

Pattern Information Processing:²² Properties of Least-Squares

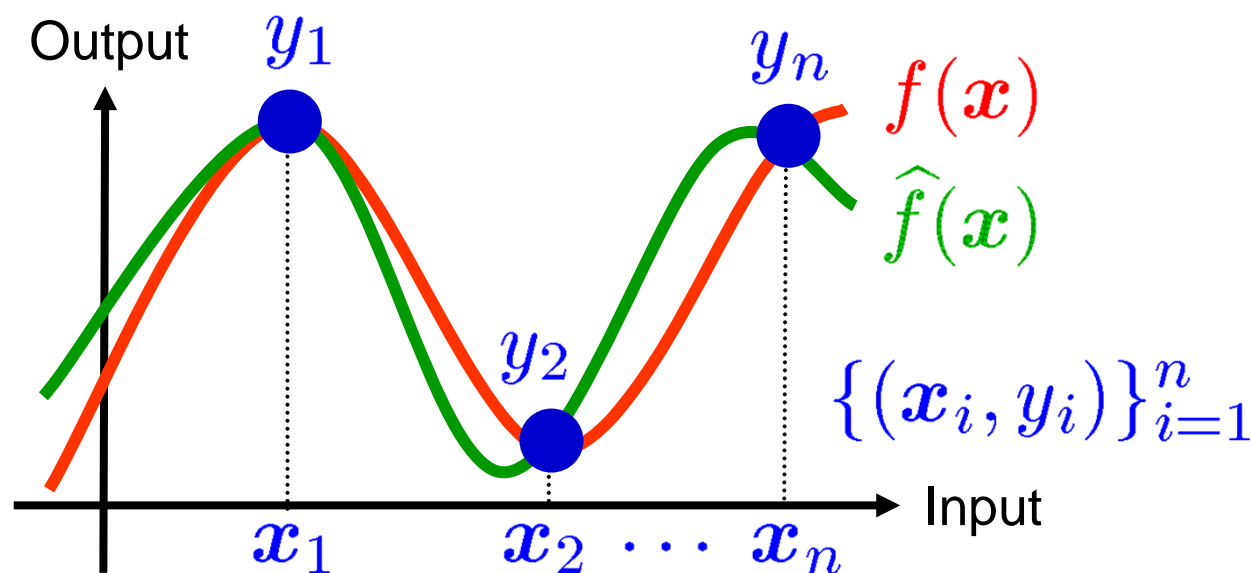
Masashi Sugiyama
(Department of Computer Science)

Contact: W8E-505

sugi@cs.titech.ac.jp

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

Supervised Learning as Function Approximation



Using training examples $\{(x_i, y_i)\}_{i=1}^n$,
find a function $\hat{f}(x)$ from a model \mathcal{M}
that well approximates the target function $f(x)$.

Assumptions

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

- Training inputs x_i : **i.i.d.** from a probability distribution with density $q(x)$
- Training outputs y_i : additive noise included

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

- Output noise ϵ_i : **i.i.d.** with mean zero

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

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

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

■ Least-squares learning for a linear model:

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

$$\hat{\boldsymbol{\alpha}}_{LS} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} J_{LS}(\boldsymbol{\alpha})$$

$$J_{LS}(\boldsymbol{\alpha}) = \sum_{i=1}^n \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2$$

■ Solution: $\hat{\boldsymbol{\alpha}}_{LS} = \mathbf{L}_{LS} \mathbf{y}$

$$\mathbf{L}_{LS} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$$

$$X_{i,j} = \varphi_j(\mathbf{x}_i)$$

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

Today's Plan

- How does LS contribute to reducing the generalization error (i.e, expected prediction error for all test input points)?

