Pattern Information Processing:²⁴ Properties of Least-Squares

Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Reviews

Linear models / kernel models:

$$\hat{f}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi_i(\boldsymbol{x}) \qquad \hat{f}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$$

Least-squares learning:

$$\hat{\boldsymbol{\alpha}}_{LS} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} J_{LS}(\boldsymbol{\alpha})$$
$$J_{LS}(\boldsymbol{\alpha}) = \sum_{i=1}^{n} \left(\hat{f}(\boldsymbol{x}_{i}) - y_{i} \right)^{2}$$

Today's Plan

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

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

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

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

Unrealizable: f(x) is not realizable

Justification in Realizable Cases²⁸

In realizable cases, generalization error is expressed as

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

$$\begin{aligned} \|\boldsymbol{\alpha}\|_{\boldsymbol{U}}^2 &= \langle \boldsymbol{U}\boldsymbol{\alpha}, \boldsymbol{\alpha} \rangle \\ U_{i,j} &= \int_{\mathcal{D}} \varphi_i(\boldsymbol{x}) \varphi_j(\boldsymbol{x}) q(\boldsymbol{x}) d\boldsymbol{x} \end{aligned}$$



Unbiasedness

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

$$(\|\mathbb{E}_{\boldsymbol{\epsilon}}\hat{\boldsymbol{\alpha}}_{LS} - \boldsymbol{\alpha}^*\|_{\boldsymbol{U}}^2 = 0$$



Proof: In realizable cases,

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

 $oldsymbol{y} = oldsymbol{X} oldsymbol{lpha}^* + oldsymbol{\epsilon}$ $oldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^ op$

Best Linear Unbiased Estimator³¹

• $\hat{\alpha}_{LS}$ is the best linear unbiased estimator (BLUE, a linear estimator which has the smallest variance among all linear unbiased estimators)

$$\begin{split} \mathbb{E}_{\boldsymbol{\epsilon}} \| \hat{\boldsymbol{\alpha}}_{LS} - \mathbb{E}_{\boldsymbol{\epsilon}} \hat{\boldsymbol{\alpha}}_{LS} \|_{\boldsymbol{U}}^2 \\ &\leq \mathbb{E}_{\boldsymbol{\epsilon}} \| \hat{\boldsymbol{\alpha}}_{LU} - \mathbb{E}_{\boldsymbol{\epsilon}} \hat{\boldsymbol{\alpha}}_{LU} \|_{\boldsymbol{U}}^2 \\ &\text{for any linear unbiased estimator } \hat{\boldsymbol{\alpha}}_{LU} \end{split}$$

Proof: Homework!

Justification of LS (Unrealizable Cases)

Decomposition: f(x) = g(x) + r(x)



Generalization Error Decomposition³³

$$\begin{aligned} G &= \int_{\mathcal{D}} \left(\hat{f}(\boldsymbol{x}) - f(\boldsymbol{x}) \right)^2 q(\boldsymbol{x}) d\boldsymbol{x} \\ &= \int_{\mathcal{D}} \left(\hat{f}(\boldsymbol{x}) - g(\boldsymbol{x}) - r(\boldsymbol{x}) \right)^2 q(\boldsymbol{x}) d\boldsymbol{x} \\ &= \int_{\mathcal{D}} \left(\hat{f}(\boldsymbol{x}) - g(\boldsymbol{x}) \right)^2 q(\boldsymbol{x}) d\boldsymbol{x} + \int_{\mathcal{D}} r(\boldsymbol{x})^2 q(\boldsymbol{x}) d\boldsymbol{x} \\ &= \|\boldsymbol{\alpha} - \boldsymbol{\alpha}^*\|_U^2 + C \end{aligned}$$

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



Asymptotic Unbiasedness

Assume $x_i \stackrel{i.i.d.}{\sim} q(x)$.
Then $\hat{\alpha}_{LS}$ is an asymptotically unbiased estimator of the optimal parameter α^* .

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

Proof

 $\boldsymbol{u} = X \alpha^* + \boldsymbol{z}_r + \boldsymbol{\epsilon}$ $\boldsymbol{z}_r = (r(\boldsymbol{x}_1), r(\boldsymbol{x}_2), \dots, r(\boldsymbol{x}_n))^\top$ $\boldsymbol{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$ $\mathbf{E}_{\boldsymbol{\epsilon}}[\hat{\boldsymbol{\alpha}}_{LS}] = \mathbb{E}_{\boldsymbol{\epsilon}}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}\boldsymbol{y}$ $= (\boldsymbol{X}^{\top}\boldsymbol{X})^{-1}\boldsymbol{X}^{\top}(\boldsymbol{X}\boldsymbol{\alpha}^{*} + \boldsymbol{z}_{r} + \mathbb{E}_{\boldsymbol{\epsilon}}\boldsymbol{\epsilon})$ $= \alpha^* + (\frac{1}{n} X^\top X)^{-1} \frac{1}{n} X^\top z_r$ $[\frac{1}{n} \boldsymbol{X}^{\top} \boldsymbol{z}_r]_k = \frac{1}{n} \sum \varphi_k(\boldsymbol{x}_i) r(\boldsymbol{x}_i)$ i=1 $\rightarrow \int \varphi_k(\boldsymbol{x}) r(\boldsymbol{x}) q(\boldsymbol{x}) d\boldsymbol{x} = 0$ (Law of large numbers)

Efficiency

- 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 Cramér-Rao bound.
 For linear model with LS and ε_i ^{i.i.d.} N(0, σ²), Cramér-Rao bound is

$$\sigma^2 \operatorname{tr}(\boldsymbol{U}(\boldsymbol{X}^{\top}\boldsymbol{X})^{-1})$$

Asymptotic Efficiency

Asymptotically efficient estimator: An unbiased estimator which asymptotically attains Cramér-Rao's lower bound. When $\epsilon_i \overset{i.i.d.}{\sim} N(0,\sigma^2)$ and $x_i \overset{i.i.d.}{\sim} q(x)$, LS estimator is asymptotically efficient. Proof: LS estimator is asymptotically unbiased and $\mathbb{E}_{\boldsymbol{\epsilon}} \| \hat{\boldsymbol{\alpha}}_{LS} - \mathbb{E}_{\boldsymbol{\epsilon}} \hat{\boldsymbol{\alpha}}_{LS} \|_{\boldsymbol{I}}^2 = \mathbb{E}_{\boldsymbol{\epsilon}} \| \boldsymbol{X} \boldsymbol{\epsilon} \|_{\boldsymbol{I}}^2$ $= \sigma^2 \operatorname{tr}(\boldsymbol{U}(\boldsymbol{X}^\top \boldsymbol{X})^{-1})$

which is Cramér-Rao's lower bound.

Homework

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

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

Hints:

• All linear unbiased estimator $\hat{\alpha}_{LU}$ satisfies $\mathbb{E}_{\epsilon}[\hat{\alpha}_{LU}] = \alpha^* \qquad \hat{\alpha}_{LU} = L_U y$

Therefore, $L_U X = I$ • By assumptions, noise satisfies $\mathbb{E}_{\epsilon}[\epsilon] = \mathbf{0}$ $\mathbb{E}_{\epsilon}[\epsilon \epsilon^{\top}] = \sigma^2 I$