

# Model Selection and Assessment

Chau-Wai Wong

Electrical & Computer Engineering  
North Carolina State University

Contact: [chauwai.wong@ncsu.edu](mailto:chauwai.wong@ncsu.edu). Updated: October 5, 2020.

# Model Selection Definition

**Model Selection:** Choose the best model out of a set of candidate models.

**Model Assessment:** Having chosen a final model, estimating its prediction/generalization error on new data.

Readings: Chapter 7 of Hastie et al.

# Model Selection Examples

(1) Time series:

$$\mathcal{S}_1 = \{AR(1), AR(2), AR(3), \dots\}$$

(2) Linear regression:

$$y_i = \beta_0 + \beta_1 x_{1i} + \dots + \beta_p x_{pi} + e_i, \quad i = 1, \dots, 50.$$

$$\mathcal{S}_2 = \{\beta_0 \neq 0, \beta_1 \neq 0, \dots, (\beta_0, \beta_1) \neq \underline{0}, (\beta_0, \beta_2) \neq \underline{0}, \\ \dots, (\beta_0, \dots, \beta_p) \neq \underline{0}\}$$

## Model Selection Examples (cont'd)

(3) Harmonic model:

$$y(n) = \sum_{i=0}^p A_i e^{j(\omega_i n + \phi)} + v(n), \quad n = 0, \dots, 999,$$

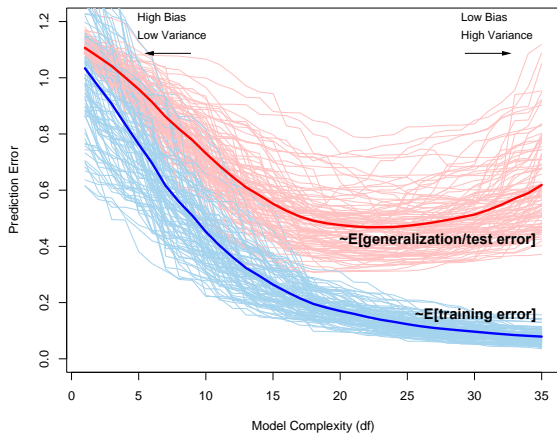
where  $v(n) \sim N(0, \sigma_v^2)$ ,  $\phi \sim \text{Uni}(0, 2\pi]$ , and  $(A_i, \omega_i)$  are fixed but unknown parameters.

$$\mathcal{S}_3 = \{A_0 \neq 0, \dots, (A_0, A_1) \neq \underline{0}, \dots, (A_0, \dots, A_p) \neq \underline{0}\}$$

Note that  $|\mathcal{S}_2| = |\mathcal{S}_3| = 2^{p+1} - 1$ .

# Model Selection Criterion: Generalization Performance

A learning method's **generalization performance** is reflected by its prediction capability assessed using **new/test data** drawn from the same population where the data used for training were drawn.



# Model Selection in Ideal, Data-Rich Scenario

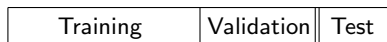
Split data into two three sets:



- 1 Fit  $K$  candidate models to the training data.
- 2 Evaluate the prediction errors using validation data for all models. Select the model with the smallest prediction error. This is called the “validation error.”
- 3 Test the selected model using the test data and evaluate the prediction error. This is called the “test/generalization error.”

# Model Selection in Ideal, Data-Rich Scenario

Split data into two three sets:



- 1 Fit  $K$  candidate models to the training data.
- 2 Evaluate the prediction errors using validation data for all models. Select the model with the smallest prediction error. This is called the “validation error.”
- 3 Test the selected model using the test data and evaluate the prediction error. This is called the “test/generalization error.”
- 4 Question: Why can't validation error be considered as the generalization error? (Hint: Test data mustn't be seen by the model selection process.)

# Model Selection in Practical, Data-Limited Scenario

<b>Strategy</b>	<b>Method</b>
Sample reuse	Crossvalidation, Bootstrap
Analytically approximate test/generalization step	AIC, BIC, MDL, etc.



**Convention:** lower vs. upper cases—deterministic vs. random;  
upper case & bold—deterministic matrix; Tilde below—vector.

