

# Model Selection under Covariate Shift



Institut  
Rechnerarchitektur  
und Softwaretechnik



**Masashi Sugiyama**

Tokyo Institute of Technology, Tokyo, Japan

**Klaus-Robert Müller**

Fraunhofer FIRST, Berlin, Germany

University of Potsdam, Potsdam, Germany

# Standard Regression Problem

■ Learning target function:  $f(\mathbf{x})$

■ Training examples:

$$\{(\mathbf{x}_i, y_i) \mid y_i = f(\mathbf{x}_i) + \epsilon_i\}_{i=1}^n$$

■ Test input:  $\{\mathbf{t}_i \mid \mathbf{t}_i \stackrel{i.i.d.}{\sim} p_t(\mathbf{x})\}_{i=1}^m$

■ Goal: Obtain approximation  $\hat{f}(\mathbf{x})$  that minimizes **expected error for test inputs** (or **generalization error**)

$$J = \int \left( \hat{f}(\mathbf{t}) - f(\mathbf{t}) \right)^2 p_t(\mathbf{t}) d\mathbf{t}$$

# Training Input Distribution

- Common assumption:

Training input  $\{\mathbf{x}_i\}_{i=1}^n$  follows the **same** distribution as test input:

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

- Here, we suppose distributions are different.

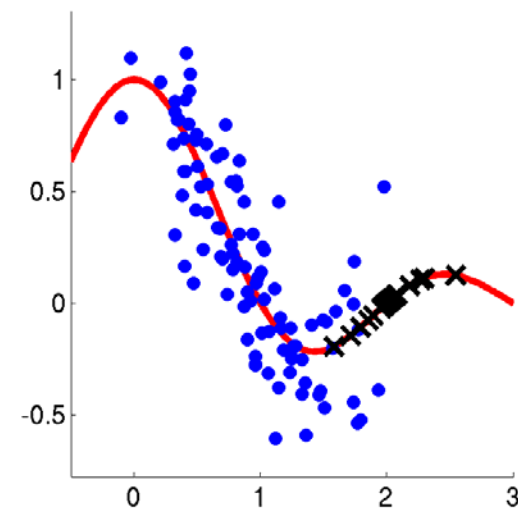
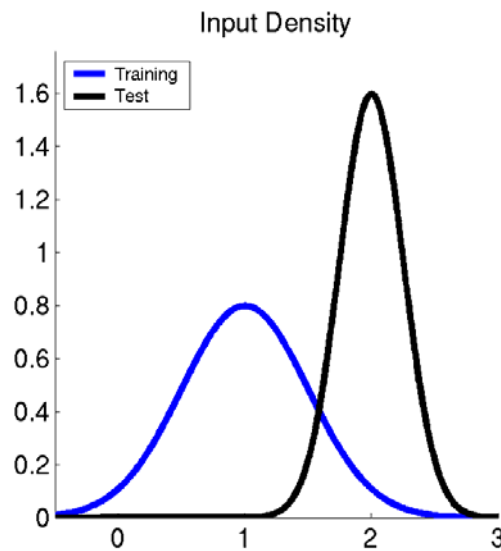
$$\begin{array}{l} \mathbf{x}_i \stackrel{i.i.d.}{\sim} p_x(\mathbf{x}) \\ \mathbf{t}_i \stackrel{i.i.d.}{\sim} p_t(\mathbf{x}) \end{array} \quad p_x(\mathbf{x}) \neq p_t(\mathbf{x})$$



**Covariate shift**

# Covariate Shift

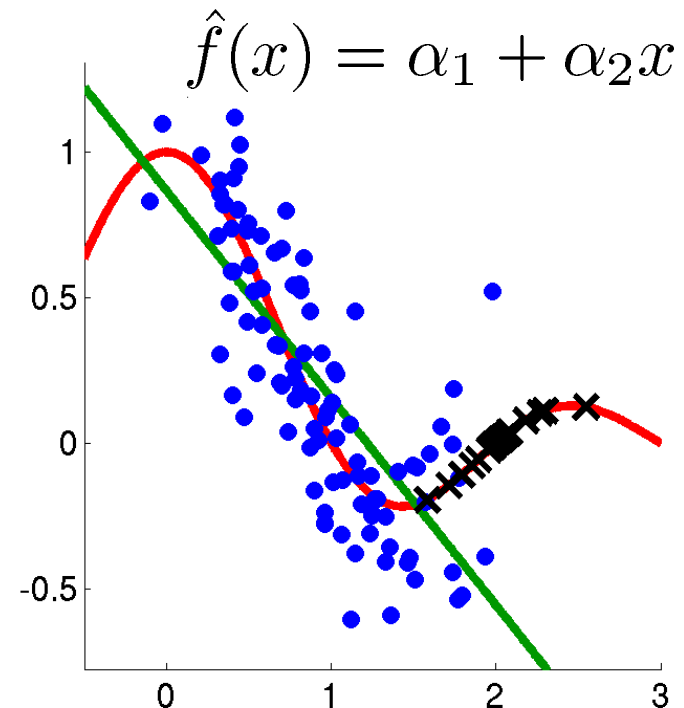
- Is covariate shift important to investigate?
- **Yes! It often happens in reality.**
  - Interpolation / **extrapolation**
  - Active learning (experimental design)
  - Classification from imbalanced data



