

# Pattern Information Processing:<sup>73</sup> Model Selection by Cross-Validation

Masashi Sugiyama  
(Department of Computer Science)

Contact: W8E-505

[sugi@cs.titech.ac.jp](mailto: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 Gaussian kernel models:
  - **Gaussian width:**  $c (> 0)$
  - **Regularization parameter:**  $\lambda (\geq 0)$

$$\hat{f}(\mathbf{x}) = \sum_{i=1}^n \alpha_i K(\mathbf{x}, \mathbf{x}_i)$$

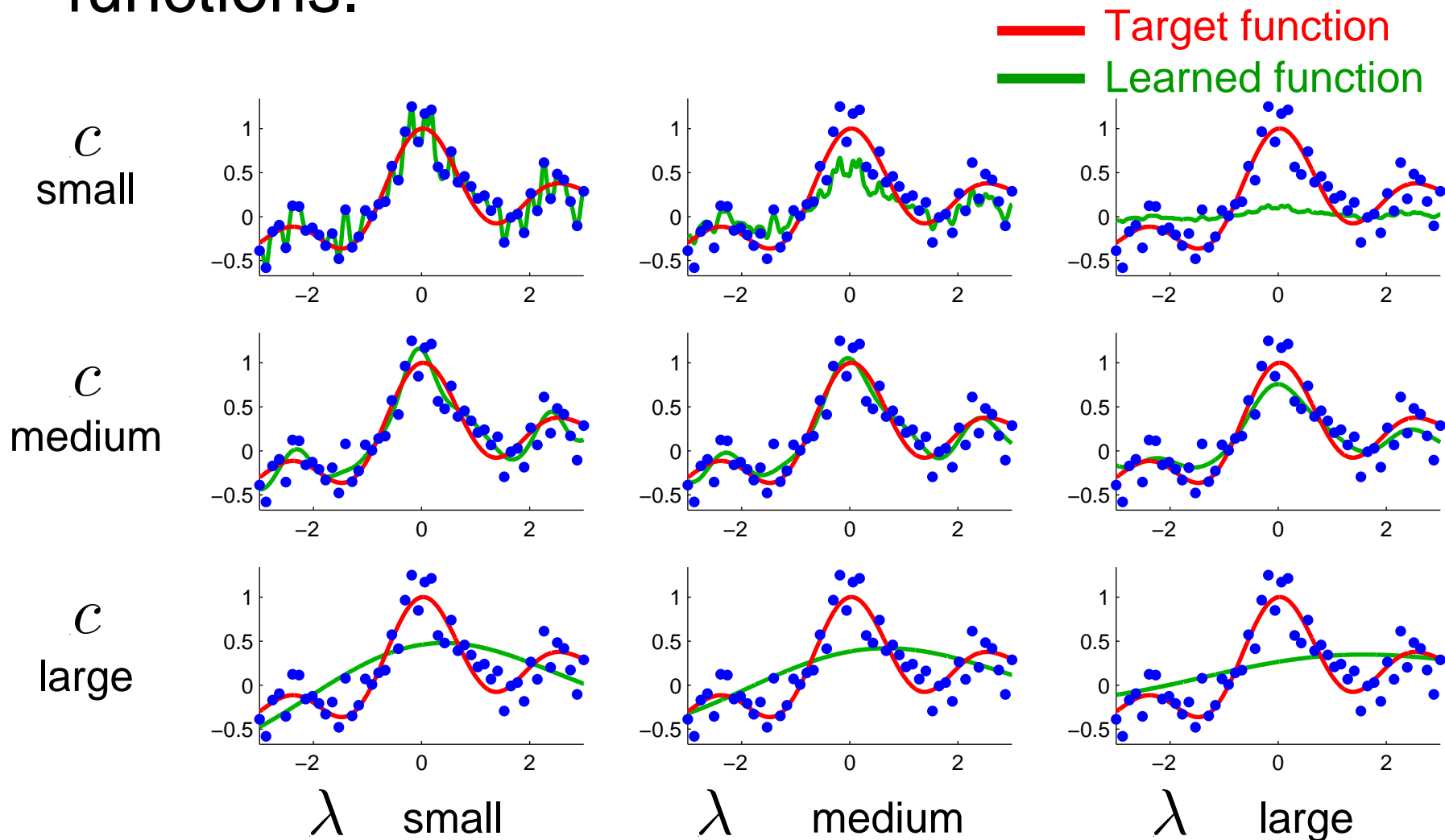
$$K(\mathbf{x}, \mathbf{x}') = \exp \left( -\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2c^2} \right)$$

