

# Pattern Information Processing:<sup>216</sup> Covariate Shift Adaptation

Masashi Sugiyama  
(Department of Computer Science)

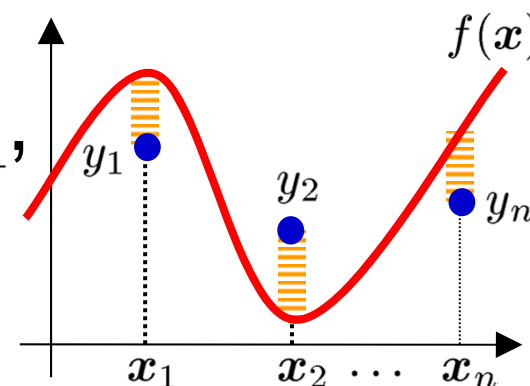
Contact: W8E-505

[sugi@cs.titech.ac.jp](mailto:sugi@cs.titech.ac.jp)

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

# Common Assumption in Supervised Learning

- Goal of supervised learning:  
From **training samples**  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ,  
predict outputs of **unseen  
test samples**



- We always assume

Training and test samples are  
drawn from the **same distribution**

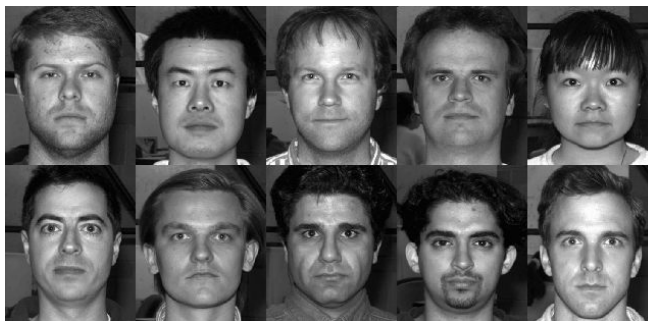
$$P_{train}(\mathbf{x}, y) = P_{test}(\mathbf{x}, y)$$

- Is this assumption really true?

# Not Always True!

- **Less women** in face dataset than reality.
- **More criticisms** in survey sampling than reality.
- **Sample generation mechanism** varies over time.

The Yale Face Database B



# Covariate Shift

- However, no chance for generalization if training and test samples have **nothing in common**.

