

On the Influence of Input Noise on a Generalization Error Estimator



Masashi Sugiyama ^(1,2)

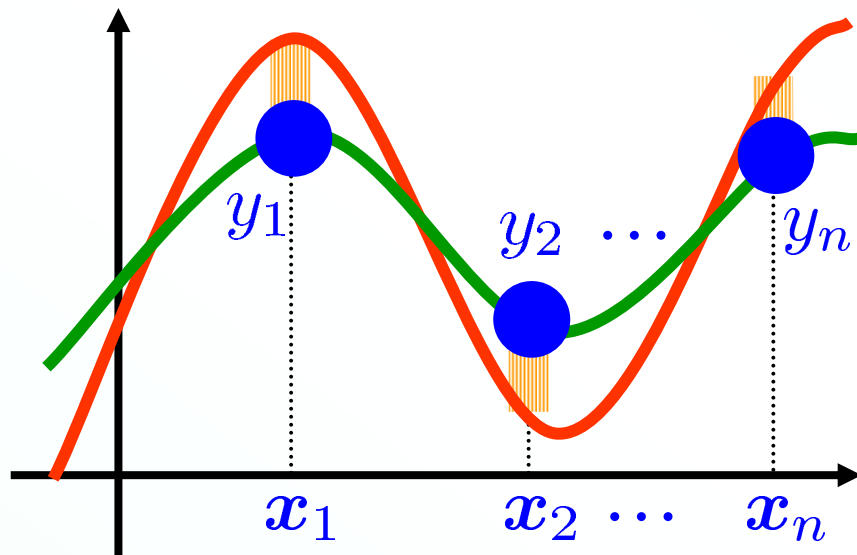
Yuta Okabe ⁽²⁾

Hidemitsu Ogawa ⁽²⁾

⁽¹⁾ Fraunhofer FIRST-IDA, Berlin, Germany

⁽²⁾ Tokyo Institute of Technology, Tokyo, Japan

Regression Problem



$f(\mathbf{x})$: Underlying function

$\hat{f}(\mathbf{x})$: Learned function

$\{(\mathbf{x}_i, y_i)\}_{i=1}^n$: Training examples

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

(noise)

$\epsilon_i \stackrel{i.i.d.}{\sim}$ mean 0, variance σ^2

From $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, obtain a good approximation $\hat{f}(\mathbf{x})$ to $f(\mathbf{x})$

Typical Method of Learning

3

■ Kernel regression model

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

α_i : Parameters to be learned

$K(\mathbf{x}, \mathbf{x}')$: Kernel function (e.g., Gaussian)

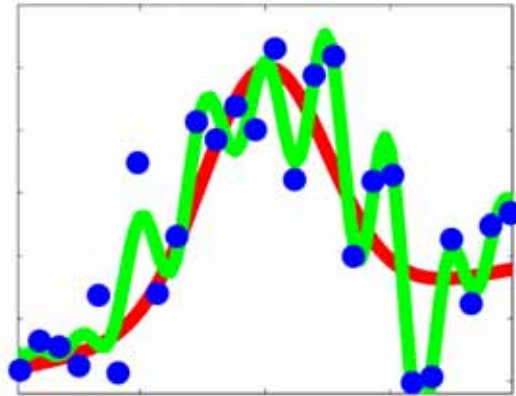
■ Ridge estimation

$$\min_{\{\alpha_i\}} \left[\sum_{i=1}^n \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 + \lambda \sum_{i=1}^n \alpha_i^2 \right]$$

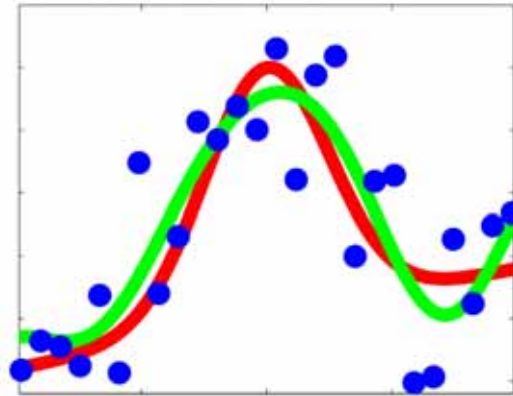
λ : Ridge parameter (model parameter)