$$J_{QCLS}(\boldsymbol{\alpha}) = J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|^2$$

# Different Model Parameters

75

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



# Determining Model Parameters<sup>76</sup>

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

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

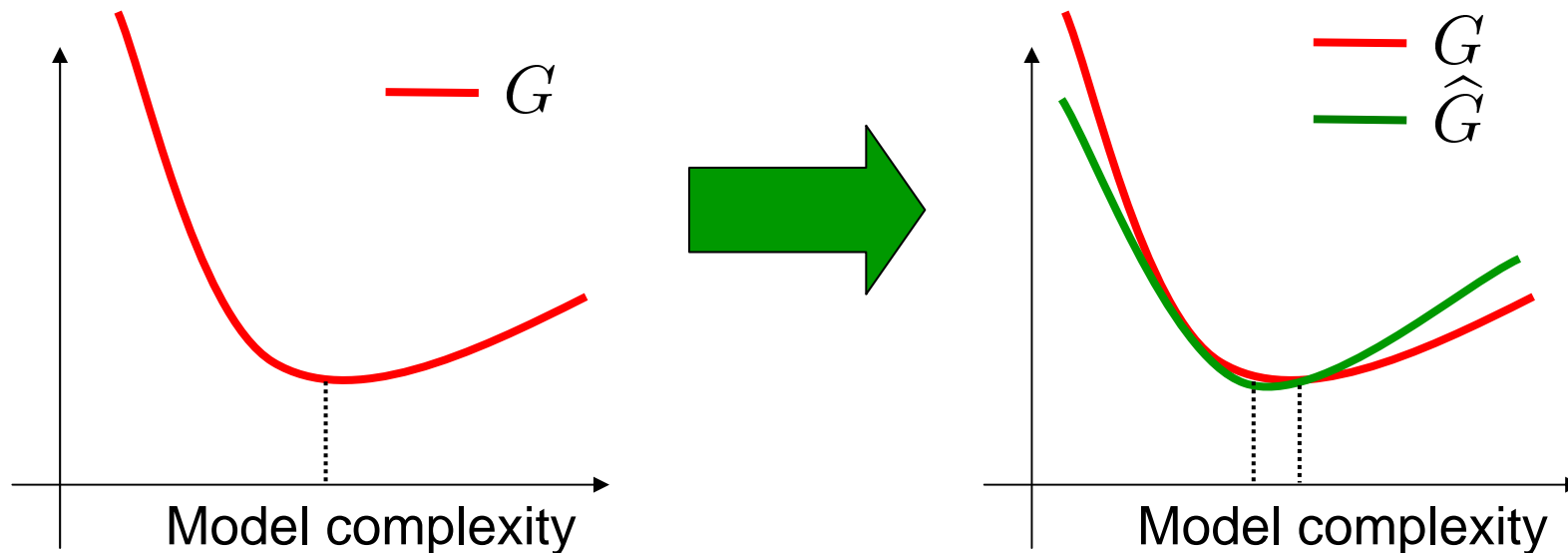
$$t \sim q(x)$$

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

# Generalization Error Estimation<sup>77</sup>

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

- Instead, we use a generalization error estimate.



# Model Selection

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

$$\{\mathcal{M}_i \mid \mathcal{M}_i = (c_i, \lambda_i)\}$$

- Estimate generalization error for each model.

