

Near-Optimal Linear Regression under Distribution Shift

Qi Lei

Princeton University

Joint work with
Wei Hu and Jason Lee.

Setup

- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
 - Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
 - Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$
-
- Covariate shift: $f_S = f_T, p_S \neq p_T$
 - Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$

Setup

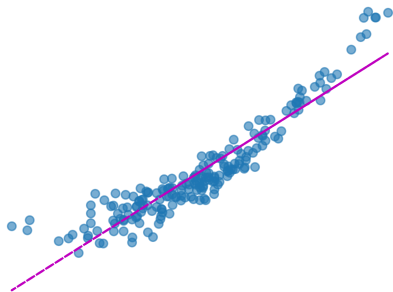
- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$

- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$

On the Problem of Distributions shift

Setup

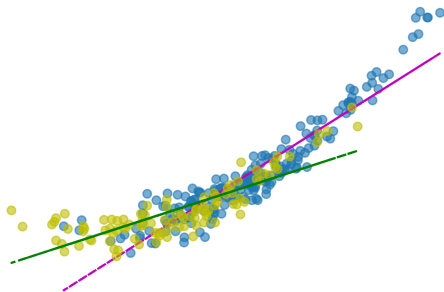
- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$
- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$



On the Problem of Distributions shift

Setup

- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$
- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$



Setup

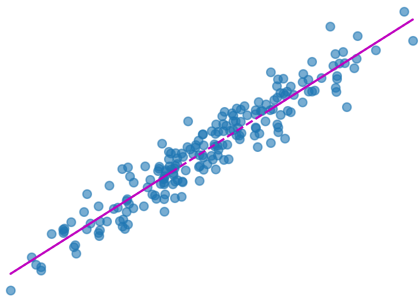
- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}$, $\mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$

- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$

On the Problem of Distributions shift

Setup

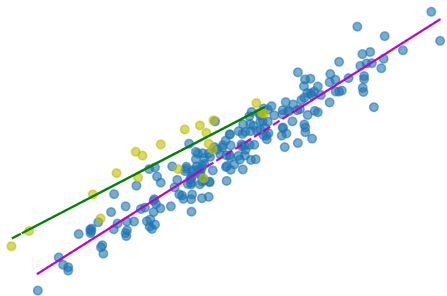
- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$
- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$



On the Problem of Distributions shift

Setup

- Observations: $y_i = f_S(\mathbf{x}_i) + \text{noise}, \mathbf{x}_i \sim p_S, i = 1, 2, \dots, n_S$
- Estimate a linear model $\hat{\beta}((\mathbf{x}_i, y_i)_{i=1}^{n_S})$
- Tested on: $\mathbb{E}_{\mathbf{x} \sim p_T} [\|f_T(\mathbf{x}) - \hat{\beta}^\top \mathbf{x}\|^2]$
- Covariate shift: $f_S = f_T, p_S \neq p_T$
- Model shift: $f_S \neq f_T$ linear, $p_S \neq p_T$



Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on source distribution $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on source distribution $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.^a

^aSee e.g. Murphy, K. P. (2012). Machine learning: a probabilistic perspective

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on source distribution $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.^a

^aSee e.g. Murphy, K. P. (2012). Machine learning: a probabilistic perspective

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on source distribution $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.^a

^aSee e.g. Murphy, K. P. (2012). Machine learning: a probabilistic perspective

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on source distribution $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.^a

^aSee e.g. Murphy, K. P. (2012). Machine learning: a probabilistic perspective

How to adapt to the target domain?

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on **source distribution** $\mathbf{x} \sim p_S$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on **target distribution** $\mathbf{x} \sim p_T$,
- with Gaussian prior $\beta^* \sim \mathcal{N}(0, r^2 I)$.

Interestingly, this does not give us different estimator.

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on **target distribution** $\mathbf{x} \sim p_T$,
- **with Gaussian prior** $\beta^* \sim \mathcal{N}(0, r^2 I)$.

Ridge Regression as MAP Inference

$$\hat{\beta}_{\text{RR}} \leftarrow \arg \min_{\beta} \frac{1}{n_S} \sum_{i=1}^{n_S} (\beta^\top \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \|\beta\|^2 \quad (1)$$

⇐ This is the optimal linear estimator

- assuming $y \sim \mathcal{N}(\mathbf{x}^\top \beta^*, \sigma^2)$,
- measured on **target distribution** $\mathbf{x} \sim p_T$,
- **for the worse-case** $\beta^*, \|\beta^*\| \leq r$.

This gives the "minimax" linear estimator. We developed a meta algorithm for this mechanism.

Step 1: Find a sufficient statistic $\hat{\beta}_{SS}$ for the optimal linear estimator given the observations

Step 2: Solve the best estimator that is linear in $\hat{\beta}_{SS}$ in the worst-case setting:

$$\hat{\beta}_{MM} \leftarrow \arg \min_{\beta \text{ linear in } \hat{\beta}_{SS}} \max_{\|\beta^*\| \leq r} \mathbb{E}_{\text{noise}, x \sim p_T} (x^\top (\beta - \beta^*))^2.$$

- Step 1:** Find a sufficient statistic $\hat{\beta}_{SS}$ for the optimal linear estimator given the observations
- Step 2:** Solve the best estimator that is linear in $\hat{\beta}_{SS}$ in the worse-case setting:

$$\hat{\beta}_{MM} \leftarrow \arg \min_{\beta \text{ linear in } \hat{\beta}_{SS}} \max_{\|\beta^*\| \leq r} \mathbb{E}_{\text{noise}, \mathbf{x} \sim p_T} (\mathbf{x}^\top (\beta - \beta^*))^2.$$

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - 1 developed algorithms with only unlabeled target data and/or
 - 2 with few labeled target data,
 - 3 show that our algorithm is optimal among all linear estimators, and
 - 4 is within constant of the best nonlinear estimators under some conditions,
 - 5 prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - 1 developed algorithm that is
 - 2 asymptotically optimal among all linear estimators,
 - 3 provide a practical way to estimate the selection bias on source distribution,
 - 4 tested on real dataset
- Small model shift under Gaussian sequence model
 - 1 developed algorithms that balance model shift and variance,
 - 2 prove optimality among all linear estimators, and
 - 3 is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Results on Different Settings

- Covariate shift under Gaussian sequence model
 - ① developed algorithms with only unlabeled target data and/or
 - ② with few labeled target data,
 - ③ show that our algorithm is optimal among all linear estimators, and
 - ④ is within constant of the best nonlinear estimators under some conditions,
 - ⑤ prove a separation result between ours and ridge regression
- Covariate shift with approximation error (nonlinear model)
 - ① developed algorithm that is
 - ② asymptotically optimal among all linear estimators,
 - ③ provide a practical way to estimate the selection bias on source distribution,
 - ④ tested on real dataset
- Small model shift under Gaussian sequence model
 - ① developed algorithms that balance model shift and variance,
 - ② prove optimality among all linear estimators, and
 - ③ is within constant of the best nonlinear estimators under some conditions

Thank you!