Model Selection

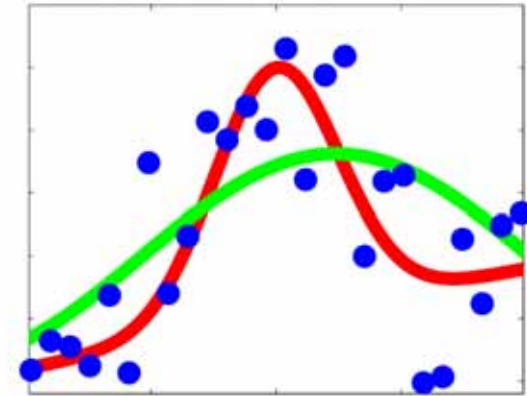
— Underlying function $f(x)$
— Learned function $\hat{f}(x)$



λ is too small



λ is appropriate



λ is too large

Choice of the model is crucial
for obtaining good learned function $\hat{f}(x)$!

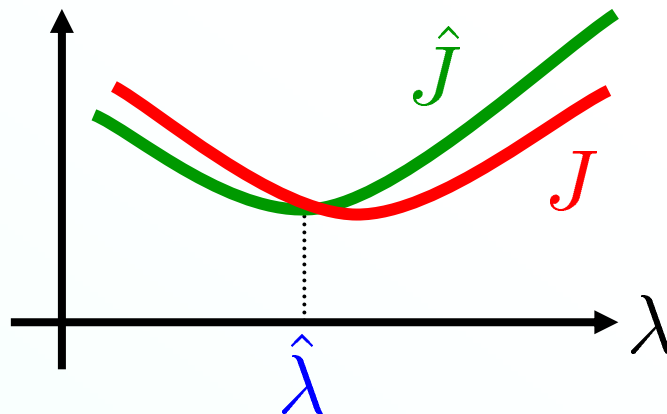
Generalization Error

For model selection, we need a criterion that measures 'closeness' between $\hat{f}(x)$ and $f(x)$:

➔ Generalization error

Determine the model λ so that an estimator \hat{J} of the unknown generalization error J is minimized.

