

Pattern Information Processing:⁷³ Model Selection by Cross-Validation

Masashi Sugiyama
(Department of Computer Science)

Contact: W8E-505

sugi@cs.titech.ac.jp

<http://sugiyama-www.cs.titech.ac.jp/~sugi/>

Model Parameters

- In the process of parameter learning, we **fixed** model parameters.
- For example, quadratically constrained least-squares with a Gaussian kernel model:
 - **Gaussian width:** $h (> 0)$
 - **Regularization parameter:** $\lambda (\geq 0)$

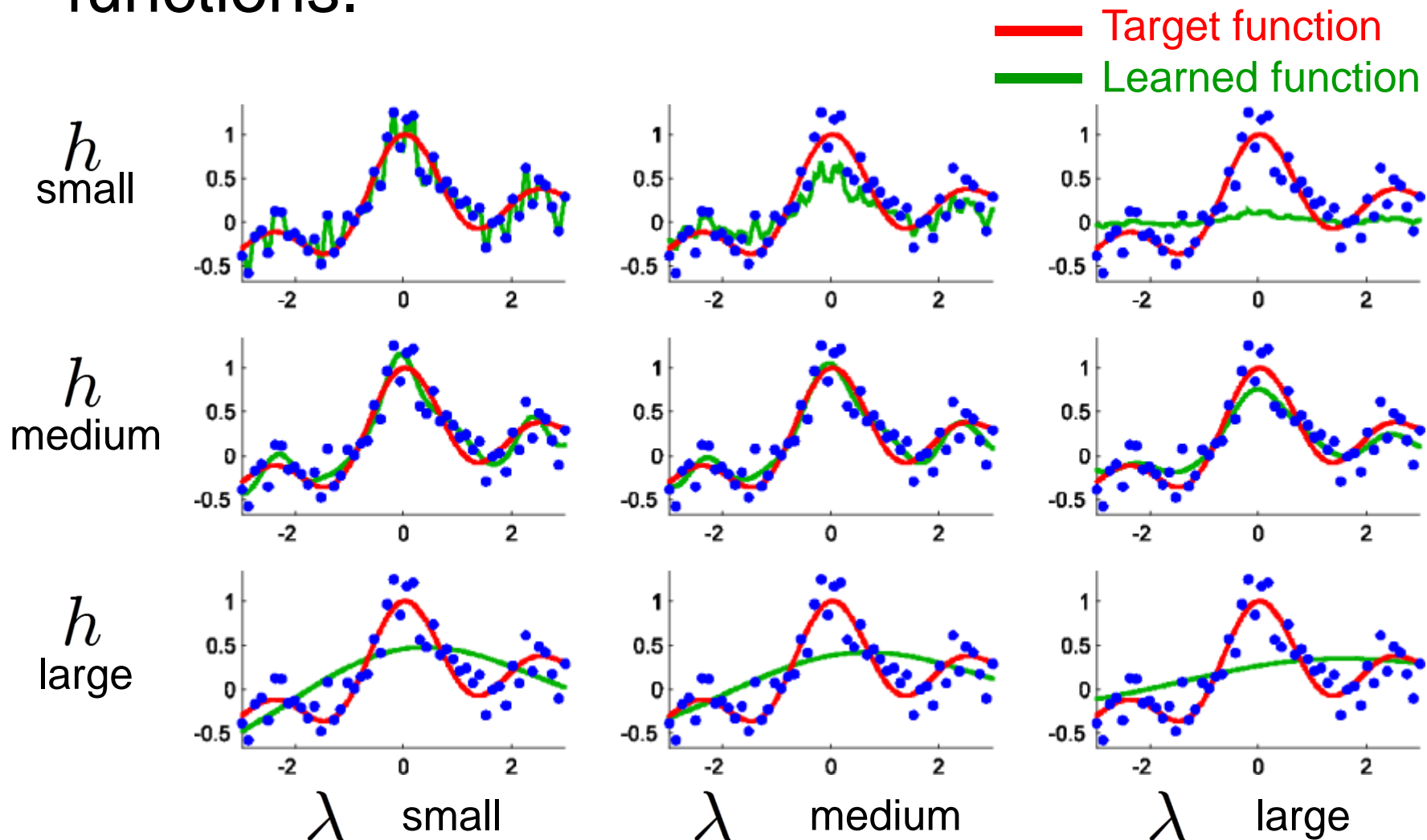
$$\min_{\boldsymbol{\alpha} \in \mathbb{R}^b} \left[\sum_{i=1}^n (f_{\boldsymbol{\alpha}}(\mathbf{x}_i) - y_i)^2 + \lambda \|\boldsymbol{\alpha}\|^2 \right]$$

$$f_{\boldsymbol{\alpha}}(\mathbf{x}) = \sum_{i=1}^n \alpha_i \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2h^2} \right)$$

Different Model Parameters

75

- Model parameters **strongly affect** learned functions.



