Pattern Information Processing<sup>171</sup> Input-Dependent Estimation of Generalization Error

> Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/ Input Dependent Estimation <sup>172</sup> of Generalization Error

CV is very general and useful.

- Its unbiasedness holds with respect to both input points and output noise.
- However, input points are known.
- Is it possible to have an unbiased estimator of the generalization error only with respect to the noise?

#### Setting

 $\mathbf{P}q(\mathbf{x})$  is known.

Linear model:

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

Regularization learning:

$$\begin{split} \min_{\boldsymbol{\alpha}} \left[ \sum_{i=1}^{n} \left( \hat{f}(\boldsymbol{x}_{i}) - y_{i} \right)^{2} + \lambda \|\boldsymbol{\alpha}\|^{2} \right] \\ \boldsymbol{L} = (\boldsymbol{X}^{\top} \boldsymbol{X} + \lambda \boldsymbol{I})^{-1} \boldsymbol{X}^{\top} \\ \hat{\boldsymbol{\alpha}} = \boldsymbol{L} \boldsymbol{y} \\ \boldsymbol{X}_{i,j} = \varphi_{j}(\boldsymbol{x}_{i}) \\ \boldsymbol{\alpha} = (\alpha_{1}, \alpha_{2}, \dots, \alpha_{b})^{\top} \\ \boldsymbol{y} = (y_{1}, y_{2}, \dots, y_{n})^{\top} \end{split}$$

Decomposition of  
Generalization Error
$$G = \int \left(\hat{f}(x) - f(x)\right)^2 q(x) dx$$

$$= \int \hat{f}(x)^2 q(x) dx \quad \text{(accessible)}$$

$$-2 \int \hat{f}(x) f(x) q(x) dx \quad \text{(to be estimated)}$$

$$+ \int f(x)^2 q(x) dx \quad \text{(constant: ignored)}$$



#### Estimation of Generalization Error

Suppose we have  $L_u, \sigma_u^2$  such that (i)  $\mathbb{E}_{\epsilon} L_u y = \alpha^*$  ( $L_u$  and L are irrelevant) (ii)  $\mathbb{E}_{\epsilon} \sigma_u^2 = \sigma^2$ 

$$\mathbb{E}_{\boldsymbol{\epsilon}} \int \hat{f}(\boldsymbol{x}) g(\boldsymbol{x}) q(\boldsymbol{x}) d\boldsymbol{x} = \mathbb{E}_{\boldsymbol{\epsilon}} \langle \boldsymbol{U} \hat{\boldsymbol{\alpha}}, \boldsymbol{\alpha}^* \rangle = \mathbb{E}_{\boldsymbol{\epsilon}} [\langle \boldsymbol{U} \boldsymbol{L} \boldsymbol{y}, \boldsymbol{L}_u \boldsymbol{y} \rangle - \sigma_u^2 \operatorname{tr}(\boldsymbol{U} \boldsymbol{L} \boldsymbol{L}_u^\top)]$$

However, such  $L_u, \sigma_u^2$  are not available in practice, so we use approximations.

# Estimation of Generalization Error (cont.)

(i) 
$$\mathbb{E}_{\boldsymbol{\epsilon}} L_u y = \boldsymbol{\alpha}^*$$

$$\widehat{\boldsymbol{L}}_u = (\boldsymbol{X}^{ op} \boldsymbol{X})^{-1} \boldsymbol{X}^{ op}$$



$$\mathbb{E}_{\boldsymbol{\epsilon}}\widehat{\boldsymbol{L}}_{u}\boldsymbol{y} = \boldsymbol{\alpha}^{*} + \mathcal{O}_{p}(\delta n^{-\frac{1}{2}})$$

