

Pseudo-Labeling for Covariate Shift Adaptation

Kaizheng Wang

IEOR & DSI
Columbia University
November 12th 2024



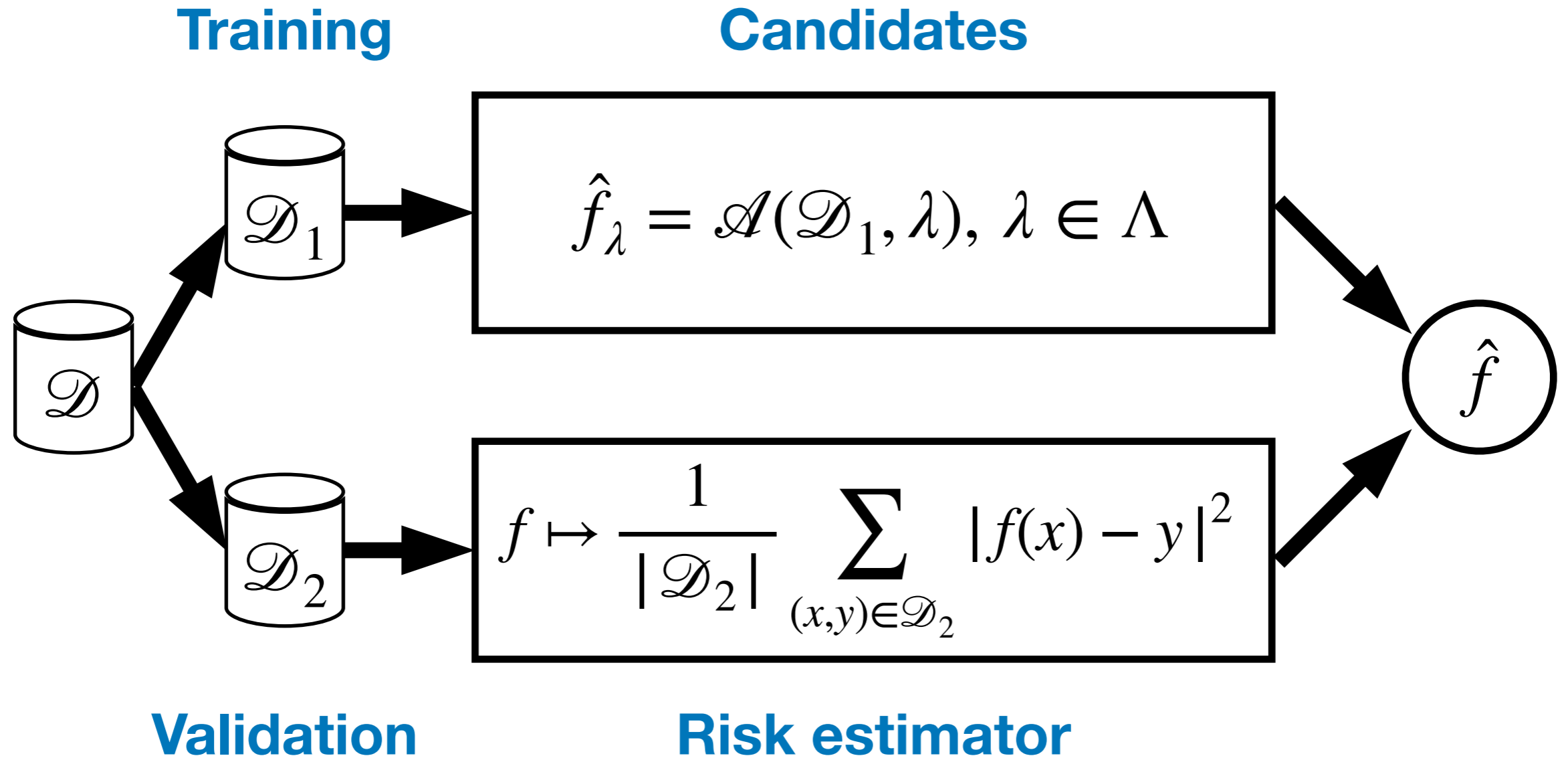
Outline

- **Introduction: covariate shift adaptation**
- Methodology: validation with pseudo-labels
- Model selection theory: general results
- Adaptivity guarantees: kernel ridge regression

Regression

- **Distribution:** $(x, y) \sim \mathcal{Q}$ over $\mathcal{X} \times \mathbb{R}$.
- **Risk (MSE):** $R(f) = \mathbb{E}_{(x,y) \sim \mathcal{Q}} |f(x) - y|^2$.
- **Data:** $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$ i.i.d. from \mathcal{Q} .
- **Goal:** learn a model \hat{f} with small risk $R(\hat{f})$.
- **Estimator:** $\hat{f}_\lambda = \mathcal{A}(\mathcal{D}, \lambda)$ with hyperparameter $\lambda \in \Lambda$.

Hold-Out Validation



Needs sufficient labeled data from \mathcal{Q} .

Covariate Shift

- **Goal:** learn a predictive model \hat{f} that works well over \mathcal{Q} .
- **Issue:** Too costly to collect labeled data from \mathcal{Q} .
- **Data:** $\mathcal{D} = \{(x_{P,i}, y_{P,i})\}_{i=1}^{n_P} \sim \mathcal{P}$ and $\{x_{Q,i}\}_{i=1}^{n_Q} \sim \mathcal{Q}_X$.

Covariate Shift

- **Goal:** learn a predictive model \hat{f} that works well over \mathcal{Q} .
- **Issue:** Too costly to collect labeled data from \mathcal{Q} .
- **Data:** $\mathcal{D} = \{(x_{P,i}, y_{P,i})\}_{i=1}^{n_P} \sim \mathcal{P}$ and $\{x_{Q,i}\}_{i=1}^{n_Q} \sim \mathcal{Q}_X$.

Examples

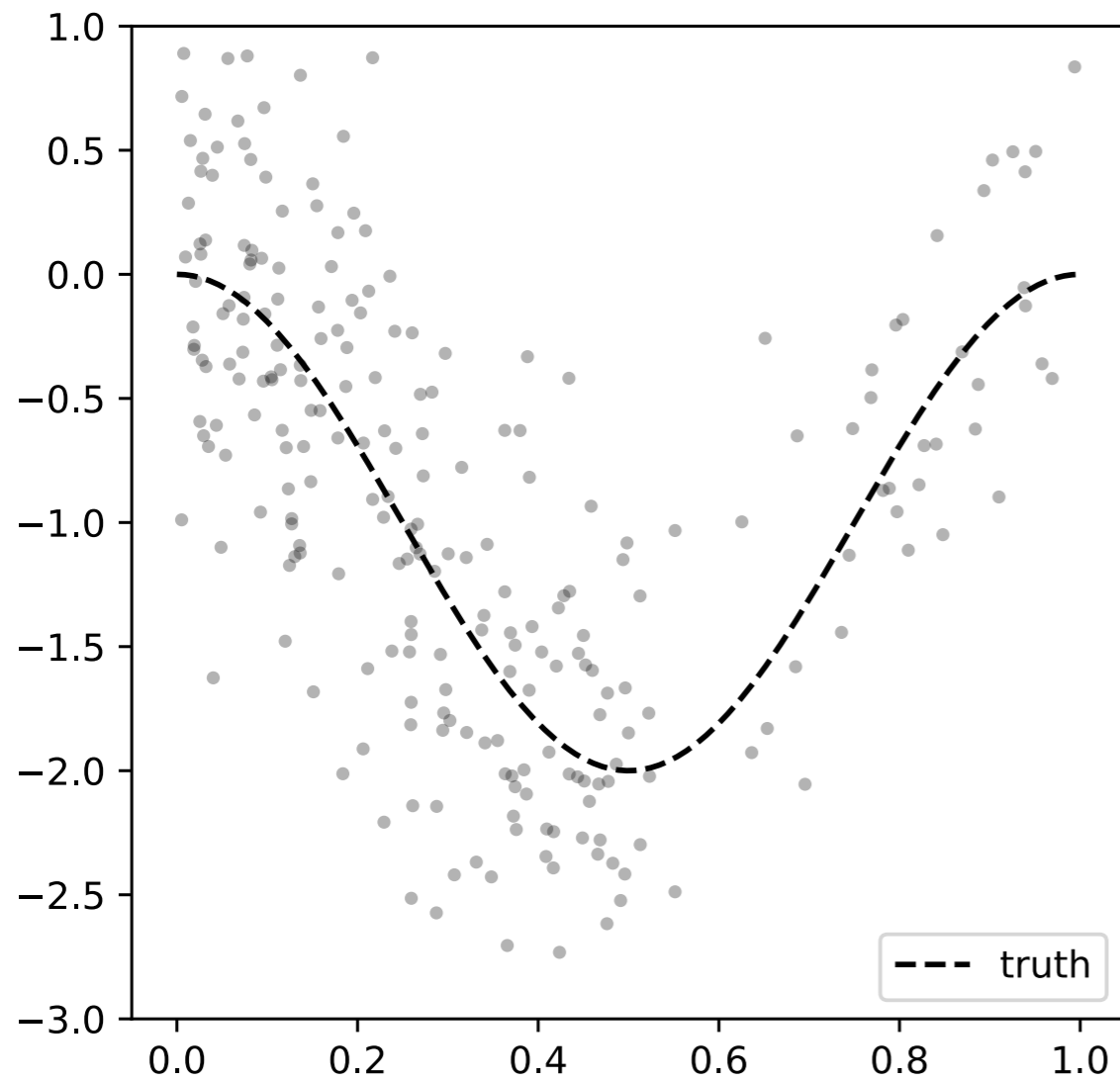
- **Chatbot:** \mathcal{Q} = Business negotiations, \mathcal{P} = News.
- **Self-driving:** \mathcal{Q} = NYC, \mathcal{P} = Berkeley.