# Ordinary Least Squares under Covariate Shift

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

- **Asymptotically unbiased** if model is correct.
- **Asymptotically biased** for misspecified models.
- Need to reduce bias.



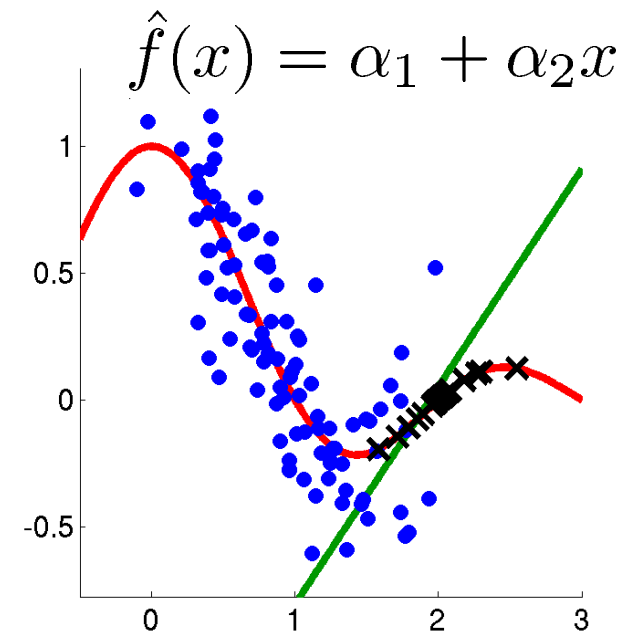
# Weighted Least Squares for Covariate Shift

(Shimodaira, 2000)

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

$p_x(\mathbf{x}), p_t(\mathbf{x})$  : Assumed known and strictly positive

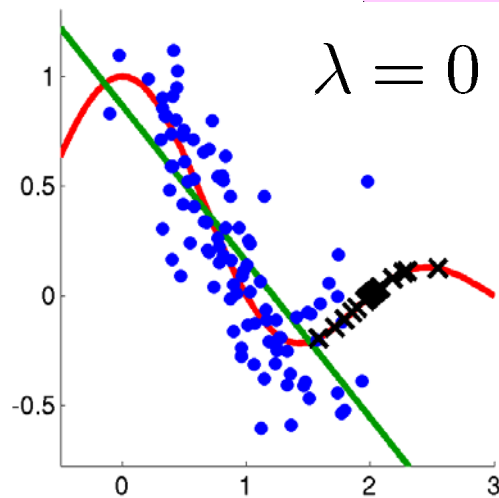
- **Asymptotically unbiased** for misspecified models.
- Can have large variance.
- Need to reduce variance.



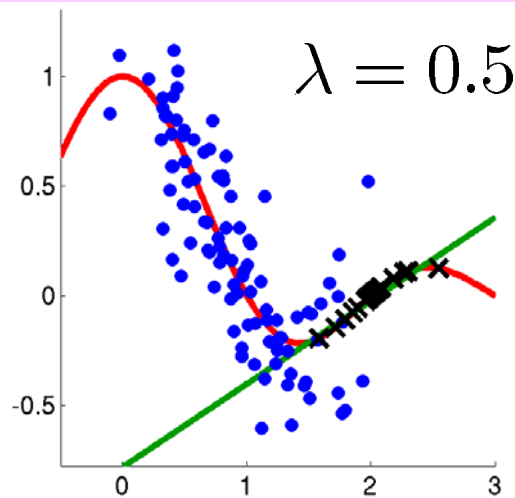
# $\lambda$ -Weighted Least Squares

(Shimodaira, 2000)

