

Prediction of a new observation : note that this is a single observation, not the regression line.

First, the variance of a generic linear combination (from Chapter 1:1.27a & b)

$$T = aW + bX + cZ$$

$$E(T) = aE(W) + bE(X) + cE(Z)$$

$$\text{Var}(T) = a^2\text{Var}(W) + b^2\text{Var}(X) + c^2\text{Var}(Z) + 2(\text{Covariances})$$

$$\begin{aligned} \text{Var}(T) = & a^2\text{Var}(W) + b^2\text{Var}(X) + c^2\text{Var}(Z) \\ & ab\text{Cov}(W,X) + bc\text{Cov}(X,Z) + ac\text{Cov}(X,Z) \end{aligned}$$

If we are able to assume that the three terms are stochastically independent, then the covariances are equal to zero.

We have already seen a series of Linear Combinations

1) First we saw,

$$b_1 = \frac{\sum(X_i - \bar{X})Y_i}{\sum(X_i - \bar{X})^2} = \sum_{i=1}^n k_i Y_i$$

$$\text{so, } \text{Var}(b_1) = k_1^2\text{Var}(Y_1) + k_2^2\text{Var}(Y_2) + k_3^2\text{Var}(Y_3) + \dots$$

since all  $Y_i$  at all values of  $X_i$  have the same variance (homogeneous), then

$$= \sum_{i=1}^n k_i^2 \text{Var}(Y_i)$$

$$\text{and recall that } \sum_{i=1}^n k_i^2 = \frac{1}{\sum(X_i - \bar{X})^2}$$

and that  $\text{Var}(Y_i)$  is estimated by the MSE, then

$$\text{Var}(b_1) = \frac{\text{MSE}}{\sum(X_i - \bar{X})^2}$$

1) Show that  $b_1$  is a linear combination of  $k_i$   $\left( = \frac{(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2} \right)$

$$a) \quad b_1 = \frac{\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{\Sigma(X_i - \bar{X})^2}$$

where

$$\begin{aligned} b) \quad \Sigma(X_i - \bar{X})(Y_i - \bar{Y}) &= \Sigma(X_i Y_i - \bar{X} Y_i - X_i \bar{Y} + \bar{X} \bar{Y}) \\ &= \Sigma(X_i - \bar{X}) Y_i - \Sigma(X_i - \bar{X}) \bar{Y} \\ &= \Sigma(X_i - \bar{X}) Y_i - \bar{Y} \Sigma(X_i - \bar{X}) \end{aligned}$$

and since

$$\Sigma(X_i - \bar{X}) = 0$$

then

$$c) \quad \Sigma(X_i - \bar{X})(Y_i - \bar{Y}) = \Sigma(X_i - \bar{X}) Y_i$$

as a result,

$$d) \quad b_1 = \frac{\Sigma(X_i - \bar{X})(Y_i - \bar{Y})}{\Sigma(X_i - \bar{X})^2} = \frac{\Sigma(X_i - \bar{X}) Y_i}{\Sigma(X_i - \bar{X})^2} = \frac{\Sigma(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2} Y_i = \sum_{i=1}^n k_i Y_i$$

where

$$e) \quad \Sigma k_i = \sum_{i=1}^n \frac{(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2} = \frac{\Sigma(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2}$$

note that  $\Sigma k_i = 0$  since  $\Sigma(X_i - \bar{X}) = 0$

now prove that  $\sum_{i=1}^n k_i^2 = \frac{1}{\Sigma(X_i - \bar{X})^2}$

$$\begin{aligned} f) \quad \Sigma k_i^2 &= \sum_{i=1}^n \left[ \frac{(X_i - \bar{X})}{\Sigma(X_i - \bar{X})^2} \right]^2 = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{[\Sigma(X_i - \bar{X})^2]^2} \\ &= \frac{1}{[\Sigma(X_i - \bar{X})^2]^2} \Sigma(X_i - \bar{X})^2 = \frac{1}{\Sigma(X_i - \bar{X})^2} \end{aligned}$$

2) Then we say that  $\hat{Y}_i = b_0 + b_1 X_i$ , also a linear combination.

$$\text{Var}(\hat{Y}_i) = 1 * \text{Var}(b_0) + X_i * \text{Var}(b_1) + 2 * 1 * X_i * \text{Cov}(b_0, b_1)$$

**note that we do NOT assume that  $b_0$  and  $b_1$  are independent.**

**The covariance is included, not equal to zero.**

Using previous definitions of  $\text{Var}(b_0)$  and  $\text{Var}(b_1)$ , and the Gaussian multipliers from the  $(X'X)^{-1}$  matrix for the covariance

$$\text{Var}(\hat{Y}_i) = \sigma^2 * 1^2 \left[ \frac{1}{n} + \frac{X^2}{\sum(X_i - \bar{X})^2} \right] + \sigma^2 * X_i^2 \left[ \frac{1}{\sum(X_i - \bar{X})^2} \right] + 2 * 1 * X_i \sigma^2 \left[ \frac{-\sum X_i}{n \sum(X_i - \bar{X})^2} \right]$$

$$\text{Var}(\hat{Y}_i) = \sigma^2 \left[ \frac{1}{n} + \frac{X^2}{\sum(X_i - \bar{X})^2} + \frac{X_i^2}{\sum(X_i - \bar{X})^2} - \frac{2X_i \bar{X}}{\sum(X_i - \bar{X})^2} \right]$$

$$\text{Var}(\hat{Y}_i) = \sigma^2 \left[ \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right]$$

3) Now we want a confidence interval for a single (new) observation.

