

[Criterion 3] Mallow's C_p

$$C_p(\mathcal{M}) = \frac{\overset{RSS}{\text{SSE}(\mathcal{M})}}{\hat{\sigma}^2} - n + 2 \times p_{\mathcal{M}}$$

- ▶ $\hat{\sigma}^2 = \text{SSE}(\mathcal{F})/df_{\mathcal{F}}$
 - ▶ \mathcal{F} denotes the fullest model
 - ▶ best estimate of σ^2
- ▶ The criterion is motivated from the Model Error (ME).
 - ▶ $\text{ME} = \|E(Y) - \hat{Y}\|^2$
 - ▶ $E(\text{ME}) = E(\text{SSE}) + \sigma^2(-n + 2p)$
 - ▶ See the derivation on Notes.

[Criterion 3] Mallow's C_p

$$C_p(\mathcal{M}) = \frac{\text{SSE}(\mathcal{M})}{\hat{\sigma}^2} - n + 2 \times p_{\mathcal{M}}$$

- ▶ $\hat{\sigma}^2 = \text{SSE}(\mathcal{F})/df_{\mathcal{F}}$
 - ▶ \mathcal{F} denotes the fullest model
 - ▶ best estimate of σ^2
- ▶ The criterion is motivated from the Model Error (ME).
 - ▶ $\text{ME} = \|\text{E}(Y) - \hat{Y}\|^2$
 - ▶ $\text{E}(\text{ME}) = \text{E}(\text{SSE}) + \sigma^2(-n + 2p)$
 - ▶ See the derivation on Notes.

[Criterion 3] Mallow's C_p

$$C_p(\mathcal{M}) = \frac{\text{SSE}(\mathcal{M})}{\hat{\sigma}^2} - n + 2 \times p_{\mathcal{M}} \quad \checkmark$$

- ▶ $\hat{\sigma}^2 = \text{SSE}(\mathcal{F})/df_{\mathcal{F}}$
 - ▶ \mathcal{F} denotes the fullest model
 - ▶ best estimate of σ^2
- ▶ The criterion is motivated from the Model Error (ME).
 - ▶ $\text{ME} = \|E(Y) - \hat{Y}\|^2$
 - ▶ $E(\text{ME}) = \underbrace{E(\text{SSE}) + \sigma^2(-n + 2p)}$
 - ▶ See the derivation on Notes.

- ▶ Let $a = E(Y) - \hat{Y}$.

$$E \|a\|^2 = \|E(a)\|^2 + \text{tr}\{\text{var}(a)\}$$

- ▶ Thus,

$$\begin{aligned} E(\text{ME}) &= \|E(Y) - E(\hat{Y})\|^2 + \text{tr}\{\text{var}(\hat{Y})\} \\ &= \text{bias}^2 + \text{variance} \end{aligned}$$

$$\begin{aligned}
 E \|a\|^2 &= E \|a - E(a) + E(a)\|^2 \\
 &= E \left\{ \underbrace{\|a - E(a)\|^2}_{(1)} + \underbrace{\|E(a)\|^2}_{(2)} + \underbrace{2 \times \{a - E(a)\}^T E(a)}_{(3)} \right\}
 \end{aligned}$$

► Let $a = \underbrace{E(Y)}_{n \times 1} - \underbrace{\hat{Y}}_{n \times 1}$.

★

$$E \|a\|^2 = \underbrace{\|E(a)\|^2}_{(2)} + \underbrace{\text{tr}\{\text{var}(a)\}}_{(1)}$$

$E(a)$ is constant

► Thus,

$$\begin{aligned}
 E(\text{ME}) &= \|E(Y) - E(\hat{Y})\|^2 + \text{tr}\{\text{var}(\hat{Y})\} \\
 &= \text{bias}^2 + \text{variance}
 \end{aligned}$$

$$(1) = E \|a - E(a)\|^2 = E \text{tr} \{ \{a - E(a)\} \{a - E(a)\}^T \} = \text{tr} \{ \text{var}(a) \}$$

$$(2) = \|E(a)\|^2$$

$$(3) = 2 \times E \{ (a - E(a))^T E(a) \} = 2 \times \underbrace{E \{ a - E(a) \}^T}_{=0} E(a)$$

$$E \{ a - E(a) \} = E(a) - E(a) = 0$$

► Let $a = E(Y) - \hat{Y}$.

