

(2) We have

$$\hat{\beta}_1 = \frac{(X^n - \overline{X_n}1^n)^T(Y^n - \overline{Y_n}1^n)}{(X^n - \overline{X_n}1^n)^T(X^n - \overline{X_n}1^n)}$$

where 1^n is the vector of all 1's, $X^n = (X_1, \dots, X_n)^T$, and $Y^n = (Y_1, \dots, Y_n)^T$. Hence, $\mathbb{V}(\hat{\beta}_1)$ is equal to

$$\frac{1}{(X^n - \overline{X_n}1^n)^T(X^n - \overline{X_n}1^n)}(X^n - \overline{X_n}1^n)^T \Sigma (X^n - \overline{X_n}1^n)$$

where Σ is the covariance matrix of $Y^n - \overline{Y_n}1^n$. Let's compute Σ .

We have $\mathbb{V}(Y_i - \overline{Y_n}) = \mathbb{V}(Y_i) + \mathbb{V}(\overline{Y_n}) - 2 \text{Cov}(Y_i, \overline{Y_n})$. We have $\mathbb{V}(Y_i) = \mathbb{V}(\beta_0 + \beta_1 X_i + \epsilon_i) = \mathbb{V}(\epsilon_i) = \sigma^2$. Similarly, $\mathbb{V}(\overline{Y_n}) = \frac{\sigma^2}{n}$. Finally, $\text{Cov}(Y_i, Y_j) = \text{Cov}(\epsilon_i, \epsilon_j) = \sigma^2 \delta_{ij}$ where δ_{ij} is the Dirac delta: $\delta_{ij} = 1$ if $i = j$ and $\delta_{ij} = 0$ otherwise. By bilinearity of Cov,

$$\text{Cov}(Y_i, \overline{Y_n}) = \frac{1}{n} \sum_{j=1}^n \text{Cov}(Y_i, Y_j) = \frac{\sigma^2}{n}.$$

Hence

$$\mathbb{V}(Y_i - \overline{Y_n}) = \sigma^2 \left(1 + \frac{1}{n} - 2\frac{1}{n} \right) = \frac{(n-1)\sigma^2}{n}.$$

For $i \neq j$ we have

$$\text{Cov}(Y_i - \overline{Y_n}, Y_j - \overline{Y_n}) = \text{Cov}(Y_i, Y_j) - \text{Cov}(Y_i, \overline{Y_n}) - \text{Cov}(Y_j, \overline{Y_n}) + \mathbb{V}(\overline{Y_n}).$$

By our previous calculations this is equal to $-\sigma^2/n$.

Finally then,

$$\Sigma = \frac{\sigma^2}{n} \begin{pmatrix} n-1 & -1 & -1 & -1 \\ -1 & n-1 & -1 & -1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \dots & n-1 \end{pmatrix}$$

and this yields

$$\mathbb{V}(\hat{\beta}_1) = \frac{\sigma^2}{n} \cdot \frac{1}{(\sum (X_i - \overline{X_n})^2)^2} \cdot \left((n-1) \sum_{i=1}^n (X_i - \overline{X_n})^2 - \sum_{i \neq j} (X_i - \overline{X_n})(X_j - \overline{X_n}) \right).$$

We have that

$$-\sum_{i=1}^n (X_i - \overline{X_n})^2 - \sum_{i \neq j} (X_i - \overline{X_n})(X_j - \overline{X_n}) = -\left(\sum_{i=1}^n (X_i - \overline{X_n}) \right)^2 = 0.$$

So $\mathbb{V}(\hat{\beta}_1) = \frac{\sigma^2}{\sum (X_i - \overline{X_n})^2} = \frac{\sigma^2}{ns_X^2}$, as claimed.