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

- Justification of LS for linear models:
 - Realizable cases
 - Unrealizable cases

Realizability

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

- **Realizable:** Learning target function $f(\mathbf{x})$ can be expressed by the model, i.e., there exists a parameter vector $\alpha^* = (\alpha_1^*, \alpha_2^*, \dots, \alpha_b^*)^\top$ such that

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

- **Unrealizable:** $f(\mathbf{x})$ is not realizable

Justification in Realizable Cases²⁸

- In realizable cases, generalization error is expressed as

$$G = \int_{\mathcal{D}} \left(\hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 q(\mathbf{x}) d\mathbf{x}$$
$$= \|\hat{\alpha} - \alpha^*\|_U^2$$

$$\|\alpha\|_U^2 = \langle U\alpha, \alpha \rangle$$

$$U_{i,j} = \int_{\mathcal{D}} \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}$$

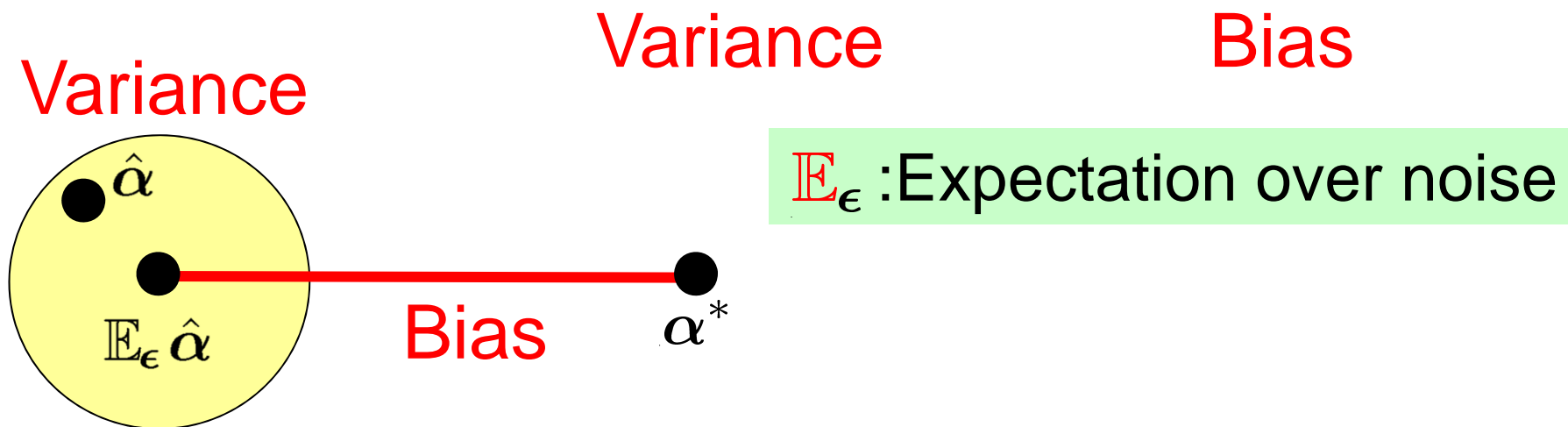
Bias/Variance Decomposition ²⁹

■ Expected generalization error:

$$\mathbb{E}_{\epsilon} [G] = \mathbb{E}_{\epsilon} \|\hat{\alpha} - \alpha^*\|_U^2$$

$$= \mathbb{E}_{\epsilon} \|\hat{\alpha} - \mathbb{E}_{\epsilon} \hat{\alpha} + \mathbb{E}_{\epsilon} \hat{\alpha} - \alpha^*\|_U^2$$

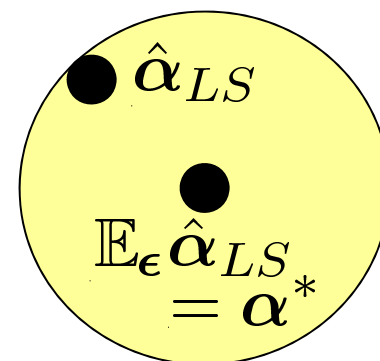
$$= \underbrace{\mathbb{E}_{\epsilon} \|\hat{\alpha} - \mathbb{E}_{\epsilon} \hat{\alpha}\|_U^2}_{\text{Variance}} + \underbrace{\|\mathbb{E}_{\epsilon} \hat{\alpha} - \alpha^*\|_U^2}_{\text{Bias}}$$