$$\hat{\lambda} = \underset{\lambda}{\operatorname{argmin}} \hat{J}(\lambda)$$

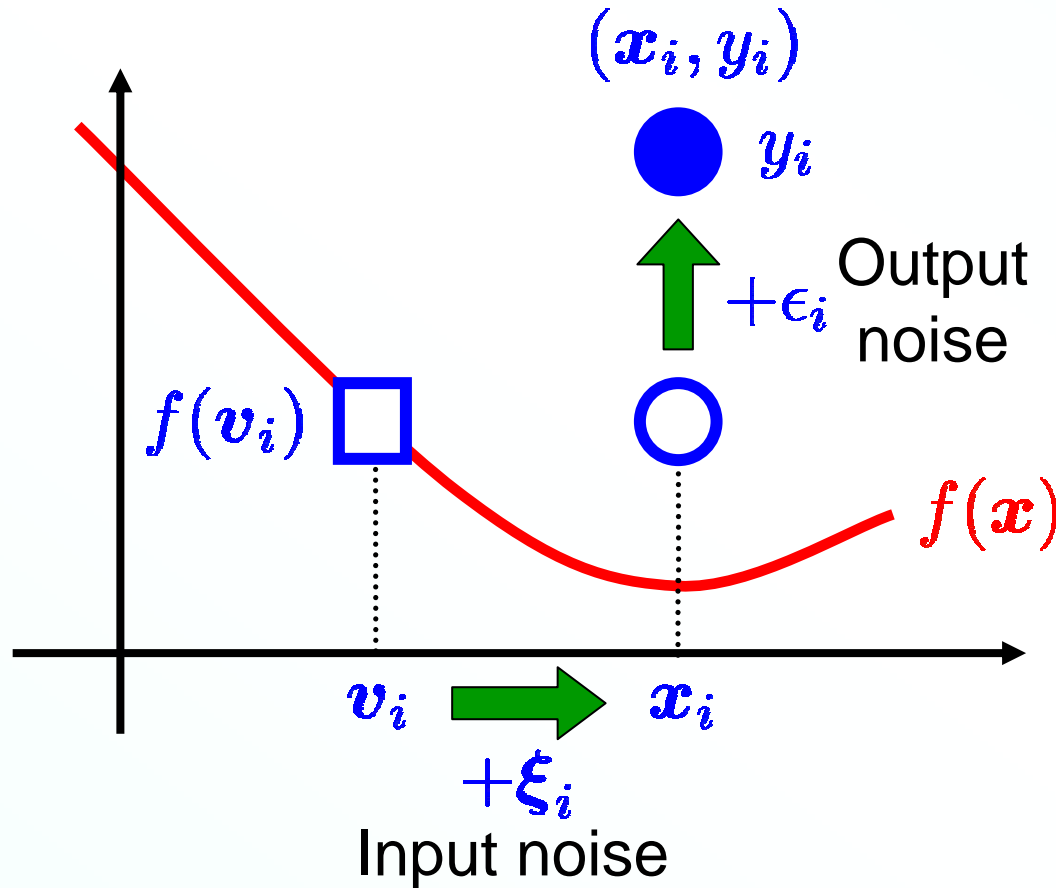


Noise in Input Points

- Previous research mainly deals with the cases where noise is included only in output values.
- However, noise is sometimes included also in input points, e.g.,
 - **Input points are measured:**
Signal/image recognition, robot motor control, and bioinformatic data analysis.
 - **Input points are estimated:**
Time series prediction of multiple-step ahead.

Noise in Input Points (cont.)

7



- We want to measure output values $f(x_i)$ at x_i
- But measurement is actually done at unknown v_i
- Output noise ϵ_i is then added

$$x_i = v_i + \xi_i$$

$$y_i = f(v_i) + \epsilon_i$$



Aim of Our Research

- So far, it seems that model selection in the presence of input noise has not been well studied yet.
- We investigate the effect of input noise on a generalization error estimator called the **subspace information criterion (SIC)**.

Sugiyama & Ogawa (Neural Computation, 2001)

Sugiyama & Müller (JMLR, 2002)

Generalization Error in RKHS

- \mathcal{H} : A reproducing kernel Hilbert space
- We assume $f, \hat{f} \in \mathcal{H}$
- We shall measure the generalization error by

$$J = \mathbb{E}_{\epsilon} \|\hat{f} - f\|^2 - \|f\|^2$$

\mathbb{E}_{ϵ} : Expectation over output noise

$\|\cdot\|$: Norm

■ Kernel regression model

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

α_i : Parameters to be learned

$K(\mathbf{x}, \mathbf{x}')$: Kernel function (e.g., Gaussian)

■ Linear estimation

$$\boldsymbol{\alpha} = \mathbf{X} \mathbf{y}$$

\mathbf{X} : Learning matrix

$$\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_n)^\top$$

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

Subspace Information Criterion

11

Sugiyama & Ogawa (Neural Computation, 2001)

Sugiyama & Müller (JMLR, 2002)

$$SIC = \langle \mathbf{K} \mathbf{X} \mathbf{y}, \mathbf{X} \mathbf{y} \rangle - 2 \langle \mathbf{K} \mathbf{X} \mathbf{y}, \mathbf{K}^\dagger \mathbf{y} \rangle + 2\sigma^2 \text{tr}(\mathbf{K}^\dagger \mathbf{K} \mathbf{X})$$

$$K_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$$

 \mathbf{K}^\dagger : Pseudo inverse of \mathbf{K} $\langle \cdot, \cdot \rangle$: Inner product

- In the absence of input noise, SIC is an unbiased estimator of J :

$$E_\epsilon SIC = J \quad J = E_\epsilon \|\hat{f} - f\|^2 - \|f\|^2$$

- We investigate how the unbiasedness of SIC is affected by input noise.

Unbiasedness of SIC in the Presence of Input Noise

- In the presence of input noise,

$$E_{\epsilon} SIC = J + \Delta J$$

$$\Delta J = \langle K^{\dagger} K X z, z_x - z \rangle$$

$$z = (f(v_1), f(v_2), \dots, f(v_n))^{\top} \quad v_i : \text{Noiseless input points}$$

$$z_x = (f(x_1), f(x_2), \dots, f(x_n))^{\top} \quad x_i : \text{Noisy input points}$$

