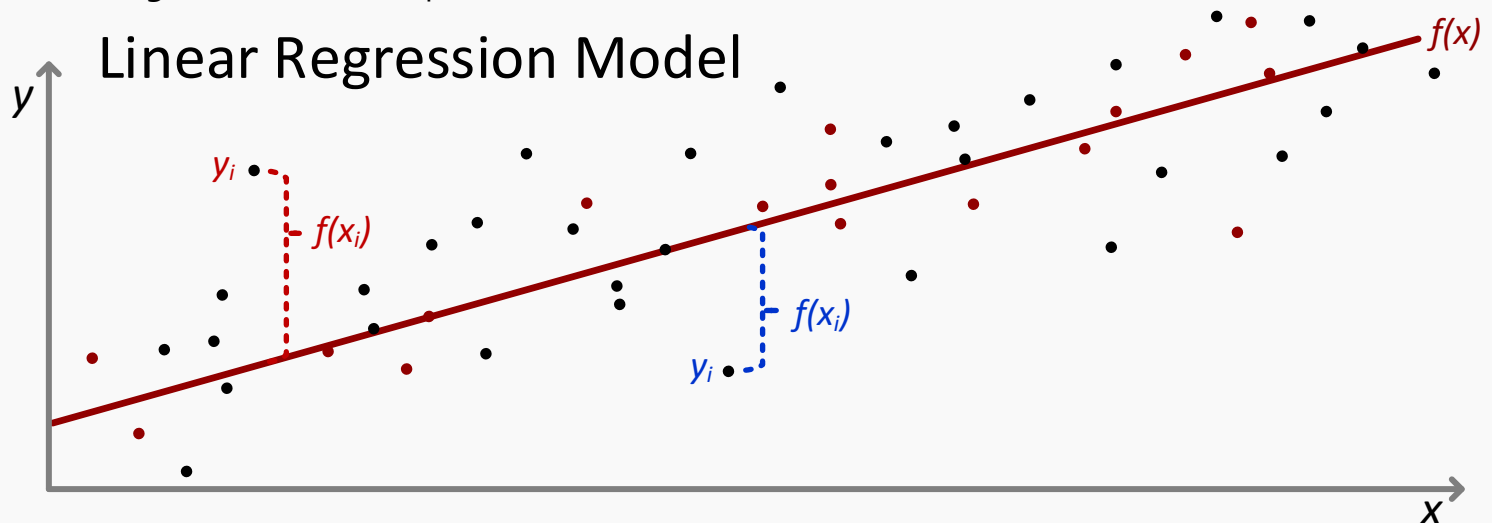


module2· regression

linear regression ↙ for machine learning

Linear Regression is used to predict real-valued outcomes (continuous variables)



The simplest form of **Linear Regression** is an equation that consists of a single linear coefficient and a constant fit the model data to a line; a function is applied in order to estimate y for each observation of $\rightarrow f(x_i) = b_0 + b_1x_i$. In order to fit the best model possible to the data, the model needs to be adjusted in response to the error that can be measured.

Error in the above model can be measured as the distance from a point to the linear model prediction. However, there are two orientations of error as seen above: $y_i - f(x_i)$ and $f(x_i) - y_i$. This application is a fallacy considering that either way will not account for the error of the other. An alternative would be to use the absolute error $|f(x_i) - y_i|$, however the latter remains inadequate. The solution to measure the model's error is the **Sum of Least Squares** error $(y_i - f(x_i))^2$.

With the **SSE** applied to a linear model, the optimization objective of minimizing error is achieved.

Single Variable Linear Regression:

$$f(x_i) = b_0 + b_1x_i$$

Choose b_0 and b_1 to minimize the total error on the training set:

$$\text{SSE}(f) = \sum_{i=1}^n (y_i - f(x_i))^2 = \text{SSE}(f) = \sum_{i=1}^n (y_i - (b_0 + b_1x_i))^2$$