The equation for that observation is

$$Y_i = b_0 + b_1 X_i + \epsilon_i$$

or

$$Y_i = \hat{Y}_i + \epsilon_i$$

We assumed independence once before (each  $Y_i$  independent of others). We are now going to assume independence again. We assume that the residuals are independent of the model (ie. assume that  $\epsilon_i$  are independent of  $\hat{Y}_i$ ).

So the variance of single observations will be

$$\text{Var}(Y_i) = \text{Var}(\hat{Y}_i) + \text{Var}(\epsilon_i) + \left[ 2 * \text{Cov}(\hat{Y}_i, \epsilon_i) = 0 \right]$$

We know from previous work that

$$\text{Var}(\hat{Y}_i) = \sigma^2 \left[ \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right]$$

$$\text{Var}(\epsilon_i) = \sigma^2$$

therefore

$$\text{Var}(Y_i) = \sigma^2 \left[ \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right] + \sigma^2$$

or

$$\text{Var}(Y_i) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right]$$

where the estimator of  $\sigma^2$  is MSE

Note that both your textbook and I have been using  $\sigma^2$  for both  $\text{Var}(Y_i)$  and for  $\text{Var}(\epsilon_i)$ . Each is “the variance”, but they are variances of different things. A better notation perhaps is  $\text{Var}(\epsilon_i) = \sigma_\epsilon^2$

There is another confidence interval of potential interest between

$$\text{Var}(\hat{Y}_i) = \sigma^2 \left[ \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right], \text{ the regression line}$$

and

$$\text{Var}(Y_i) = \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right], \text{ a new observation}$$

This is the confidence interval for the mean of a new sample taken at some particular value of  $X_i$ , where  $m$  is the size of the new sample. This cannot be as narrow as the confidence interval for the regression, but should be narrower than the confidence interval for a single sample. This CI is given by,

$$\text{Var}(Y_i) = \sigma^2 \left[ \frac{1}{m} + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum(X_i - \bar{X})^2} \right], \bar{X} \text{ for a new sample}$$

Example : From *vial breakage regressed on number of airline transfers* example  
Place a confidence interval on the breakage for 3 transfers for a single new observation.

$$s_{\hat{Y}_i}^2 = \text{MSE} \left[ 1 + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right]$$

$$= 2.2 \left[ 1 + \frac{1}{10} + \frac{(3-1)^2}{20 - \frac{10^2}{10}} \right] = 2.2 \left[ 1 + \frac{1}{10} + \frac{4}{10} \right] = 2.2 * 1.5 = 3.3$$

we previously calculated the variance of the regression line at  $s_{\hat{Y}_i}^2 = 1.1$ . Note that the variance of a single point is  $s_{\hat{Y}_i}^2 + s^2 = 1.1 + 2.2 = 3.3$

$$s_{Y_i} = \sqrt{3.3} = 1.816$$

since  $t_{\frac{\alpha}{2}, 8 \text{ df}} = 2.306$ , then

$$P(\hat{Y}_{X=3} - t_{1-\frac{\alpha}{2}, n-2} s_{Y_i} \leq E(\hat{Y}) \leq \hat{Y}_{X=3} + t_{1-\frac{\alpha}{2}, n-2} s_{Y_i}) = 1-\alpha$$

$$P(22.2 - 2.306 * 1.816 \leq E(\hat{Y}) \leq 22.2 + 2.306 * 1.816) = 1-\alpha$$

$$P(18.011 \leq E(\hat{Y}) \leq 26.389) = 0.95$$

SAS will calculate confidence intervals for either the regression line (option CLM) or for individual points (option CLI). But not for a new sample.

Check this against the SAS output

Suppose we were to ship 4 cases through 3 transfers. What is the confidence interval for the mean breakage of 4 cases?

$$s_{\bar{Y}_i}^2 = \text{MSE} \left( \frac{1}{m} + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2} \right) = 2.2 \left( \frac{1}{4} + \frac{1}{10} + \frac{(3-1)^2}{20 - \frac{10^2}{10}} \right)$$

$$= 2.2 \left( \frac{1}{4} + \frac{1}{10} + \frac{4}{10} \right) = 2.2 * 0.75 = 1.65$$

$$s_{\bar{Y}_i} = \sqrt{1.65} = 1.2845$$

since  $t_{\frac{\alpha}{2}, 8 \text{ df}} = 2.306$ , then

$$P(\bar{Y}_{X=3} - t_{1-\frac{\alpha}{2}, n-2} s_{\bar{Y}_i} \leq E(\bar{Y}) \leq \bar{Y}_{X=3} + t_{1-\frac{\alpha}{2}, n-2} s_{\bar{Y}_i}) = 1-\alpha$$

$$P(22.2 - 2.306 * 1.2845 \leq E(\bar{Y}) \leq 22.2 + 2.306 * 1.2845) = 1-\alpha$$

$$P(19.238 \leq E(\bar{Y}) \leq 25.162) = 0.95$$

**The MEAN of the 4 cases falls in this range.**

The CI for the regression line is narrower

The CI for individual points is higher