$$P_{train}(\mathbf{x}, y) \neq P_{test}(\mathbf{x}, y)$$

- **Covariate shift:**

- Input distribution changes

$$P_{train}(\mathbf{x}) \neq P_{test}(\mathbf{x})$$

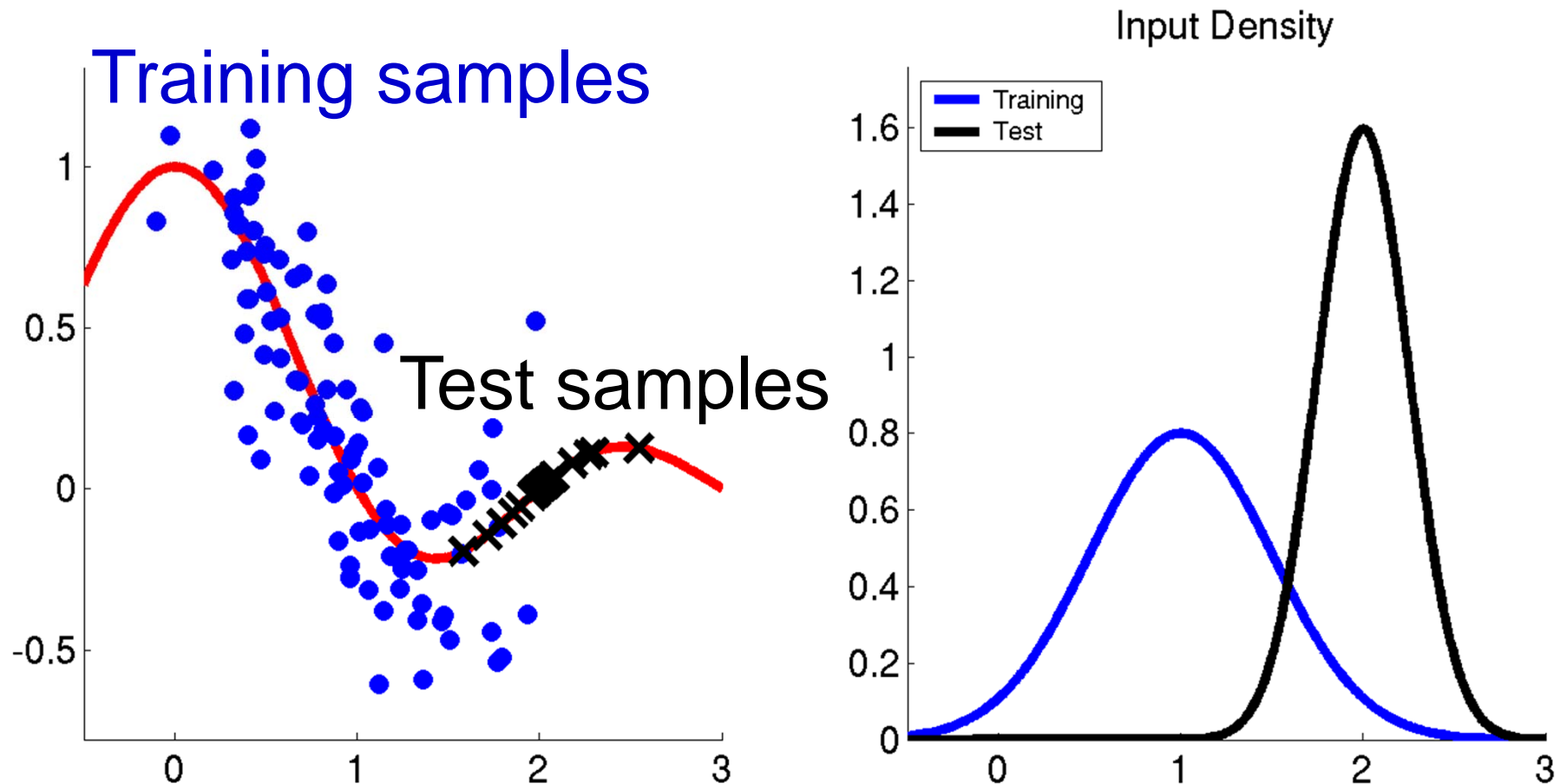
- Functional relation remains unchanged

$$P_{train}(y|\mathbf{x}) = P_{test}(y|\mathbf{x})$$

# Examples of Covariate Shift <sup>220</sup>

(Weak) extrapolation:

Predict output values outside training region



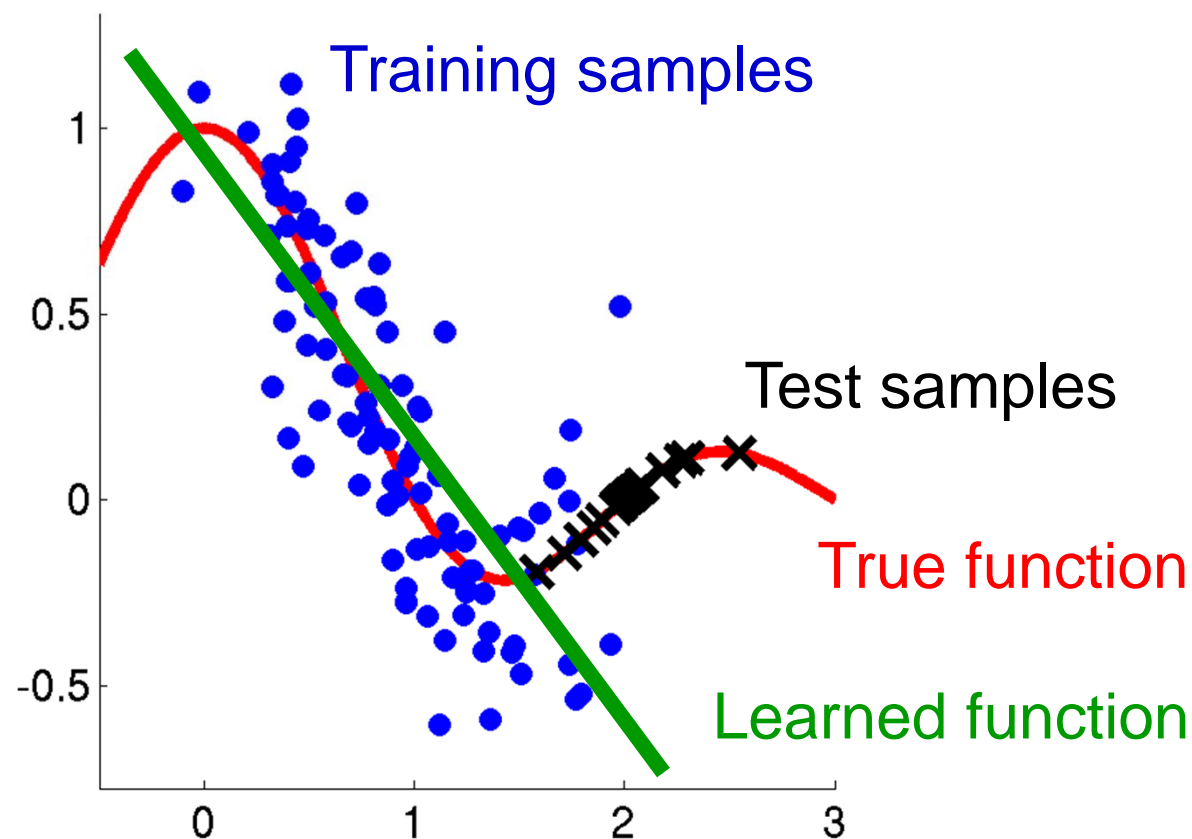
# Organization

1. Linear regression under covariate shift
2. **Parameter learning**
3. Importance estimation
4. Model selection



# Covariate Shift

- To illustrate the effect of covariate shift, let's focus on **linear extrapolation**