Unbiasedness

- When $f(x)$ is realizable, $\hat{\alpha}_{LS}$ is an **unbiased estimator**:

$$\mathbb{E}_{\epsilon}[\hat{\alpha}_{LS}] = \alpha^*$$



- **Proof:** In realizable cases,

$$y = X\alpha^* + \epsilon$$

Then

$$\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$$

$$\mathbb{E}_{\epsilon}[\hat{\alpha}_{LS}] = \mathbb{E}_{\epsilon}(X^\top X)^{-1} X^\top y$$

$$= (X^\top X)^{-1} X^\top (X\alpha^* + \mathbb{E}_{\epsilon}[\epsilon])$$

$$= \alpha^*$$

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

Best Linear Unbiased Estimator³¹

- $\hat{\alpha}_{LS}$ is the **best linear unbiased estimator** (BLUE; a linear estimator that has the smallest variance among all linear unbiased estimators).

$$\mathbb{E}_{\epsilon} \|\hat{\alpha}_{LS} - \mathbb{E}_{\epsilon} \hat{\alpha}_{LS}\|_U^2$$

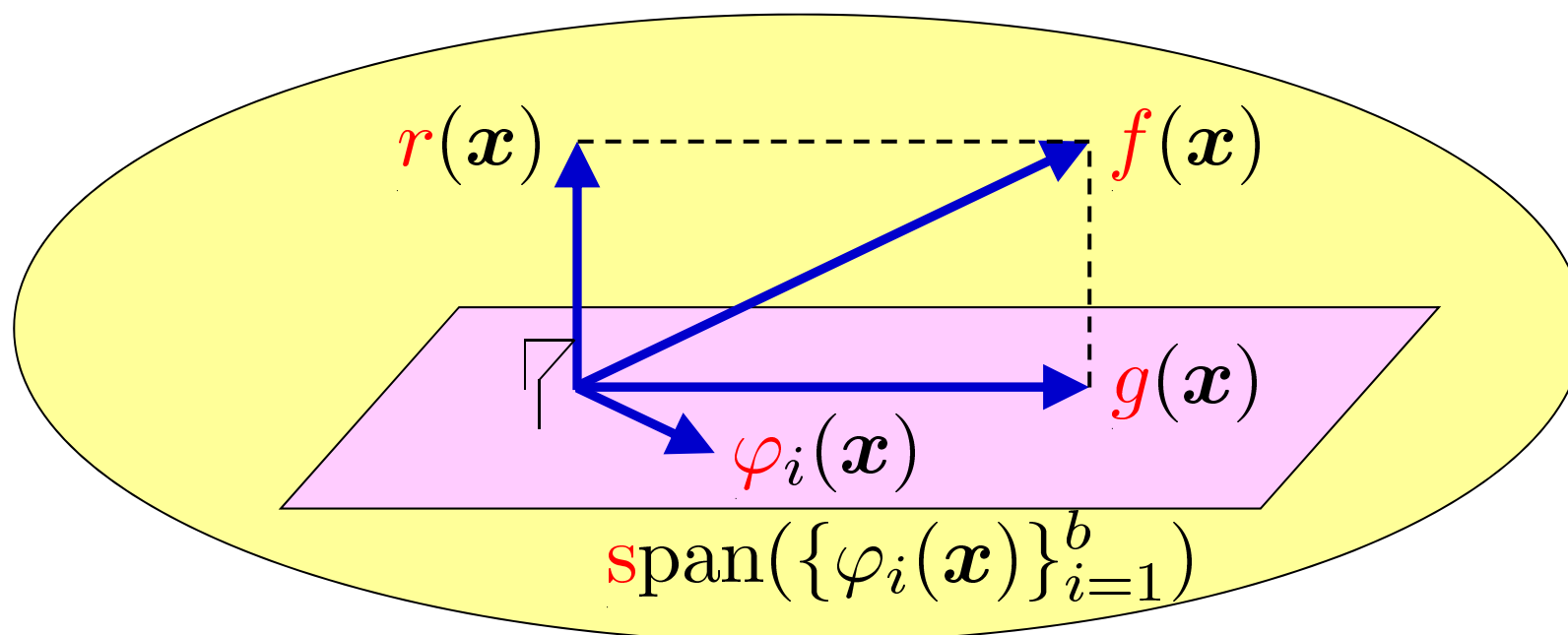
$$\leq \mathbb{E}_{\epsilon} \|\hat{\alpha}_{LU} - \mathbb{E}_{\epsilon} \hat{\alpha}_{LU}\|_U^2$$

for any linear unbiased estimator $\hat{\alpha}_{LU}$

- **Proof:** Homework!