Covariate Shift

- **Goal:** learn a predictive model \hat{f} that works well over \mathcal{Q} .
- **Issue:** Too costly to collect labeled data from \mathcal{Q} .
- **Data:** $\mathcal{D} = \{(x_{P,i}, y_{P,i})\}_{i=1}^{n_P} \sim \mathcal{P}$ and $\{x_{Q,i}\}_{i=1}^{n_Q} \sim \mathcal{Q}_X$.
- $\mathcal{P}_X \neq \mathcal{Q}_X$ but $\mathcal{P}_{Y|X} = \mathcal{Q}_{Y|X}$.
- \mathcal{P} , \mathcal{Q} are called **source** and **target**.

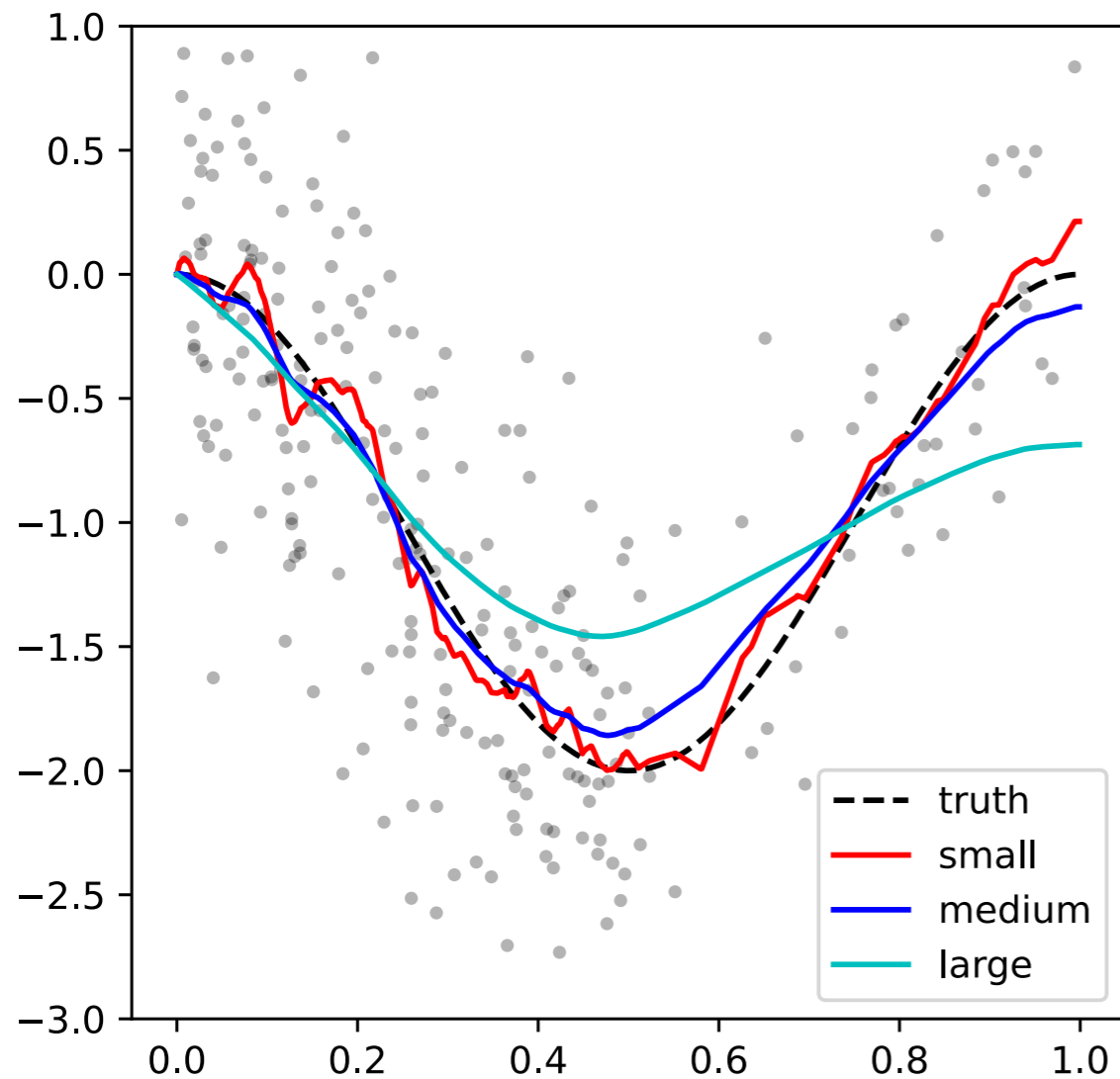
Pan and Yang (2010), Sugiyama and Kawanabe (2012)

