

Can SGD Select Good Fishermen?

Alkis Kalavasis, Anay Mehrotra, Felix Zhou

Yale University

Image generated by Sora.

Collaborators



Alkis Kalavasis
(Yale University)



Anay Mehrotra
(Yale University)

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Hunters vs Fishermen

In a small village, two mutually exclusive occupations are available: hunting and fishing.

Hunters vs Fishermen

In a small village, two mutually exclusive occupations are available: hunting and fishing.

Simple question:

What makes a good fisherman and what makes a good hunter?

Hunters vs Fishermen

In a small village, two mutually exclusive occupations are available: hunting and fishing.

Simple question:

What makes a good fisherman and what makes a good hunter?

1. Collect random sample of hunters and fishermen from the village.

Hunters vs Fishermen

In a small village, two mutually exclusive occupations are available: hunting and fishing.

Simple question:

What makes a good fisherman and what makes a good hunter?

1. Collect random sample of hunters and fishermen from the village.
2. Record relevant features and income.

Hunters vs Fishermen

In a small village, two mutually exclusive occupations are available: hunting and fishing.

Simple question:

What makes a good fisherman and what makes a good hunter?

1. Collect random sample of hunters and fishermen from the village.
2. Record relevant features and income.
3. Estimate parameters of 2 linear models, one per occupation.

Hunters vs Fishermen

The resulting linear fits can be biased due to [selection bias](#).

Hunters vs Fishermen

The resulting linear fits can be biased due to [selection bias](#).

- ▶ Better hunters will opt to hunt, and vice-versa.

Hunters vs Fishermen

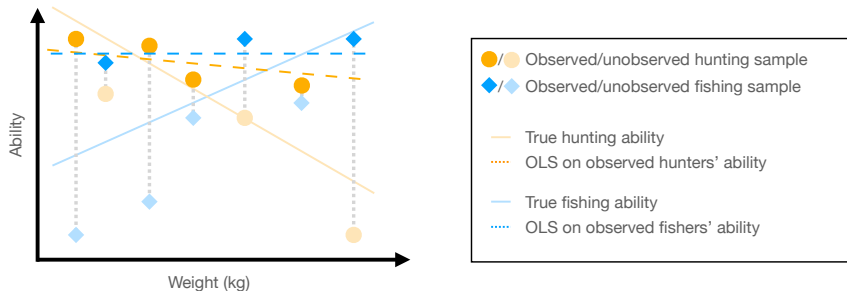
The resulting linear fits can be biased due to **selection bias**.

- ▶ Better hunters will opt to hunt, and vice-versa.
- ▶ Never observe fishing earnings of an individual better at hunting!

Hunters vs Fishermen

The resulting linear fits can be biased due to **selection bias**.

- ▶ Better hunters will opt to hunt, and vice-versa.
- ▶ Never observe fishing earnings of an individual better at hunting!



Courtesy of Cherapanamjeri, Daskalakis, Ilyas, Zampetakis [CDIZ23].

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],
- ▶ Learning from auction data [AH02; AH07; CDIZ22].

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],
- ▶ Learning from auction data [AH02; AH07; CDIZ22].
- ▶ Studies of participation in the labor force [Hec74; Han76; Nel77; Hec79; Cog80; Han80]

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],
- ▶ Learning from auction data [AH02; AH07; CDIZ22].
- ▶ Studies of participation in the labor force [Hec74; Han76; Nel77; Hec79; Cog80; Han80]
- ▶ Studies of migration and income [NZ80; Bor87]

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],
- ▶ Learning from auction data [AH02; AH07; CDIZ22].
- ▶ Studies of participation in the labor force [Hec74; Han76; Nel77; Hec79; Cog80; Han80]
- ▶ Studies of migration and income [NZ80; Bor87]
- ▶ Studies of the effect of unions on wages [Lee78; AF82]

History: Inference under Selection Bias

Rich history in Statistics and Econometrics, starting with foundational works of Roy [Roy51], Heckman [Hec79], Willis and Rosen [WR79], Fair and Jaffe [FJ72], and has since found many applications:

- ▶ Causal inference and imitation learning [Hec90]
- ▶ Learning from strategically reported data [HMPW16; DRSW+18; KR20],
- ▶ Learning from auction data [AH02; AH07; CDIZ22].
- ▶ Studies of participation in the labor force [Hec74; Han76; Nel77; Hec79; Cog80; Han80]
- ▶ Studies of migration and income [NZ80; Bor87]
- ▶ Studies of the effect of unions on wages [Lee78; AF82]
- ▶ Studies of returns on education [GHH78; KLMT79; WR79]

Model: Regression under Self-Selection Biases

Goal: recover unknown regressors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ to error $\varepsilon > 0$ given observations $(\mathbf{x}_1, y_1^{\max}), \dots, (\mathbf{x}_n, y_n^{\max})$.

