

Pattern Information Processing:¹⁸⁹ Active Learning

Masashi Sugiyama
(Department of Computer Science)

Contact: W8E-505

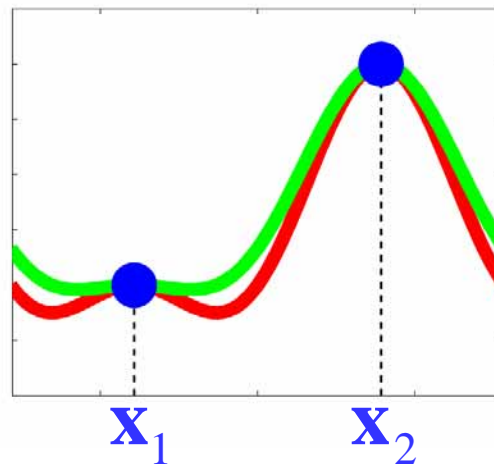
sugi@cs.titech.ac.jp

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

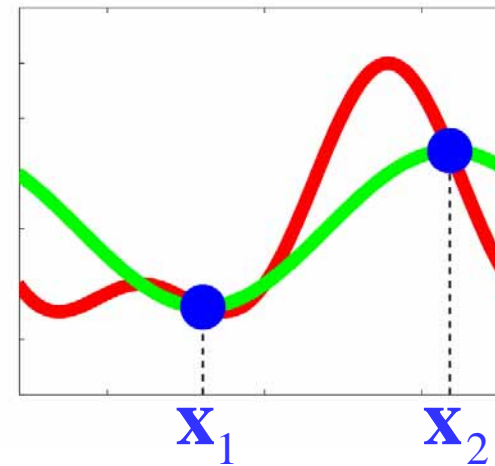
Active Learning

190

For obtaining good learning results, training input points should be determined appropriately.



Good questions

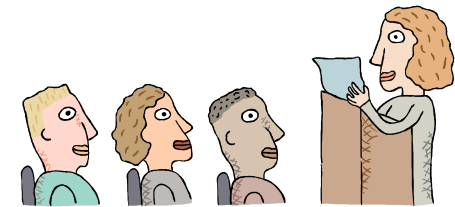


Bad questions

— Target function
— Learned function

Active Learning: Analogy to Real Life¹⁹¹

- It is not interesting to **passively** attend the lecture.



- It is more effective to **actively** ask questions in the lecture.



Formal Description

192

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

- Determine training input points so that

$$\min_{\{\mathbf{x}_i\}_{i=1}^n} G$$

Setting

■ $q(\mathbf{x})$ is known.

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

■ Least-squares learning:

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

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

$$\hat{\alpha} = \mathbf{L} \mathbf{y}$$

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

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_b)^\top$$

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

Estimating Generalization Error¹⁹⁴

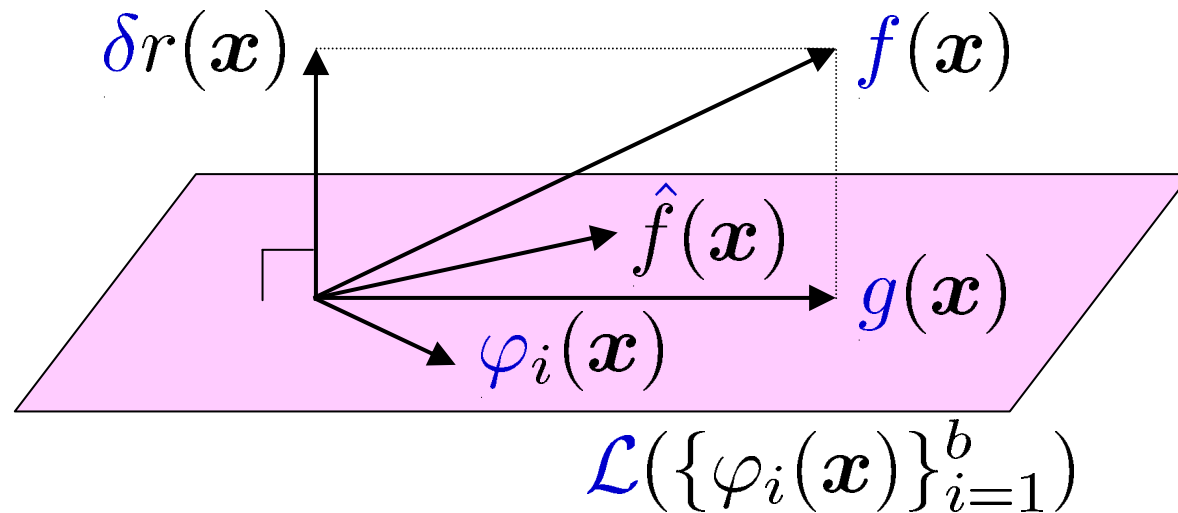
$$\min_{\{\mathbf{x}_i\}_{i=1}^n} G \quad G = \int_{\mathcal{D}} \left(\hat{f}(\mathbf{t}) - f(\mathbf{t}) \right)^2 q(\mathbf{t}) d\mathbf{t}$$

- We have to estimate unknown generalization error.
- This is similar to model selection.
- We do not have training output values $\{y_i\}_{i=1}^n$ in active learning!