Example: Kernel Ridge Regression



$$y = f^*(x) + \varepsilon$$

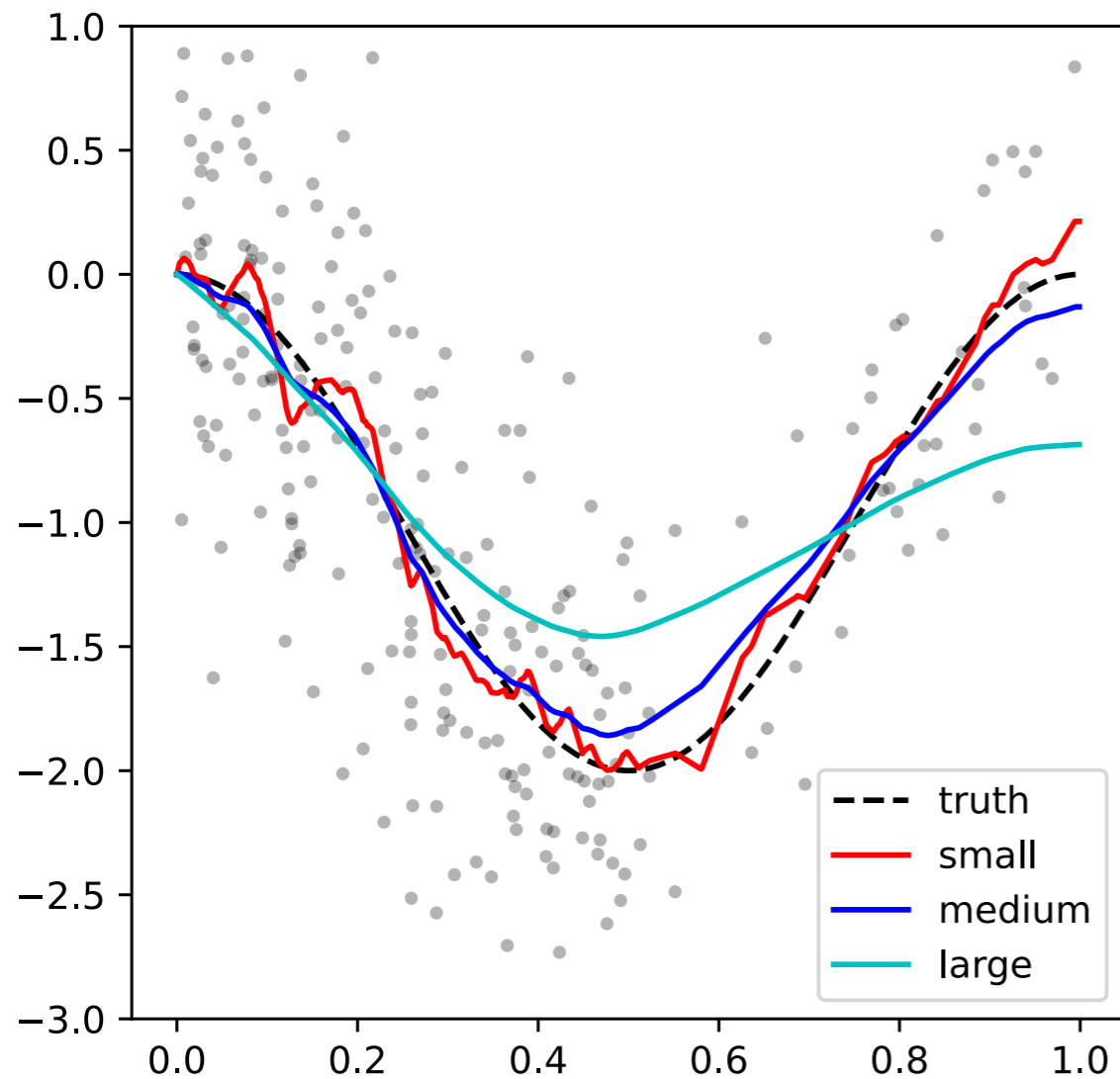
Example: Kernel Ridge Regression



KRR with **S**, **M**, **L** penalties.

Which one is the best?

Example: Kernel Ridge Regression



- If Q_X concentrates on $[0,0.5]$, **blue** is the best;
- If Q_X concentrates on $[0.5,1]$, **red** is the best.

The “best” model depends on Q .

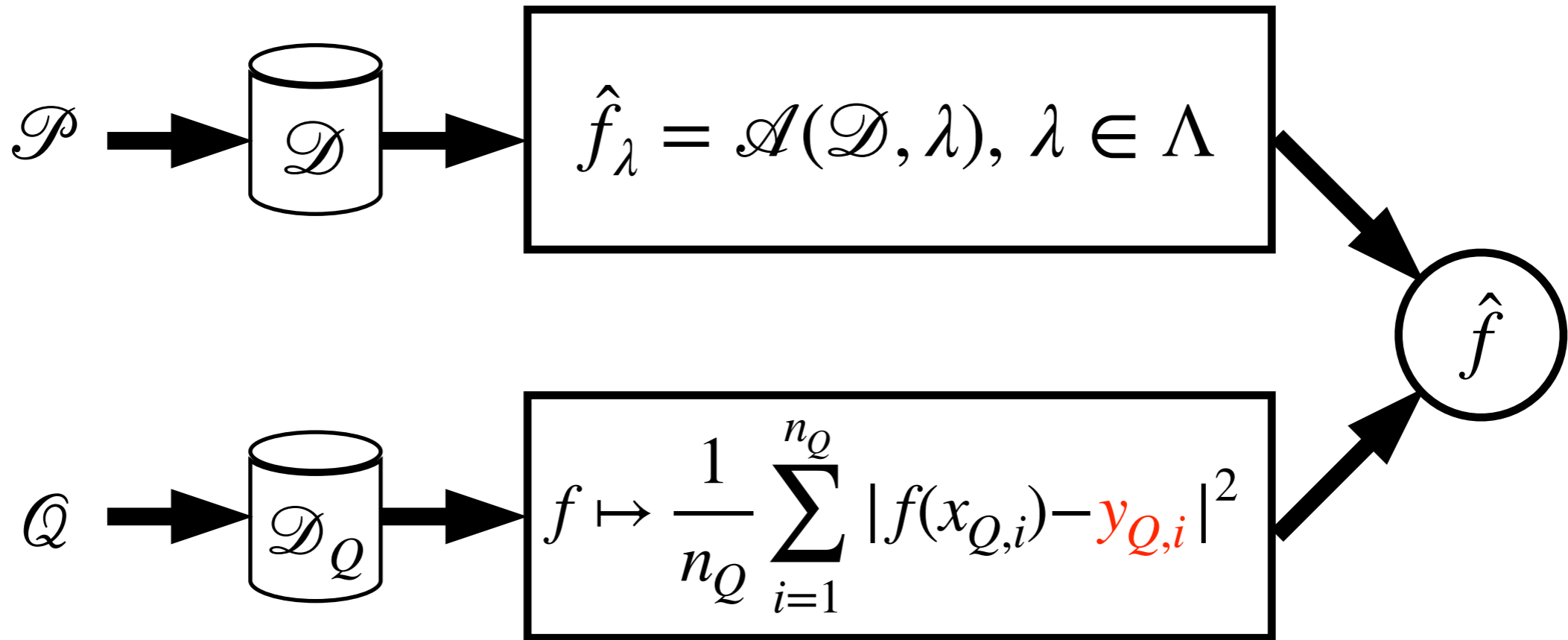
Outline

- Introduction: covariate shift adaptation
- **Methodology: validation with pseudo-labels**
- Model selection theory: general results
- Adaptivity guarantees: kernel ridge regression

If Target Data were Labeled...

Source

Candidates



Target

Risk estimator

Issue: $\{y_{Q,i}\}_{i=1}^{n_Q}$ are missing.

Regression Imputation

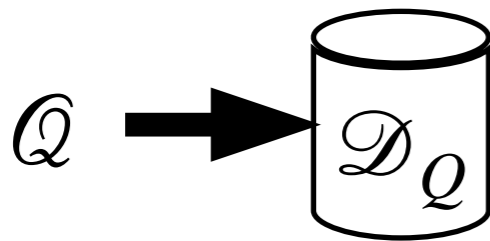
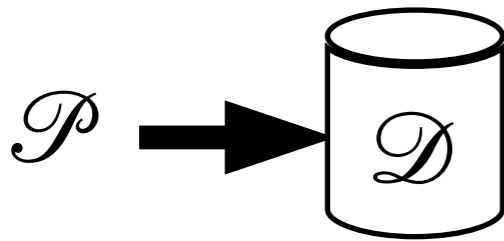
Idea: fill the missing labels $\{y_{Q,i}\}_{i=1}^{n_Q}$.

Train an **imputation model** \tilde{f} to label the data: $\tilde{y}_{Q,i} = \tilde{f}(x_{Q,i})$.