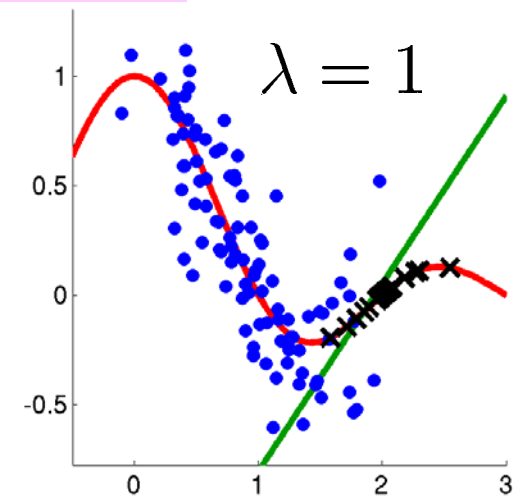
$$\min_{\alpha} \left[ \sum_{i=1}^n \left( \frac{p_t(\mathbf{x}_i)}{p_x(\mathbf{x}_i)} \right)^{\lambda} \left( \hat{f}(\mathbf{x}_i) - y_i \right)^2 \right]$$



Large bias  
Small variance



(Intermediate)

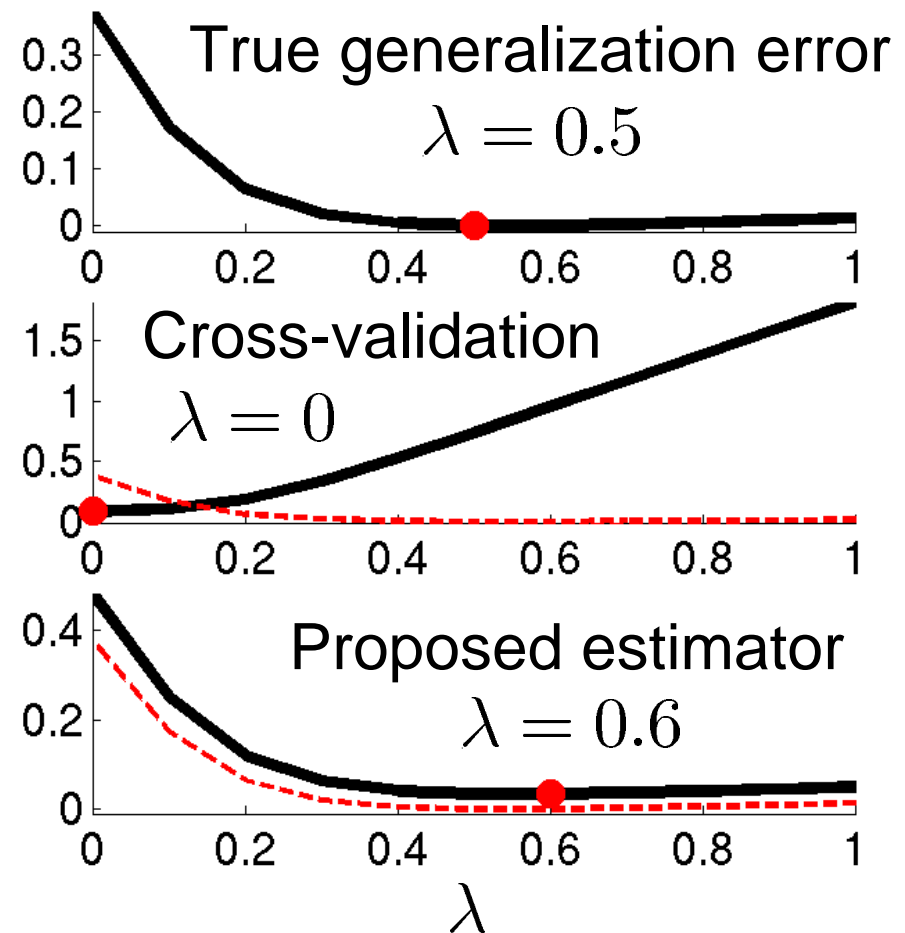


Small bias  
Large variance

$\lambda$  should be chosen appropriately!  
(Model Selection)

# Generalization Error Estimation <sup>8</sup> under Covariate Shift

- $\lambda$  is determined so that (estimated) generalization error is minimized.
- However, standard methods such as cross-validation is heavily biased.
- **Goal:** Derive better estimator





# Setting

- I.i.d. noise with mean 0 and variance  $\sigma^2$
- Linear regression model:

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

- $\lambda$ -weighted least squares:

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

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

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

$$\alpha = (\alpha_1, \alpha_2, \dots, \alpha_p)^{\top}$$

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

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

$$D = \text{diag} \left( \frac{p_t(\mathbf{x}_i)}{p_x(\mathbf{x}_i)} \right)$$

# Decomposition of Generalization Error

$$J = \int \left( \hat{f}(\mathbf{x}) - f(\mathbf{x}) \right)^2 p_t(\mathbf{x}) d\mathbf{x}$$

$$= \|\hat{f} - f\|_{L_2(p_t)}^2$$

$$= \|\hat{f}\|_{L_2(p_t)}^2 - 2\langle \hat{f}, f \rangle_{L_2(p_t)} + \|f\|_{L_2(p_t)}^2$$

Accessible

Estimated

Constant  
(ignored)

■ We estimate  $Z \equiv \langle \hat{f}, f \rangle_{L_2(p_t)}$

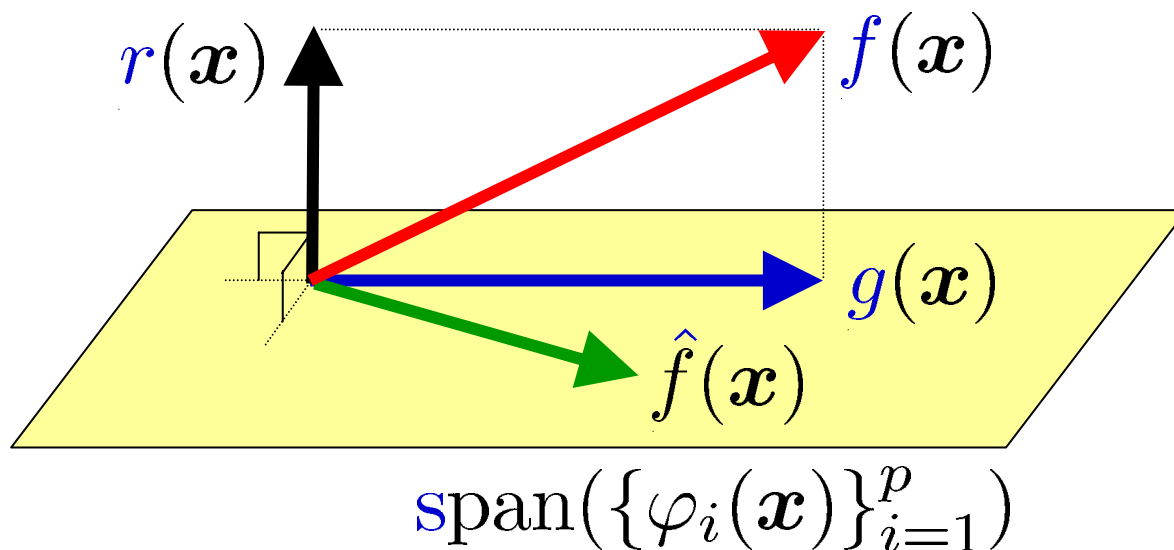
# Orthogonal Decomposition of Learning Target Function <sup>11</sup>

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

$$\langle \varphi_i, r \rangle_{L_2(p_t)} = 0$$

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

$\alpha^*$  : Optimal parameter



$$\begin{aligned} Z &= \langle \hat{f}, f \rangle_{L_2(p_t)} \\ &= \langle \hat{f}, g \rangle_{L_2(p_t)} \\ &= \langle \mathbf{U} \hat{\alpha}, \alpha^* \rangle \end{aligned}$$

$$U_{i,j} = \langle \varphi_i, \varphi_j \rangle_{L_2(p_t)}$$

# Unbiased Estimation of $\mathbb{E}_\epsilon Z$ 12

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

■ Suppose we have

- $L_u$ , which gives linear unbiased estimator of  $\alpha^*$

$$\mathbb{E}_\epsilon L_u \mathbf{y} = \alpha^*$$

- $\sigma_u^2$ : Unbiased estimator of noise variance

$$\mathbb{E}_\epsilon \sigma_u^2 = \sigma^2$$

■ Then we have an **unbiased estimator** of  $\mathbb{E}_\epsilon Z$  :

$$\hat{Z} \equiv \langle U L \mathbf{y}, L_u \mathbf{y} \rangle - \sigma_u^2 \text{tr}(U L L_u^\top)$$

■ But  $L_u, \sigma_u^2$  are not always available.