**Notations:**  $y_i$  response,  $\underline{x}_i$  collection of predictors for  $y_i$ ,  
 $\mathcal{T} = \{(\underline{x}_i, y_i), i = 1, \dots, N\}$  deterministic data set,  
 $\hat{f}_{\mathcal{T}}(\cdot)$  or  $\hat{y}_{\mathcal{T}}(\cdot)$  prediction function based on/conditioned on  $\mathcal{T}$ ,  
 $L(\cdot, \cdot)$  loss function, e.g.,  $L(a, b) = (a - b)^2$  or  $L(a, b) = |a - b|$ .

Examples when the prediction function is linear:

$$\underbrace{\begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}}_{\underline{y}} = \underbrace{\begin{bmatrix} \underline{x}_1^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}}_{\mathbf{X}} \beta + \underline{\epsilon},$$

$$\begin{aligned} \hat{f}_{\mathcal{T}}(\underline{x}_0) &= \underline{x}_0^T \hat{\beta}_{\mathcal{T}} = \underline{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \underline{y}, \\ \text{or} &= \underline{x}_0^T \tilde{\beta}_{\mathcal{T}} = \underline{x}_0^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \underline{y}. \end{aligned}$$

# Definitions of Test and Training Errors

## Generalization/Test error

$$\text{Err}_{\mathcal{T}} = \mathbb{E}[L(Y^0, \hat{f}_{\mathcal{T}}(\underline{X}^0)) | \mathcal{T}] \text{ (extra-sample error).}$$

## Expected generalization/test error

$$\text{Err} = \mathbb{E}[\text{Err}_{\mathcal{T}}] = \mathbb{E}\left[\mathbb{E}[L(Y^0, \hat{f}_{\mathcal{T}}(\underline{X}^0)) | \mathcal{T}]\right] = \mathbb{E}[L(Y^0, \hat{f}_{\mathcal{T}}(\underline{X}^0))].$$

## Training error

$$\bar{\text{err}} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}_{\mathcal{T}}(\underline{x}_i)).$$

Question: How can you modify the definition of training error to define validation error?

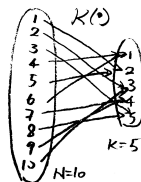
## Cross-Validation Motivation & Example

Cross-Validation (CV), sometimes called rotation estimation, or out-of-sample testing.

**Data Reuse:** Each segment will act as the validation set once, while data in the remaining  $K - 1$  segments are used to calculate a prediction model.

$K$ -Fold CV, typical choice  $K = 5$  or  $10$ . A random partition example when  $K = 5$ :

Data index:	4, 6	1, 5	2, 10	7, 9	3, 8
Segment index:	1	2	3	4	5
A random partition when $K = 5$	Train	Train	Train	Validation	Train



A training-validation split when the 4th segment is acting as the validation set.

# Cross-Validation Error

## Cross-Validation error

$$CV(\hat{f}) = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{f}^{-\kappa(i)}(\underline{x}_i)),$$

where  $\kappa : \{1, \dots, N\} \rightarrow \{1, \dots, K\}$  is a random partition function.

All data points,  $(\underline{x}_i, y_i)$ ,  $i = 1, \dots, N$ , or all segments, contribute to the CV error.

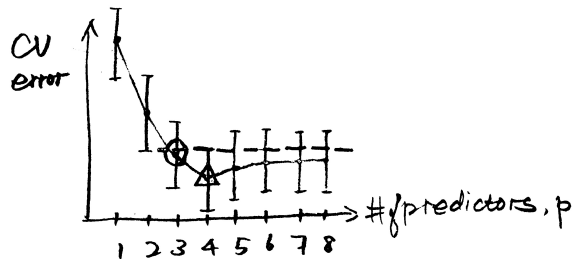
CV error is used to approximate the generalization error.

Note:  $CV(\hat{f})$  estimates the expected generalization error,  $\text{Err}$ , better than the conditional generalization error,  $\text{Err}_{\mathcal{T}}$ . (See Section 7.12 for more details.)

## LOOCV and One SE Rule

**Leave-One-Out Cross-Validation (LOOCV):** A special case of CV when  $K = N$ . Approximately unbiased but has large variance as the training datasets are almost the same.

“One standard error rule”: Choose the most parsimonious model.  
Example: CV error for linear regression on polynomials



$\hat{p}_{\text{lowest}} = 4$  and  $\hat{p}_{\text{one-std-rule}} = 3$ .