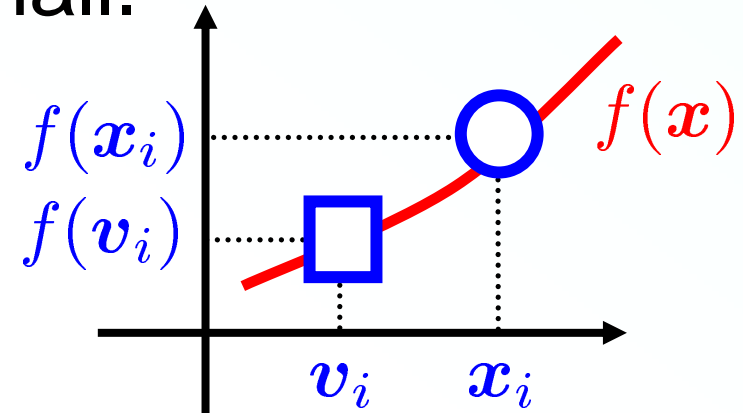
Unbiasedness of SIC does not generally hold in the presence of input noise.

Effect of Small Input Noise

- When $f(\mathbf{x})$ is continuous, small input noise varies the output value only slightly, i.e., $|f(\mathbf{x}_i) - f(\mathbf{v}_i)|$ is small.

\mathbf{v}_i : Noiseless input points

\mathbf{x}_i : Noisy input points



- Therefore, we expect that the unbiasedness of SIC is not severely affected (ΔJ is small) by small input noise.

$$E_{\epsilon} SIC = J + \Delta J$$

Effect of Small Input Noise (cont.)¹⁴

- However, we can show that, for some learning matrix X , it holds that

$$|\Delta J| \not\rightarrow 0 \text{ as } \|\xi_i\| \rightarrow 0 \text{ for all } i .$$

ξ_i : Input noise

- This implies that, for some X , the unbiasedness of SIC is heavily affected even when input noise is very small.

Theorem

- Let $\|\mathbf{X}\|$ be the matrix norm defined by

$$\|\mathbf{X}\| = \sup_{\mathbf{z} \neq 0} \frac{\|\mathbf{X}\mathbf{z}\|}{\|\mathbf{z}\|}$$

- If the learning matrix \mathbf{X} satisfies

$$\|\mathbf{X}\| = o(1/\delta) \quad \delta = \|\mathbf{z}_x - \mathbf{z}\|$$

then $|\Delta J| \rightarrow 0$ as $\|\xi_i\| \rightarrow 0$ for all i .

$$\begin{aligned} \mathbf{z} &= (f(\mathbf{v}_1), f(\mathbf{v}_2), \dots, f(\mathbf{v}_n))^{\top} & \mathbf{v}_i &: \text{Noiseless input points} \\ \mathbf{z}_x &= (f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_n))^{\top} & \mathbf{x}_i &: \text{Noisy input points} \end{aligned}$$

Ridge Estimation

■ Ridge estimation

λ : Ridge parameter

$$\min_{\{\alpha_i\}} \left[\sum_{i=1}^n \left(\hat{f}(\mathbf{x}) - y_i \right)^2 + \lambda \sum_{i=1}^n \alpha_i^2 \right]$$

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

$\mathbf{K}_{i,j} = K(\mathbf{x}_i, \mathbf{x}_j)$
 \mathbf{I} : Identity matrix

■ We can prove that ridge estimation satisfies

$$\|\mathbf{X}\| = O(1) = o(1/\delta)$$

■ Therefore, SIC with ridge estimation is robust against small input noise.

Simulation

- H : Gaussian RKHS

$$K(x, x') = \exp \left(-(x - x')^2 / 2 \right)$$

- Learning target function $f(x)$: sinc function

- Training examples $\{(x_i, y_i)\}_{i=1}^n$:

$$v_i \stackrel{i.i.d.}{\sim} U(-\pi, \pi)$$

$$x_i = v_i + \xi_i, \quad \xi_i \stackrel{i.i.d.}{\sim} N(0, \sigma_x^2)$$

$$y_i = f(x_i) + \epsilon_i, \quad \epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