From $\hat{\beta}_0 = \overline{Y_n} - \hat{\beta}_1 \overline{X_n}$, we see

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \text{Cov}(\overline{Y_n}, \hat{\beta}_1) - \overline{X_n} \text{Cov}(\hat{\beta}_1, \hat{\beta}_1).$$

By our previous calculations,

$$\text{Cov}(\bar{Y}_n, \hat{\beta}_1) = \frac{1}{ns_X^2} \sum (X_i - \bar{X}_n) \text{Cov}(\bar{Y}_n, Y_i - \bar{Y}_n) = \frac{1}{ns_X^2} \sum (X_i - \bar{X}_n) \left(\frac{\sigma^2}{n} - \frac{\sigma^2}{n} \right) = 0$$

and this leaves $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = -\bar{X}_n \mathbb{V}(\hat{\beta}_1)$, as claimed.

Finally,

$$\mathbb{V}(\hat{\beta}_0) = \mathbb{V}(\bar{Y}_n) - 2\bar{X}_n \text{Cov}(\bar{Y}_n, \hat{\beta}_1) + \bar{X}_n^2 \mathbb{V}(\hat{\beta}_1) = \frac{\sigma^2}{n} + \frac{\sigma^2 \bar{X}_n^2}{ns_X^2} = \frac{\sigma^2(s_X^2 + \bar{X}_n^2)}{ns_X^2}$$

and we calculate that $s_X^2 + \bar{X}_n^2 = \frac{1}{n} \sum X_i^2$. This gives the (1, 1) entry of the covariance matrix.

- (3) We assume $Y_i = \beta X_i + \epsilon_i$ where the ϵ_i are iid with variance σ^2 . Our model is $\hat{r}(X) = \hat{\beta}X$. For data points $(Y_1, X_1), \dots, (Y_n, X_n)$, the sum of squared residuals is $\sum \hat{\epsilon}_i^2 = \sum (Y_i - \hat{\beta}X_i)^2$. This function of $\hat{\beta}$ is convex with a single local minimum. We have

$$\frac{\partial}{\partial \hat{\beta}} \left(\sum \hat{\epsilon}_i^2 \right) = 2 \sum \left(-X_i Y_i + \hat{\beta} X_i^2 \right).$$

Setting this equal to zero and solving for $\hat{\beta}$ yields

$$\hat{\beta} = \frac{\sum X_i Y_i}{\sum X_i^2} = \frac{X^n \cdot Y^n}{X^n \cdot X^n}.$$

where $X^n = (X_1, \dots, X_n)^T$ and similarly for Y^n . Considering X^n to be constant, $\mathbb{V}(\hat{\beta}) = \frac{1}{(X^n \cdot X^n)^2} (X^n)^T \Sigma (X^n)$ where Σ is the covariance matrix of Y^n . As in the previous exercise, $\text{Cov}(Y_i, Y_j) = \sigma^2 \delta_{ij}$ where δ_{ij} is the Dirac delta. So $\Sigma = \sigma^2 I$ where I is the $n \times n$ identity matrix and

$$\mathbb{V}(\hat{\beta}) = \frac{\sigma^2 X^n \cdot X^n}{(X^n \cdot X^n)^2} = \frac{\sigma^2}{\sum X_i^2}.$$

As long as $\sum X_i^2 \rightarrow \infty$ we have $\mathbb{E}(\hat{\beta}) \rightarrow \beta$ and $\mathbb{V}(\hat{\beta}) \rightarrow 0$. Thus $\hat{\beta}$ converges to β in quadratic mean and therefore also in probability (see e.g. exercise 5.2).

- (6) See the Jupyter Notebook 6.ipynb.
 (7) See the Jupyter Notebook 7.ipynb.
 (8) In this case, up to addition of a constant not depending on $\hat{\beta}$, the AIC is

$$\text{AIC} = \ell_S - |S| = -n \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^n (Y_i - (X\hat{\beta})_i)^2 - |S|.$$

Mallow's C_p statistic is

$$\hat{R}_{\text{tr}}(S) + 2|S|\sigma^2 = \sum_{i=1}^n (Y_i - (X\hat{\beta})_i)^2 + 2|S|\sigma^2.$$

Thus, up to adding a constant to AIC, which doesn't affect where the maximum AIC is achieved, we have $C_p = -2\sigma^2 \text{AIC}$. So the $\hat{\beta}$ which maximizes AIC is the same as the $\hat{\beta}$ which minimizes C_p .

(11) See the Jupyter Notebook 11.ipynb.