# Linear Regression

With linear regression, you are trying to reduce a set of data into a line. The equation for a line is

$$y = a + bx$$

Intercept    Slope

Slope

Sample Standard Deviation of y

$$b = r * \frac{s_y}{s_x}$$

Pearson's Correlation

Sample Standard Deviation of x

The first thing to calculate is the slope of the regression line. That value is the correlation of the two data sets, multiplied by the ratio of their standard deviations. Note that those are sample standard deviations, not population

Correlation shows how much two sets of data change together. Correlation is always between -1 and 1, and is unit-less. Correlation is frequently around the average x, average y, but if you want to force the line through a specific point, you can use that point instead of average x, y in all the

Sum Over All Data Points

x & y values of each point minus x & y mean values

$$r = \frac{\sum ( (x - \bar{x}) * (y - \bar{y}) )}{(n - 1) * s_x * s_y}$$

Pearson's Correlation

\# of Data Points

Standard Deviation of x & y

$$s_x = \sqrt{\frac{\sum (x - \bar{x})^2}{(n - 1)}} \qquad s_y = \sqrt{\frac{\sum (y - \bar{y})^2}{(n - 1)}}$$

The sample standard deviation of x and y measure how spread out the x and y values are around their mean. They have units

Once you have done those calculations you have the slope. With the slope and a point the line goes through you can calculate the intercept. If you used the average x, y before use it here again. Otherwise use the same x, y that you specified before

Intercept      Slope

$$a = \bar{y} - b\bar{x}$$

Y & X Average

Sum Squared Regression Error

$$R^2 = 1 - \frac{SS_{Regression}}{SS_{Total}}$$

Sum Squared Total Error

When you have done the regression, one way of evaluating it's quality is R-squared. R-squared is a measure of the summed squared error in the regression vs the error if you didn't do a regression. Summed Squared error is shown below

Sum Over All Data Points

Square The Result

$$SS_{Total} = \sum (y_i - \bar{y})^2$$

Sum Squared Total Error    Each Data Point    Mean Value

Sum Over All Data Points

Square The Result

$$SS_{Regression} = \sum (y_i - y_{Regression})^2$$

Sum Squared Regression Error    Each Data Point    Regression Value

**Visit www.FairlyNerdy.com for more FREE Engineering Cheat Sheets**
**Print This, Save It, Or Share With A Friend!**