$$E(a) = E(Y) - E(\hat{Y})$$

✓ $E \|a\|^2 = \|E(a)\|^2 + \text{tr}\{\text{var}(a)\}$

► Thus,

$$\begin{aligned} E(\text{ME}) &= \|E(Y) - E(\hat{Y})\|^2 + \text{tr}\{\text{var}(\hat{Y})\} \\ &= \text{bias}^2 + \text{variance} \end{aligned}$$

$$\text{var}(a) = \text{var}\{E(Y) - \hat{Y}\} = \text{var}(-\hat{Y}) = \text{var}(\hat{Y})$$

Want both bias and variance to be small.

[Criterion 4] AIC: Akaike Information Criterion

$$Y = X\beta + \epsilon$$

- Definition: For a general model \mathcal{M} with parameter θ .

$$AIC(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + 2 \times p_{\mathcal{M}} \Rightarrow \text{penalty}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE. $\hat{\theta}$

negative $L \Rightarrow$ small
complexity $p \Rightarrow$ small

- Motivation: Kullback-Leibler discrepancy

$$KL(f, g) = \int \log \frac{f(\mathbf{y})}{g(\mathbf{y}; \theta)} f(\mathbf{y}) d\mathbf{y}$$

- A measure of difference between a true fixed f and various competing models g depending on parameter θ .
 - non-symmetric $KL(f, g) \neq KL(g, f)$.
 - $KL(f, g) \geq KL(f, f) = 0$.
 - $KL(f, g) = - \int \log g(\mathbf{y}; \theta) f(\mathbf{y}) d\mathbf{y} + \text{constant}$

[Criterion 4] AIC: Akaike Information Criterion

- Definition: For a general model \mathcal{M} with parameter θ ,

$$\text{AIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + 2 \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- Motivation: Kullback-Leibler discrepancy

KL divergence

$$KL(f, g) = \int \log \frac{f(\mathbf{y})}{g(\mathbf{y}; \theta)} f(\mathbf{y}) d\mathbf{y}$$

*f: \mathbf{Y} ↘
g: $\mathbf{x} \beta \leftarrow$
parameters*

- A measure of difference between a true fixed f and various competing models g depending on parameter θ .

- non-symmetric $KL(f, g) \neq KL(g, f)$.
- $KL(f, g) \geq KL(f, f) = 0$.
- $KL(f, g) = - \int \log g(\mathbf{y}; \theta) f(\mathbf{y}) d\mathbf{y} + \text{constant}$

[Criterion 4] AIC: Akaike Information Criterion

- Definition: For a general model \mathcal{M} with parameter θ ,

$$\text{AIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + 2 \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- Motivation: Kullback-Leibler discrepancy

$$\int \log\left(-\frac{g}{f}\right) \times \underline{g} \, d\cdot$$

$$KL(f, g) = \int \log \frac{\underline{f}(\mathbf{y})}{\underline{g}(\mathbf{y}; \theta)} \underline{f}(\mathbf{y}) d\mathbf{y}$$

- A measure of difference between a true fixed f and various competing models g depending on parameter θ .
 - non-symmetric $KL(f, g) \neq KL(g, f)$.
 - $KL(f, g) \geq KL(f, f) = 0$.
 - $KL(f, g) = - \int \log g(\mathbf{y}; \theta) f(\mathbf{y}) d\mathbf{y} + \text{constant}$

[Criterion 4] AIC: Akaike Information Criterion

- Definition: For a general model \mathcal{M} with parameter θ ,

$$\text{AIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + 2 \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- Motivation: Kullback-Leibler discrepancy

$$\text{KL}(f, g) = \int \log \frac{f(\mathbf{y})}{g(\mathbf{y}; \theta)} f(\mathbf{y}) d\mathbf{y}$$

Handwritten notes in orange:

$$\begin{aligned} & \text{KL}(f, f) \\ &= \int \log \frac{f(y)}{f(y)} f(y) dy \\ &= \int \log(1) \times f(y) dy \\ &= 0 \end{aligned}$$