# Generalization Error = Bias + Variance

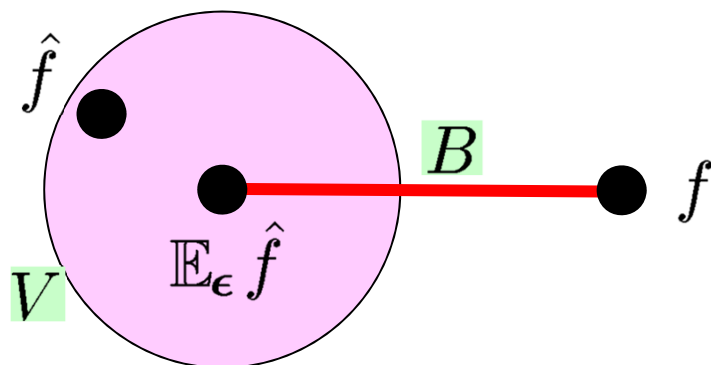
$$\mathbb{E}_{\epsilon} \int \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 p_{test}(\mathbf{x}) d\mathbf{x}$$

$$= \int \left( \mathbb{E}_{\epsilon} \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 p_{test}(\mathbf{x}) d\mathbf{x}$$

Bias

$$+ \mathbb{E}_{\epsilon} \int \left( \mathbb{E}_{\epsilon} \hat{f}(\mathbf{x}) - \hat{f}(\mathbf{x}) \right)^2 p_{test}(\mathbf{x}) d\mathbf{x}$$

Variance



$\mathbb{E}_{\epsilon}$  : expectation over noise



# Model Specification

- Model is said to be **correctly specified** if

$$\exists \alpha^*, \hat{f}(x; \alpha^*) = f(x)$$

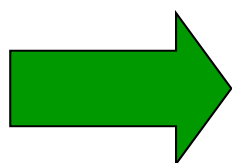
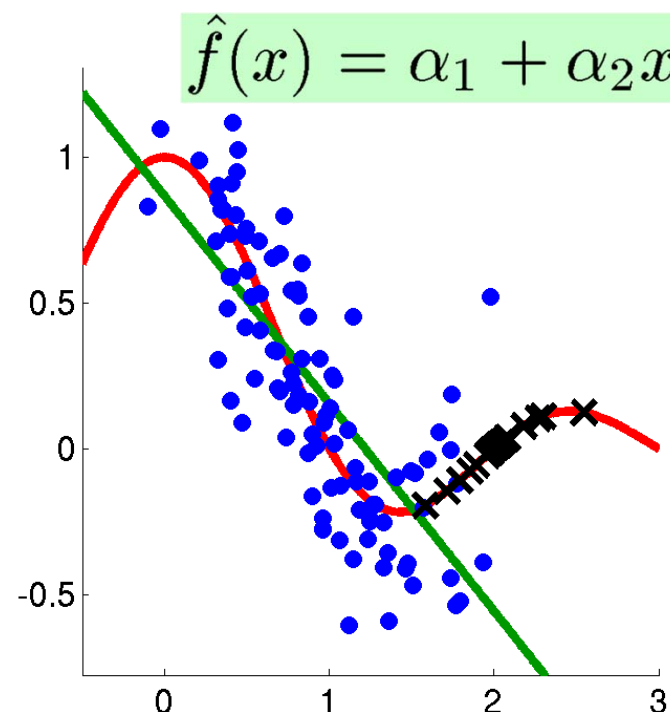
- In practice, our model may not **be correct**.
- Therefore, we need a theory for **misspecified models**!

# Ordinary Least-Squares

225

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

- If model is correct:
  - OLS minimizes bias asymptotically
- If model is misspecified:
  - OLS does **not minimize bias even asymptotically**.



We want to reduce bias!

# Law of Large Numbers

- Sample average converges to the population mean:

$$\frac{1}{n} \sum_{i=1}^n A(\mathbf{x}_i) \longrightarrow \int A(\mathbf{x}) p_{train}(\mathbf{x}) d\mathbf{x}$$

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} p_{train}(\mathbf{x})$$

- We want to estimate the expectation over **test input points** only using **training input points**  $\{\mathbf{x}_i\}_{i=1}^n$  .

$$\int A(\mathbf{t}) p_{test}(\mathbf{t}) d\mathbf{t}$$

$$\mathbf{t} \sim p_{test}(\mathbf{x})$$

# Key Trick:

## Importance-Weighted Average

- **Importance**: Ratio of test and training input densities

$$\frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})}$$

- **Importance-weighted average**:

$$\frac{1}{n} \sum_{i=1}^n \frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)} A(\mathbf{x}_i) \longrightarrow \int \frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})} A(\mathbf{x}) p_{train}(\mathbf{x}) d\mathbf{x}$$

$$\mathbf{x}_i \stackrel{i.i.d.}{\sim} p_{train}(\mathbf{x}) \quad = \int A(\mathbf{x}) p_{test}(\mathbf{x}) d\mathbf{x}$$