- **Pseudo-labeling and self-training** (Lee, 2013, Kumar et al., 2020, Cai et al., 2021, Liu et al., 2021)
- **Regression imputation** (Little and Rubin, 2019, Hirshberg et al., 2019, Kallus, 2020; Hirshberg and Wager, 2021; Mou et al., 2022)
- Different approaches: **propensity weighting and matching** (Shimodaira, 2000; Huang et al., 2006; Cortes et al., 2010; Ma et al., 2022)

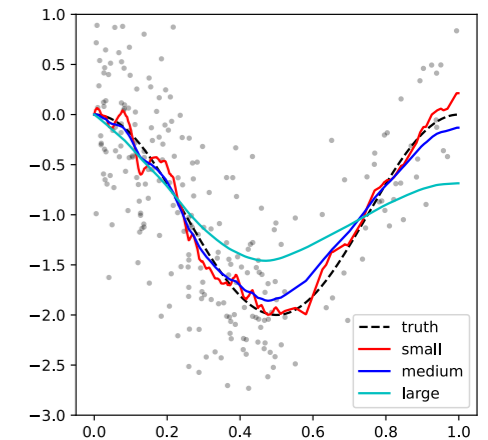
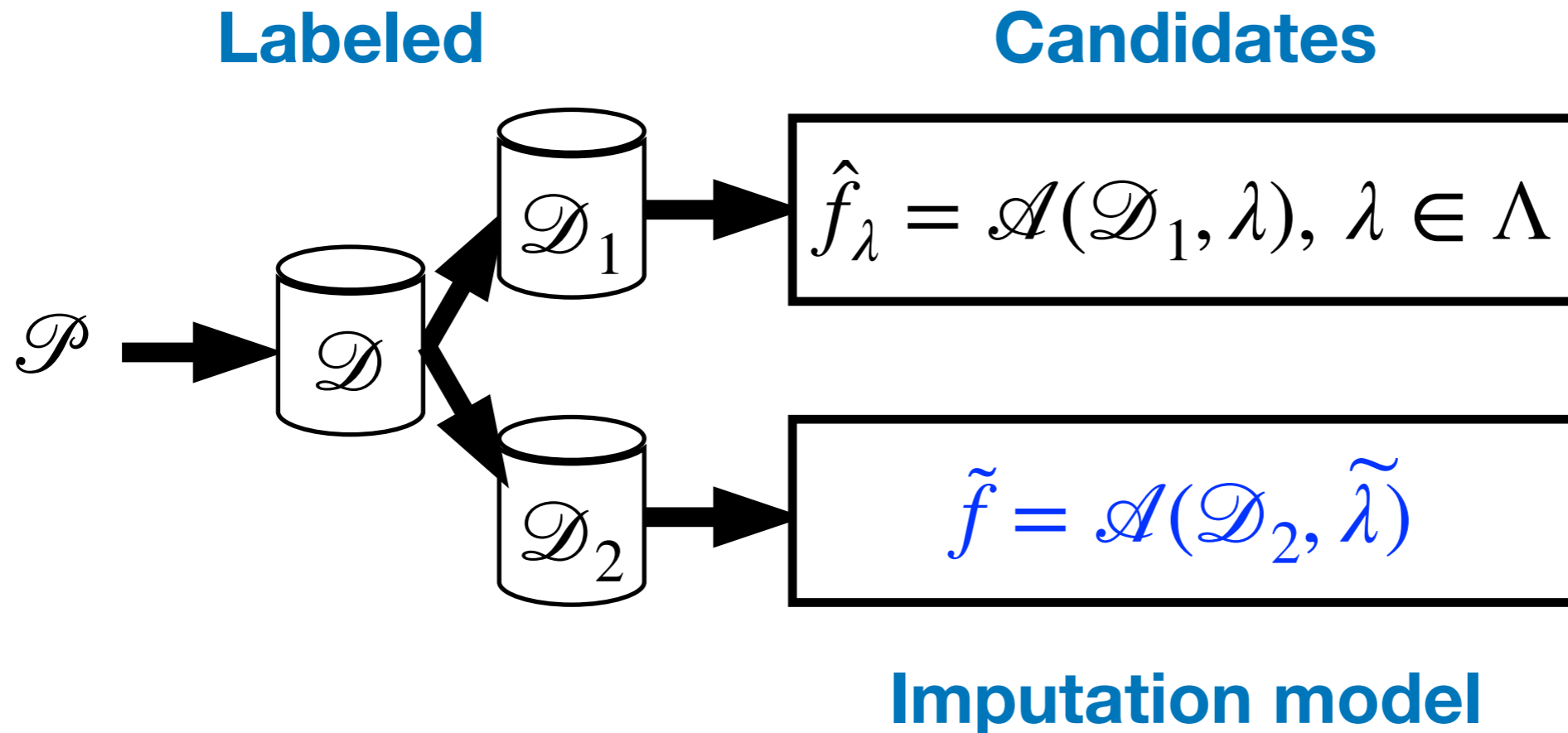
Methodology: Start

Labeled

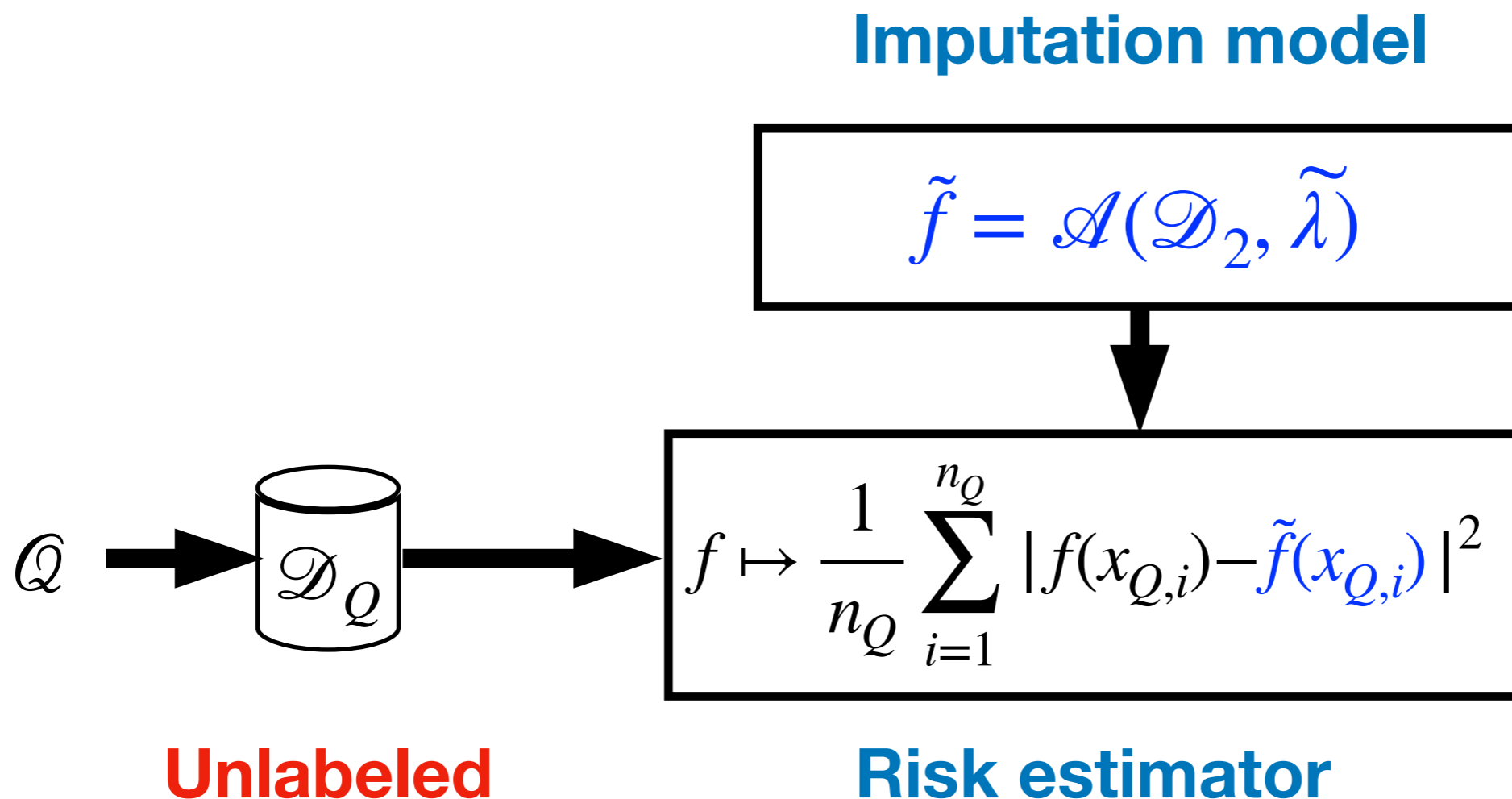


Unlabeled

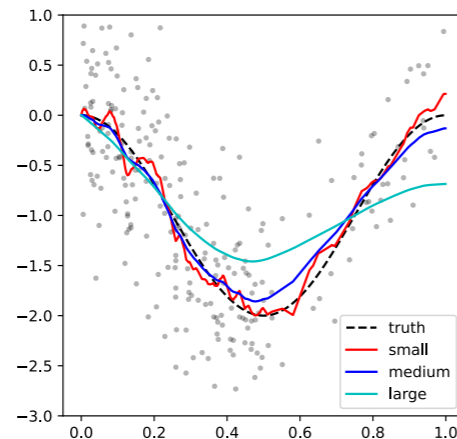
Training



Pseudo-Labeling (Imputation)