- A measure of difference between a true fixed f and various competing models g depending on parameter θ .
 - non-symmetric $\text{KL}(f, g) \neq \text{KL}(g, f)$.
 - $\text{KL}(f, g) \geq \text{KL}(f, f) = 0$.
 - $\text{KL}(f, g) = - \int \log g(\mathbf{y}; \theta) f(\mathbf{y}) d\mathbf{y} + \text{constant}$

[Criterion 4] AIC: Akaike Information Criterion

- ▶ Definition: For a general model \mathcal{M} with parameter θ ,

$$\text{AIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + 2 \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- ▶ Motivation: Kullback-Leibler discrepancy

$$KL(f, g) = \int \log \frac{f(\mathbf{y})}{g(\mathbf{y}; \theta)} f(\mathbf{y}) d\mathbf{y}$$

- ▶ A measure of difference between a true fixed f and various competing models g depending on parameter θ .
 - ▶ non-symmetric $KL(f, g) \neq KL(g, f)$.
 - ▶ $KL(f, g) \geq KL(f, f) = 0$.
 - ▶ $KL(f, g) = - \int \log g(\mathbf{y}; \theta) f(\mathbf{y}) d\mathbf{y} + \text{constant}$

Modification from KL

① Because θ is unknown plug in MLE $\hat{\theta}(\gamma)$ from the observed data $\gamma \Rightarrow$

$$\Delta := - \int \log g(z; \hat{\theta}(\gamma)) f(z) dz$$

$$= -E_z \log g(z; \hat{\theta}(\gamma)) \quad (z \text{ independent with } \gamma)$$

can be used to approximate KL-divergence

② In the MLR with σ^2 known, next prove

$$E_\gamma(\text{AIC}) = E_\gamma(z \times \Delta) + \text{constant}$$

Proof:

Step 1. Form of AIC

(1.1) Likelihood of Y under the model $E(Y) = X\beta$ is

$$L_M(\beta) = -\frac{1}{2\sigma^2} \|Y - X\beta\|^2 - \frac{n}{2} \log(2\pi\sigma^2) \quad \text{known}$$

$$(1.2) \text{ MLE } \hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\begin{aligned} (1.3) \quad \text{AIC} &= -2L_M(\hat{\beta}) + 2 \times p \\ &= \|Y - X\hat{\beta}\|^2 + \text{constant} + 2 \times p \\ &= \text{RSS}_M + \text{constant} + 2 \times p \end{aligned}$$

Step 2: $E(AIC)$

By analysis of C_p (Step 3)

We have calculated

$$E(RSS_M) = \mu^T (I - P_p) \mu + (n-p) \sigma^2$$

with $\mu = E(Y)$

$$P_p = X_p (X_p^T X_p)^{-1} X_p \quad \text{using } p \text{ covariates}$$

Thus $E(AIC)$

$$= \frac{1}{\sigma^2} E(RSS_M) + 2p + \text{constant}$$

$$= \frac{1}{\sigma^2} \mu^T (I - P_p) \mu + p + n$$

Step 3: Calculate Δ

- \mathbf{Z} denotes a random matrix independent with \mathbf{Y} but follows the same distribution as \mathbf{Y}

- So we have $\mathbf{Z} = \mu + \epsilon_{\mathbf{Z}}$

$$\begin{aligned}\Delta &= -E_{\mathbf{Z}} \log q(\mathbf{Z}; \hat{\theta}(\mathbf{Y})) \\ &= \frac{1}{2\sigma^2} E_{\mathbf{Z}} \|\mathbf{Z} - \mathbf{X} \hat{\beta}(\mathbf{Y})\|^2 + \text{constant}\end{aligned}$$

$$= \frac{1}{2\sigma^2} E_{\mathbf{Z}} \|\mu + \epsilon_{\mathbf{Z}} - \mathbf{X} \hat{\beta}(\mathbf{Y})\|^2 + \text{constant}$$

$$= \frac{1}{2\sigma^2} \left\{ \|\mu - \mathbf{X} \hat{\beta}(\mathbf{Y})\|^2 + E(\epsilon_{\mathbf{Z}}^T \epsilon_{\mathbf{Z}}) \right\}$$

$$= \frac{1}{2\sigma^2} (\|\mu - \mathbf{X} \hat{\beta}(\mathbf{Y})\|^2 + n\sigma^2)$$

