L2. Correction : $Var(e_i) = \sigma^2 [ 1 - \frac{1}{n} - \frac{(x_i - \bar{x})^2}{S_{xx}} ]$

- For a given data $(X_1, Y_1), \cdots, (X_n, Y_n)$

Fit a SLR :

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

where $\hat{\beta}_0$ and $\hat{\beta}_1$ are the OLSE based on $(X_1, Y_1), \cdots, (X_n, Y_n)$

The residual $e_i$ is the estimation bias of $\boxed{\text{in sample}}$ $y_i$

$$e_i = y_i - \hat{y}_i$$

- $Var(e_i) = Var(y_i - \hat{y}_i) = Var[ (\beta_0 + \beta_1 X_i + \varepsilon_i) - (\hat{\beta}_0 + \hat{\beta}_1 X_i) ]$

$= Var[ (\beta_0 + \beta_1 X_i) + \varepsilon_i - (\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 X_i ]$

$= Var(\varepsilon_i) + Var(\bar{y}) + (X_i - \bar{x})^2 Var(\hat{\beta}_1)$

$\qquad - 2 cov(\varepsilon_i, \bar{y})$

$\qquad - 2(X_i - \bar{x}) cov(\varepsilon_i, \hat{\beta}_1)$

$\qquad + 2(X_i - \bar{x}) cov(\bar{y}, \hat{\beta}_1)$

- $\text{Cov}(\bar{Y}, \hat{\beta_1}) = \text{Cov}\left(\sum \frac{Y_i}{n}, \sum k_j y_j\right) = \sum_i \sum_j \text{Cov}\left(\frac{y_i}{n}, k_j y_j\right)$

$$= \sum_i \sum_j \frac{k_j}{n} \text{Cov}(y_i, y_j)$$

Because $\text{Cov}(y_i, y_j) = 0$ for $i \neq j$

$\text{Cov}(y_j, y_j) = \text{Var}(y_i) = \sigma^2$

$$= \sum_j \frac{k_j}{n} \sigma^2 = \frac{\sigma^2}{n} \underset{\underset{0}{\shortparallel}}{\sum k_j} = 0$$

- $\text{Cov}(\bar{Y}, \varepsilon_i) = \text{Cov}\left(\sum_{j=1}^{n} \frac{y_j}{n}, \varepsilon_i\right)$

$$= \sum_{j=1}^{n} \frac{1}{n} \text{Cov}(y_j, \varepsilon_i)$$

when $i \neq j$, $y_j = \beta_0 + \beta_1 X_j + \varepsilon_j$, $\text{Cov}(y_j, \varepsilon_i) = 0$

when $i = j$, $\text{Cov}(y_j, \varepsilon_i) = \text{Cov}(\beta_0 + \beta_1 x_i + \varepsilon_i, \varepsilon_i) = \sigma^2$

$$= \frac{1}{n} \text{Cov}(y_i, \varepsilon_i) = \frac{\sigma^2}{n}$$

- $\text{cov}(\varepsilon_i, \hat{\beta}_1)$     Similarly

$$= \text{cov}(\varepsilon_i, \sum_j k_j y_j)$$

$$= k_i \, \text{cov}(\varepsilon_i, y_i)$$

$$= k_i \sigma^2 = \frac{(X_i - \bar{X})}{\sum (X_j - \bar{X})^2} \sigma^2 = \frac{(X_i - \bar{X})}{S_{xx}} \sigma^2$$

so  $Var(e_i) = Var(\bar{Y}) + (X_i - \bar{X})^2 \, Var(\hat{\beta}_1) + Var(\varepsilon_i)$

$$\overline{\ } \, 2(X_i - \bar{X}) \, \text{cov}(\varepsilon_i, \hat{\beta}_1) - 2 \, \text{cov}(\bar{Y}, \varepsilon_i)$$

$$= \frac{\sigma^2}{n} + (X_i - \bar{X})^2 \cdot \frac{\sigma^2}{S_{xx}} + \sigma^2$$

$$- \frac{2(X_i - \bar{X})^2}{S_{xx}} \sigma^2 - 2 \frac{\sigma^2}{n}$$

$$= \sigma^2 - \frac{\sigma^2}{n} - \frac{(X_i - \bar{X})^2}{S_{xx}} \sigma^2$$

$$\boxed{= \sigma^2 \left[ 1 - \frac{1}{n} - \frac{(X_i - \bar{X})^2}{S_{xx}} \right]}$$

In L3. when giving a new $X_0$, and look at
the $\boxed{\text{out of sample}}$ prediction bias :

$$e_0 = y_0 - \hat{y}_0$$

$$Var(e_0) = Var(y_0 - \hat{y}_0)$$

$$= Var(\varepsilon_0) + Var(\bar{y}) + (X_0 - \bar{X}) Var(\hat{\beta}_1)$$

$$- 2 cov(\varepsilon_0, \bar{y})$$

$$- 2 (X_0 - \bar{X}) cov(\varepsilon_0, \hat{\beta}_1)$$

$$+ 2 (X_i - \bar{X}) cov(\bar{y}, \hat{\beta}_1)$$

Since $\varepsilon_0$ is out of the sample, so $\varepsilon_0$ is uncorrelated
with $\varepsilon_i$ for all $i = 1, \cdots, n$, therefore

- $cov(\varepsilon_0, \bar{y}) = 0$

- $cov(\varepsilon_0, \hat{\beta}_1) = 0$

- $cov(\bar{y}, \hat{\beta}_1) = 0$ is proved the same
  as before.

$$Var(e_0) = Var(\bar{y}) + (x_0 - \bar{x})^2 \, Var(\hat{\beta_1}) + Var(\varepsilon_0)$$

$$= \frac{\sigma^2}{n} + (x_0 - \bar{x})^2 \frac{\sigma^2}{S_{xx}} + \sigma^2$$

$$= \sigma^2 \left[ 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}} \right]$$

- When the new $x_0$ is far away from $\bar{x}$, variance of the prediction bias is large. $\rightarrow$ unstable prediction.