Decomposition of Target Function¹⁹⁵

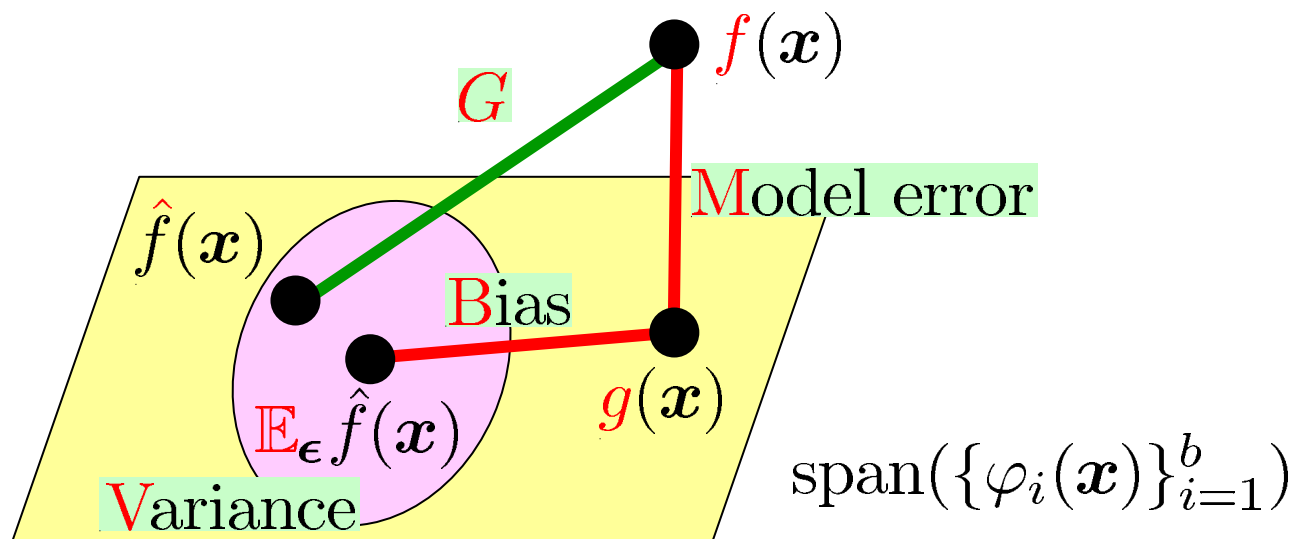
$$\begin{aligned} f(\mathbf{x}) &= g(\mathbf{x}) + \delta r(\mathbf{x}) & \int r(\mathbf{x})^2 q(\mathbf{x}) d\mathbf{x} &= 1 \\ g(\mathbf{x}) &= \sum_{i=1}^b \alpha_i^* \varphi_i(\mathbf{x}) & \int \varphi_i(\mathbf{x}) r(\mathbf{x}) q(\mathbf{x}) d\mathbf{x} &= 0 \end{aligned}$$

$$\delta \geq 0$$



Bias/Variance Decomposition¹⁹⁶

$$\begin{aligned}\mathbb{E}_\epsilon G &= \mathbb{E}_\epsilon \int \left(\hat{f}(x) - g(x) - \delta r(x) \right)^2 q(x) dx \\ &= \delta^2 \int r(x)^2 q(x) dx && \text{(model error)} \\ &+ \int \left(\mathbb{E}_\epsilon \hat{f}(x) - g(x) \right)^2 q(x) dx && \text{(bias)} \\ &+ \mathbb{E}_\epsilon \int \left(\hat{f}(x) - \mathbb{E}_\epsilon \hat{f}(x) \right)^2 q(x) dx && \text{(variance)}\end{aligned}$$



Assumption

- We assume that **model is correct**
 - $\delta = 0$: model error vanishes
 - Least squares is unbiased: bias vanishes
- Only variance remains!