Step 4: Calculate $E_Y(z\Delta)$

$$E_Y(z\Delta) = \frac{1}{\sigma^2} E_Y [\| \mu - X \hat{\beta}(Y) \|^2] + n$$

$$= \frac{1}{\sigma^2} E_Y [\| \mu - PY \|^2] + n$$

$$= \frac{1}{\sigma^2} \mu^T (I - P) \mu + p\sigma^2 \quad \begin{array}{l} \leadsto \text{(same as ME)} \\ \text{in } C_p \end{array}$$

In summary,

$$E_Y(z\Delta) = \frac{1}{\sigma^2} \mu^T (I - P) \mu + p$$

$$E_Y(\text{AIC}) = \frac{\mu^T (I - P) \mu}{\sigma^2} + (n - p) + 2p$$

n is fixed, does not influence model selection

AIC under multiple linear models

- ▶ if σ^2 is known,

$$\text{AIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + 2p_{\mathcal{M}}.$$

- Similar to C_p if replace σ^2 by $\hat{\sigma}^2$ (only differ by $-n$)

- ▶ if σ^2 is unknown,

$$\text{AIC} = n \log(\text{SSE}_{\mathcal{M}} / n) + 2p_{\mathcal{M}}$$

AIC under multiple linear models

- ▶ if σ^2 is known,

$$\text{AIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + 2p_{\mathcal{M}}.$$

- Similar to C_p if replace σ^2 by $\hat{\sigma}^2$ (only differ by $-n$)

- ▶ if σ^2 is unknown,

$$\text{AIC} = n \log(\text{SSE}_{\mathcal{M}} / n) + 2p_{\mathcal{M}}$$

[Criterion 5] BIC: Bayesian Information Criterion

For a general model \mathcal{M} with parameter θ ,

$$\text{BIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + \log(n) \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- BIC penalizes larger models more heavily and so will tend to prefer smaller models in comparison to AIC.

[Criterion 5] BIC: Bayesian Information Criterion

For a general model \mathcal{M} with parameter θ ,

$$\text{BIC}(\mathcal{M}) = -2 \log L_{\mathcal{M}}(\hat{\theta}) + \log(n) \times p_{\mathcal{M}}$$

where $L_{\mathcal{M}}(\hat{\theta})$ denotes the likelihood function of the parameters in the model \mathcal{M} evaluated at the MLE.

- BIC penalizes larger models more heavily and so will tend to prefer smaller models in comparison to AIC.

BIC is derived under the Bayesian perspective

- ▶ Consider the multiple linear model with σ^2 known.
- ▶ Suppose we fit a submodel with $\mathbf{X}_p\beta_p$
 - ▶ p can be smaller than the total number of covariates
 - ▶ Assume β has prior distribution $N_p(\mathbf{m}, \sigma^2 V)$
 - ▶ The **log posterior distribution** of β_p is proportional to

$$\text{BIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + \log(n)p_{\mathcal{M}}$$

- ▶ Detailed proof: read Linear Regression Analysis (Lee and Seber) Section 12.3.4

BIC is derived under the Bayesian perspective

- ▶ Consider the multiple linear model with σ^2 known.
- ▶ Suppose we fit a submodel with $\mathbf{X}_p\beta_p$
 - ▶ p can be smaller than the total number of covariates
 - ▶ Assume β has prior distribution $N_p(\mathbf{m}, \sigma^2 V)$
 - ▶ The **log posterior distribution** of β_p is proportional to

$$\text{BIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + \log(n)p_{\mathcal{M}}$$

- ▶ Detailed proof: read Linear Regression Analysis (Lee and Seber) Section 12.3.4

BIC is derived under the Bayesian perspective

- ▶ Consider the multiple linear model with σ^2 known.
- ▶ Suppose we fit a submodel with $\mathbf{X}_p\beta_p$
 - ▶ p can be smaller than the total number of covariates
 - ▶ Assume β has prior distribution $N_p(\mathbf{m}, \sigma^2 V)$
 - ▶ The **log posterior distribution** of β_p is proportional to

$$\text{BIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + \log(n)p_{\mathcal{M}}$$

- ▶ Detailed proof: read Linear Regression Analysis (Lee and Seber) Section 12.3.4