## Analytic Approximations

**Observation:** Training error  $\bar{\text{err}} < \text{Err}_{\mathcal{T}}$ , because the fitted model  $\hat{f}_{\mathcal{T}}$  has adapted to data  $\mathcal{T}$ .

Can we find an correction term and add it to the training error to approximate the generalization error, i.e.,  $\bar{\text{err}} + \square = \text{Err}_{\mathcal{T}}$ ?

### In-sample prediction error

$$\text{Err}_{\text{in}} = \frac{1}{N} \sum_{k=1}^N \mathbb{E}[L(Y_k^0, \hat{f}_{\mathcal{T}}(\underline{x}_k)) | \mathcal{T}],$$

which is defined similarly to  $\text{Err}_{\mathcal{T}}$  but uses  $\{(\underline{x}_i, Y_i^0)\}_{i=1}^N$  instead of  $\{(X_i^0, Y_i^0)\}_{i=1}^{\infty}$ .

$\text{Err}_{\text{in}} \approx \text{Err}_{\mathcal{T}}$  if (1)  $\underline{x}_i$  is uniformly sampled from population, and (2)  $N$  is large.

# The Correction Term: Optimism

## Optimism

$$\text{op} \stackrel{\text{def}}{=} \text{Err}_{\text{in}} - \bar{\text{err}}.$$

## Expected optimism

$$\omega \stackrel{\text{def}}{=} \mathbb{E}[\text{op} | \{\mathbf{x}_i\}_{i=1}^N].$$

Example:  $\omega = \frac{2}{N} \sum_{i=1}^N \text{cov}(\hat{y}_i, y_i)$ . The harder we fit, the greater the covariance, and the more op.

## Analytic Form of Optimism

$$\mathbb{E}[\text{Err}_{\text{in}}|\{\underline{x}_i\}] = \bar{\text{err}} + \frac{2}{N} \sum_{i=1}^N \text{cov}(\hat{y}_i, y_i).$$

If  $\hat{y}_i$  is from linear model with  $d$  predictors, we have

$$\mathbb{E}[\text{Err}_{\text{in}}|\{\underline{x}_i\}] = \bar{\text{err}} + 2 \cdot \frac{d}{N} \cdot \sigma_e^2.$$

Try to validate the above expression for parameters  $d$ ,  $N$ , and  $\sigma_e^2$  using a linear regression model as a special case.



# Analytic Approximations

**Analytic Models:** Akaike information criterion (AIC), Bayesian information criterion (BIC), Minimum description length (MDL).

★ One way to estimate the in-sample prediction error  $\text{Err}_{\text{in}}$  is to estimate the optimism and then add it to the training error  $\bar{\text{err}}$ :

$$AIC \text{ or } C_p = \bar{\text{err}} + 2 \cdot \frac{d}{N} \cdot \hat{\sigma}_e^2$$

$$BIC = \frac{N}{\hat{\sigma}_e^2} \left[ \bar{\text{err}} + (\log N) \cdot \frac{d}{N} \cdot \hat{\sigma}_e^2 \right]$$

# Detailed Derivations

Evaluating  $\mathbb{E}[\text{Err}_{\text{in}}|\{\underline{x}_i\}]$ 

$$\begin{aligned}\mathbb{E}[\text{Err}_{\text{in}}|\{\underline{x}_i\}] &= \mathbb{E}\left[\frac{1}{N} \sum_{k=1}^N \mathbb{E}[L(Y_k^0, \hat{f}_{\mathcal{T}}(\underline{x}_k))|\mathcal{T}]|\{\underline{x}_i\}\right] \\ &= \frac{1}{N} \sum_{k=1}^N \mathbb{E}\left[\mathbb{E}[L(Y_k^0, \hat{f}_{\mathcal{T}}(\underline{x}_k))|\{\underline{x}_i\}, \{y_i\}]|\{\underline{x}_i\}\right] \\ &= \frac{1}{N} \sum_{k=1}^N \mathbb{E}[L(Y_k^0, \hat{f}_{\mathcal{T}}(\underline{x}_k))|\{\underline{x}_i\}] \\ &\stackrel{\text{def}}{=} \frac{1}{N} \sum_{k=1}^N \text{Err}(\underline{x}_k)\end{aligned}$$

# The Bias-Variance Decomposition for $\text{Err}(\underline{x}_k)$

...

## Special Case for the Linear Regression Model

Linear model  $\underline{y} = \mathbf{X}\beta + \underline{\epsilon}$

...