$$t \sim p_{test}(\mathbf{x}) \quad (\text{cf. importance sampling})$$

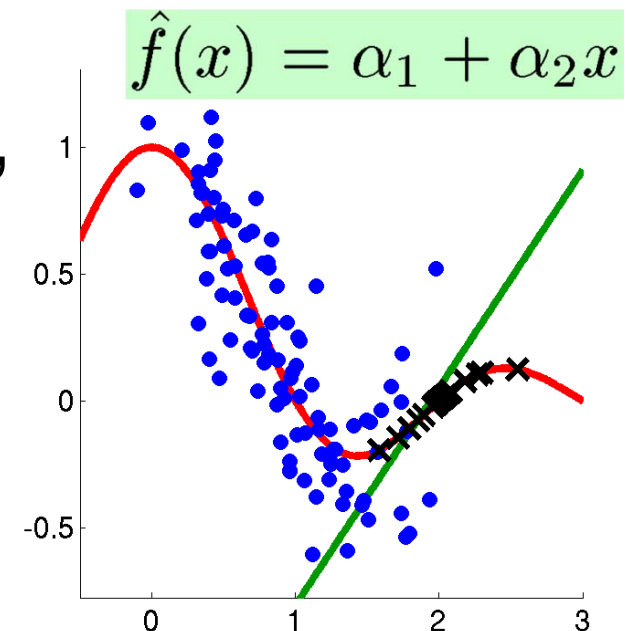
# Importance-Weighted LS

228

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

$p_{train}(\mathbf{x}), p_{test}(\mathbf{x})$  : Assumed strictly positive

- Even for misspecified models, IWLS **minimizes bias asymptotically**.
- We need to estimate importance in practice.



# Organization

1. Linear regression under covariate shift
2. Parameter learning
3. Importance estimation
4. Model selection



# Importance Estimation

$$w(\mathbf{x}_i) = \frac{p_{test}(\mathbf{x}_i)}{p_{train}(\mathbf{x}_i)}$$

- **Assumption:** We have training inputs  $\{\mathbf{x}_i^{train}\}_{i=1}^{n_{train}}$  and test inputs  $\{\mathbf{x}_i^{test}\}_{i=1}^{n_{test}}$ .
- **Naïve approach:** Estimate  $p_{train}(\mathbf{x})$  and  $p_{test}(\mathbf{x})$  separately, and take the ratio of the density estimates
- This does not work well since density estimation is hard in high dimensions.

# Modeling Importance Function<sup>231</sup>

$$w(\mathbf{x}) = \frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})}$$

- We use a linear model:

$$\hat{w}(\mathbf{x}) = \sum_{i=1}^t \theta_i \phi_i(\mathbf{x}) \quad \theta_i, \phi_i(\mathbf{x}) \geq 0$$

- Test density is approximated by

$$\hat{p}_{test}(\mathbf{x}) = \hat{w}(\mathbf{x}) p_{train}(\mathbf{x})$$

- **Idea:** Learn  $\{\theta_i\}_{i=1}^t$  so that  $\hat{p}_{test}(\mathbf{x})$  well approximates  $p_{test}(\mathbf{x})$ .



# Kullback-Leibler Divergence <sup>232</sup>

$$\min_{\{\theta_i\}_{i=1}^t} KL[p_{test}(\mathbf{x}) || \hat{p}_{test}(\mathbf{x})]$$

$$\hat{p}_{test}(\mathbf{x}) = \hat{w}(\mathbf{x}) p_{train}(\mathbf{x})$$

■  $KL[p_{test}(\mathbf{x}) || \hat{w}(\mathbf{x}) p_{train}(\mathbf{x})]$

$$= \int p_{test}(\mathbf{x}) \log \frac{p_{test}(\mathbf{x})}{\hat{w}(\mathbf{x}) p_{train}(\mathbf{x})} d\mathbf{x}$$

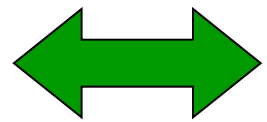
$$= \int p_{test}(\mathbf{x}) \log \frac{p_{test}(\mathbf{x})}{p_{train}(\mathbf{x})} d\mathbf{x} \quad \text{(constant)}$$

$$- \int p_{test}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x} \quad \text{(relevant)}$$

# Learning Importance Function<sup>233</sup>

■ Thus

$$\min_{\{\theta_i\}_{i=1}^t} KL[\hat{w}(\mathbf{x})p_{train}(\mathbf{x})||p_{test}(\mathbf{x})]$$



$$\max_{\{\theta_i\}_{i=1}^t} \int p_{test}(\mathbf{x}) \log \hat{w}(\mathbf{x}) d\mathbf{x}$$

(objective function)

■ Since  $\hat{p}_{test}(\mathbf{x}) = \hat{w}(\mathbf{x})p_{train}(\mathbf{x})$  is density,