Determining Model Parameters⁷⁶

- We want to determine the model parameters so that the **generalization error (expected test error)** is minimized.

$$G = \int_{\mathcal{D}} \left(\hat{f}(t) - f(t) \right)^2 q(t) dt$$

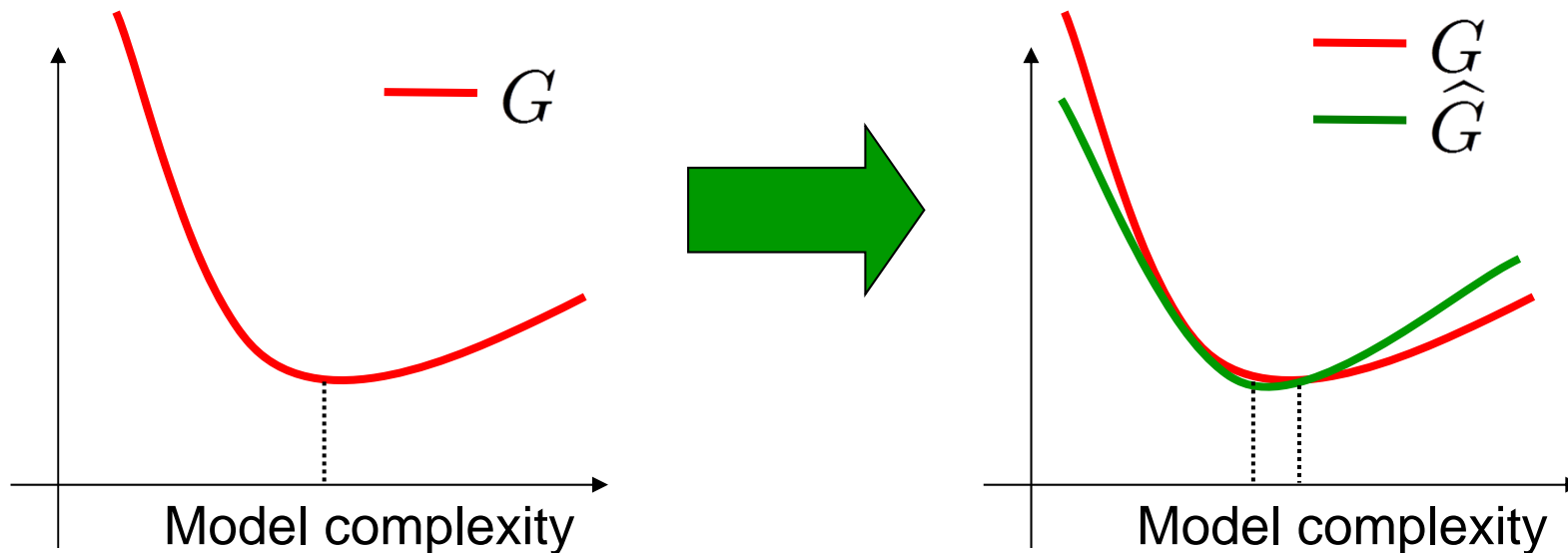
$t \sim q(x)$

- However, $f(x)$ is unknown so the generalization error is not accessible.
- $q(x)$ may also be unknown.

Generalization Error Estimation⁷⁷

$$G = \int_{\mathcal{D}} \left(\hat{f}(\mathbf{t}) - f(\mathbf{t}) \right)^2 q(\mathbf{t}) d\mathbf{t}$$

- Instead, we use a generalization error estimate.



