

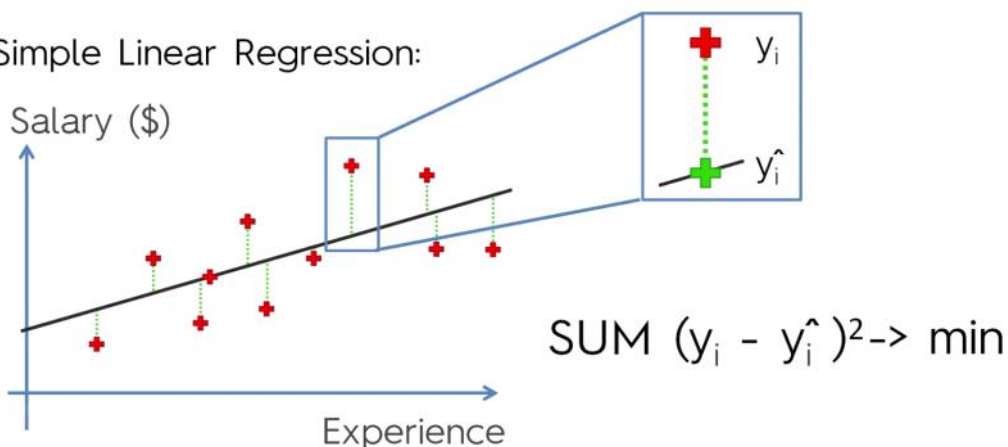
R Squared and Adjusted R Squared Intuition

Saturday, January 21, 2017 9:10 AM

R-Squared

Consider a Simple Linear Regression Model

Simple Linear Regression:



Here, the line represents the Simple Linear Regression Line.

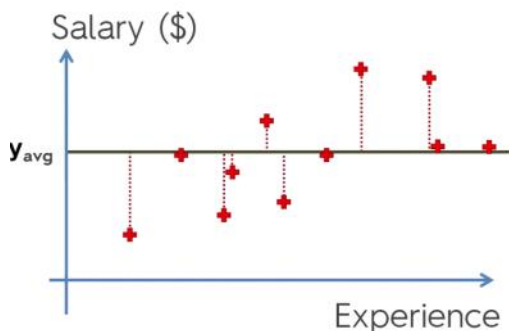
The root mean squared error is given by the formulae stated above and we try to minimize it when predicting parameters.

This is called the Sum of Squares of Residuals. This represents $SS(\text{res})$.

Now, instead let us draw a line which pertains to the average y value. This is a line which will be parallel to the X -axis.

Now, the sum of squares of this value difference is called the total sum of squares.

Simple Linear Regression:



$$SS_{\text{res}} = \text{SUM } (y_i - \hat{y}_i)^2$$

$$SS_{\text{tot}} = \text{SUM } (y_i - y_{\text{avg}})^2$$

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

Based on this, what R^2 is telling us is that how good our regression line is compared to the average line.

Because, we are trying to minimize $SS(\text{res})$ and $SS(\text{tot})$ is a constant for a given set of values.

So ideally, if our $SS(\text{res})$ is 0, then R^2 value = 1 (Best possible line)

R^2 can be negative if $SS(\text{res})$ fits the data worse than the average line.

Hence, R^2 is an assessment of how good our prediction is.

Adjusted R-Squared

According to R-squared, the best fitting model for prediction is the one which minimizes the Sum of Squares of Residuals.

R^2 - goodness of fit - bigger the better

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

R^2 – Goodness of fit
(greater is better)

$$y = b_0 + b_1 * x_1$$

$$y = b_0 + b_1 * x_1 + b_2 * x_2$$

Problem:

$$+ b_3 * x_3$$

$SS_{res} \rightarrow \text{Min}$

R^2 will never decrease

Consider a multi variable regression. Here, suppose we add another variable to the equation x_3 . Now, the regression function is to decrease $SS(res)$. Hence, the model will give a value to b_3 which actually decreases the $SS(res)$ of the model thereby increasing R^2 . Or, the model will make $b_3 = 0$, making sure that the $SS(res)$ value will not increase, thereby making the value of R^2 stay constant.

Note - Most often, there is some sort of basic correlation even between two random entities.

The problem with this is by adding new variables, we will not know if the variables are actually impacting the model or not and to what extent.

For this reason, we use adjusted R -squared.

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$$

$$\text{Adj } R^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1}$$

p - number of regressors

n - sample size

R^2 will never decrease !

The basic conception over here is that, as we add more regressors (variables) to the model, the value of R^2 generally tends to increase thereby leading to increase in the value of Adjusted R -squared.

But on the other hand, as we increase the number of regressors, the value of p increases thereby leading to lowering the value of Adjusted R -squared. Hence, there is a battle on our hands between adding the variables or not adding the variables. Because an insignificant increase in R^2 is penalized by a relatively significant increase in p .

Thus, adjusted R^2 helps us identify whether we are adding good variables to our model or not. It acts as a tradeoff metric of sorts which strives to capture the essence of using minimum number of features to make the best possible prediction. (All features used by the model have a good impact on the prediction the model makes)