Use approximations instead

# Approximations of $L_u, \sigma_u^2$

■  $\hat{L}_u = (X^\top D X)^{-1} X^\top D$

■  $\hat{\sigma}_u^2 = \frac{\|y - H y\|^2}{n - p} \quad H = X(X^\top X)^{-1} X^\top$

$$\mathbb{E}_\epsilon L_u y = \alpha^*$$

$$\mathbb{E}_\epsilon \sigma_u^2 = \sigma^2$$

- If model is correct,

$$\mathbb{E}_\epsilon \hat{L}_u y = \alpha^* \quad \mathbb{E}_\epsilon \hat{\sigma}_u^2 = \sigma^2$$

- If model is misspecified,

$$\mathbb{E}_\epsilon \hat{L}_u y \rightarrow \alpha^* \quad \mathbb{E}_\epsilon \hat{\sigma}_u^2 \not\rightarrow \sigma^2 \quad (n \rightarrow \infty)$$

# New Generalization Error Estimator<sup>14</sup>

$$\hat{J} = \langle ULy, Ly \rangle - 2\langle ULy, \hat{L}_u y \rangle + 2\hat{\sigma}_u^2 \text{tr}(ULL\hat{L}_u^\top)$$

$$\text{Bias: } B_\epsilon = \mathbb{E}_\epsilon[\hat{J} - J] + C \quad C = \|f\|_{L_2(p_t)}^2$$

- If model is correct,

$$B_\epsilon = 0$$

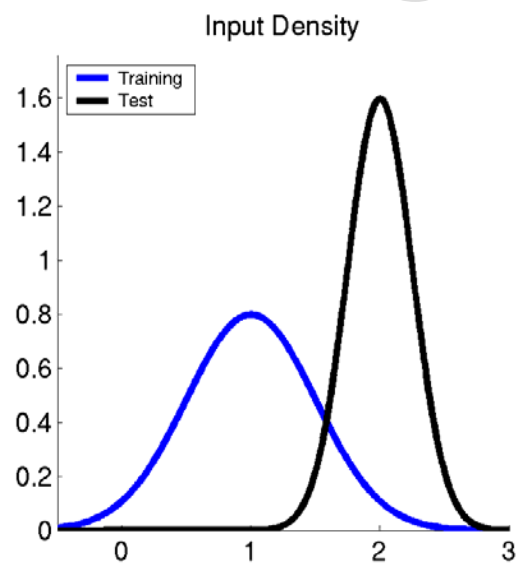
- If model is almost correct,

$$B_\epsilon = \mathcal{O}(\delta) \quad \delta = \max\{r(\mathbf{x}_i)\}$$

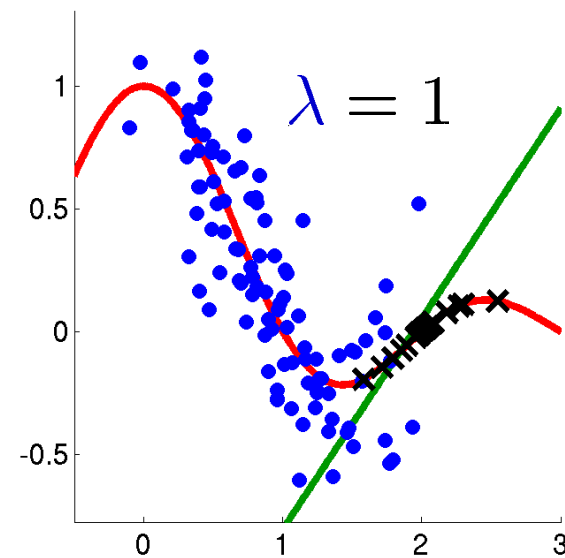
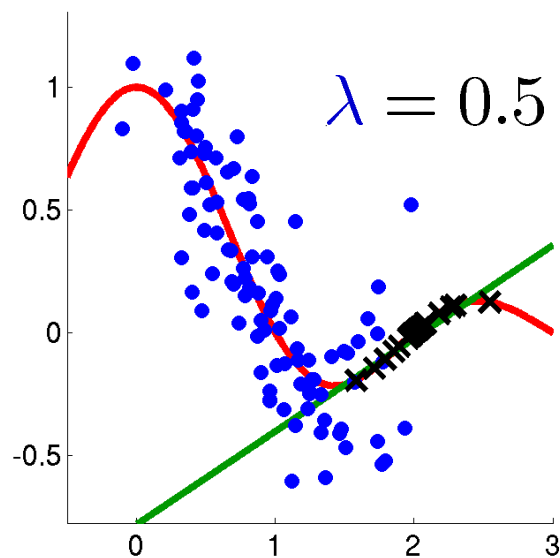
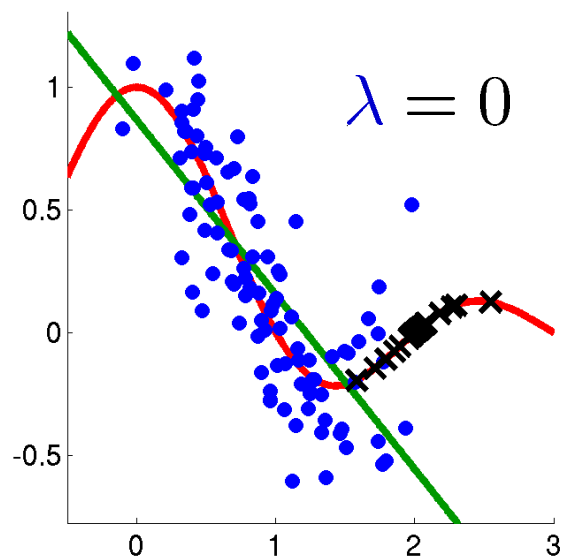
- If model is misspecified,