$$\int \hat{w}(\mathbf{x})p_{train}(\mathbf{x})d\mathbf{x} = 1$$

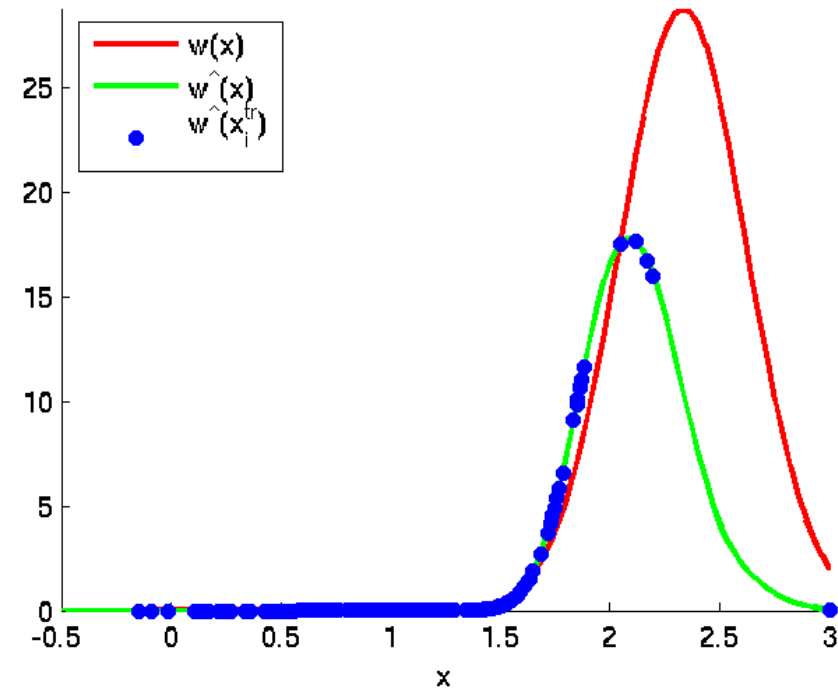
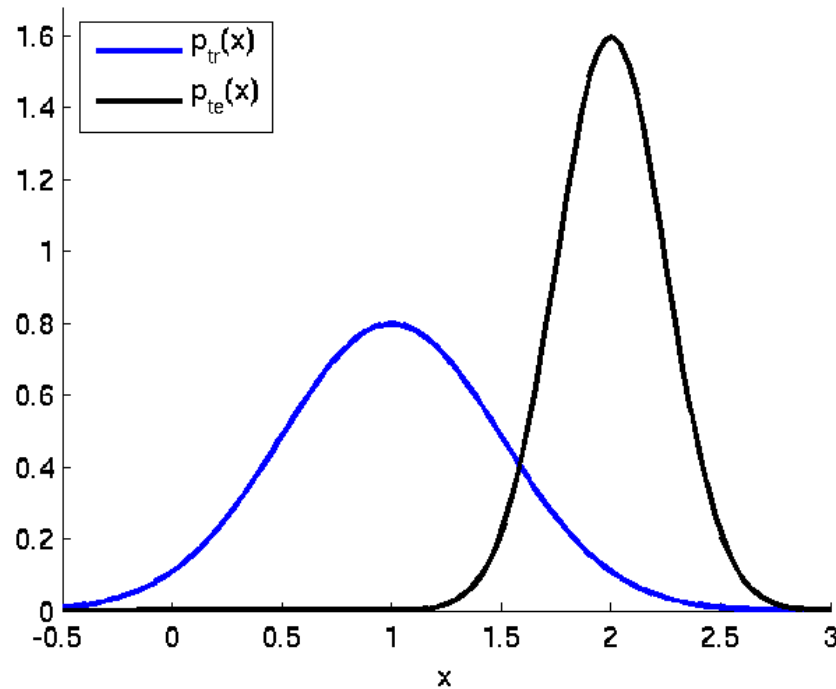
(constraint)

# KLIEP (Kullback-Leibler Importance Estimation Procedure)

$$\begin{aligned} \max_{\{\theta_i\}_{i=1}^t} & \left[ \sum_{i=1}^{n_{test}} \log \hat{w}(\mathbf{x}_i^{test}) \right] \\ \text{subject to} & \sum_{i=1}^{n_{train}} \hat{w}(\mathbf{x}_i^{train}) = n_{train} \\ & \theta_1, \theta_2, \dots, \theta_t \geq 0 \end{aligned}$$
$$\hat{w}(\mathbf{x}) = \sum_{i=1}^t \theta_i \phi_i(\mathbf{x})$$

- **Convexity:** unique global solution is available
- **Sparse solution:** prediction is fast!

# Examples



$$\hat{w}(x) = \sum_{i=1}^{n_{test}} \theta_i K(x, x_i^{test})$$

$$K(x, x') = \exp \left( -\frac{\|x - x'\|^2}{2\sigma^2} \right)$$

# Model Selection of KLIEP

236

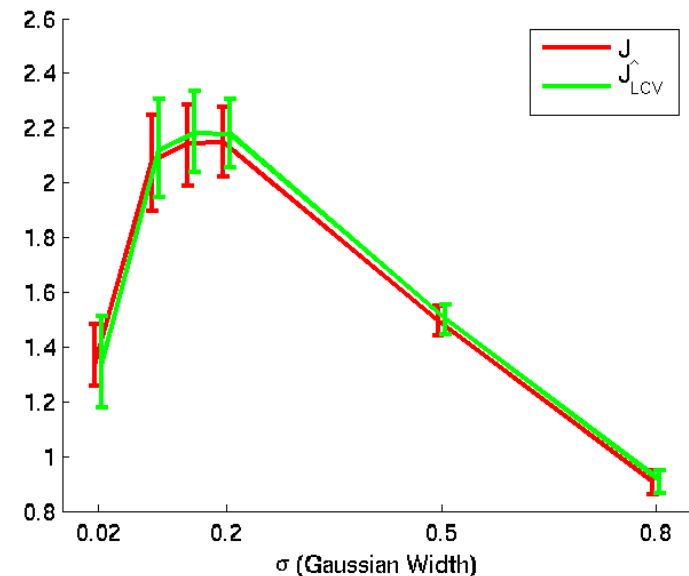
■ How to choose tuning parameters (such as Gaussian width)?