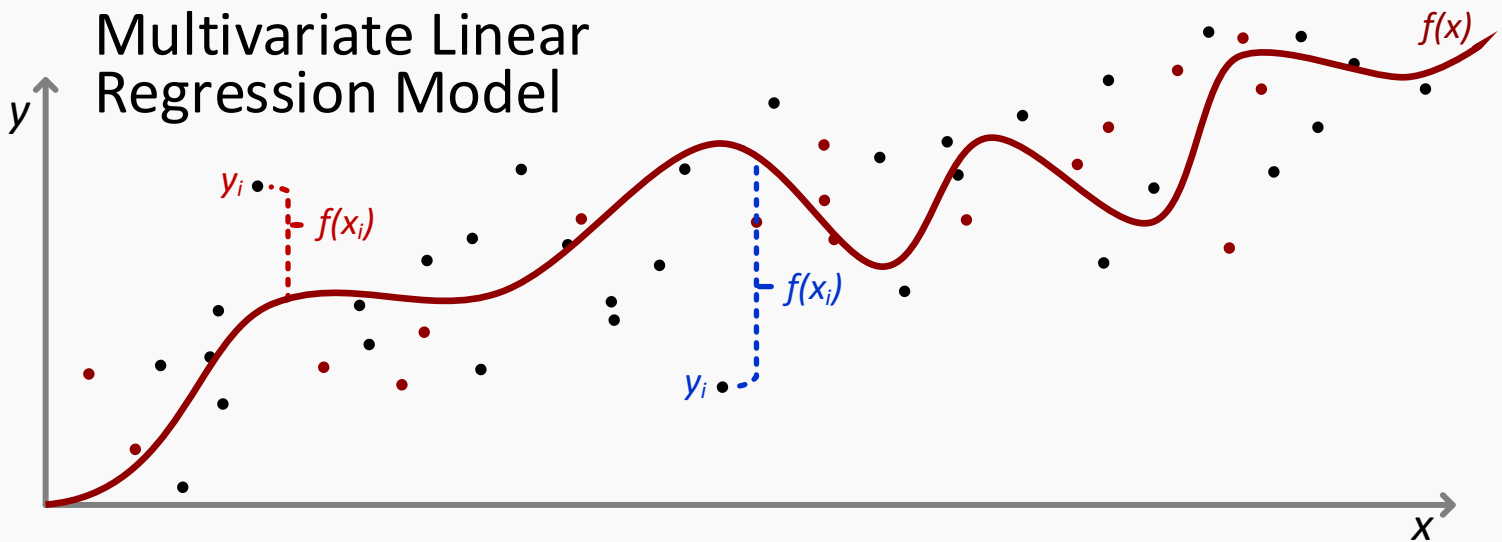
multivariate linear regression for machine learning

Multivariate Linear Regression expands upon simple linear regression when multiple weighted coefficients are added in perpetuity to the algorithm.

$$f(x_i) = b_0 + b_1x_i \rightarrow \mathbf{f(x_i) = b_0 + b_1x_{i,1} + b_2x_{i,2} + \cdots + b_px_{i,p}}$$

Method of Least Squares choosing coefficients to minimize:

$$\text{SSE}(f) = \sum_{i=1}^n (y_i - f(x_i))^2 = \text{SSE}(f) = \sum_{i=1}^n \left(y_i - (b_0 + b_1x_{i,1} + b_2x_{i,2} + \cdots + b_px_{i,p}) \right)^2$$



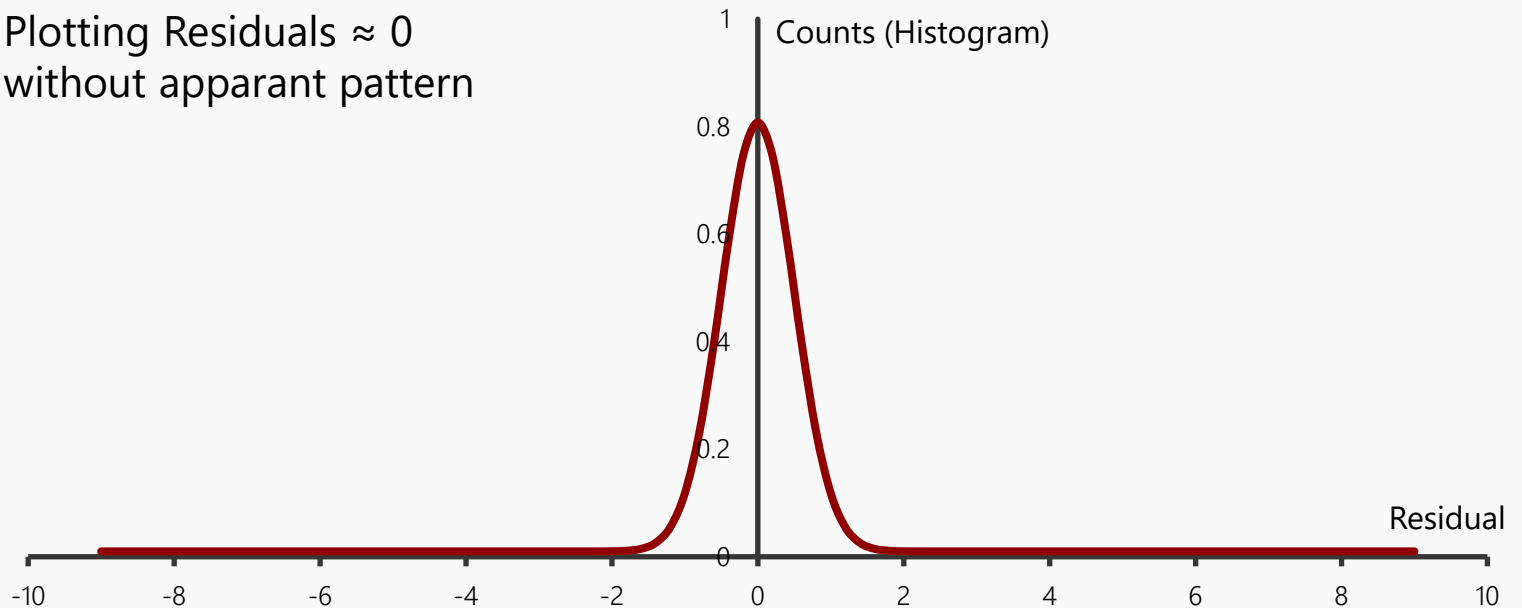
Multivariate Linear Regressive models can fit data in a much more dynamic manner as seen above. This is achieved through the use of polynomial, quadratic, trigonometric, and other various functions.

evaluation ↗ regression models

The **Residual** or **Error** is the difference between the predicted and actual values $f(x_i) - y_i$.