Model: Regression under Self-Selection Biases

Goal: recover unknown regressors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ to error $\varepsilon > 0$ given observations $(\mathbf{x}_1, y_1^{\max}), \dots, (\mathbf{x}_n, y_n^{\max})$.

DEFINITION 1 (MAXIMUM SELF-SELECTION MODEL [CDIZ23])

An observation (\mathbf{x}, y^{\max}) is generated as follows:

Model: Regression under Self-Selection Biases

Goal: recover unknown regressors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ to error $\varepsilon > 0$ given observations $(\mathbf{x}_1, y_1^{\max}), \dots, (\mathbf{x}_n, y_n^{\max})$.

DEFINITION 1 (MAXIMUM SELF-SELECTION MODEL [CDIZ23])

An observation (\mathbf{x}, y^{\max}) is generated as follows:

1. $\mathbf{x} \sim \mathcal{N}(0, I_d)$.

Model: Regression under Self-Selection Biases

Goal: recover unknown regressors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ to error $\varepsilon > 0$ given observations $(\mathbf{x}_1, y_1^{\max}), \dots, (\mathbf{x}_n, y_n^{\max})$.

DEFINITION 1 (MAXIMUM SELF-SELECTION MODEL [CDIZ23])

An observation (\mathbf{x}, y^{\max}) is generated as follows:

1. $\mathbf{x} \sim \mathcal{N}(0, I_d)$.
2. $y_i = \mathbf{w}_i^\top \mathbf{x} + \xi_i$ where $\xi_i \sim_{i.i.d.} \mathcal{N}(0, 1)$.

Model: Regression under Self-Selection Biases

Goal: recover unknown regressors $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ to error $\varepsilon > 0$ given observations $(\mathbf{x}_1, y_1^{\max}), \dots, (\mathbf{x}_n, y_n^{\max})$.

DEFINITION 1 (MAXIMUM SELF-SELECTION MODEL [CDIZ23])

An observation (\mathbf{x}, y^{\max}) is generated as follows:

1. $\mathbf{x} \sim \mathcal{N}(0, I_d)$.
2. $y_i = \mathbf{w}_i^\top \mathbf{x} + \xi_i$ where $\xi_i \sim_{i.i.d.} \mathcal{N}(0, 1)$.
3. Observe (\mathbf{x}, y^{\max}) where $y^{\max} = \max\{y_1, \dots, y_k\}$.

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Prior Works: Regression under Self-Selection Biases

Efficient algorithms for *finite samples* were not known until recently.

Prior Works: Regression under Self-Selection Biases

Efficient algorithms for *finite samples* were not known until recently.

[CDIZ23] Cherapanamjeri, Daskalakis, Ilyas, Zampetakis [STOC'23] designed a moment-based algorithm with $\text{poly}(d) \cdot \exp(k/\epsilon)$ sample complexity and running time.

Prior Works: Regression under Self-Selection Biases

Efficient algorithms for *finite samples* were not known until recently.

[CDIZ23] Cherapanamjeri, Daskalakis, Ilyas, Zampetakis [STOC'23] designed a moment-based algorithm with $\text{poly}(d) \cdot \exp(k/\epsilon)$ sample complexity and running time.

[GM24] Gaitonde and Mossel also used moments to design an algorithm with $\text{poly}(d, k, 1/\epsilon)$ sample complexity but $\text{poly}(d) + (1/\epsilon)^{\tilde{O}(k)}$ running time.

Prior Works: Regression under Self-Selection Biases

Efficient algorithms for *finite samples* were not known until recently.

[CDIZ23] Cherapanamjeri, Daskalakis, Ilyas, Zampetakis [STOC'23] designed a moment-based algorithm with $\text{poly}(d) \cdot \exp(k/\epsilon)$ sample complexity and running time.