#### Proof

## Estimation of Generalization Error (ii) $\mathbb{E}_{\epsilon} \sigma_u^2 = \sigma^2$ (cont.)



#### Proof

$$G^{2} = G$$

$$GX\alpha^{*} = 0$$

$$\|Gz_{r}\|^{2} = \mathcal{O}(\delta^{2})$$

$$\mathbb{E}_{\epsilon}\sigma_{u}^{2} = \frac{\mathbb{E}_{\epsilon}\|Gy\|^{2}}{|G|^{2}}$$

$$\mathbb{E}_{\boldsymbol{\epsilon}} \sigma_{u}^{2} = \frac{\mathbb{E}_{\boldsymbol{\epsilon}} \|\boldsymbol{G}\boldsymbol{y}\|^{2}}{\operatorname{tr}(\boldsymbol{G})}$$
$$= \frac{\|\boldsymbol{G}\boldsymbol{X}\boldsymbol{\alpha}^{*} + \boldsymbol{G}\boldsymbol{z}_{r}\|^{2} + \mathbb{E}_{\boldsymbol{\epsilon}} \|\boldsymbol{G}\boldsymbol{\epsilon}\|^{2}}{\operatorname{tr}(\boldsymbol{G})}$$
$$= \sigma^{2} + \mathcal{O}(\delta^{2})$$

$$\widehat{G} = \langle \boldsymbol{U}\boldsymbol{L}\boldsymbol{y}, \boldsymbol{L}\boldsymbol{y} \rangle - 2\langle \boldsymbol{U}\boldsymbol{L}\boldsymbol{y}, \widehat{\boldsymbol{L}}_{u}\boldsymbol{y} \rangle + 2\widehat{\sigma_{u}^{2}} \mathrm{tr}(\boldsymbol{U}\boldsymbol{L}\widehat{\boldsymbol{L}}_{u}^{\top})$$

Bias: 
$$B_{\epsilon} = \mathcal{O}_p(\delta n^{-\frac{1}{2}})$$

$$B_{\epsilon} = \mathbb{E}_{\epsilon}[\widehat{G} - G] + C$$

$$C = \int f(\boldsymbol{x})^2 q(\boldsymbol{x}) d\boldsymbol{x}$$

A purpose of estimating generalization error is model selection.

We want to know whether  $\widehat{G}$  can distinguish good models from poor ones.

 $\mathcal{M} = \{\{\varphi_i(\boldsymbol{x})\}_{i=1}^b, \lambda\}$ 

183



If  $\operatorname{sgn}(\mathbb{E}_{\epsilon}\Delta G) = \operatorname{sgn}(\mathbb{E}_{\epsilon}\Delta \widehat{G})$ , better model can be selected on average.

However, checking the sign is not easy, so we simplify the criterion.



Good" if  $0 < \mathbb{E}_{\epsilon} \Delta \widehat{G} < 2\mathbb{E}_{\epsilon} \Delta G$  ( $\mathbb{E}_{\epsilon} \Delta G > 0$ )  $0 > \mathbb{E}_{\epsilon} \Delta \widehat{G} > 2\mathbb{E}_{\epsilon} \Delta G$  ( $\mathbb{E}_{\epsilon} \Delta G < 0$ )

#### Difference in the bias $B_{\epsilon}$ : $\Delta B_{\epsilon} = \mathbb{E}_{\epsilon} [\Delta \widehat{G} - \Delta G]$

Effective in model comparison:  $|\Delta B_{\epsilon}| < |\mathbb{E}_{\epsilon} \Delta G|$ 

Asymptotically effective in model comparison:

$$\Delta B_{\epsilon} = o_p(n^{-t}), \quad \mathbb{E}_{\epsilon} \Delta G \neq o_p(n^{-t})$$

### Effectiveness in Model Comparis<sup>186</sup>

#### $\widehat{G}$ is

 Effective in model comparison (if f(x) is realizable)
 Asymptotically effective in model comparison (o.w.)

#### Schedule

July 11<sup>th</sup>

: regular lecture (active learning) workshop registration

July 18<sup>th</sup>July 25<sup>th</sup>

: no class

: mini-workshop