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

- Choose the one with the **minimum estimated generalization error**.

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

# Assumptions

- Training input points:  $\mathbf{x}_i \stackrel{i.i.d.}{\sim} q(\mathbf{x})$
- Training output values:  $y_i = f(\mathbf{x}_i) + \epsilon_i$
- Noise  $\epsilon_i$  : i.i.d., mean 0, variance  $\sigma^2$

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

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

# Extra-Sample Method

- Suppose we have an **extra example**  $(x', y')$  in addition to  $\{(x_i, y_i)\}_{i=1}^n$ .

$$x' \sim q(x) \quad y' = f(x') + \epsilon' \quad \mathbb{E}_{\epsilon}[\epsilon'] = 0$$

$$\mathbb{E}_{\epsilon}[\epsilon'^2] = \sigma^2$$

$$\mathbb{E}_{\epsilon}[\epsilon' \epsilon_i] = 0, \quad \forall i$$

- Test the prediction performance of the learned function using the extra example.

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

$$\hat{f} \leftarrow \{(x_i, y_i)\}_{i=1}^n$$



# Extra-Sample Method (cont.) 81

$\hat{G}_{extra}$  is **unbiased** w.r.t.  $x'$  and  $\epsilon'$  (up to  $\sigma^2$ )

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

## ■ Proof:

$$\begin{aligned} & \mathbb{E}_{x'} \mathbb{E}_{\epsilon'} \left( \hat{f}(x') - f(x') - \epsilon' \right)^2 \\ &= \mathbb{E}_{x'} \mathbb{E}_{\epsilon'} \left[ (\hat{f}(x') - f(x'))^2 - 2\epsilon'(\hat{f}(x') - f(x')) + \epsilon'^2 \right] \\ &= G + \sigma^2 \end{aligned}$$

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

# Holdout Method

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

1. Divide training set  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  into  $\{(\mathbf{x}_i, y_i)\}_{i \neq j}$  and  $(\mathbf{x}_j, y_j)$ .
2. Train a learning machine using  $\{(\mathbf{x}_i, y_i)\}_{i \neq j}$   
$$\hat{f}_j(\mathbf{x}) \leftarrow \{(\mathbf{x}_i, y_i)\}_{i \neq j}$$
3. Test its prediction performance using the holdout sample  $(\mathbf{x}_j, y_j)$ :

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

# Almost Unbiasedness of Holdout<sup>83</sup>

- Holdout method is almost unbiased w.r.t.

$\mathbf{x}_j, \epsilon_j$  :

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

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

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

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

# Leave-One-Out Cross-Validation<sup>84</sup>

- 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.  $\{\mathbf{x}_i, \epsilon_i\}_{i=1}^n$

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

# k-Fold Cross-Validation

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

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

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

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

- k-fold is easier to compute and more stable.

# Advantages of CV

86

- **Wide applicability:** Almost unbiasedness of LOOCV holds for (virtually) any learning methods
- **Practical usefulness:** CV is 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

- Linear model

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

- Quadratically constrained least-squares

$$J_{QCLS}(\boldsymbol{\alpha}) = J_{LS}(\boldsymbol{\alpha}) + \lambda \|\boldsymbol{\alpha}\|^2$$

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

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

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



# Homework (cont.)

1. (Try to) prove the closed-form expression of leave-one-out cross-validation score for 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 fact:**

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

# Homework (cont.)

90

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

# Suggestions

- Please look for software which can solve
  - Linearly constrained quadratic programming

$$\min_{\beta} \left[ \frac{1}{2} \langle Q\beta, \beta \rangle + \langle \beta, q \rangle \right]$$

subject to  $V\beta \leq v$  and  $G\beta = g$

- Linearly constrained linear programming

$$\min_{\beta} \langle \beta, q \rangle \quad \text{subject to } V\beta \leq v \text{ and } G\beta = g$$

- For example, MOSEK, LOQO, or SeDuMi.
- The software does not have to be sophisticated; just an elementary one is enough.