[GM24] Gaitonde and Mossel also used moments to design an algorithm with $\text{poly}(d, k, 1/\epsilon)$ sample complexity but $\text{poly}(d) + (1/\epsilon)^{\tilde{O}(k)}$ running time.

Question: Polynomial number of samples is sufficient. *Can we design an algorithm with polynomial running time?*

Table of Contents

Regression under Self-Selection Biases

Prior Works

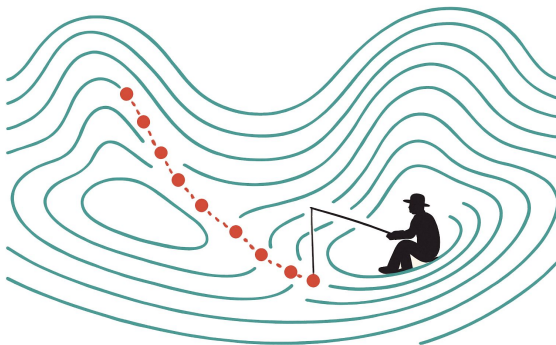
Our Contributions

Technical Overview

Our Results

THEOREM 1 (KALAVASIS, MEHROTRA, Z. '25)

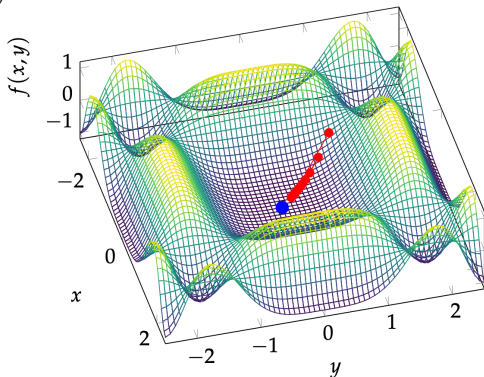
There is an algorithm for regression under self-selection bias with $\text{poly}(d, 1/\epsilon, k)$ sample complexity and $\text{poly}(d, 1/\epsilon) + 2^{\tilde{O}(k)}$ running time.



Our Results

THEOREM 2 (KALAVASIS, MEHROTRA, Z. '25)

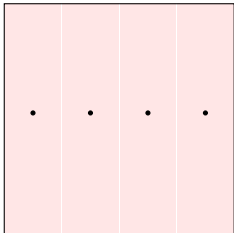
There is an SGD-based local convergence algorithm for regression under self-selection bias with $\text{poly}(d, 1/\epsilon, k)$ sample complexity and running time, given a $\text{poly}(1/k)$ -warm start.



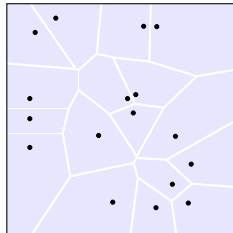
1. Unexpected connection to learning with “*coarse* observations”
[Fotakis, Kalavasis, Kontonis, Tzamos, 2021].

1. Unexpected connection to learning with “*coarse* observations”
[Fotakis, Kalavasis, Kontonis, Tzamos, 2021].
 - Simplifying example: instead of observing $z \sim \mathcal{N}(\mu^*, I)$, we observe a set from some given partition containing z . Can we recover μ^* ?

Key Idea



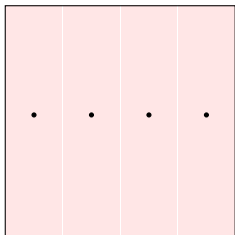
(a) Non-Identifiable Case



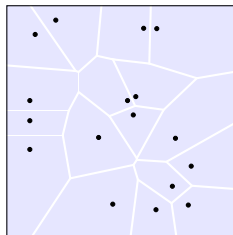
(b) Convex Partition Case

1. Unexpected connection to learning with “*coarse* observations” [Fotakis, Kalavasis, Kontonis, Tzamos, 2021].
 - Simplifying example: instead of observing $z \sim \mathcal{N}(\mu^*, I)$, we observe a set from some given partition containing z . Can we recover μ^* ?