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

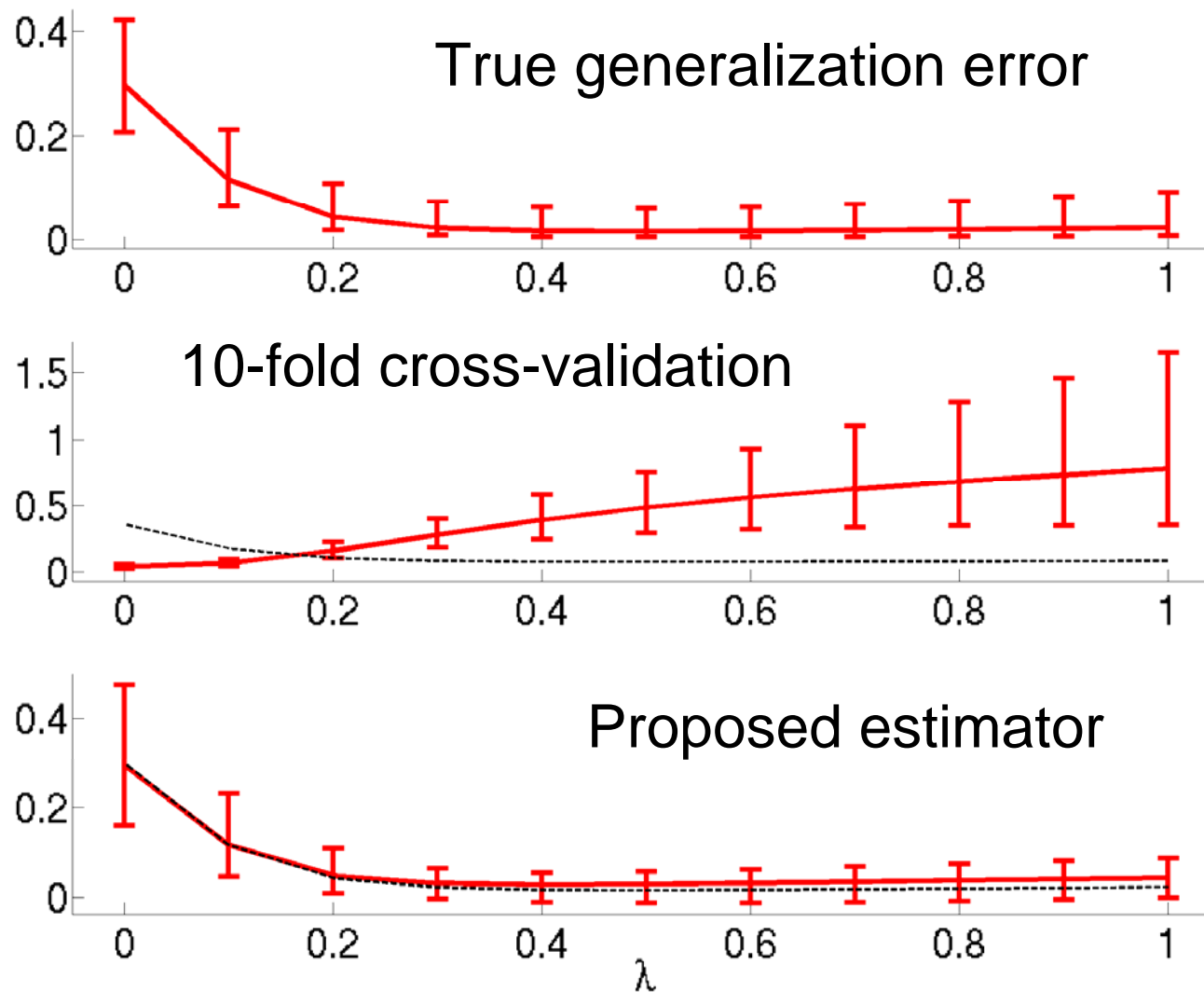
# Simulation (Toy)