Model Selection

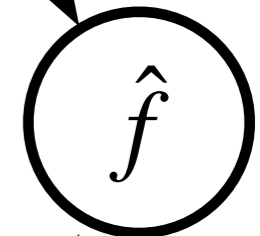


Candidates

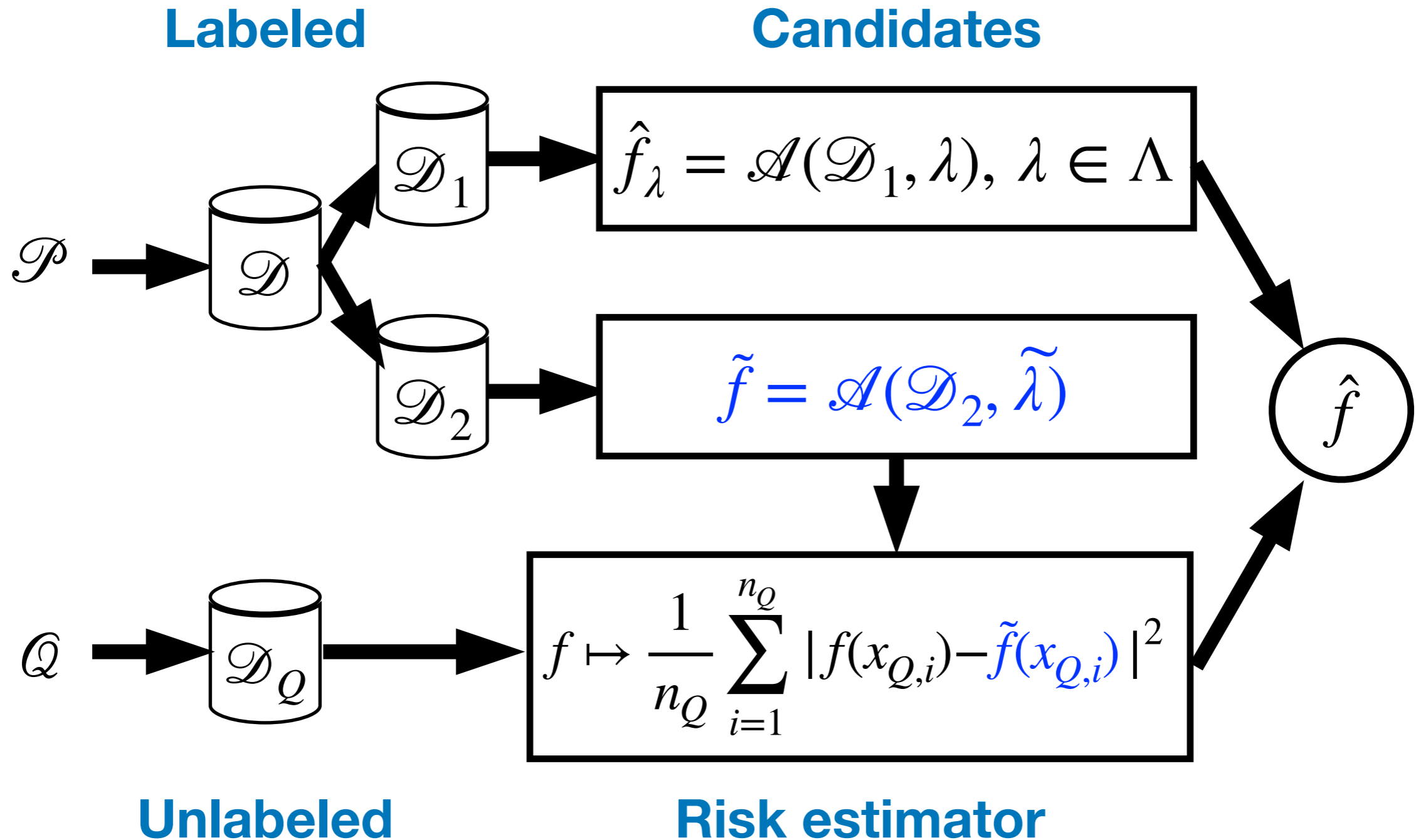
$$\hat{f}_\lambda = \mathcal{A}(\mathcal{D}_1, \lambda), \lambda \in \Lambda$$

$$f \mapsto \frac{1}{n_Q} \sum_{i=1}^{n_Q} |f(x_{Q,i}) - \tilde{f}(x_{Q,i})|^2$$

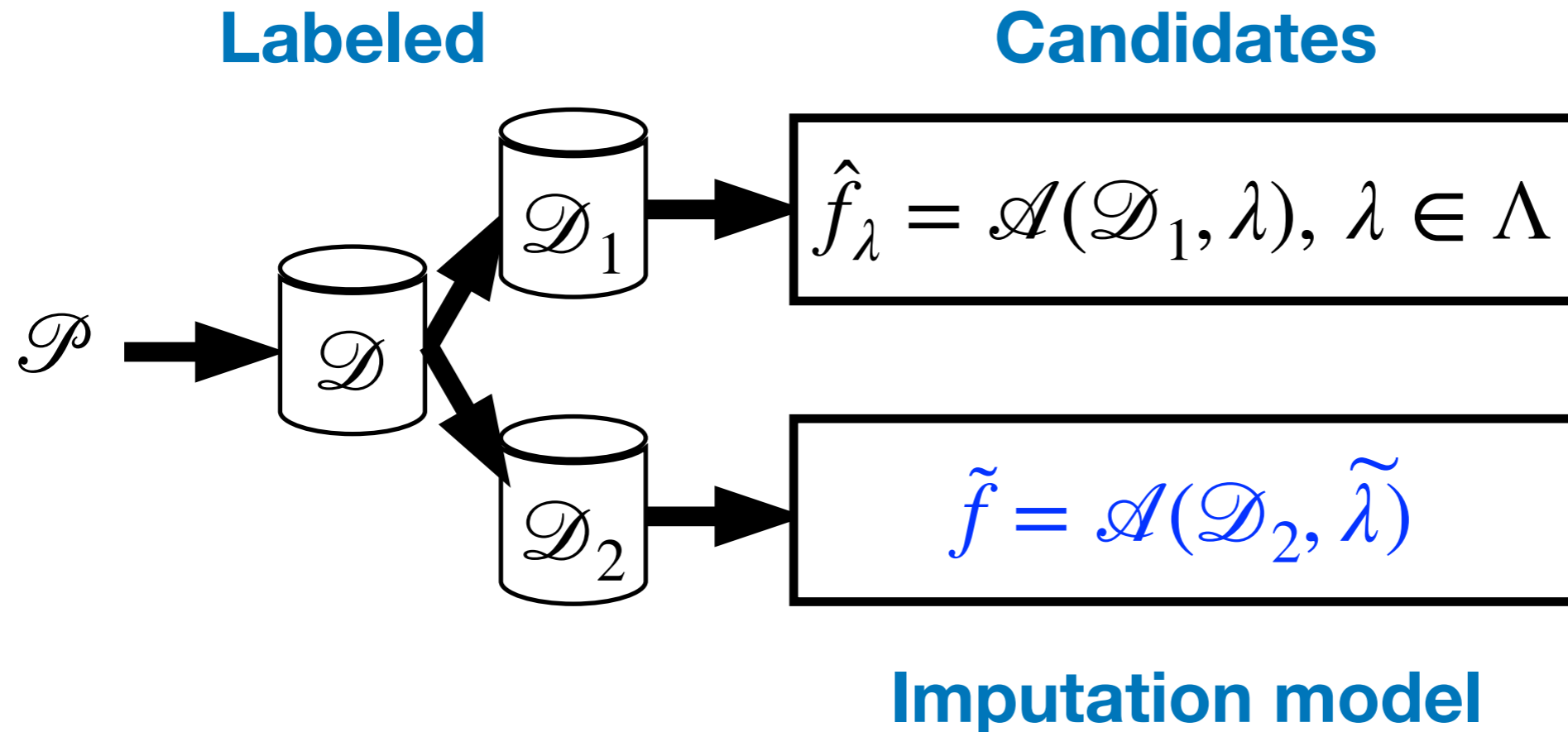
Risk estimator



Methodology



A Chicken-and-Egg Dilemma?