Key Idea



(a) Non-Identifiable Case



(b) Convex Partition Case

1. Unexpected connection to learning with “*coarse* observations” [Fotakis, Kalavasis, Kontonis, Tzamos, 2021].
 - Simplifying example: instead of observing $z \sim \mathcal{N}(\mu^*, I)$, we observe a set from some given partition containing z . Can we recover μ^* ?
2. Run stochastic gradient descent (SGD) on the “coarse negative log-likelihood function.”

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

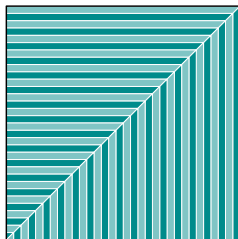
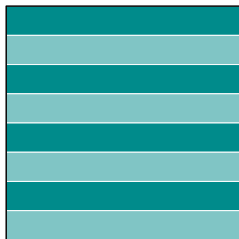
Step I: Coarse Learning

Step II: Optimizing Coarse Likelihood

Challenges

Learning with Coarse Observations

Goal: recover unknown parameter θ^* given coarsened observations P_1, \dots, P_n from a given partition \mathcal{P} of \mathbb{R}^d .

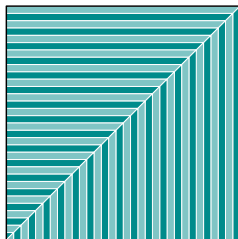


Learning with Coarse Observations

Goal: recover unknown parameter θ^* given coarsened observations P_1, \dots, P_n from a given partition \mathcal{P} of \mathbb{R}^d .

DEFINITION 2 (COARSE LEARNING MODEL [FKKT21])

An observation $P \in \mathcal{P}$ is generated as follows:



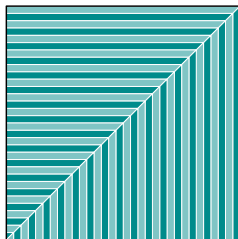
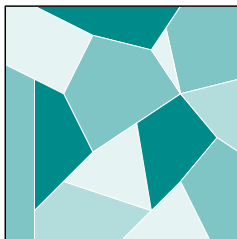
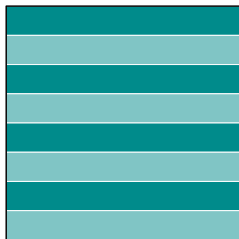
Learning with Coarse Observations

Goal: recover unknown parameter θ^* given coarsened observations P_1, \dots, P_n from a given partition \mathcal{P} of \mathbb{R}^d .

DEFINITION 2 (COARSE LEARNING MODEL [FKKT21])

An observation $P \in \mathcal{P}$ is generated as follows:

1. $\mathbf{z} \sim q_{\theta^*}$.



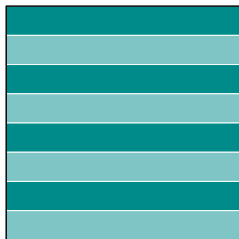
Learning with Coarse Observations

Goal: recover unknown parameter θ^* given coarsened observations P_1, \dots, P_n from a given partition \mathcal{P} of \mathbb{R}^d .

DEFINITION 2 (COARSE LEARNING MODEL [FKKT21])

An observation $P \in \mathcal{P}$ is generated as follows:

1. $\mathbf{z} \sim q_{\theta^*}$.
2. Observe unique $P \in \mathcal{P}$ s.t. $P \ni \mathbf{z}$.



Self-Selection Partition

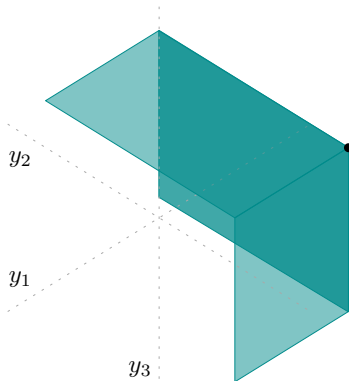
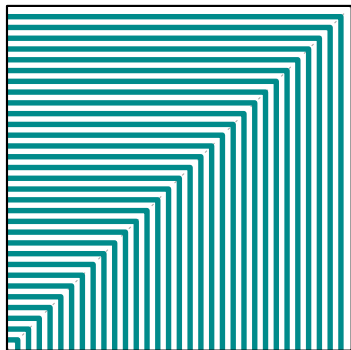
- ▶ $\theta^* = (\mathbf{w}_1, \dots, \mathbf{w}_k)$ and q_{θ^*} is distribution of $\mathbf{z} = (\mathbf{x}, y^{\max})$.