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



# Results

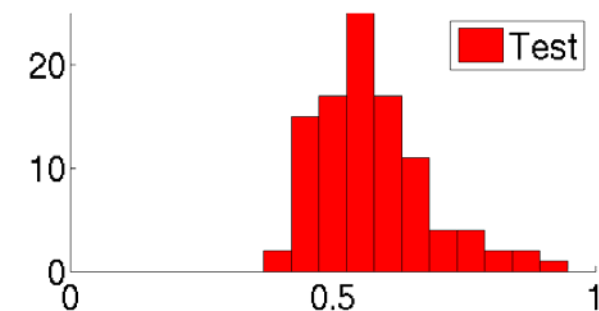
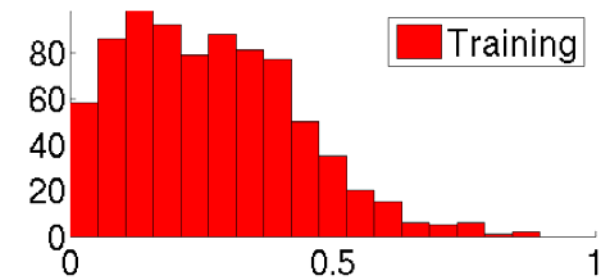
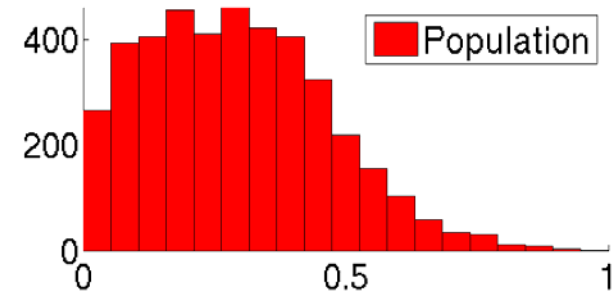




# Simulation (Abalone from DELVE)<sup>17</sup>

- Estimate the age of abalones from 7 physical measurements.
- We add bias to 4<sup>th</sup> attribute (weight of abalones)
- Training and test input densities are estimated by standard kernel density estimator.

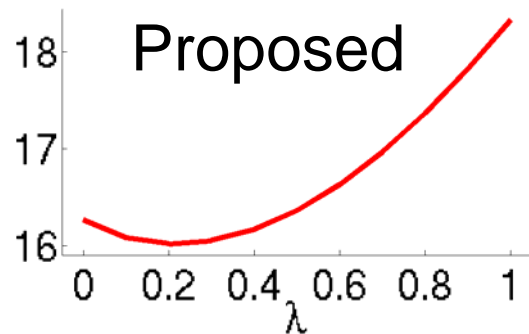
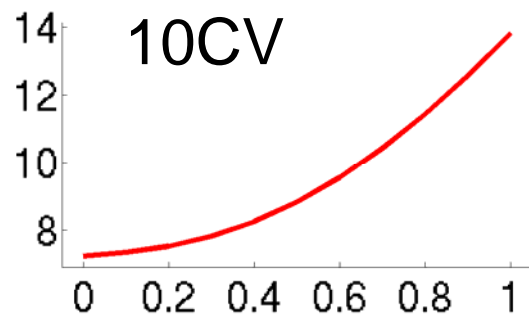
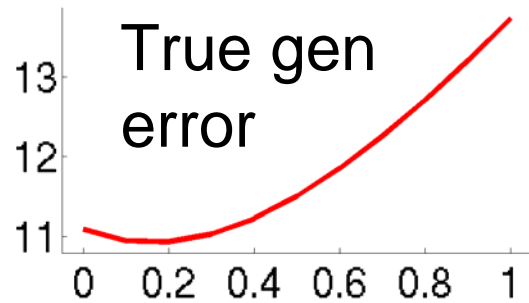
- $$\hat{f}(\mathbf{x}) = \alpha_1 + \sum_{i=1}^7 \alpha_{i+1} x^{(i)}$$