■ Likelihood cross-validation:

- Divide test samples  $\{x_i^{test}\}_{i=1}^{n_{test}}$  into  $\mathcal{X}$  and  $\mathcal{X}'$ .
- Learn importance from  $\mathcal{X}$ .
- Estimate the likelihood using  $\mathcal{X}'$ .

$$\frac{1}{|\mathcal{X}'|} \sum_{x \in \mathcal{X}'} \log \hat{w}_{\mathcal{X}}(x)$$

■ This gives **an unbiased estimate of KL** (up to an irrelevant constant).



# Organization

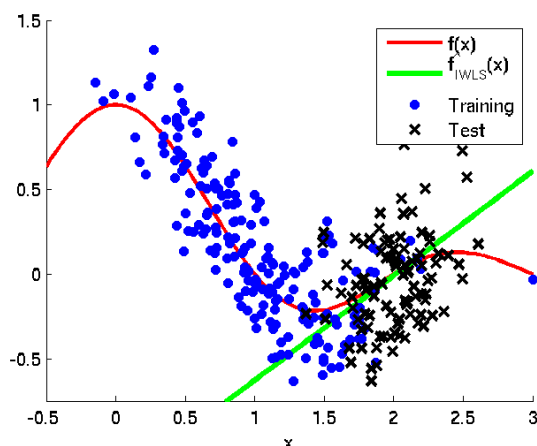
1. Linear regression under covariate shift
2. Parameter learning
3. Importance estimation
4. **Model selection**



# Model Selection

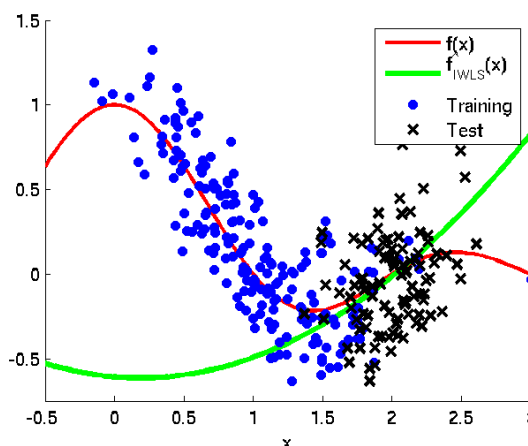
Choice of models is crucial:

Polynomial of order 1



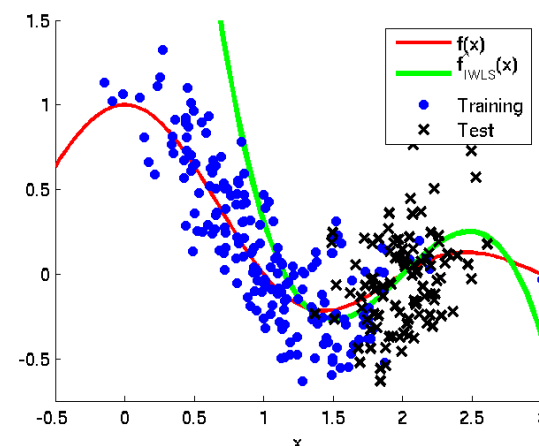
$$\hat{f}(x) = \alpha_1 + \alpha_2 x$$

Polynomial of order 2



$$\hat{f}(x) = \alpha_1 + \alpha_2 x + \alpha_3 x^2$$

Polynomial of order 3



$$\hat{f}(x) = \alpha_1 + \alpha_2 x + \alpha_3 x^2 + \alpha_4 x^3$$

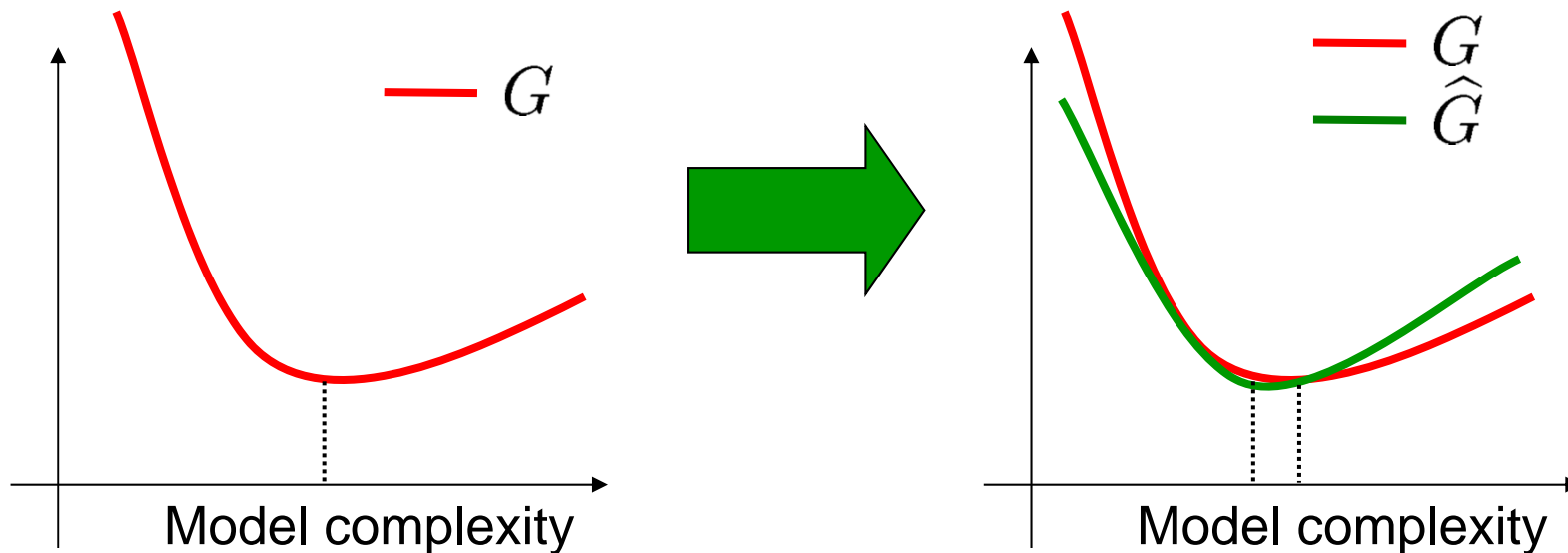
We want to determine the model so that **generalization error is minimized**:

$$G = \int \left( \hat{f}(x) - f(x) \right)^2 p_{test}(x) dx$$

# Generalization Error Estimation<sup>239</sup>

$$G = \int \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 p_{test}(\mathbf{x}) d\mathbf{x}$$

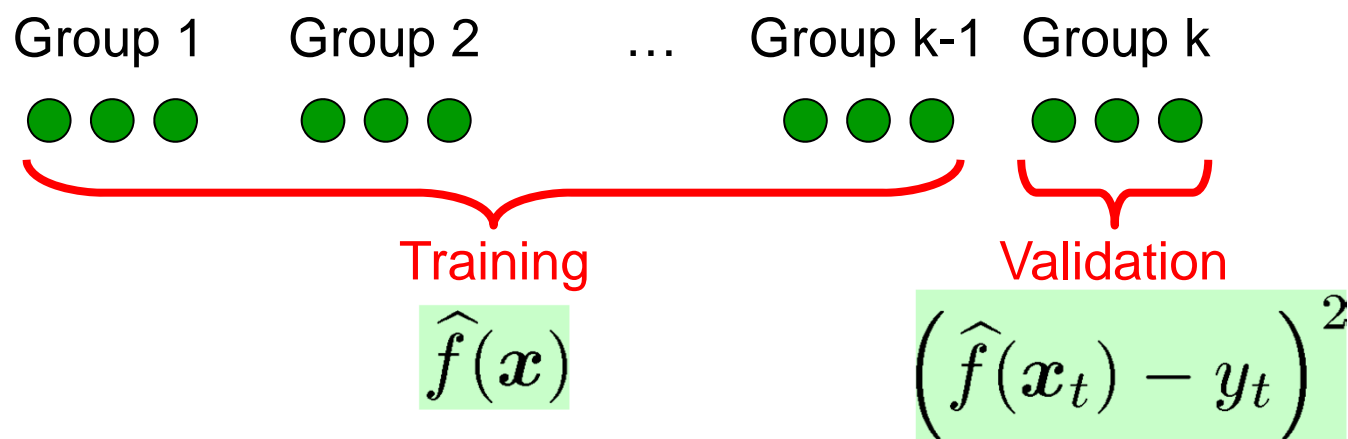
- Generalization error is not accessible since the target function  $f(\mathbf{x})$  is unknown.
- Instead, we use **a generalization error estimate**.





# Cross-Validation

- Divide training samples into  $k$  groups.
- Train a learning machine with  $k - 1$  groups.
- Validate the trained machine using the rest.
- Repeat this for all combinations and output the mean validation error.



- CV is almost unbiased without covariate shift.
- But, **CV is heavily biased under covariate shift!**

# Importance-Weighted CV (IWCV)<sup>241</sup>

- When testing the classifier in CV process, we also **importance-weight the test error**.