- To select a good λ , need a good \tilde{f} , which needs a good $\tilde{\lambda}$.

Outline

- Introduction: covariate shift adaptation
- Methodology: validation with pseudo-labels
- **Model selection theory: general results**
- Adaptivity guarantees: kernel ridge regression

Validation with Pseudo-Labels

$\{x_i\}_{i=1}^n$ deterministic, $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$.

Given $\{f_j\}_{j=1}^m$, select $\hat{j} \in \arg \min_{j \in [m]} \left\{ \frac{1}{n} \sum_{i=1}^n |f_j(x_i) - \tilde{y}_i|^2 \right\}$.

Validation with Pseudo-Labels

$\{x_i\}_{i=1}^n$ deterministic, $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$.

Given $\{f_j\}_{j=1}^m$, select $\hat{j} \in \arg \min_{j \in [m]} \left\{ \frac{1}{n} \sum_{i=1}^n |f_j(x_i) - \tilde{y}_i|^2 \right\}$.

Define $\mathcal{E}(f) = \frac{1}{n} \sum_{i=1}^n |f(x_i) - f^*(x_i)|^2$.

How does $\hat{f} = f_{\hat{j}}$ compare to the best in class?

Validation with Pseudo-Labels

$\{x_i\}_{i=1}^n$ deterministic, $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$.

Given $\{f_j\}_{j=1}^m$, select $\hat{j} \in \arg \min_{j \in [m]} \left\{ \frac{1}{n} \sum_{i=1}^n |f_j(x_i) - \tilde{y}_i|^2 \right\}$.

Define $\mathcal{E}(f) = \frac{1}{n} \sum_{i=1}^n |f(x_i) - f^*(x_i)|^2$.

Theorem

$$\mathcal{E}(\hat{f}) \lesssim \min_{j \in [m]} \mathcal{E}(f_j) + \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^*(x_i)|^2 + \frac{\|\mathbf{S}\| \log m}{n}.$$

Bias-Variance Decomposition

Impact of $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$ on model selection:

$$\mathcal{E}(\hat{f}) \lesssim \min_{j \in [m]} \mathcal{E}(f_j) + \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^*(x_i)|^2 + \frac{\|\mathbf{S}\| \log m}{n}.$$

Bias-Variance Decomposition

Impact of $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$ on model selection:

$$\mathcal{E}(\hat{f}) \lesssim \min_{j \in [m]} \mathcal{E}(f_j) + \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^\star(x_i)|^2 + \frac{\|\mathbf{S}\| \log m}{n}.$$

MSE of $(\tilde{y}_1, \dots, \tilde{y}_n)^\top$:

$$\mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n |\tilde{y}_i - f^\star(x_i)|^2 \right) = \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^\star(x_i)|^2 + \frac{\text{Tr}(\mathbf{S})}{n}.$$

- True labels: $\bar{y}_i = f^\star(x_i)$, $\mathbf{S} = \sigma^2 \mathbf{I}$.

Bias-Variance Decomposition

Impact of $(\tilde{y}_1, \dots, \tilde{y}_n)^\top \sim N(\bar{y}, \mathbf{S})$ on model selection:

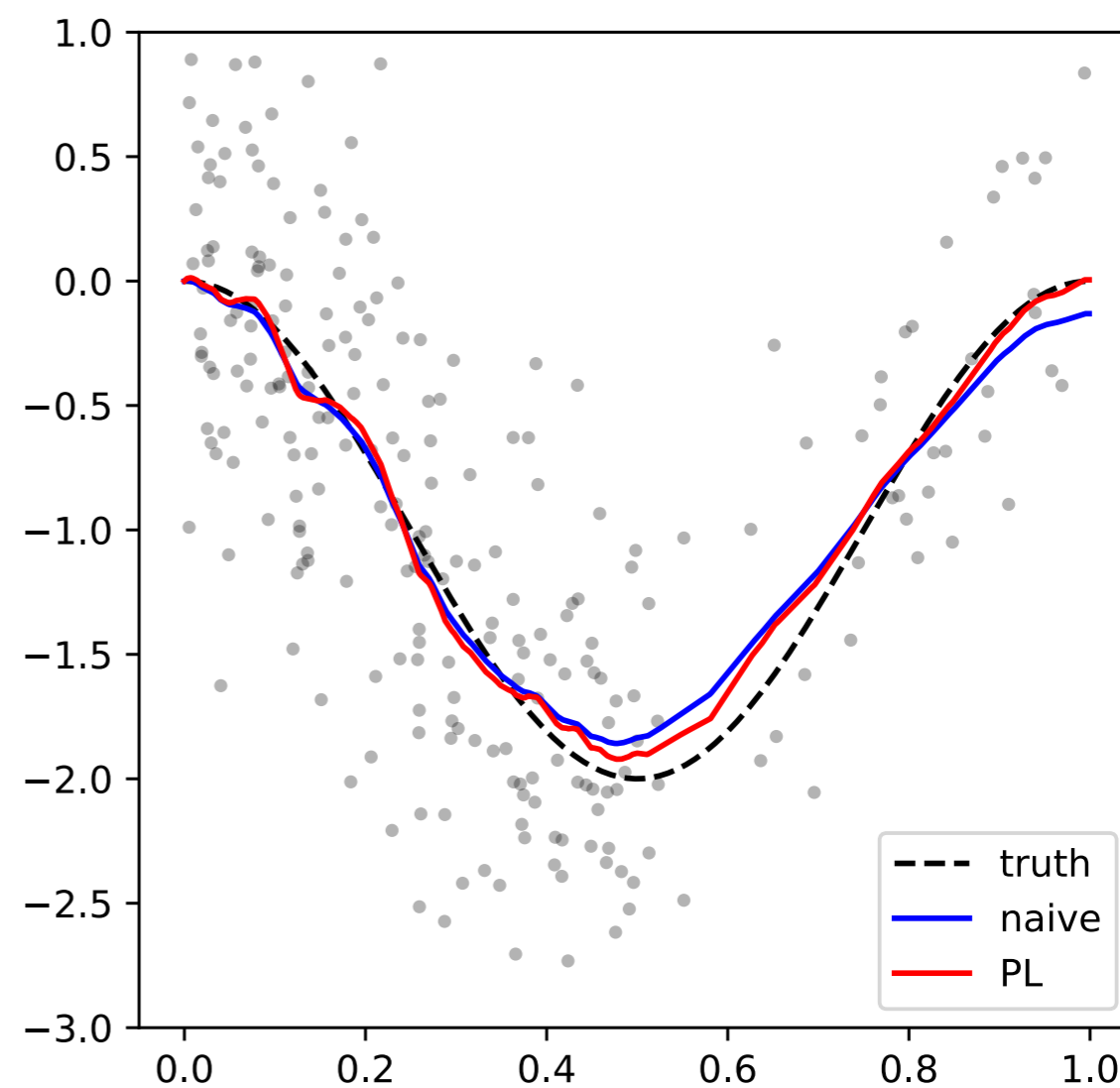
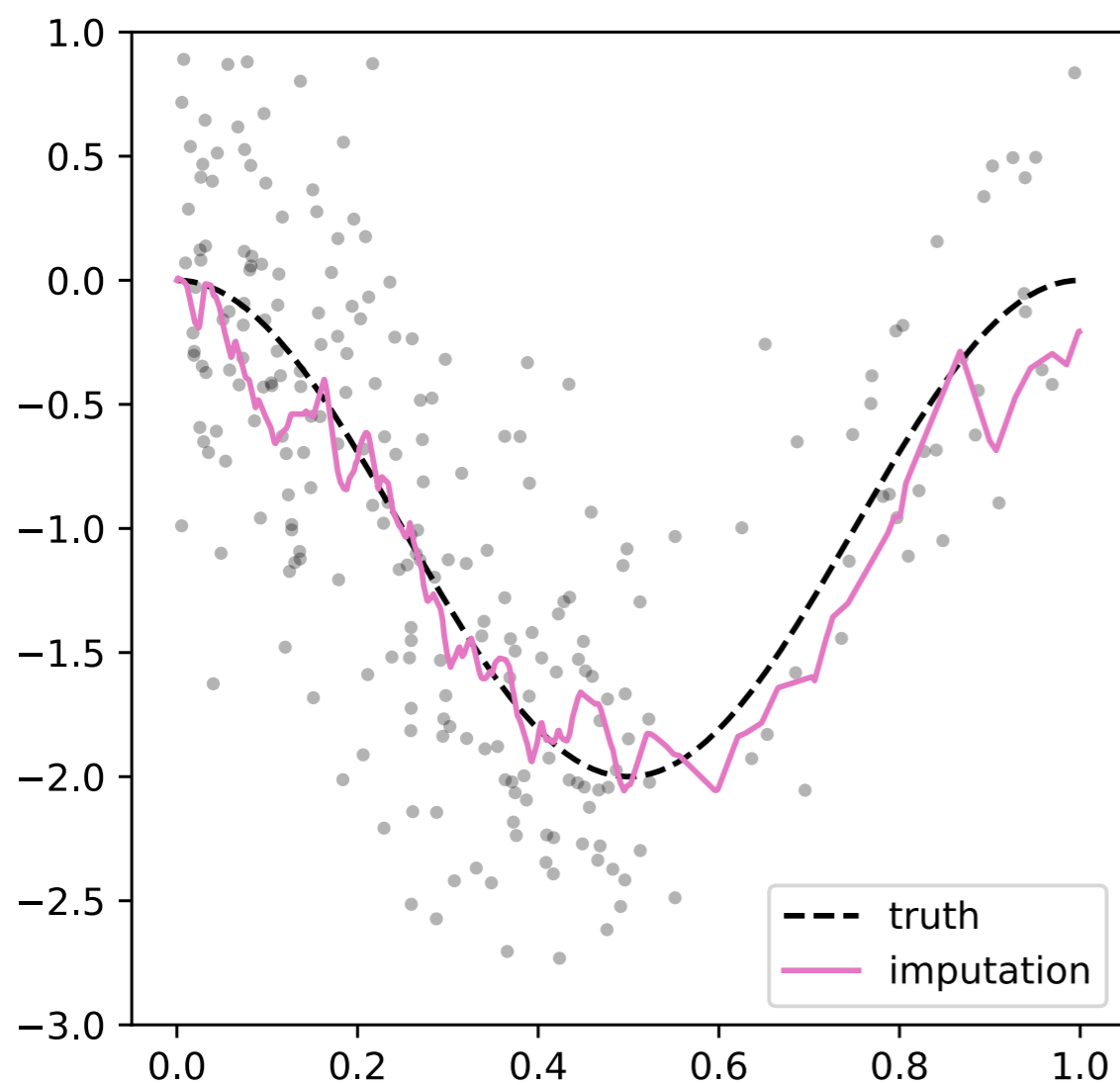
$$\mathcal{E}(\hat{f}) \lesssim \min_{j \in [m]} \mathcal{E}(f_j) + \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^\star(x_i)|^2 + \frac{\|\mathbf{S}\| \log m}{n}.$$