Self-Selection Partition

- ▶ $\theta^* = (\mathbf{w}_1, \dots, \mathbf{w}_k)$ and q_{θ^*} is distribution of $\mathbf{z} = (\mathbf{x}, y^{\max})$.
- ▶ Observing $(\mathbf{x}, y^{\max}) \equiv \{\mathbf{x}\} \times P_{y^{\max}}$ where $P_{y^{\max}} \in \mathcal{P}_{\max}$ below.

Self-Selection Partition

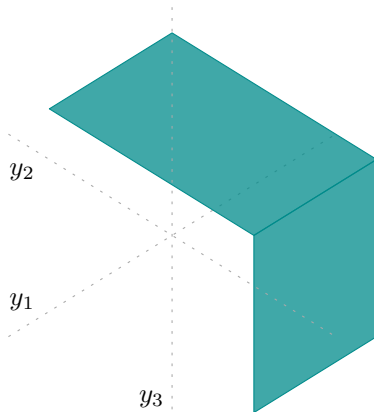
- ▶ $\theta^* = (\mathbf{w}_1, \dots, \mathbf{w}_k)$ and q_{θ^*} is distribution of $\mathbf{z} = (\mathbf{x}, y^{\max})$.
- ▶ Observing $(\mathbf{x}, y^{\max}) \equiv \{\mathbf{x}\} \times P_{y^{\max}}$ where $P_{y^{\max}} \in \mathcal{P}_{\max}$ below.



Self-selection partition for $k = 2$ and a single observation for $k = 3$.

Remark

Reduction to coarse learning is general and captures other problems such as regression with “second-price auction data”.



A single observation of second-price auction data for $k = 3$.

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Step I: Coarse Learning

Step II: Optimizing Coarse Likelihood

Challenges

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta})$$

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = -\mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} [\log q_{\boldsymbol{\theta}}^{\mathcal{P}}(P)]$$

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = -\mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} [\log q_{\boldsymbol{\theta}}^{\mathcal{P}}(P)] = - \sum_{P \in \mathcal{P}} q_{\boldsymbol{\theta}^*}(P) \cdot \log \int_{\mathbf{z} \in P} q_{\boldsymbol{\theta}}(\mathbf{z}) .$$

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = -\mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} [\log q_{\boldsymbol{\theta}}^{\mathcal{P}}(P)] = - \sum_{P \in \mathcal{P}} q_{\boldsymbol{\theta}^*}(P) \cdot \log \int_{\mathbf{z} \in P} q_{\boldsymbol{\theta}}(\mathbf{z}) .$$

- ▶ $\boldsymbol{\theta}^*$ is a stationary point of $\mathcal{L}_{\mathcal{P}}$

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = -\mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} [\log q_{\boldsymbol{\theta}}^{\mathcal{P}}(P)] = - \sum_{P \in \mathcal{P}} q_{\boldsymbol{\theta}^*}(P) \cdot \log \int_{\mathbf{z} \in P} q_{\boldsymbol{\theta}}(\mathbf{z}) .$$

► $\boldsymbol{\theta}^*$ is a stationary point of $\mathcal{L}_{\mathcal{P}}$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}|P}[\mathbf{z}] .$$

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = -\mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} [\log q_{\boldsymbol{\theta}}^{\mathcal{P}}(P)] = - \sum_{P \in \mathcal{P}} q_{\boldsymbol{\theta}^*}(P) \cdot \log \int_{\mathbf{z} \in P} q_{\boldsymbol{\theta}}(\mathbf{z}) .$$

- ▶ $\boldsymbol{\theta}^*$ is a stationary point of $\mathcal{L}_{\mathcal{P}}$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{P}}(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\boldsymbol{\theta}^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\theta}}|P}[\mathbf{z}] .$$

- ▶ $\mathcal{L}_{\mathcal{P}}$ is convex* if each $P \in \mathcal{P}$ is convex (Brascamp–Lieb Inequality)

