

## 6. Model Selection

[MS 1] Why?

# [MS 1] Why?

## 1. Model Interpretability.

- ▶ Often not all predictors are associated with the response.

$$Y = \beta_0 + \beta_A X_A + \beta_I X_I + \epsilon \quad (\beta_A \neq 0, \beta_I = 0)$$

- ▶ Discover associated predictors  $X_A$ .
  - ▶ Remove irrelevant predictors  $X_I$ ; reduce unnecessary complexity.
- ▶ Occam's razor: Prefer models easier to interpret.

## 2. Prediction Accuracy.

- ▶ Models are evaluated by prediction accuracy.

# [MS 1] Why?

body fat ← measurements

## 1. Model Interpretability.

- ▶ Often not all predictors are associated with the response.

$$Y = \beta_0 + \beta_A \underline{X_A} + \beta_I \underline{X_I} + \epsilon \quad (\underline{\beta_A \neq 0}, \underline{\beta_I = 0})$$

- ▶ Discover associated predictors  $X_A$ . ✓
- ▶ Remove irrelevant predictors  $X_I$ ; reduce unnecessary complexity.
- ▶ Occam's razor: Prefer models easier to interpret.

## 2. Prediction Accuracy.

- ▶ Models are evaluated by prediction accuracy.

# [MS 1] Why?

## 1. Model Interpretability.

- ▶ Often not all predictors are associated with the response.

$$Y = \beta_0 + \beta_A X_A + \beta_I X_I + \epsilon \quad (\beta_A \neq 0, \beta_I = 0)$$

- ▶ Discover associated predictors  $X_A$ .
    - ▶ Remove irrelevant predictors  $X_I$ ; reduce unnecessary complexity.
  - ▶ Occam's razor: Prefer models easier to interpret.
- 

## 2. Prediction Accuracy.

- ▶ Models are evaluated by prediction accuracy.

# [MS 1] Why?

## 1. Model Interpretability.

- ▶ Often not all predictors are associated with the response.

$$Y = \beta_0 + \beta_A X_A + \beta_I X_I + \epsilon \quad (\beta_A \neq 0, \beta_I = 0)$$

- ▶ Discover associated predictors  $X_A$ .
  - ▶ Remove irrelevant predictors  $X_I$ ; reduce unnecessary complexity.
- ▶ Occam's razor: Prefer models easier to interpret.

## 2. Prediction Accuracy.

- ▶ Models are evaluated by prediction accuracy.

[MS 2] Criteria

# [Criterion 1] Coefficient of multiple determination $R^2$

*goodness-of-fit*

Given a particular model  $\mathcal{M}$ ,

*corr( $Y, \hat{Y}$ )*

$$R^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M})}{\text{SST}}$$

- ▶  $\text{SSE}(\mathcal{M}) = \|\underline{Y} - \underline{\hat{Y}}_{\mathcal{M}}\|^2$ : SSE of the model  $\mathcal{M}$
- ▶  $\text{SST} = \|\underline{Y} - \underline{\bar{Y}}\|^2$ : sum of squares total

- ▶ Always increases as the model size increases.
  - ▶ Tends to prefer a larger model.
- ▶ Can be used to compare 2 models of the same size.



# [Criterion 1] Coefficient of multiple determination $R^2$

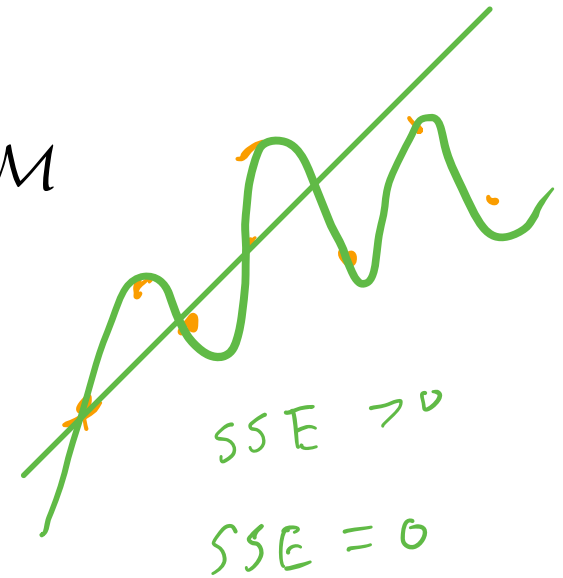
Given a particular model  $\mathcal{M}$ ,

$$R^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M})}{\text{SST}}$$

- ▶  $\text{SSE}(\mathcal{M}) = \|Y - \hat{Y}_{\mathcal{M}}\|^2$ : SSE of the model  $\mathcal{M}$
- ▶  $\text{SST} = \|Y - \bar{Y}\|^2$ : sum of squares total

SSE ↓

- ▶ Always increases as the model size increases.
  - ▶ Tends to prefer a larger model.
- ▶ Can be used to compare 2 models of the same size.