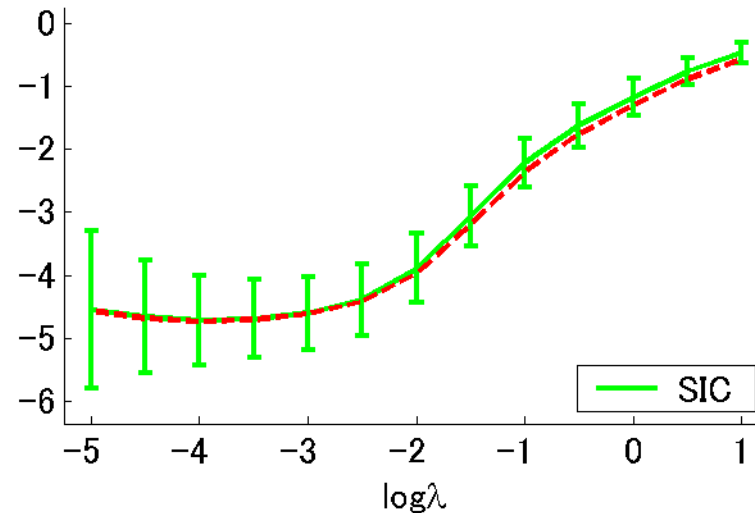
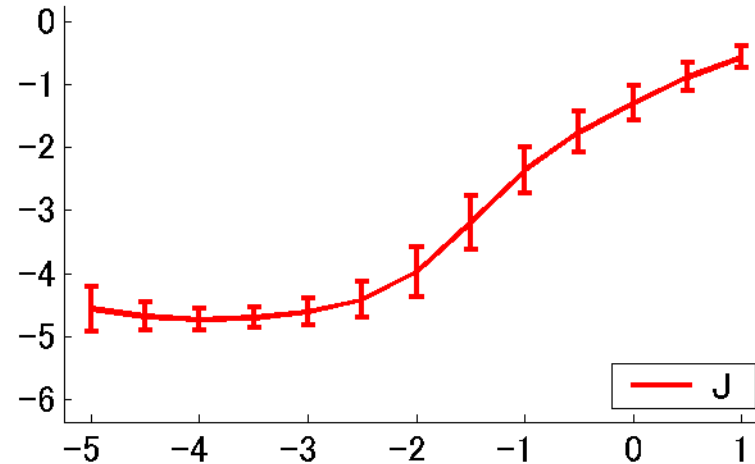
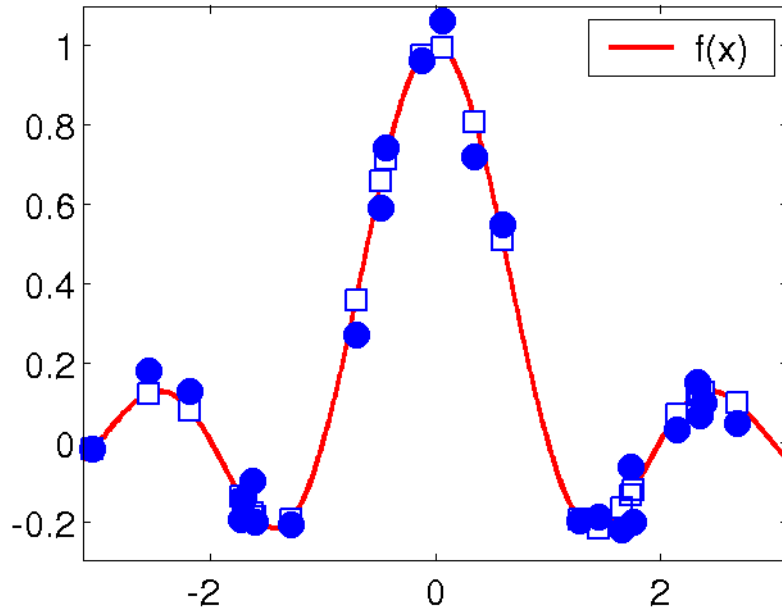
- $n = 25$, $\sigma = 0.05$, $\sigma_x = 0, 0.1, 0.2$

- Ridge estimation is used for learning.