Coarse Negative Log-Likelihood

Facts about the general coarse NLL

$$\mathcal{L}_{\mathcal{P}}(\theta) = -\mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} [\log q_{\theta}^{\mathcal{P}}(P)] = - \sum_{P \in \mathcal{P}} q_{\theta^*}(P) \cdot \log \int_{\mathbf{z} \in P} q_{\theta}(\mathbf{z}) .$$

- θ^* is a stationary point of $\mathcal{L}_{\mathcal{P}}$

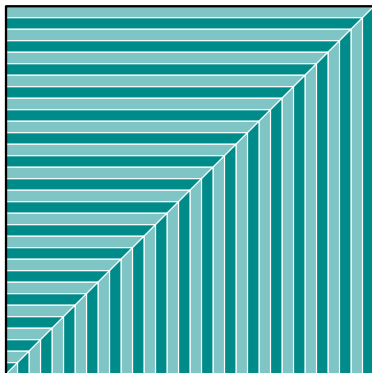
$$\nabla_{\theta} \mathcal{L}_{\mathcal{P}}(\theta) = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}] .$$

- $\mathcal{L}_{\mathcal{P}}$ is convex* if each $P \in \mathcal{P}$ is convex (Brascamp–Lieb Inequality)

$$\nabla_{\theta}^2 \mathcal{L}_{\mathcal{P}}(\theta) = \text{Cov}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \text{Cov}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}] .$$

Remark

If we also observe *index* i_{\max} of the regressor attaining y^{\max} , the partition becomes convex, and we can straightforwardly recover the efficient algorithm of [CDIZ23] for the *known-index* variant of self-selection.



Known-index self-selection partition for $k = 2$.

Table of Contents

Regression under Self-Selection Biases

Prior Works

Our Contributions

Technical Overview

Step I: Coarse Learning

Step II: Optimizing Coarse Likelihood

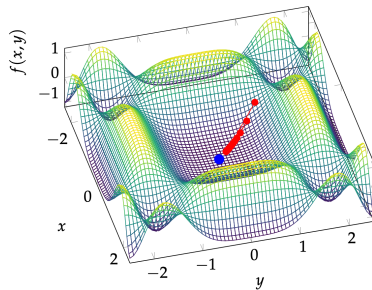
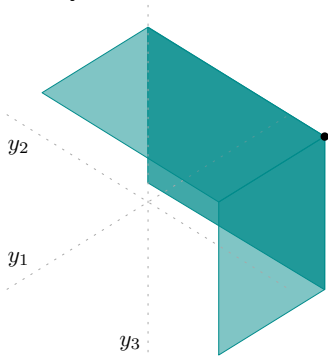
Challenges

Challenges

1. How can we compute $\nabla \mathcal{L}_{\mathcal{P}_{\max}} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}]$?

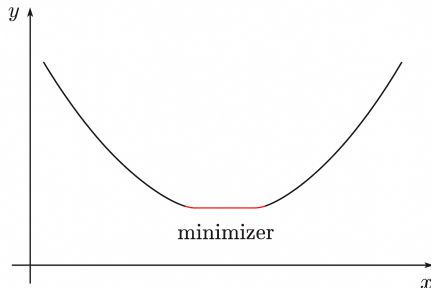
Challenges

1. How can we compute $\nabla \mathcal{L}_{\mathcal{P}_{\max}} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta} | P}[\mathbf{z}]$?
2. Self-selection partition \mathcal{P}_{\max} is not convex, hence $\mathcal{L}_{\mathcal{P}_{\max}}$ potentially has many local minima.



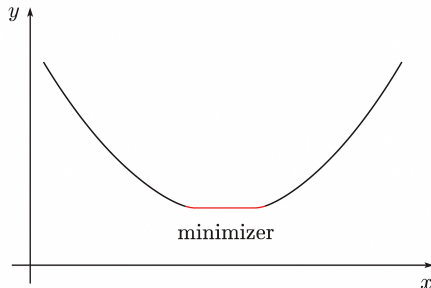
Challenges

1. How can we compute $\nabla \mathcal{L}_{\mathcal{P}_{\max}} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}]$?
2. Self-selection partition \mathcal{P}_{\max} is not convex, hence $\mathcal{L}_{\mathcal{P}_{\max}}$ potentially has many local minima.
3. $\mathcal{L}_{\mathcal{P}_{\max}}$ can be “flat” near θ^* , and SGD may be unable to recover θ^* .