Model Selection

$$\min_{\mathcal{M}} G$$

$$G = \int_{\mathcal{D}} \left(\hat{f}(\mathbf{t}) - f(\mathbf{t}) \right)^2 q(\mathbf{t}) d\mathbf{t}$$

1. Prepare a set of **model candidates**.

$$\{\mathcal{M} \mid \mathcal{M} = (h, \lambda)\}$$

2. Estimate generalization error for each model.

$$\hat{G}(\mathcal{M})$$

3. Choose the one that minimizes the **estimated generalization error**.

$$\hat{\mathcal{M}} = \operatorname{argmin}_{\mathcal{M}} \hat{G}(\mathcal{M})$$

Extra-Sample Method

- Suppose we have an **extra example** (\mathbf{x}', y') in addition to $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$.
- **Idea**: Test the prediction performance of the learned function using the extra example.

$$\hat{G}_{extra} = \left(\hat{f}(\mathbf{x}') - y' \right)^2$$

$$\hat{f} \longleftarrow \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

Extra-Sample Method (cont.) ⁸⁰

- Suppose (\mathbf{x}', y') satisfies:
- | | |
|---|---|
| $\mathbb{E}_{\epsilon'}[\epsilon'] = 0$ | |
| $\mathbf{x}' \sim q(\mathbf{x})$ | $\mathbb{E}_{\epsilon'}[\epsilon'^2] = \sigma^2$ |
| $y' = f(\mathbf{x}') + \epsilon'$ | $\mathbb{E}_{\epsilon'}[\epsilon' \epsilon_i] = 0, \quad \forall i$ |
- $\mathbb{E}_{\epsilon'}:$ Expectation over noise ϵ'

\hat{G}_{extra} is **unbiased** w.r.t. \mathbf{x}' and ϵ' (up to σ^2):

$$\mathbb{E}_{\mathbf{x}'} \mathbb{E}_{\epsilon'} [\hat{G}_{extra}] = G + \sigma^2$$

- **Proof:** $\mathbb{E}_{\mathbf{x}'} \mathbb{E}_{\epsilon'} \left(\hat{f}(\mathbf{x}') - f(\mathbf{x}') - \epsilon' \right)^2$
- $$= \mathbb{E}_{\mathbf{x}'} \mathbb{E}_{\epsilon'} \left[(\hat{f}(\mathbf{x}') - f(\mathbf{x}'))^2 - 2\epsilon'(\hat{f}(\mathbf{x}') - f(\mathbf{x}')) + \epsilon'^2 \right]$$
- $$= G + \sigma^2$$

Extra-Sample Method (cont.) ⁸¹

$$\hat{G}_{extra} = \left(\hat{f}(\mathbf{x}') - y' \right)^2$$

$$\hat{f} \longleftarrow \{(\mathbf{x}_i, y_i)\}_{i=1}^n$$

- \hat{G}_{extra} may be used for model selection.
- However, in practice, such an extra example is not available (or if we have it, it should be included in the original training set!).

Holdout Method

■ **Idea:** Use one of the training samples as an extra sample

- Train a learning machine using $\{(\mathbf{x}_i, y_i)\}_{i \neq j}$

$$\hat{f}_j(\mathbf{x}) \leftarrow \{(\mathbf{x}_i, y_i)\}_{i \neq j}$$

- Test its prediction performance using the holdout sample (\mathbf{x}_j, y_j) :

$$\hat{G}_j = \left(\hat{f}_j(\mathbf{x}_j) - y_j \right)^2$$

Holdout Method (cont.)

- Suppose $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ satisfies:

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} q(\mathbf{x})$$

$$\mathbb{E}_{\epsilon_j} [\epsilon_i] = 0$$

$$y_i = f(\mathbf{x}_i) + \epsilon_i$$

$$\mathbb{E}_{\epsilon_i} \mathbb{E}_{\epsilon_j} [\epsilon_i \epsilon_j] = \begin{cases} \sigma^2 & (i = j) \\ 0 & (i \neq j) \end{cases}$$

- Holdout method is almost unbiased w.r.t. \mathbf{x}_j, ϵ_j :

$$\mathbb{E}_{\mathbf{x}_j} \mathbb{E}_{\epsilon_j} [\hat{G}_j] = G_j + \sigma^2 \approx G + \sigma^2$$

$$G_j = \int_{\mathcal{D}} \left(\hat{f}_j(\mathbf{x}) - f(\mathbf{x}) \right)^2 q(\mathbf{x}) d\mathbf{x}$$

$$\hat{f}_j(\mathbf{x}) \approx \hat{f}(\mathbf{x}) \text{ if } n \text{ is large}$$

- However, \hat{G}_j is heavily affected by the choice of the holdout sample (\mathbf{x}_j, y_j) .