## [Criterion 2] Adjusted $R^2$

$$\underline{R_{adj}^2(\mathcal{M})} = 1 - \frac{\text{SSE}(\mathcal{M}) / \underbrace{(n - p_{\mathcal{M}})}_{\substack{\text{df. of } \mathcal{M}}} }{\text{SST} / \underbrace{(n - 1)}_{\substack{\text{df. of} \\ \text{null model} \\ \text{(no covariates)}}}} = 1 - \frac{\hat{\sigma}_{\mathcal{M}}^2}{\hat{\sigma}_{\text{null}}^2}$$

where  $\underbrace{p_{\mathcal{M}}}$  is the number of parameters in the model  $\mathcal{M}$

- ▶ When  $p_{\mathcal{M}}$  increases,
  - ▶  $\text{SSE}(\mathcal{M})$  always decreases,
  - ▶  $\hat{\sigma}_{\mathcal{M}}^2$  could increase or decrease.
  - ▶  $R_{adj}^2(\mathcal{M})$  could increase or decrease.
- ▶ Prefer models with larger  $R_{adj}^2(\mathcal{M})$ 
  - ▶ A goodness-of-fit measure.

## [Criterion 2] Adjusted $R^2$

$$R_{adj}^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M})/(n - p_{\mathcal{M}})}{\text{SST}/(n - 1)} = 1 - \frac{\hat{\sigma}_{\mathcal{M}}^2}{\hat{\sigma}_{\text{null}}^2}$$

where  $p_{\mathcal{M}}$  is the number of parameters in the model  $\mathcal{M}$

- ▶ When  $p_{\mathcal{M}}$  increases,
  - ▶  $\text{SSE}(\mathcal{M})$  always decreases,
  - ▶  $\hat{\sigma}_{\mathcal{M}}^2$  could increase or decrease.
  - ▶  $R_{adj}^2(\mathcal{M})$  could increase or decrease.
- ▶ Prefer models with larger  $R_{adj}^2(\mathcal{M})$ 
  - ▶ A goodness-of-fit measure.

## [Criterion 2] Adjusted $R^2$

$$R_{adj}^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M}) / (n - p_{\mathcal{M}})}{\text{SST} / (n - 1)} = 1 - \frac{\hat{\sigma}_{\mathcal{M}}^2}{\hat{\sigma}_{\text{null}}^2}$$

where  $p_{\mathcal{M}}$  is the number of parameters in the model  $\mathcal{M}$

- ▶ When  $p_{\mathcal{M}}$  increases,
  - ▶  $\text{SSE}(\mathcal{M})$  always decreases,
  - ▶  $\hat{\sigma}_{\mathcal{M}}^2$  could increase or decrease.
  - ▶  $R_{adj}^2(\mathcal{M})$  could increase or decrease.
- ▶ Prefer models with larger  $R_{adj}^2(\mathcal{M})$ 
  - ▶ A goodness-of-fit measure.

## [Criterion 2] Adjusted $R^2$

$$R_{adj}^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M}) / (n - p_{\mathcal{M}})}{\text{SST} / (n - 1)} = 1 - \frac{\hat{\sigma}_{\mathcal{M}}^2}{\hat{\sigma}_{\text{null}}^2}$$

where  $p_{\mathcal{M}}$  is the number of parameters in the model  $\mathcal{M}$

- ▶ When  $p_{\mathcal{M}}$  increases,
  - ▶  $\text{SSE}(\mathcal{M})$  always decreases,
  - ▶  $\hat{\sigma}_{\mathcal{M}}^2$  could increase or decrease.
  - ▶  $R_{adj}^2(\mathcal{M})$  could increase or decrease.

- ▶ Prefer models with larger  $R_{adj}^2(\mathcal{M})$ 
  - ▶ A goodness-of-fit measure.

$$1 - \frac{n-1}{n-p_{\mathcal{M}}} \times \frac{\text{SSE}}{\text{SST}}$$

↓  
correction  
penalty on model  
complexity

$p_{\mathcal{M}} \uparrow$

## [Criterion 2] Adjusted $R^2$

$$R_{adj}^2(\mathcal{M}) = 1 - \frac{\text{SSE}(\mathcal{M})/(n - p_{\mathcal{M}})}{\text{SST}/(n - 1)} = 1 - \frac{\hat{\sigma}_{\mathcal{M}}^2}{\hat{\sigma}_{\text{null}}^2}$$

where  $p_{\mathcal{M}}$  is the number of parameters in the model  $\mathcal{M}$

- ▶ When  $p_{\mathcal{M}}$  increases,
  - ▶  $\text{SSE}(\mathcal{M})$  always decreases,
  - ▶  $\hat{\sigma}_{\mathcal{M}}^2$  could increase or decrease.
  - ▶  $R_{adj}^2(\mathcal{M})$  could increase or decrease.
- ▶ Prefer models with larger  $R_{adj}^2(\mathcal{M})$ 
  - ▶ A goodness-of-fit measure.