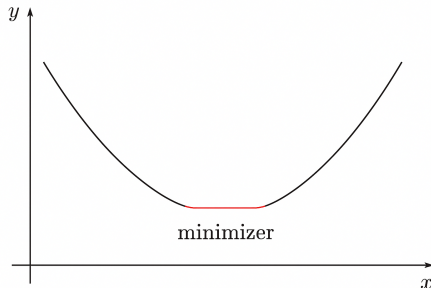
Challenges

1. How can we compute $\nabla \mathcal{L}_{\mathcal{P}_{\max}} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}]$?
 - Unbiased estimates given samples $x, P_{y^{\max}}$. ✓
2. Self-selection partition \mathcal{P}_{\max} is not convex, hence $\mathcal{L}_{\mathcal{P}_{\max}}$ potentially has many local minima.
3. $\mathcal{L}_{\mathcal{P}_{\max}}$ can be “flat” near θ^* , and SGD may be unable to recover θ^* .



Challenges

1. How can we compute $\nabla \mathcal{L}_{\mathcal{P}_{\max}} = \mathbb{E}_{\mathbf{z} \sim q_{\theta}}[\mathbf{z}] - \mathbb{E}_{P \sim q_{\theta^*}^{\mathcal{P}}} \mathbb{E}_{\mathbf{z} \sim q_{\theta}|P}[\mathbf{z}]$?
 - Unbiased estimates given samples $x, P_{y^{\max}}$. ✓
2. Self-selection partition \mathcal{P}_{\max} is not convex, hence $\mathcal{L}_{\mathcal{P}_{\max}}$ potentially has many local minima.
 - We show it is *locally* convex about θ^* . ✓ (with warm start)
3. $\mathcal{L}_{\mathcal{P}_{\max}}$ can be “flat” near θ^* , and SGD may be unable to recover θ^* .



Quadratic Growth from “Information Preservation”

Claim: It suffices to show $\text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}}) \geq \Omega(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|)$.

Quadratic Growth from “Information Preservation”

Claim: It suffices to show $\text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}}) \geq \Omega(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|)$.

$$\alpha \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2 \leq \text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}})^2$$

Quadratic Growth from “Information Preservation”

Claim: It suffices to show $\text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}}) \geq \Omega(\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|)$.

$$\begin{aligned}\alpha \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|^2 &\leq \text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}})^2 \\ &\leq \text{KL}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}} \| q_{\boldsymbol{\theta}^*}^{\mathcal{P}_{\max}}) \quad \text{ Pinsker's}\end{aligned}$$

Quadratic Growth from “Information Preservation”

Claim: It suffices to show $\text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^{\star}}^{\mathcal{P}_{\max}}) \geq \Omega(\|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|)$.

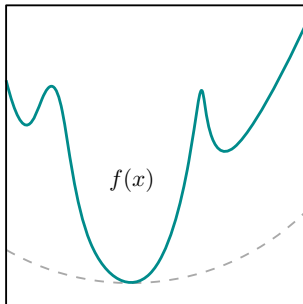
$$\begin{aligned}\alpha \cdot \|\boldsymbol{\theta} - \boldsymbol{\theta}^{\star}\|^2 &\leq \text{TV}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}}, q_{\boldsymbol{\theta}^{\star}}^{\mathcal{P}_{\max}})^2 \\ &\leq \text{KL}(q_{\boldsymbol{\theta}}^{\mathcal{P}_{\max}} \| q_{\boldsymbol{\theta}^{\star}}^{\mathcal{P}_{\max}}) && \text{Pinsker's} \\ &= \mathcal{L}_{\mathcal{P}_{\max}}(\boldsymbol{\theta}). && \text{by definition}\end{aligned}$$

Quadratic Growth from “Information Preservation”

Claim: It suffices to show $\text{TV}(q_{\theta}^{\mathcal{P}_{\max}}, q_{\theta^*}^{\mathcal{P}_{\max}}) \geq \Omega(\|\theta - \theta^*\|)$.