Leave-One-Out Cross-Validation⁸⁴

- **Idea:** Repeat the holdout procedure for all combinations and output the average.

$$\hat{G}_{LOOCV} = \frac{1}{n} \sum_{j=1}^n \hat{G}_j$$

$$\hat{G}_j = \left(\hat{f}_j(\mathbf{x}_j) - y_j \right)^2$$

- LOOCV is almost unbiased w.r.t. $\{\mathbf{x}_i, \epsilon_i\}_{i=1}^n$:

$$\begin{aligned} \mathbb{E}_{\{\mathbf{x}_i\}_{i=1}^n} \mathbb{E}_{\{\epsilon_i\}_{i=1}^n} [\hat{G}_{LOOCV}] \\ \approx \mathbb{E}_{\{\mathbf{x}_i\}_{i=1}^n} \mathbb{E}_{\{\epsilon_i\}_{i=1}^n} [G] + \sigma^2 \end{aligned}$$

k-Fold Cross-Validation

- **Idea:** Randomly split training set into k disjoint subsets $\{\mathcal{T}_j\}_{j=1}^k$.

$$\hat{G}_{kCV} = \frac{1}{k} \sum_{j=1}^k \hat{G}_{\mathcal{T}_j}$$

$$\hat{G}_{\mathcal{T}_j} = \frac{1}{|\mathcal{T}_j|} \sum_{i \in \mathcal{T}_j} \left(\hat{f}_{\mathcal{T}_j}(\mathbf{x}_i) - y_i \right)^2$$

$$\hat{f}_{\mathcal{T}_j}(\mathbf{x}) \longleftarrow \{(\mathbf{x}_i, y_i) \mid i \notin \mathcal{T}_j\}$$

- k-fold is easier to compute and more stable than leave-one-out.

Advantages of CV

86

- **Wide applicability:** Almost unbiasedness of LOOCV holds for (virtually) any learning methods
- **Practical usefulness:** CV has been shown to work very well in many practical applications

Disadvantages of CV

- **Computationally expensive:**

It requires repeating training of models with different subsets of training samples

- **Number of folds:**

It is often recommended to use $k = 5, 10$. However, how to optimally choose k is still open.

Closed Form of LOOCV

88

$$f_{\alpha}(\mathbf{x}) = \sum_{i=1}^b \alpha_i \varphi_i(\mathbf{x})$$

$$\min_{\alpha \in \mathbb{R}^b} \left[\sum_{i=1}^n (f_{\alpha}(\mathbf{x}_i) - y_i)^2 + \lambda \|\alpha\|^2 \right]$$

- For a linear model trained by quadratically constrained least-squares, the LOOCV score can be expressed as

$$\hat{G}_{LOOCV} = \frac{1}{n} \|\widetilde{\mathbf{H}}^{-1} \mathbf{H} \mathbf{y}\|^2$$

$$\mathbf{H} = \mathbf{I} - \mathbf{X}(\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^{\top}$$

$\widetilde{\mathbf{H}}$: same diagonal as \mathbf{H} but zero for off-diagonal

Homework

1. Prove the closed-form expression of leave-one-out cross-validation score for a linear model with quadratically constraint least-squares:

$$\hat{G}_{LOOCV} = \frac{1}{n} \|\widetilde{\mathbf{H}}^{-1} \mathbf{H} \mathbf{y}\|^2$$

Hint: Express $\hat{\alpha}_j$ in terms of $\hat{\alpha}$

- $\hat{\alpha}_j$: Learned parameter without the j-th sample
- $\hat{\alpha}$: Learned parameter with all samples.
- **Key formula:**

$$(\mathbf{U} - \mathbf{u}\mathbf{u}^\top)^{-1} = \mathbf{U}^{-1} + \frac{\mathbf{U}^{-1} \mathbf{u} \mathbf{u}^\top \mathbf{U}^{-1}}{1 - \mathbf{u}^\top \mathbf{U}^{-1} \mathbf{u}}$$

Homework (cont.)

90

2. For your own toy 1-dimensional data, perform simulations using
- Gaussian kernel models
 - Quadratically-constrained least-squares learning
- and optimize
- Width of Gaussian kernel
 - Regularization parameter
- based on cross-validation. Analyze the results when changing
- Target function
 - Number of samples
 - Noise level