VARIANCE OF RESIDUALS IN SIMPLE LINEAR REGRESSION

ALLEN BACK

Suppose we use the usual denominator (n-1) in defining the sample variance and sample covariance for samples of size n:

$$Var(x) = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})^2$$
$$Cov(x, y) = \frac{1}{n-1} \sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})$$

Of course the correlation coefficient r is related to this covariance by

$$r = \frac{1}{s_x s_y} \left(\operatorname{Cov}(x, y) \right).$$

Then since $(a + b)^2 = a^2 + 2ab + b^2$, it follows that Var(a + b) = Var(a) + Var(b) + 2Cab

$$Var(a + b) = Var(a) + Var(b) + 2Cov(a, b).$$

If we apply this to the usual simple linear regression setup, we obtain:

Proposition: The sample variance of the residuals d_i in a simple linear regression satisfies

$$\operatorname{Var}(d_i) = (1 - r^2)\operatorname{Var}(y_i)$$

where $Var(y_i)$ is the sample variance of the original response variable.

Proof: The line of regression may be written as

$$\hat{y} - \bar{y} = b_1(x - \bar{x})$$

 $y - y = o_1(x - x)$ where $b_1 = \frac{rs_y}{s_x}$. The residual of a point (x_i, y_i) is $d_i = y_i - \hat{y}_i$, so:

$$Var(d_i) = Var(y_i - \hat{y}_i)$$

= $Var(y_i - (\bar{y} + b_1(x_i - \bar{x})))$
= $Var((y_i - \bar{y}) - b_1(x_i - \bar{x}))$
= $s_y^2 + b_1^2 s_x^2 - 2b_1 Cov(y_i - \bar{y}, x_i - \bar{x})$
= $s_y^2 + r^2 \frac{s_y^2}{s_x^2} s_x^2 - 2r \frac{s_y}{s_x} (rs_x s_y)$
= $s_y^2 - r^2 s_y^2$
= $(1 - r^2) s_y^2$.