$$\begin{aligned}\alpha \cdot \|\theta - \theta^*\|^2 &\leq \text{TV}(q_{\theta}^{\mathcal{P}_{\max}}, q_{\theta^*}^{\mathcal{P}_{\max}})^2 \\ &\leq \text{KL}(q_{\theta}^{\mathcal{P}_{\max}} \| q_{\theta^*}^{\mathcal{P}_{\max}}) \\ &= \mathcal{L}_{\mathcal{P}_{\max}}(\theta) .\end{aligned}$$

Pinsker's
by definition

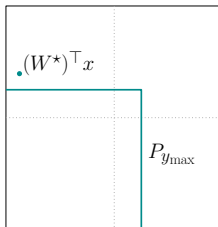


Proof Intuition of “Information Preservation”

- Identify an event \mathcal{E} such that $\mathbf{w}_{i_{\max}}^\top \mathbf{x} \gg \mathbf{w}_j^\top \mathbf{x}$ for all $j \neq i_{\max}$.

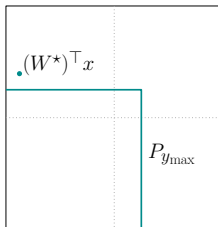
Proof Intuition of “Information Preservation”

- ▶ Identity an event \mathcal{E} such that $\mathbf{w}_{i_{\max}}^\top \mathbf{x} \gg \mathbf{w}_j^\top \mathbf{x}$ for all $j \neq i_{\max}$.
- ▶ Conditional on \mathcal{E} , $P_{y_{\max}}$ “looks” like a convex set under q_{θ^*} .



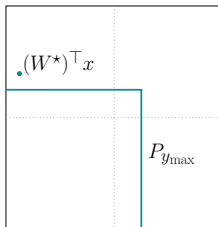
Proof Intuition of “Information Preservation”

- ▶ Identity an event \mathcal{E} such that $\mathbf{w}_{i_{\max}}^\top \mathbf{x} \gg \mathbf{w}_j^\top \mathbf{x}$ for all $j \neq i_{\max}$.
- ▶ Conditional on \mathcal{E} , $P_{y_{\max}}$ “looks” like a convex set under q_{θ^*} .
- ▶ Conditional on \mathcal{E} and under regularity conditions, prove that $\text{TV}(q_{\theta}^{\mathcal{P}_{\max}} \mid \mathcal{E}, q_{\theta^*}^{\mathcal{P}_{\max}} \mid \mathcal{E}) \geq \Omega(\|\theta - \theta^*\|)$.



Proof Intuition of “Information Preservation”

- ▶ Identity an event \mathcal{E} such that $\mathbf{w}_{i_{\max}}^\top \mathbf{x} \gg \mathbf{w}_j^\top \mathbf{x}$ for all $j \neq i_{\max}$.
- ▶ Conditional on \mathcal{E} , $P_{y_{\max}}$ “looks” like a convex set under q_{θ^\star} .
- ▶ Conditional on \mathcal{E} and under regularity conditions, prove that $\text{TV}(q_{\theta^{\mathcal{P}_{\max}}}^{\mathcal{P}_{\max}} \mid \mathcal{E}, q_{\theta^\star}^{\mathcal{P}_{\max}} \mid \mathcal{E}) \geq \Omega(\|\theta - \theta^\star\|)$.
- ▶ Conclude using the fact $\text{TV}(q_{\theta^{\mathcal{P}_{\max}}}^{\mathcal{P}_{\max}}, q_{\theta^\star}^{\mathcal{P}_{\max}}) \geq \Pr[\mathcal{E}] \cdot \text{TV}(q_{\theta^{\mathcal{P}_{\max}}}^{\mathcal{P}_{\max}} \mid \mathcal{E}, q_{\theta^\star}^{\mathcal{P}_{\max}} \mid \mathcal{E})$.



Conclusion

- ▶ Study the geometry of regression with self-selection bias through the lenses of the coarse learning framework.

Conclusion

- ▶ Study the geometry of regression with self-selection bias through the lenses of the coarse learning framework.
- ▶ Leads to an SGD-based local convergence algorithm, which improves on the running time of [\[GM24\]](#).

- ▶ Study the geometry of regression with self-selection bias through the lenses of the coarse learning framework.
- ▶ Leads to an SGD-based local convergence algorithm, which improves on the running time of [GM24].

*Is there a **fully-polynomial** (SGD-based) algorithm for regression with self-selection bias?*

That's All!



felix-zhou.com

felix.zhou@yale.edu