

Linear Regression

Machine Learning Context

- ▶ Given a set of data with associated measurements
- ▶ Predict the results of future measurements given a set of known results

Data could be a collection of images, measurements say “this is a duck”.

Data could be numerical (such as time intervals) and measurements could be numerical (such as speed of an object or a stock price).

Simplest case is finding a linear relationship.

Basic Problem

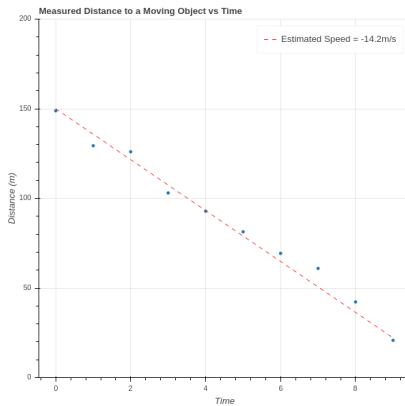


Figure 1: Physics Experiment

Engine size and MPG

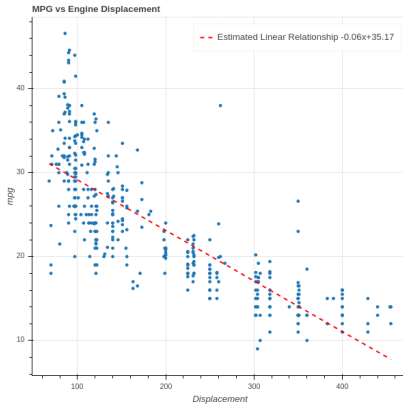


Figure 2: MPG vs Displacement

Mean Squared Error

Data consists of pairs $\{(x_i, y_i)\}$.

$$MSE(m, b) = \frac{1}{N} \sum_{i=1}^n (y_i - mx_i - b)^2$$

Minimize MSE

Write E instead of MSE for simplicity.

$$\frac{\partial E}{\partial m} = \frac{1}{N} \sum_1^N -2x_i(y_i - mx_i - b)$$

$$\frac{\partial E}{\partial b} = \frac{1}{N} \sum_1^N -2(y_i - mx_i - b)$$

Compute the derivatives

$$\frac{1}{N} \left(\sum_{i=1}^N x_i^2 \right) m + \frac{1}{N} \left(\sum_{i=1}^N x_i \right) b = \frac{1}{N} \sum_{i=1}^N x_i y_i$$
$$\frac{1}{N} \left(\sum_{i=1}^N x_i \right) m + b = \frac{1}{N} \sum_{i=1}^N y_i$$

- ▶ $\bar{x} = \frac{1}{N} \sum x_i$
- ▶ $\bar{y} = \frac{1}{N} \sum y_i$
- ▶ S_{xx} , S_{xy} , and S_{yy} are $\frac{1}{N} \sum x_i^2$, $\frac{1}{N} \sum x_i y_i$, $\frac{1}{N} \sum y_i^2$ respectively.

Solve to find the minima

$$S_{xx}m + \bar{x}b = S_{xy}$$

$$\bar{x}m + b = \bar{y}$$

Solution

$$m = \frac{S_{xy} - \bar{x}\bar{y}}{S_{xx} - \bar{x}^2}$$

$$b = \frac{S_{xx}\bar{y} - S_{xy}\bar{x}}{S_{xx} - \bar{x}^2}$$