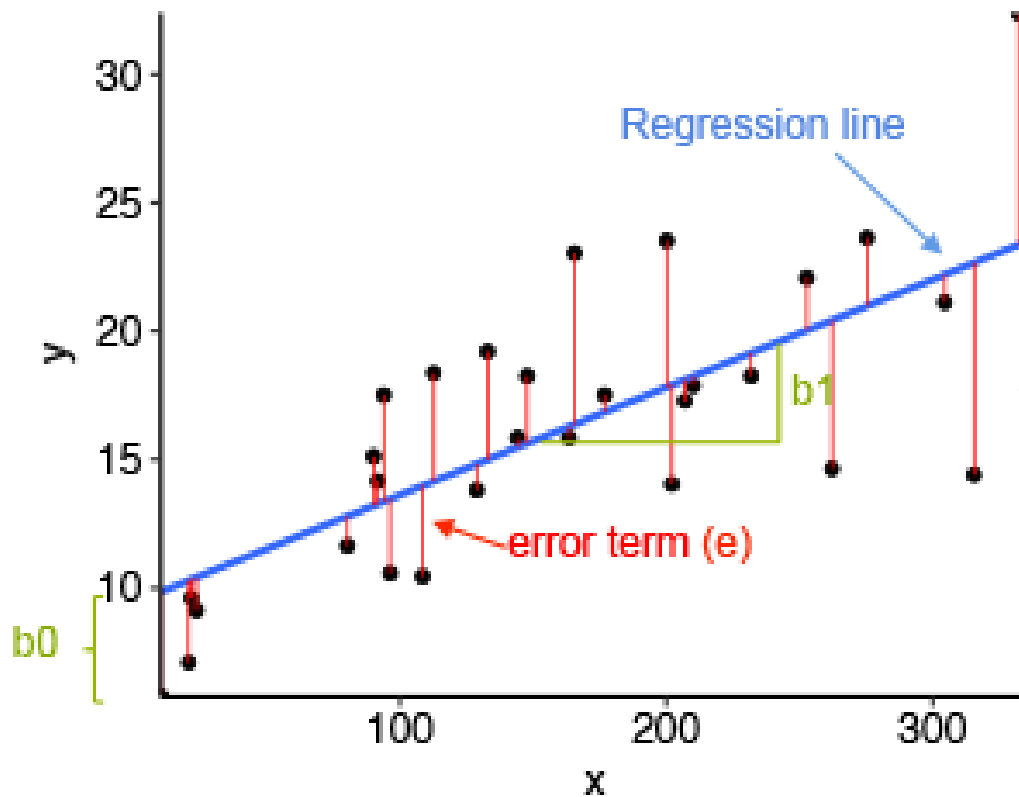
Author: Stanley Sayianka

Date: 10/02/2021

LinkedIn: [linkedin.com/in/stanley-sayianka-8a6450170](https://www.linkedin.com/in/stanley-sayianka-8a6450170)

GitHub: [github.com/stanleyrazor](https://github.com/stanleyrazor)

I assume that the reader has working knowledge on the basic simple linear regression model, deriving the coefficients and interpreting them, the reader should also have some knowledge on Expectations and variances and their properties as they relate to random variables.



LinkedIn: [linkedin.com/in/stanley-sayianka-8a6450170](https://www.linkedin.com/in/stanley-sayianka-8a6450170)

GitHub: [github.com/stanleyrazor](https://github.com/stanleyrazor)

The simple linear regression model is given as:

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \text{ where } i = 1, \dots, n$$

Where

$y_i$  – is the response variable

$\beta_0$  – is the intercept

$\beta_1$  – is the slope of the regression line

$x_i$  – is an independent variable

$\varepsilon_i$  – is the random error term

We also know that:

$$E(\varepsilon_i) = 0$$

$$\text{Var}(\varepsilon_i) = \sigma^2$$

*The error terms are uncorrelated, and independent, and follow a normal distribution with mean 0 and variance  $\sigma^2$*

Based on my idea of the noise-added linear regression model, (I will stick to the x-y notation, where x is the independent variable and y is the dependent variable)

For a new test input x, we use the linear regression model to predict its y value, and then for every y value predicted, I will add an error term which is obtained by averaging the error terms of the nearest neighbor of our new x instance. Remember the nearest neighbors of the test instance are all obtained from the training set that was used to fit the linear regression model.

Therefore:

*Let  $\hat{y}_i$  be the predicted value of the new test instance  $x_i$ ,*

*and let  $\theta_i$  be the error component that we desire to add to  $\hat{y}_i$*

Assuming we obtain K nearest neighbors to the new test instance and we obtain their residuals from the linear regression model as shown below:

$$\varepsilon_1, \varepsilon_2, \varepsilon_3, \varepsilon_4, \dots, \varepsilon_k$$

Then:

$$\theta_i = \frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k}{k}$$

Thus for the new test instance, its predicted value will be given by:

LinkedIn: [linkedin.com/in/stanley-sayianka-8a6450170](https://www.linkedin.com/in/stanley-sayianka-8a6450170)

GitHub: [github.com/stanleyrazor](https://github.com/stanleyrazor)

$$\hat{y}_i = \beta_0 + \beta_1 x_i + \frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k}{k}$$

The desired error term has the following properties:

$$\begin{aligned} E(\theta_i) &= E\left(\frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k}{k}\right) \\ &= \frac{1}{k} E(\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k) \\ &= \frac{1}{k} \{E(\varepsilon_1) + E(\varepsilon_2) + E(\varepsilon_3) + \dots + E(\varepsilon_k)\} \end{aligned}$$

*But  $E(\varepsilon_i) = 0$ , therefore*

$$E(\theta_i) = \frac{1}{k} * 0 = 0$$

$$\begin{aligned} Var(\theta_i) &= Var\left(\frac{\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k}{k}\right) \\ &= \frac{1}{k^2} * Var(\varepsilon_1 + \varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_k) \end{aligned}$$

$$= \frac{1}{k^2} \{Var(\varepsilon_1) + Var(\varepsilon_2) + Var(\varepsilon_3) + \dots + Var(\varepsilon_k) + 2 * Covariance(\varepsilon_1, \varepsilon_2, \varepsilon_3, \dots, \varepsilon_k)\}$$

*Since the error terms are uncorrelated and independent, their covariance vanishes*

$$\text{But } Var(\varepsilon_i) = \sigma^2$$

$$= \frac{1}{k^2} * k\sigma^2$$

$$\text{Thus: } Var(\theta_i) = \frac{\sigma^2}{k}$$

We see that then the desired noise follows a Normal distribution with mean 0 and variance equal to  $\frac{\sigma^2}{k}$

$$\theta_i \sim N(\mu = 0, \sigma^2 = \frac{\sigma^2}{k})$$