

Here  $Y = \begin{bmatrix} Y_{11} & \dots & Y_{1q} \\ \vdots & & \vdots \\ Y_{n1} & \dots & Y_{nq} \end{bmatrix} = \begin{bmatrix} -y_1^T - \\ \vdots \\ -y_n^T - \end{bmatrix}_{n \times q}$   $y_i \in \mathbb{R}^q$

$X = \begin{bmatrix} - & x_1^T - \\ \vdots & \\ - & x_n^T - \end{bmatrix}_{n \times K}$ ,  $x_i \in \mathbb{R}^K$

$B = \begin{bmatrix} 1 & & 1 \\ \beta_1 & \dots & \beta_q \\ 1 & & 1 \end{bmatrix}_{K \times q}$   $\beta_i \in \mathbb{R}^K$

$\Sigma = \begin{bmatrix} - & \Sigma_1^T - \\ \vdots & \\ - & \Sigma_n^T - \end{bmatrix}_{n \times q}$   $\Sigma_i \in \mathbb{R}^q$

$\Sigma_i \stackrel{i.i.d}{\sim} N(0, \underset{q \times q}{\Sigma})$  and  $y_i \sim N(\underset{q \times K}{B} \underset{K \times 1}{X_i}, \Sigma)$

$\forall i=1, \dots, n$ ,  $y_i$  are i.i.d.

$$4 \quad AIC = -2 \log(\text{likelihood}) + 2 (\# \text{ parameters})$$

$$= nq \log 2\pi + n \log |\hat{\Sigma}_k| + nq + 2kq + q(q+1)$$

we can ignore the constant w.r.t "k".

$$\Rightarrow AIC = n \log |\hat{\Sigma}_k| + 2kq$$

$$\text{OR } AIC_k = \log |\hat{\Sigma}_k| + \frac{2kq}{n}$$


---

$$BIC = -2 (\log \text{likelihood}) + \log(n) \times (\# \text{ parameters})$$

$$= \underbrace{nq \log 2\pi}_{\text{No } k} + n \log |\hat{\Sigma}_k| + \underbrace{nq}_{\text{No } k} + \log(n) kq$$

$$+ \frac{q(q+1) \log(n)}{2}$$

$$\underbrace{\hspace{10em}}_{\text{No } k}$$

NO k

$$\Rightarrow BIC_k = \log |\hat{\Sigma}_k| + \frac{\log(n) kq}{n}$$

5:

BIC tends to select small order model and therefore are less prone to overfitting.

## Question 2

$$X = \begin{bmatrix} - & x_1^T & - \\ & \vdots & \\ - & x_n^T & - \end{bmatrix} \rightarrow i^{\text{th}} \text{ row} = x_i^T. \quad n \times p.$$

2):

$$AIC = -2 \times \log(\text{likelihood}) + 2 \times (\# \text{ parameters})$$

$$-2 \log(\text{likelihood}) = n \log(2\pi) + n \log \sigma^2 + \frac{1}{\sigma^2} \left| \sum_{i=1}^n (y_i - x_i^T \beta)^2 \right|$$

$$-2 \log(\text{likelihood}) \Big|_{\theta = \hat{\theta}_{MLE}} = n \log(2\pi) + n \log \hat{\sigma}_{MLE}^2 + n$$

$$\begin{aligned} \text{where } \hat{\sigma}_{MLE}^2 &= \frac{RSS}{n} = \frac{(y - \hat{y})^T (y - \hat{y})}{n} \\ &= \frac{(y - X \hat{\beta})^T (y - X \hat{\beta})}{n} \\ &= \frac{\sum_{i=1}^n (y_i - x_i^T \hat{\beta})^2}{n} \end{aligned}$$

Since the constant  $n \log 2\pi + n$  play no role in model selection and can be ignored.

$$AIC = n \log \hat{\sigma}_k^2 + 2 (\# \text{ parameters})$$

$$= n \log \hat{\sigma}_k^2 + 2(k+1)$$

$$\Rightarrow AIC_k = \frac{n \log \hat{\sigma}_k^2}{n} + \frac{2(k+1)}{n}$$

2): Suppose the # of parameters in the true model is  $k_0$ , but we fit a candidate model of # parameter is  $k_0 + L$  ( $L > 0$ ).

overfitting measure by the AIC means:  $AIC_{k_0+L} < AIC_{k_0}$

$$Pr\{AIC_{k_0+L} < AIC_{k_0}\}$$

$$= Pr\left\{\log(\hat{\sigma}_{k_0+L}^2) + \frac{2(k_0+L+1)}{n} < \log(\hat{\sigma}_{k_0}^2) + \frac{2(k_0+1)}{n}\right\}$$

where  $\hat{\sigma}_{k_0+L}^2 = \frac{RSS_{k_0+L}}{n} \Rightarrow \log \hat{\sigma}_{k_0+L}^2 = \log(RSS_{k_0+L}) - \log n$