Result (No Input Noise)

18

$\sigma_x = 0$



■ SIC is surely unbiased without input noise

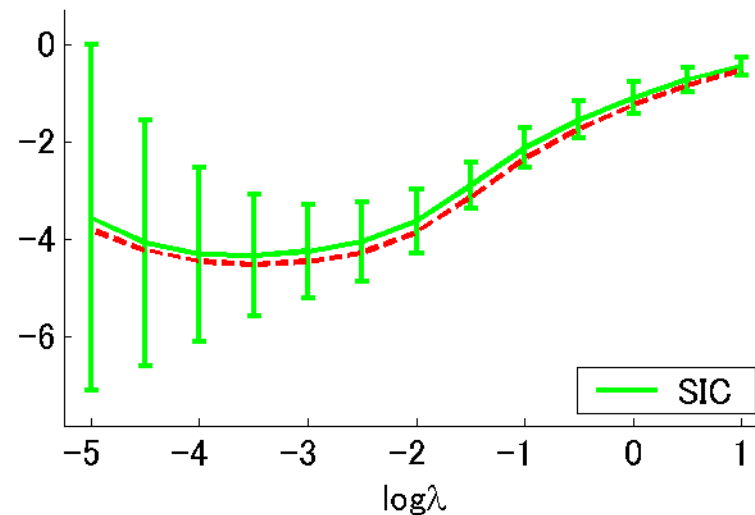
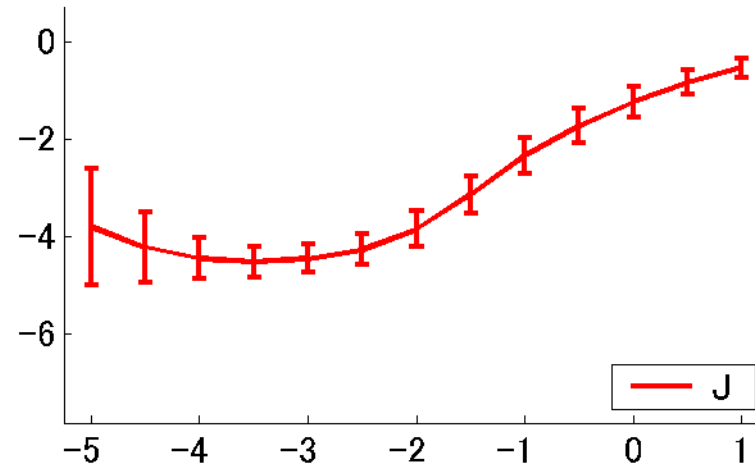
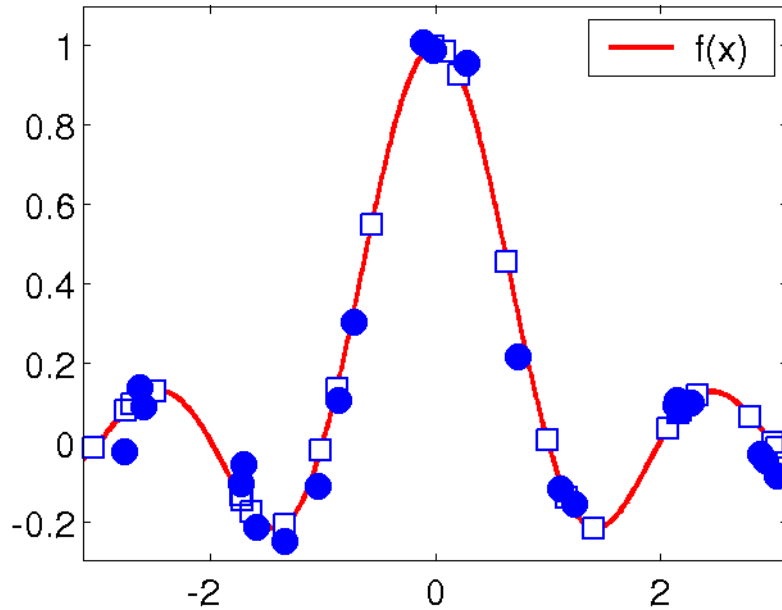
$$J = \mathbb{E}_{\epsilon} \|\hat{f} - f\|^2 - \|f\|^2$$

λ : Ridge parameter

Result (Small Input Noise)