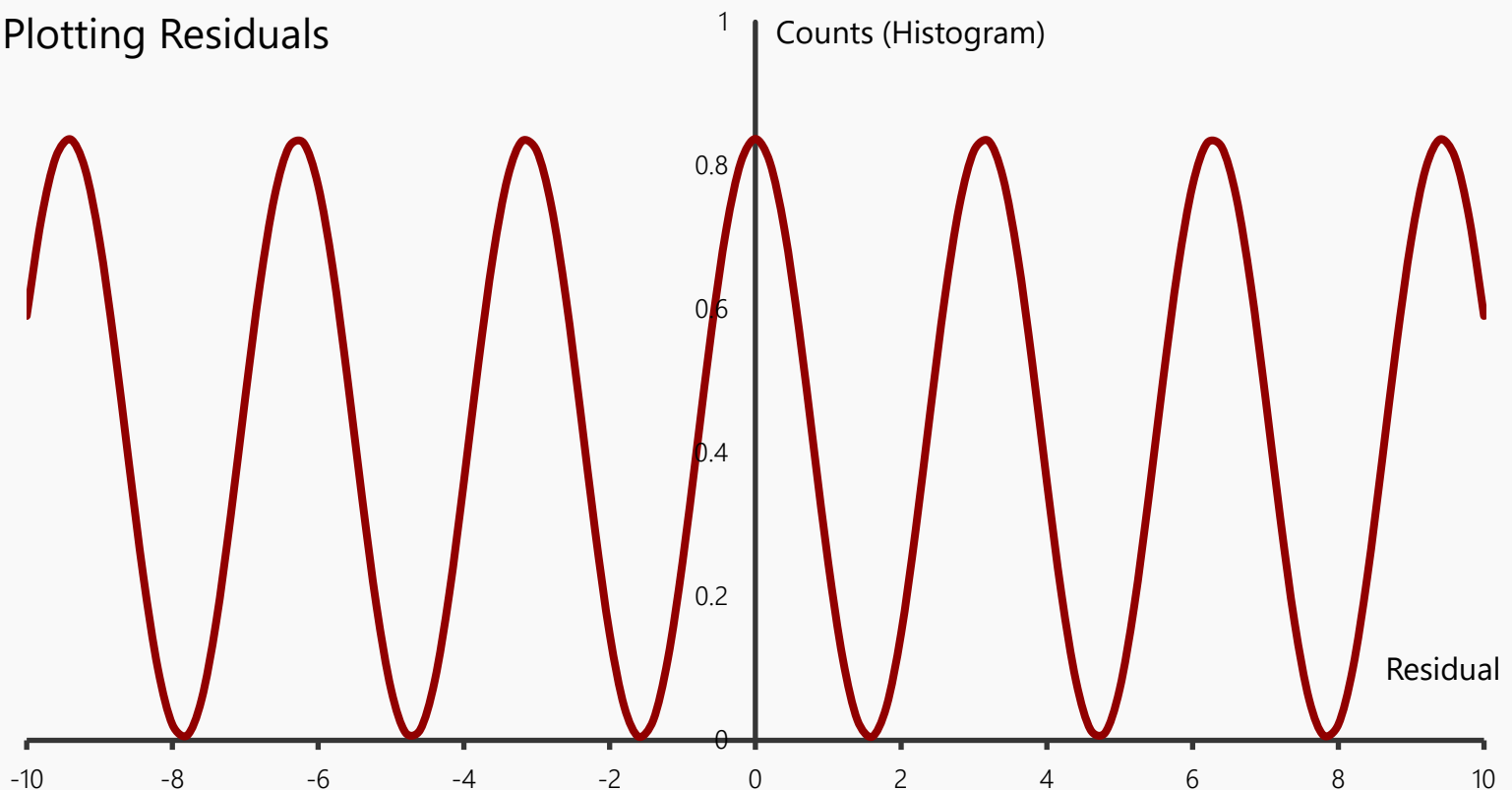
The reported Residuals are ideally ≈ 0 and will lack any specific structure or pattern amongst them:

Plotting Residuals ≈ 0
without apparant pattern



Plotted Residuals displaying an apparent pattern indicates missed properties in modeling the data:

Plotting Residuals



The above model could result from multiple linear features modeled inappropriately causing signal.

Plotted examples displaying a skew in either direction indicates a generally poorly performing model:

Plotting Residuals