$\hat{\sigma}_{k_0}^2 = \frac{RSS_{k_0}}{n} \Rightarrow \log \hat{\sigma}_{k_0}^2 = \log(RSS_{k_0}) - \log n$

$$Pr\left\{\log RSS_{k_0+L} + \frac{2(k_0+L+1)}{n} < \log RSS_{k_0} + \frac{2(k_0+1)}{n}\right\}$$

$$= Pr\left\{\log \frac{RSS_{k_0+L}}{RSS_{k_0}} < \frac{-2L}{n}\right\}$$

$$= Pr\left\{\log \frac{RSS_{k_0}}{RSS_{k_0+L}} > \frac{2L}{n}\right\}$$

$$\log \frac{RSS_{k_0+L}}{RSS_{k_0}} = -\log \frac{RSS_{k_0}}{RSS_{k_0+L}}$$

$$= Pr\left\{\frac{RSS_{k_0}}{RSS_{k_0+L}} > e^{\frac{2L}{n}}\right\}$$

$$= \Pr \left\{ \frac{RSS_{k_0} - RSS_{k_0+L}}{RSS_{k_0+L}} > e^{\frac{2L}{n}} - 1 \right\}$$

where  $\frac{RSS_{k_0} - RSS_{k_0+L}}{\sigma^2} \sim \chi^2_L$  idpt

$\frac{RSS_{k_0+L}}{\sigma^2} \sim \chi^2_{n-k_0-L}$

$$\Rightarrow \Pr \left\{ \frac{\chi^2_L / L}{\chi^2_{n-k_0-L} / (n-k_0-L)} > \frac{(n-k_0-L)}{L} \left( e^{\frac{2L}{n}} - 1 \right) \right\}$$

$$= \Pr \left\{ F_{L, n-k_0-L} > \frac{(n-k_0-L)}{L} \left( e^{\frac{2L}{n}} - 1 \right) \right\}$$

why  $RSS_{k_0} - RSS_{k_0+L} \perp RSS_{k_0+L}$  ?

$$\text{let } X_{k_0}^* = \begin{bmatrix} X_{k_0} & | & \mathbf{0} \end{bmatrix}_{n \times (k_0+L)} \quad X_{k_0} = \begin{bmatrix} x_1 & \dots & x_{k_0} \\ | & & | \end{bmatrix}$$

$$X_{k_0+L} = \begin{bmatrix} X_{k_0} & | & X_L \end{bmatrix} \quad X_L = \begin{bmatrix} x_{k_0+1} & \dots & x_{k_0+L} \\ | & & | \end{bmatrix}$$

$$H_{k_0} = X_{k_0}^* (X_{k_0}^{*T} X_{k_0}^*)^{-1} X_{k_0}^{*T} \quad H_{k_0+L} = X_{k_0+L} (X_{k_0+L}^T X_{k_0+L})^{-1} X_{k_0+L}^T$$

$\underbrace{\quad}_{\in \mathbb{R}^{k_0+L} \times \mathbb{R}^{k_0+L}}$

Fact:  $\text{colspace}(X_{k_0}^*) \subseteq \text{colspace}(X_{k_0+L})$

$$\text{thus } H_{k_0} H_{k_0+L} = H_{k_0+L} H_{k_0} = H_{k_0}.$$

$$RSS_{k_0} = y^T (I - H_{k_0}) y \quad RSS_{k_0+L} = y^T (I - H_{k_0+L}) y.$$

$$RSS_{k_0} - RSS_{k_0+L} = y^T (H_{k_0+L} - H_{k_0}) y.$$

$$\text{cov} \{ (H_{k_0+L} - H_{k_0}) y, (I - H_{k_0+L}) y \}$$

$$= \sigma^2 (H_{k_0+L} - H_{k_0}) (I - H_{k_0+L}) = \sigma^2 (H_{k_0+L} - H_{k_0+L} - H_{k_0} + H_{k_0}) = \mathbf{0}$$

$$3): \frac{n-k_0-L}{L} \left( e^{\frac{2L}{n}} - 1 \right)$$

$$= \frac{n-k_0-L}{L} \left( \frac{2L}{n} + o\left(\frac{1}{n^2}\right) \right)$$

$$= 2 \frac{n-k_0-L}{n} + o\left(\frac{1}{n}\right) \xrightarrow{\text{as } n \rightarrow \infty} 2$$

$$\text{and } F_{L, n-k_0-L} \longrightarrow \frac{\chi_L^2}{L} \quad (\text{why? since } \frac{\chi_{cn}^2}{n} = \frac{1}{n} \sum_{i=1}^n z_i^2 \xrightarrow{LNN} \mathbb{E}\{\chi^2(1)\} = 1)$$

thus, as  $n \rightarrow \infty$ ,

$$P_n \{AIC_{k_0+L} < AIC_{k_0}\} = P_n \{\chi_L^2 > 2L\}.$$

$$z_i \stackrel{iid}{\sim} N(0,1)$$

$$z_i^2 \stackrel{iid}{\sim} \chi^2(1)$$