MSE of $(\tilde{y}_1, \dots, \tilde{y}_n)^\top$:

$$\mathbb{E} \left(\frac{1}{n} \sum_{i=1}^n |\tilde{y}_i - f^\star(x_i)|^2 \right) = \frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^\star(x_i)|^2 + \frac{\text{Tr}(\mathbf{S})}{n}.$$

- For model selection, variance has a much weaker effect.
- **Good pseudo-labels: low bias and reasonable variance.**

Example: Kernel Ridge Regression



Target: $Q_X = U[0.5, 1]$

Outline

- Introduction: covariate shift adaptation
- Methodology: validation with pseudo-labels
- Model selection theory: general results
- **Adaptivity guarantees: kernel ridge regression**

Kernel Ridge Regression

Kernel ridge regression:

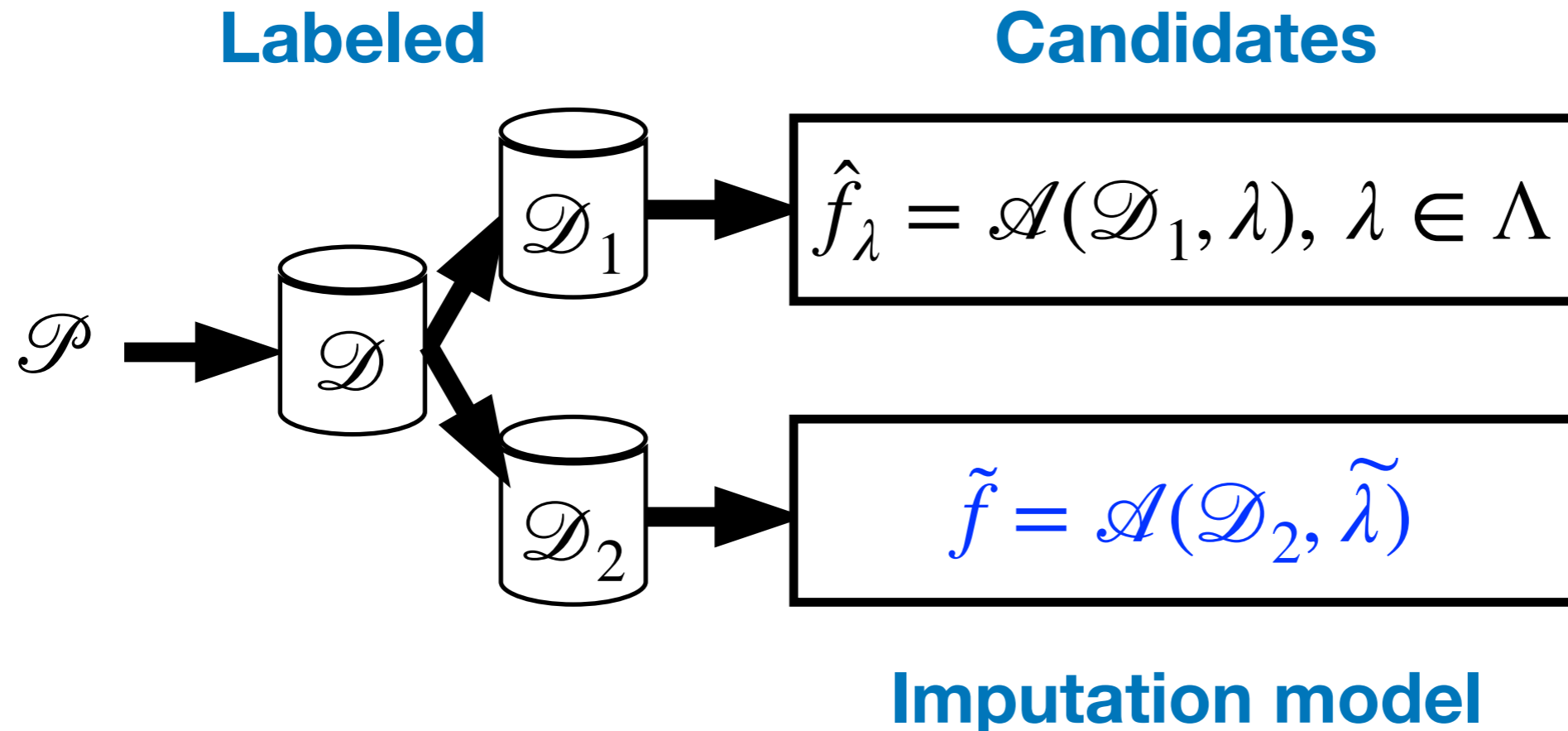
$$\mathcal{A}(\mathcal{D}, \lambda) = \operatorname{argmin}_{f \in \mathcal{F}} \left\{ \frac{1}{|\mathcal{D}|} \sum_{(x,y) \in \mathcal{D}} |f(x) - y|^2 + \lambda \|f\|_{\mathcal{F}}^2 \right\}.$$

\mathcal{F} : **reproducing kernel Hilbert space (RKHS)** containing f^\star .