$$\begin{aligned}\mathbb{E}_{\epsilon} G &= \mathbb{E}_{\epsilon} \int \left(\hat{f}(\mathbf{x}) - \mathbb{E}_{\epsilon} \hat{f}(\mathbf{x}) \right)^2 q(\mathbf{x}) d\mathbf{x} \\ &= \sigma^2 \text{tr}(\mathbf{U} \mathbf{L} \mathbf{L}^{\top}) \\ &= \sigma^2 \text{tr}(\mathbf{U} (\mathbf{X}^{\top} \mathbf{X})^{-1}) \\ &\propto \text{tr}(\mathbf{U} (\mathbf{X}^{\top} \mathbf{X})^{-1})\end{aligned}$$

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

Active Learning with LS

198

- Determine $\{\mathbf{x}_i\}_{i=1}^n$ so that

$$\operatorname{argmin}_{\{\mathbf{x}_i\}_{i=1}^n} \left[\operatorname{tr}(\mathbf{U}(\mathbf{X}^\top \mathbf{X})^{-1}) \right]$$

- In active learning, we can not use training output values $\{y_i\}_{i=1}^n$ for estimating generalization error.
- We considered zero-bias cases and evaluated the variance!

How to Optimize

199

- Determine $\{\mathbf{x}_i\}_{i=1}^n$ so that

$$\operatorname{argmin}_{\{\mathbf{x}_i\}_{i=1}^n} \left[\operatorname{tr}(\mathbf{U}(\mathbf{X}^\top \mathbf{X})^{-1}) \right]$$

- For trigonometric polynomial models, the solution can be analytically obtained.
- However, in general, simultaneously optimizing n points is not tractable.

How to Optimize (cont.)

200

- Major approaches to avoid intractability:
 - Optimize points one by one in a greedy manner
 - Optimize probability distribution from which training input points are drawn.

$$\{\mathbf{x}_i^{(k)}\}_{i=1}^n \stackrel{i.i.d.}{\sim} p^{(k)}(\mathbf{x})$$

$$\operatorname{argmin}_k \left[\operatorname{tr}(\mathbf{U}(\mathbf{X}^{(k)\top} \mathbf{X}^{(k)})^{-1}) \right]$$

$$\mathbf{X}_{i,j}^{(k)} = \varphi_j(\mathbf{x}_i^{(k)})$$

When Model Is Not Correct 201

- When model is not correct, least-squares is no longer unbiased (even asymptotically).
- Instead, the following **importance-weighted LS** is asymptotically unbiased.

$$\min_{\alpha} \left[\sum_{i=1}^n \frac{q(\mathbf{x}_i)}{p(\mathbf{x}_i)} \left(\hat{f}(\mathbf{x}_i) - y_i \right)^2 \right]$$

$$\{\mathbf{x}_i\}_{i=1}^n \stackrel{i.i.d.}{\sim} p(\mathbf{x})$$

$$t \sim q(\mathbf{x})$$

IWLS

- IWLS learning result is given by

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

$$\hat{\alpha} = L_W \mathbf{y}$$

$$L_W = (\mathbf{X}^\top \mathbf{D} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{D}$$

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

$$\mathbf{D} = \text{diag} \left(\frac{q(\mathbf{x}_1)}{p(\mathbf{x}_1)}, \frac{q(\mathbf{x}_2)}{p(\mathbf{x}_2)}, \dots, \frac{q(\mathbf{x}_n)}{p(\mathbf{x}_n)} \right)$$

Asymptotic Unbiasedness of IWLS²⁰³

$$\blacksquare \textcolor{red}{y} = X\alpha^* + \delta \mathbf{z}_r + \boldsymbol{\epsilon}$$

$$\textcolor{red}{z}_r = (r(\mathbf{x}_1), r(\mathbf{x}_2), \dots, r(\mathbf{x}_n))^\top$$

$$\textcolor{red}{\epsilon} = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$$

$$\blacksquare \textcolor{red}{L}_W X\alpha^* = \alpha^*$$

$$\blacksquare \left[\frac{1}{n} X^\top D \mathbf{z}_r \right]_k = \frac{1}{n} \sum_{i=1}^n \varphi_k(\mathbf{x}_i) r(\mathbf{x}_i) \frac{q(\mathbf{x}_i)}{p(\mathbf{x}_i)}$$

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

(Law of large numbers)

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

Active Learning with IWLS

204

- Variance of IWLS is

$$\sigma^2 \text{tr}(\mathbf{U} \mathbf{L}_W \mathbf{L}_W^\top)$$

- Determine $p(\mathbf{x})$ so that

$$\underset{p(\mathbf{x})}{\text{argmin}} \left[\text{tr}(\mathbf{U} \mathbf{L}_W \mathbf{L}_W^\top) \right]$$

Notification of Final Assignment

1. Apply supervised learning techniques to your data set and analyze it.
 2. Write your opinion about this course
- Final report deadline: **Aug 11th (Fri.)**
 - Only **e-mail submission** is accepted!
sugi@cs.titech.ac.jp