

Selecting Ridge Parameters in Infinite Dimensional Hypothesis Spaces

Masashi Sugiyama

Tokyo Institute of Technology, Tokyo, Japan

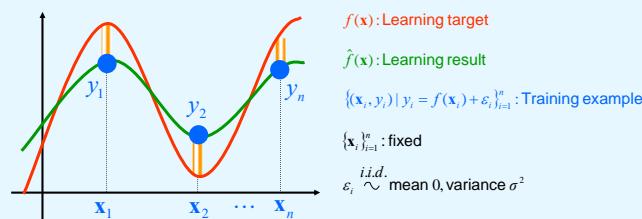
Klaus-Robert Müller

Fraunhofer FIRST and University of Potsdam, Berlin, Germany

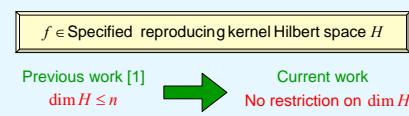
■ Introduction

- Properly tuning ridge parameter is crucial for better generalization
- Usually generalization error estimator is first derived, and ridge parameter is tuned so that estimated generalization error is minimized
- Most of generalization error estimators proposed so far are derived within asymptotic setting (i.e., large sample assumption)
- However, small sample case is of high practical importance
- We derive an exact unbiased estimator of generalization error: **SIC**
- Unbiasedness of SIC is guaranteed even with small sample cases

■ Function Approximation Problem



■ Assumption



[1] Sugiyama, M. & Ogawa, H.: Subspace information criterion for model selection. *Neural Computation*, vol.13, no.8, pp.1863-1889, 2001.

■ Generalization Measure

$$J_G = E \|\hat{f} - f\|^2$$

E : Expectation over noise
 $\|\cdot\|$: Norm in RKHS H

■ Kernel Ridge Regression

$$\text{Regression model: } \hat{f}(x) = \sum_{i=1}^n \alpha_i K(x, x_i)$$

$$\min_{\alpha} \left[\sum_{i=1}^n \left(\sum_{j=1}^n \alpha_j K(x_i, x_j) - y_i \right)^2 + \lambda \sum_{i=1}^n \alpha_i^2 \right]$$

λ : Ridge parameter

$$\mathbf{a} = \mathbf{Xy}$$

$$\mathbf{X} = (\mathbf{K}^2 + \lambda \mathbf{I})^{-1} \mathbf{K}$$

$$\mathbf{a} = (\alpha_1, \alpha_2, \dots, \alpha_n)^T$$

$$\mathbf{y} = (y_1, y_2, \dots, y_n)^T$$

$$\mathbf{K}_y = K(x_i, x_j)$$

$$\mathbf{I}$$
: Identity matrix

■ Extracting Essential Part of Generalization Error

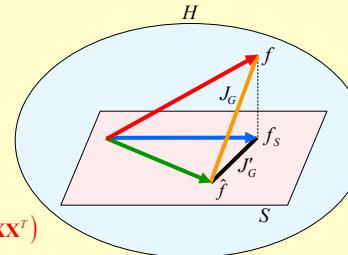
$$J_G = E \|\hat{f} - f\|^2$$

$$= E \underbrace{\|\hat{f} - f_s\|^2}_{\text{Essential}} + \underbrace{\|f_s - f\|^2}_{\text{Constant}} = J'_G$$

$$J'_G = E \|\hat{f} - f_s\|^2$$

$$= \underbrace{\|E\hat{f} - f_s\|^2}_{\text{Bias}} + E \underbrace{\|\hat{f} - E\hat{f}\|^2}_{\text{Variance}} = \sigma^2 \text{trace}(\mathbf{KXX}^T)$$

S : Subspace spanned by $\{K(x, x_i)\}_{i=1}^n$
 f_s : Orthogonal projection of f onto S



■ Estimation of Bias

If unbiased estimate \hat{f}_s^u of f_s is available, (i.e., $E\hat{f}_s^u = f_s$)

$$\text{Bias} = \|\hat{f} - \hat{f}_s^u\|^2 = \|g\|^2 + 2\langle Eg, Eg - g \rangle - \|Eg - g\|^2$$

$$\widehat{\text{Bias}} = \|g\|^2 - 0 - E\|Eg - g\|^2$$

Theorem

Unbiased estimate \hat{f}_s^u of f_s is given by

$$\hat{f}_s^u(x) = \sum_{i=1}^n \alpha_i^u K(x, x_i) \quad \mathbf{a}^u = (\alpha_1^u, \alpha_2^u, \dots, \alpha_n^u)^T = \mathbf{K}^+ \mathbf{y}$$

\mathbf{K}^+ : Generalized Inverse

$$\widehat{\text{Bias}} = \mathbf{y}^T (\mathbf{X} - \mathbf{K}^+) \mathbf{K} (\mathbf{X} - \mathbf{K}^+) \mathbf{y} + \sigma^2 \text{trace}(\mathbf{K} (\mathbf{X} - \mathbf{K}^+)^2)$$

■ Subspace Information Criterion (SIC)

$$\text{SIC} = \widehat{\text{Bias}} + \text{Variance} = \mathbf{y}^T (\mathbf{X} - \mathbf{K}^+) \mathbf{K} (\mathbf{X} - \mathbf{K}^+) \mathbf{y} + \sigma^2 \text{trace}(2\mathbf{K}^+ \mathbf{K} \mathbf{X} - \mathbf{K}^+)$$

→ SIC is an unbiased estimator of J'_G with finite samples

$$ESIC = J'_G$$

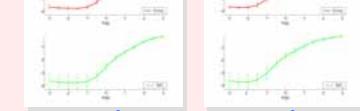
NOTE: Unbiasedness holds even without taking expectation over training input points $\{x_i\}_{i=1}^n$

■ Computer Simulations (Toy datasets)

- Gaussian RKHS: $K(x, x) = \exp(-(x - x')^2 / 2)$

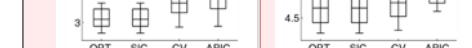
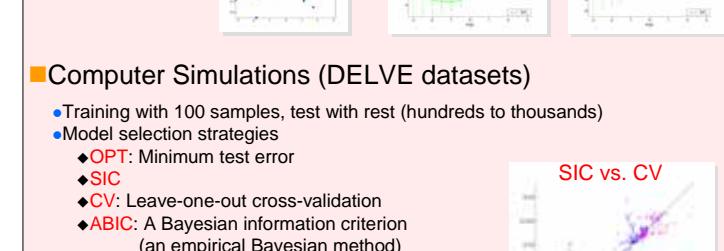
$$n = 100, \sigma^2 = 0.01$$

$$n = 50, \sigma^2 = 0.01$$



■ Computer Simulations (DELVE datasets)

- Training with 100 samples, test with rest (hundreds to thousands)



- SIC works very well for most datasets
- However its performance can be degraded with very large noise (datasets specified by "fh" or "nh")