For some Hilbert space \mathbb{H} and mapping $\phi : \mathcal{X} \rightarrow \mathbb{H}$,

KRR = ridge regression in $\{f_\theta(\cdot) = \langle \phi(\cdot), \theta \rangle : \theta \in \mathbb{H}\}$.

Undersmoothing for Imputation



Requirement on \tilde{f} : small bias, reasonable variance.

- $\tilde{\lambda} \asymp n_P^{-1}, \Lambda \asymp \{n_P^{-1}, 2n_P^{-1}, 2^2n_P^{-1}, \dots, 1\}.$

Main Result

Theorem (informal)

$\hat{f} \approx$ KRR on n_e labeled samples from \mathcal{Q} , with optimal λ .

- n_e : effective sample size, gauges transfer from \mathcal{P} to \mathcal{Q} .
- Adaptation to the spectrum of \mathcal{Q} and the covariate shift.

Excess Risk Bounds

Theorem

Let $\Sigma_Q = \mathbb{E}_Q[\phi(x) \otimes \phi(x)]$, and $\{\mu_j\}_{j=1}^{\infty}$ be its eigenvalues.

- If $\text{rank}(\Sigma_Q) \leq D$, then $R(\hat{f}) - R(f^*) \lesssim Dn_e^{-1} + n_Q^{-1}$;
- If $\mu_j \lesssim e^{-cj}$, then $R(\hat{f}) - R(f^*) \lesssim n_e^{-1} + n_Q^{-1}$;
- If $\mu_j \lesssim j^{-2\alpha}$, then $R(\hat{f}) - R(f^*) \lesssim n_e^{-\frac{2\alpha}{2\alpha+1}} + n_Q^{-1}$.

Excess Risk Bounds

Theorem

Let $\Sigma_Q = \mathbb{E}_Q[\phi(x) \otimes \phi(x)]$, and $\{\mu_j\}_{j=1}^{\infty}$ be its eigenvalues.

- If $\text{rank}(\Sigma_Q) \leq D$, then $R(\hat{f}) - R(f^*) \lesssim Dn_e^{-1} + n_Q^{-1}$;
- If $\mu_j \lesssim e^{-cj}$, then $R(\hat{f}) - R(f^*) \lesssim n_e^{-1} + n_Q^{-1}$;
- If $\mu_j \lesssim j^{-2\alpha}$, then $R(\hat{f}) - R(f^*) \lesssim n_e^{-\frac{2\alpha}{2\alpha+1}} + n_Q^{-1}$.

- First term: **optimal** rate for n_e labeled samples from \mathcal{Q} .
- Second term: **minor** cost of using limited samples from \mathcal{Q} .

Effective Sample Size

Definition

$$n_e = \sup\{t \leq n_P : t\boldsymbol{\Sigma}_Q \leq n_P\boldsymbol{\Sigma}_P + \mathbf{I}\}.$$

The followings have the same excess risk bound:

- KRR on n_P **labeled source samples**, with \mathcal{Q} -optimal λ ;
- KRR on n_e **labeled target samples**, with \mathcal{Q} -optimal λ .

Pseudo-labeling helps select a near-optimal λ .

Adaptivity

Theorem (informal)

$\hat{f} \approx$ optimal KRR on n_e labeled samples from \mathcal{Q} .

$$n_e = \sup\{t \leq n_P : t\mathbf{\Sigma}_Q \leq n_P\mathbf{\Sigma}_P + \mathbf{I}\}.$$

Example: $\frac{d\mathcal{Q}_X}{d\mathcal{P}_X} \leq B \Rightarrow \mathbf{\Sigma}_Q \leq B\mathbf{\Sigma}_P \Rightarrow n_e \geq n_P/B.$

Adaptivity

Theorem (informal)

$\hat{f} \approx$ optimal KRR on n_e labeled samples from \mathcal{Q} .

$$n_e = \sup\{t \leq n_P : t\mathbf{\Sigma}_Q \leq n_P\mathbf{\Sigma}_P + \mathbf{I}\}.$$

Example: $\frac{d\mathcal{Q}_X}{d\mathcal{P}_X} \leq B \Rightarrow \mathbf{\Sigma}_Q \leq B\mathbf{\Sigma}_P \Rightarrow n_e \geq n_P/B.$

- Optimal excess risk bound. **Adaptive to B & spectral decay.**
- Matching lower bound: Ma, Pathak and Wainwright, 2023.

Adaptivity

Theorem (informal)

$\hat{f} \approx$ optimal KRR on n_e labeled samples from \mathcal{Q} .

$$n_e = \sup \{ t \leq n_P : t \Sigma_Q \leq n_P \Sigma_P + \mathbf{I} \}.$$

Example: $\mathcal{P}_X = U[0,1]$, $\mathcal{Q}_X = \delta_{x_0}$, 1st-order Sobolev space

- $n_e \asymp ?$

Adaptivity

Theorem (informal)

$\hat{f} \approx$ optimal KRR on n_e labeled samples from \mathcal{Q} .

$$n_e = \sup\{t \leq n_P : t\mathbf{\Sigma}_Q \leq n_P\mathbf{\Sigma}_P + \mathbf{I}\}.$$

Example: $\mathcal{P}_X = U[0,1]$, $\mathcal{Q}_X = \delta_{x_0}$, 1st-order Sobolev space

- $n_e \asymp n_P^{1/2} \Rightarrow R(\hat{f}) - R(f^*) \lesssim n_e^{-1} \asymp n_P^{-1/2}$, optimal.
- **Adaptive to singularity!**
- Each \hat{f}_λ is **sub-optimal** for either \mathcal{P}_X or \mathcal{Q}_X .

Summary

- Covariate shift adaptation given unlabeled target data
- **Method: hold-out validation with pseudo-labels**
- General requirement: low bias, reasonable variance
- Analysis of kernel ridge regression: Adaptivity

Wang. Pseudo-Labeling for Kernel Ridge Regression under Covariate Shift. arXiv:2302.10160, 2023.

Supported by NSF DMS-2210907, Columbia University DSI SF-181.

Q & A

Thank you!

Validation with Pseudo-Labels

$\{(x_i, y_i)\}_{i=1}^n$, x_i deterministic, $\mathbf{y} = (y_1, \dots, y_n)^\top \sim N(\bar{\mathbf{y}}, \mathbf{S})$.

Given $\{f_j\}_{j=1}^m$, select $\hat{j} \in \arg \min_{j \in [m]} \left\{ \frac{1}{n} \sum_{i=1}^n |f(x_i) - y_i|^2 \right\}$.

Define $\mathcal{E}(f) = \frac{1}{n} \sum_{i=1}^n |f(x_i) - f^*(x_i)|^2$.

Theorem

$$\sqrt{\mathcal{E}(\hat{f})} - \min_{j \in [m]} \sqrt{\mathcal{E}(f_j)} \lesssim \sqrt{\frac{1}{n} \sum_{i=1}^n |\bar{y}_i - f^*(x_i)|^2} + \sqrt{\frac{\|\mathbf{S}\| \log m}{n}}.$$

Challenges

Example: $\mathcal{X} = [0,1]$, $\mathcal{P}_X = U[0,1]$, Sobolev kernel.

- $Q_X = \mathcal{P}_X$ (no shift): $\lambda \asymp n^{-2/3}$ is minimax optimal.
- $Q_X = \delta_{x_0}$ (e.g., conditional ATE): $\lambda \asymp n^{-1}$.

Speckman (1979), Li (1982), Donoho (1994), Wainwright (2019).