BIC is derived under the Bayesian perspective

- ▶ Consider the multiple linear model with σ^2 known.
- ▶ Suppose we fit a submodel with $\mathbf{X}_p\beta_p$
 - ▶ p can be smaller than the total number of covariates
 - ▶ Assume β has prior distribution $N_p(\mathbf{m}, \sigma^2 V)$
 - ▶ The **log posterior distribution** of β_p is proportional to

$$\text{BIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + \log(n)p_{\mathcal{M}}$$

- ▶ Detailed proof: read Linear Regression Analysis (Lee and Seber) Section 12.3.4

BIC is derived under the Bayesian perspective

- ▶ Consider the multiple linear model with σ^2 known.
- ▶ Suppose we fit a submodel with $\mathbf{X}_p\beta_p$
 - ▶ p can be smaller than the total number of covariates
 - ▶ Assume β has prior distribution $N_p(\mathbf{m}, \sigma^2 V)$
 - ▶ The **log posterior distribution** of β_p is proportional to

$$\text{BIC} = \frac{\text{SSE}_{\mathcal{M}}}{\sigma^2} + \log(n)p_{\mathcal{M}}$$

- ▶ Detailed proof: read Linear Regression Analysis (Lee and Seber) Section 12.3.4

Summary

- ▶ R^2 : motivated from $\text{corr}^2(\hat{Y}, Y)$ (prefer larger)
- ▶ Adjusted R^2 : penalizes model complexity (prefer larger)
- ▶ C_p : motivated from $\|E(Y) - \hat{Y}\|^2$ (prefer smaller)
- ▶ AIC and BIC: motivated from KL divergence (prefer smaller)

Summary

- ▶ R^2 : motivated from $\text{corr}^2(\hat{Y}, Y)$ (prefer larger)
- ▶ Adjusted R^2 : penalizes model complexity (prefer larger)
- ▶ C_p : motivated from $\|E(Y) - \hat{Y}\|^2$ (prefer smaller)
- ▶ AIC and BIC: motivated from KL divergence (prefer smaller)

Summary


- ▶ R^2 : motivated from $\text{corr}^2(\hat{Y}, Y)$ (prefer larger)
- ▶ Adjusted R^2 : penalizes model complexity (prefer larger)
- ▶ C_p : motivated from $\|E(Y) - \hat{Y}\|^2$ (prefer smaller)
- ▶ AIC and BIC: motivated from KL divergence (prefer smaller)

Summary

- ▶ R^2 : motivated from $\text{corr}^2(\hat{Y}, Y)$ (prefer larger)
- ▶ Adjusted R^2 : penalizes model complexity (prefer larger)
- ▶ C_p : motivated from $\|E(Y) - \hat{Y}\|^2$ (prefer smaller)
- ▶ AIC and BIC: motivated from KL divergence (prefer smaller)

[MS 3] Criterion-based Procedures

3.1 Best Subset Selection

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.

2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here best is defined as having the smallest $\text{RSS} = \text{SSE}$, or equivalently largest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.1 Best Subset Selection

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here best is defined as having the smallest $\text{RSS} = \text{SSE}$, or equivalently largest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.1 Best Subset Selection

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here best is defined as having the smallest $\text{RSS} = \text{SSE}$, or equivalently largest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.1 Best Subset Selection

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 1, 2, \dots, p$:
 - (a) Fit all $\binom{p}{k}$ models that contain exactly k predictors.
 - (b) Pick the best among these $\binom{p}{k}$ models, and call it \mathcal{M}_k . Here best is defined as having the smallest $\text{RSS} = \text{SSE}$, or equivalently largest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.2 Stepwise Selection: Forward

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 0, 1, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the best among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here best is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.2 Stepwise Selection: Forward

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 0, 1, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the best among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here best is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.2 Stepwise Selection: Forward

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 0, 1, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the best among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here best is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .

3.2 Stepwise Selection: Forward

1. Let \mathcal{M}_0 denote the null model, which contains no predictors. This model simply predicts the sample mean for each observation.
2. For $k = 0, 1, \dots, p - 1$:
 - (a) Consider all $p - k$ models that augment the predictors in \mathcal{M}_k with one additional predictor.
 - (b) Choose the best among these $p - k$ models, and call it \mathcal{M}_{k+1} . Here best is defined as having smallest RSS or highest R^2 .
3. Select a single best model from among $\mathcal{M}_0, \dots, \mathcal{M}_p$ using C_p , AIC, BIC, or adjusted R^2 .