## [Criterion 3] Mallow's $C_p$

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

►  $\hat{\sigma}^2 = \frac{\text{SSE}(\mathcal{F})}{df_{\mathcal{F}}}$

►  $\mathcal{F}$  denotes the fullest model

► best estimate of  $\sigma^2$

$X_1 \cdot \dots \cdot X_{10}$

► The criterion is motivated from Model Error (ME)

►  $\text{ME} = \|E(Y | X) - \hat{Y}\|^2$

## [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  $\sigma^2$

$$E(Y|X) = X\beta$$
$$\hat{Y} = X\hat{\beta}$$

- ▶ The criterion is motivated from Model Error (ME)

- ▶  $\text{ME} = \|E(Y|X) - \hat{Y}\|^2$

prediction to mean value  
of the design matrix  $X$



Step 1. ME form

$$ME = \| \mu - P\mu + P(Y - \mu) \|^2 \quad P = X(X^T X)^{-1} X^T$$

$$\star = \| (I - P)\mu \|^2 + \| P(Y - \mu) \|^2 + 0 \quad \Rightarrow \text{cross term} = 0 \\ \text{by } P(I - P) = 0$$

$$= \| (I - P)\mu \|^2 + \| P\epsilon \|^2 \quad (\text{by } Y = \mu + \epsilon)$$

$$= \mu^T (I - P)^T (I - P) \mu + \epsilon^T P^T P \epsilon$$

$$= \mu^T (I - P) \mu + \epsilon^T P \epsilon$$

Step 2: Calculate expectation of ME

$$\text{By } E(\epsilon^T P \epsilon) = E \{ \text{tr}(P \epsilon \epsilon^T) \} = \text{tr} \{ P E(\epsilon \epsilon^T) \}$$

$$= \text{tr} \{ P \times \sigma^2 I_n \} = \sigma^2 \text{tr}(P) = p$$

$$(\text{By } \text{tr}(P) = p \quad \text{Notes Sep 26})$$

$$\begin{aligned}
 E(ME) &= \mu^T (I - P) \mu + E(\epsilon^T P \epsilon) \\
 &= \mu^T (I - P) \mu + p \sigma^2
 \end{aligned}$$

### Step 3: Expectation of SSE

If we fit a model with  $p$  parameters, let

$$P_p = X_p (X_p^T X_p)^{-1} X_p \quad \text{corresponding hat matrix}$$

$RSS_p$  denote the results of residual sum of squares

$$\begin{aligned}
 E[RSS_p] &= E[Y^T (I - P_p) Y] \\
 &= E[(\mu + \epsilon)^T (I - P_p) (\mu + \epsilon)]
 \end{aligned}$$

$$\begin{aligned}
&= \mu^T (I - P_p) \mu + \cancel{E(\epsilon^T (I - P_p) \mu)} \\
&\quad + \cancel{E\{\mu^T (I - P_p) \mu\}} + E\{\epsilon^T (I - P_p) \epsilon\} \\
&= \mu^T (I - P_p) \mu + \text{tr}(I - P_p) \times \sigma^2 \\
&= \mu^T (I - P_p) \mu + (n - p) \sigma^2
\end{aligned}$$

Then  $\frac{E(ME)}{\sigma^2} = \frac{E(RSS_p)}{\sigma^2} + 2p - n$


Plug in  $\hat{\sigma}^2$  obtained from fullest model

$$\Rightarrow C_p = \frac{RSS_p}{\hat{\sigma}^2} + 2p - n$$

Comment 1 on  $C_p$ :

$$\text{By } E(C_p) \approx \frac{E(RSS_p)}{\sigma^2} + 2p - n$$

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

  
If  $\sigma^2 \rightarrow 0$

$$\approx 0 + n - p + 2p - n = p$$

Thus select model with small  $C_p$  and  $C_p \approx p$ .

## Comment 2 on ME:

Note  $\mu = E(Y)$  &  $P\mu = P E(Y) = E(PY) = E(\hat{Y})$

In  $\star$ , we obtain

$$\begin{aligned} ME &= \| \mu - P\mu \|^2 + \| PY - P\mu \|^2 \\ &= \| E(Y) - E(\hat{Y}) \|^2 + \| \hat{Y} - E(\hat{Y}) \|^2 \end{aligned}$$

$\Downarrow$

①

(Fixed)

$\Downarrow$

②

(Random)

Let  $a = \hat{Y} - E(\hat{Y})$  satisfying  $E(a) = 0$

$$E(\textcircled{2}) = E(a^T a) = E \{ \text{tr}(a a^T) \} = \text{tr} \{ E(a a^T) \} = \text{tr} \{ \text{Var}(a) \}$$

$$E(ME) = \textcircled{1} + E(\textcircled{2}) = \text{Total Bias} + \text{Total Variance}$$