19

$\sigma_x = 0.1$



■ SIC is still almost unbiased with small input noise

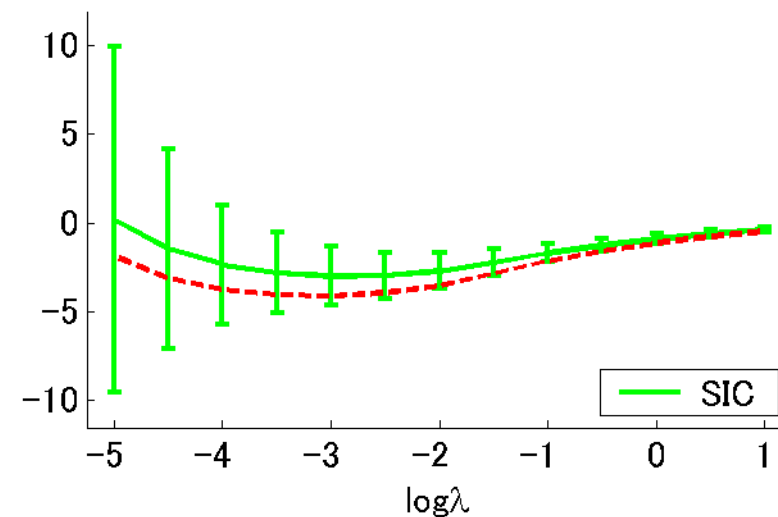
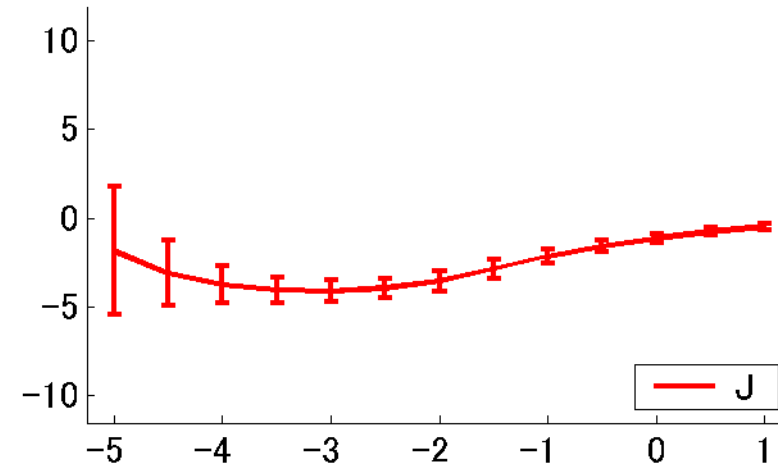
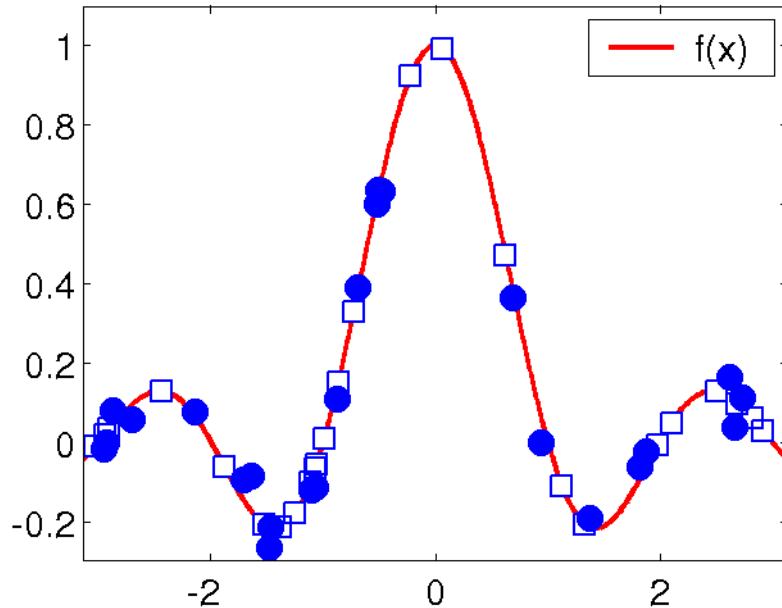
$$J = \mathbb{E}_\epsilon \|\hat{f} - f\|^2 - \|f\|^2$$

λ : Ridge parameter

Result (Large Input Noise)

20

$$\sigma_x = 0.2$$



- SIC is no longer reliable with large input noise

$$J = \mathbb{E}_\epsilon \|\hat{f} - f\|^2 - \|f\|^2$$

λ : Ridge parameter



Conclusions

- Effect of input noise on SIC.
- In some cases, the unbiasedness of SIC is heavily affected even by small input noise.
- A sufficient condition for unbiasedness.
- Ridge estimation satisfies this condition.
- Experiments: SIC is still almost unbiased for small input noise.
- Future work: Accurately estimate the generalization error when input noise is large.