Justification of LS (Unrealizable Cases)

■ Decomposition: $f(\mathbf{x}) = g(\mathbf{x}) + r(\mathbf{x})$



$$\int_{\mathcal{D}} \varphi_i(\mathbf{x}) r(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} = 0$$

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

Generalization Error Decomposition³³

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

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

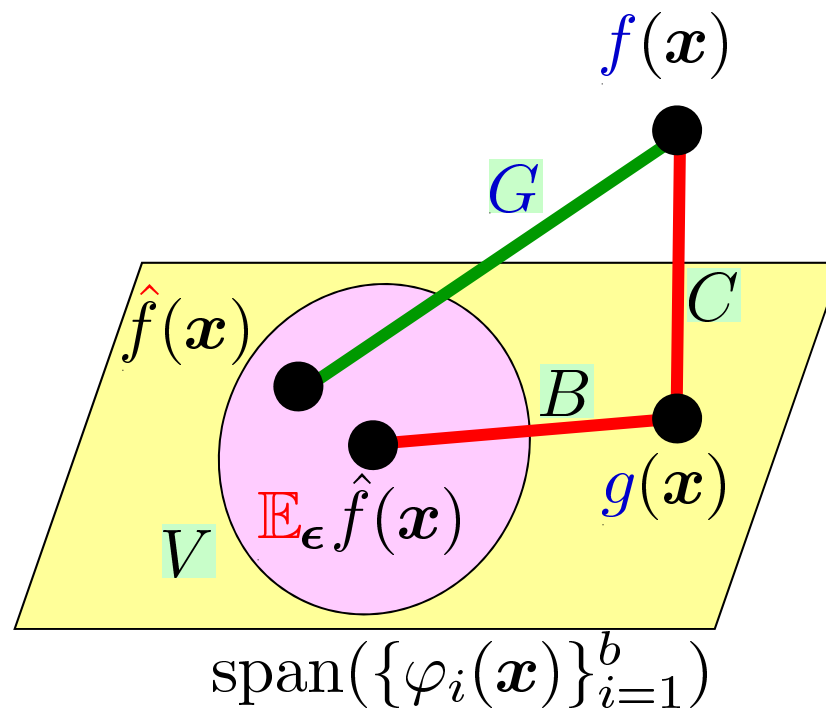
$$= \int_{\mathcal{D}} \left(\hat{f}(\mathbf{x}) - g(\mathbf{x}) \right)^2 q(\mathbf{x}) d\mathbf{x} + \int_{\mathcal{D}} r(\mathbf{x})^2 q(\mathbf{x}) d\mathbf{x}$$

$$= \|\hat{\boldsymbol{\alpha}} - \boldsymbol{\alpha}^*\|_U^2 + C$$

$$C = \int_{\mathcal{D}} r(\mathbf{x})^2 q(\mathbf{x}) d\mathbf{x}$$

Bias/Variance Decomposition ³⁴

$$\mathbb{E}_\epsilon [G] = \underbrace{\mathbb{E}_\epsilon \|\hat{\alpha} - \mathbb{E}_\epsilon \hat{\alpha}\|_U^2}_{\text{Variance}} + \underbrace{\|\mathbb{E}_\epsilon \hat{\alpha} - \alpha^*\|_U^2}_{\text{Bias}} + \underbrace{C}_{\text{Model error}}$$



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

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

Asymptotic Unbiasedness

- $\hat{\alpha}_{LS}$ is an **asymptotically unbiased estimator** of the optimal parameter α^* :

$$\mathbb{E}_{\epsilon}[\hat{\alpha}_{LS}] \rightarrow \alpha^* \text{ as } n \rightarrow \infty$$

- Proof:

- $y = X\alpha^* + z_r + \epsilon$ $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$

$$z_r = (r(x_1), r(x_2), \dots, r(x_n))^\top$$

- $$\begin{aligned}\mathbb{E}_{\epsilon}[\hat{\alpha}_{LS}] &= \mathbb{E}_{\epsilon}(X^\top X)^{-1} X^\top y \\ &= (X^\top X)^{-1} X^\top (X\alpha^* + z_r + \mathbb{E}_{\epsilon}\epsilon) \\ &= \alpha^* + \left(\frac{1}{n} X^\top X\right)^{-1} \frac{1}{n} X^\top z_r\end{aligned}$$

Proof (cont.)

- By the law of large numbers,

- $$\left[\frac{1}{n} \mathbf{X}^\top \mathbf{X}\right]_{i,j} = \frac{1}{n} \sum_{k=1}^n \varphi_i(\mathbf{x}_k) \varphi_j(\mathbf{x}_k)$$

$$\rightarrow \int_{\mathcal{D}} \varphi_i(\mathbf{x}) \varphi_j(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} = U_{i,j}$$

- $$\left[\frac{1}{n} \mathbf{X}^\top \mathbf{z}_r\right]_i = \frac{1}{n} \sum_{k=1}^n \varphi_i(\mathbf{x}_k) r(\mathbf{x}_k)$$

$$\rightarrow \int_{\mathcal{D}} \varphi_k(\mathbf{x}) r(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} = 0$$

- Thus, $\mathbb{E}_{\epsilon} [\hat{\boldsymbol{\alpha}}_{LS}] \rightarrow \boldsymbol{\alpha}^*$ as $n \rightarrow \infty$

(Q.E.D.)

Efficiency

37

- **The Cramér-Rao lower bound:** Lower bound of the variance of all (possibly non-linear) unbiased estimators.
- **Efficient estimator:** An unbiased estimator whose variance attains the Cramér-Rao bound.
- For linear model with LS and $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, the Cramér-Rao bound is

$$\sigma^2 \text{tr}(\mathbf{U}(\mathbf{X}^\top \mathbf{X})^{-1})$$

Asymptotic Efficiency

- **Asymptotically efficient estimator:** An asymptotically unbiased estimator that attains the Cramér-Rao lower bound asymptotically.
- When $\epsilon_i \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$, LS estimator is asymptotically efficient.
 - **Proof:** LS estimator is asymptotically unbiased and

$$\begin{aligned}\mathbb{E}_{\epsilon} \|\hat{\alpha}_{LS} - \mathbb{E}_{\epsilon} \hat{\alpha}_{LS}\|_U^2 &= \mathbb{E}_{\epsilon} \|\mathbf{L}_{LS} \epsilon\|_U^2 \\ &= \sigma^2 \text{tr}(\mathbf{U}(\mathbf{X}^\top \mathbf{X})^{-1})\end{aligned}$$

which is the Cramér-Rao lower bound.

Summary

- LS is unbiased in realizable cases.
- LS has the smallest variance among all linear unbiased estimators in realizable cases.
- However, the generalization error (i.e., the sum of bias and variance) is not necessarily minimized.

min variance
subject to bias=0

min bias+variance

- Theoretical guarantees in unrealizable cases hold only asymptotically.

Homework

- Prove $\hat{\alpha}_{LS}$ is BLUE in realizable cases, i.e.,

$$\mathbb{E}_{\epsilon} \|\hat{\alpha}_{LS} - \mathbb{E}_{\epsilon} \hat{\alpha}_{LS}\|^2 \leq \mathbb{E}_{\epsilon} \|\hat{\alpha}_{LU} - \mathbb{E}_{\epsilon} \hat{\alpha}_{LU}\|^2$$

Hints:

- All linear unbiased estimator $\hat{\alpha}_{LU} = L_U y$ satisfies

$$\mathbb{E}_{\epsilon} [\hat{\alpha}_{LU}] = \alpha^*$$

Therefore, $L_U X = I$.

- By assumptions, noise satisfies

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

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

Announcement

41

- There will be no class next week (May 1st).
- The next class will be on May 8th.