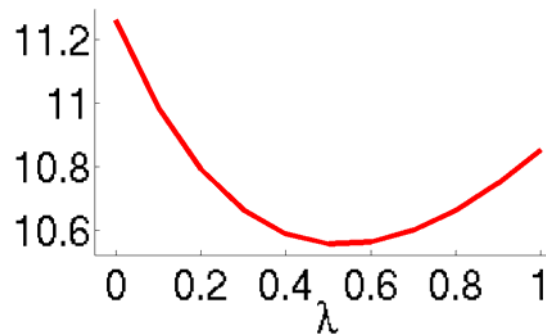
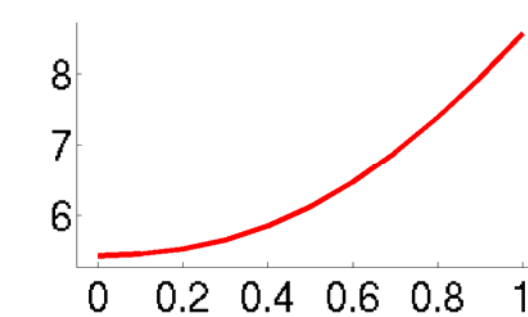
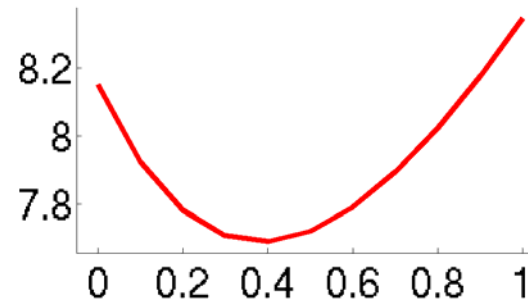
# Generalization Error Estimation<sup>18</sup>

Mean over 300 trials

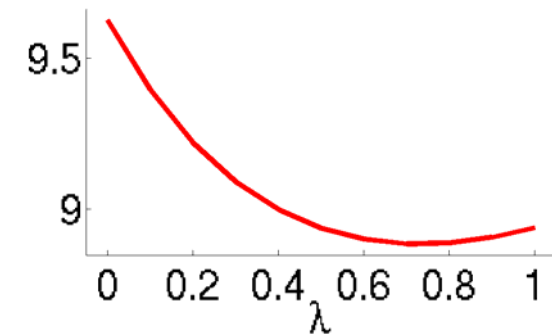
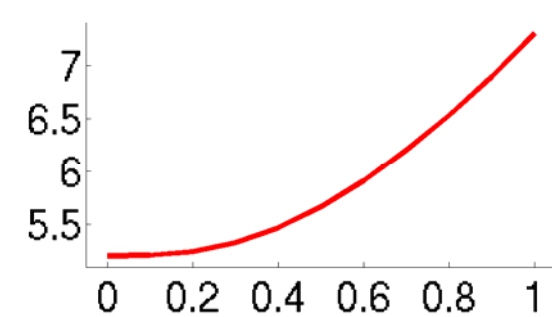
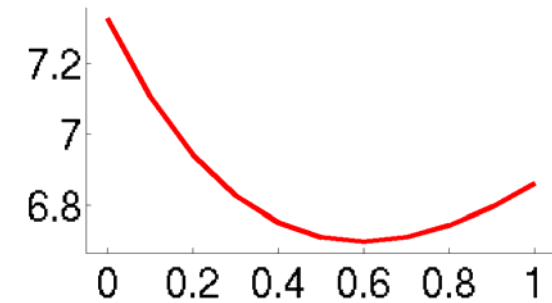
$n = 50$



$n = 200$



$n = 800$



# Test Error After Model Selection<sup>19</sup>

Extrapolation in 4<sup>th</sup> attribute

n	50	200	800
OPT	9.86 ± 4.27	7.40 ± 1.77	6.54 ± 1.34
Proposed	11.67 ± 5.74	7.95 ± 2.15	6.77 ± 1.40
10CV	10.88 ± 5.05	8.06 ± 1.91	7.24 ± 1.37

T-test (5%)

Extrapolation in 6<sup>th</sup> attribute

n	50	200	800
OPT	9.04 ± 4.04	6.76 ± 1.68	6.05 ± 1.25
Proposed	10.67 ± 6.19	7.31 ± 2.24	6.20 ± 1.33
10CV	10.15 ± 4.95	7.42 ± 1.81	6.68 ± 1.25

# Conclusions

- **Covariate shift**: Training and test input distributions are different
- Ordinary LS: Biased
- Weighted LS: Unbiased but large variance.
- $\lambda$ -WLS: Model selection needed.
- Cross-validation: Biased
- Proposed generalization error estimator:
  - Exactly unbiased (correct models)
  - Asymptotically unbiased (misspecified models)