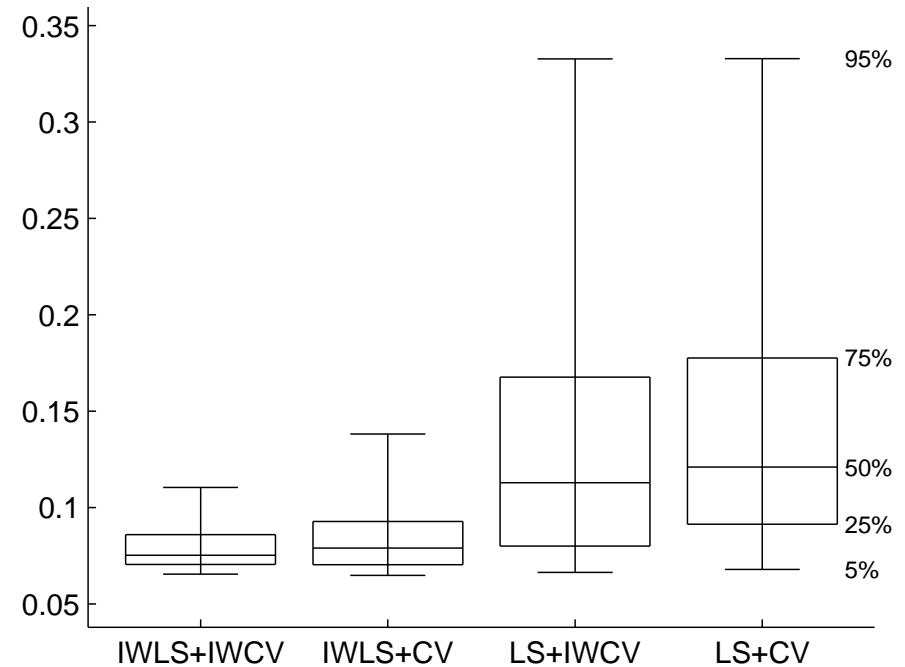
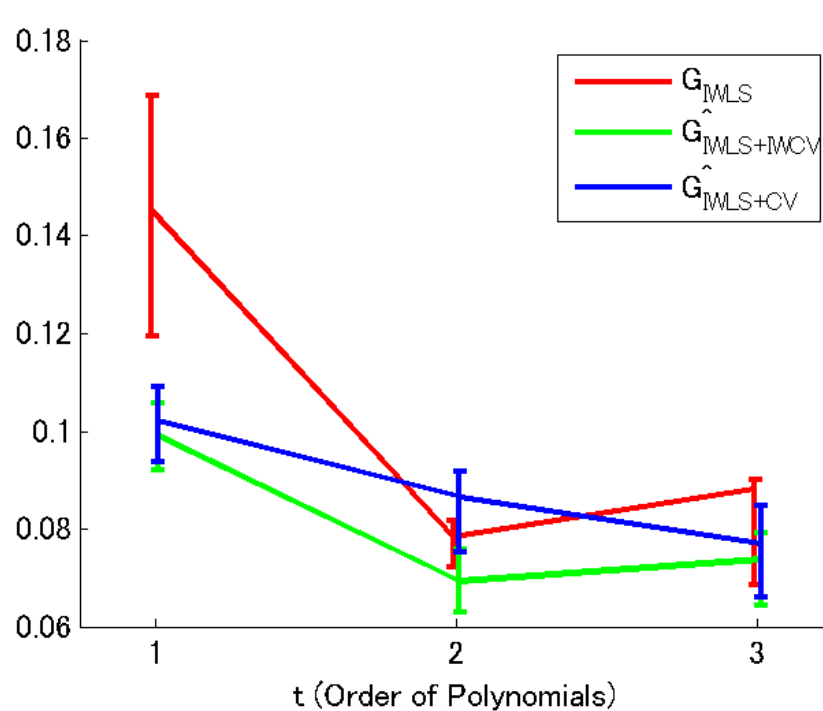
$$\hat{f}(\mathbf{x})$$

$$\frac{p_{test}(\mathbf{x}_t)}{p_{train}(\mathbf{x}_t)} \left( \hat{f}(\mathbf{x}_t) - y_t \right)^2$$

IWCV gives almost unbiased estimates of generalization error even under covariate shift

# Example of IWCV

242

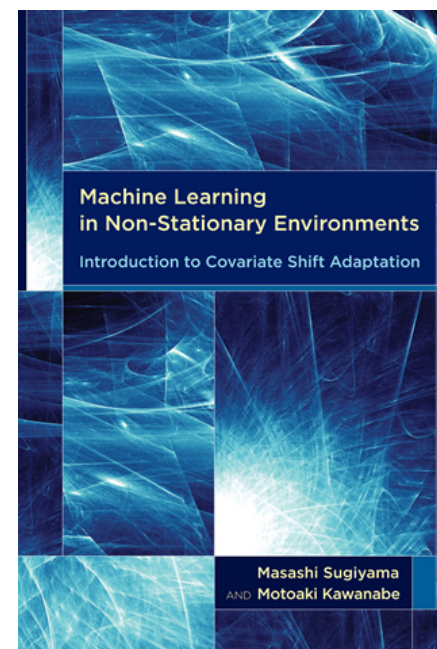


- IWCV gives better estimates of generalization error than CV.
- Model selection by IWCV outperforms CV!

# Summary

- **Covariate shift**: input distribution varies but functional relation remains unchanged
- **Importance weighting** for adaptation.
  - **IW least-squares**: consistent
  - **KLIEP**: direct importance estimation
  - **IW cross-validation**: unbiased
- **Further reading**:

Sugiyama & Kawanabe  
**Machine Learning  
in Non-Stationary Environments**,  
MIT Press, 2012



# 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 3<sup>rd</sup> (Fri.)
  - E-mail submission is also accepted!  
*sugi@cs.titech.ac.jp*