

Towards a Statistical Theory of Data Selection Under Weak Supervision

Germain Kolossov, Andrea Montanari, Pulkit Tandon

Data Selection Setup

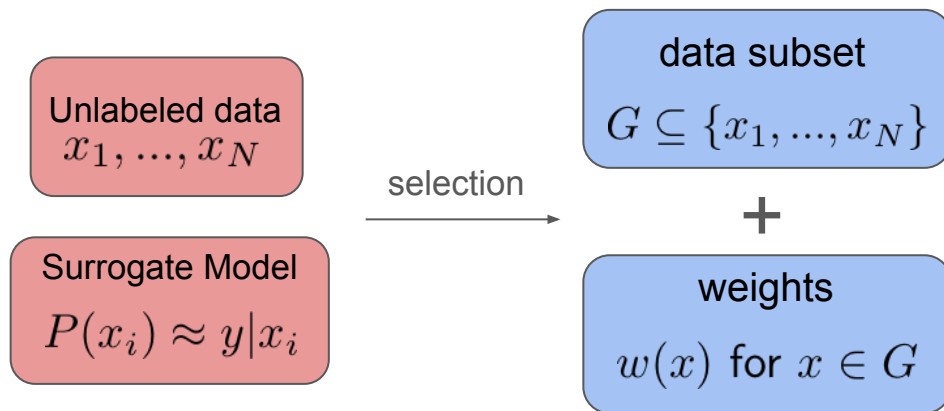
Unlabeled data

$$x_1, \dots, x_N$$

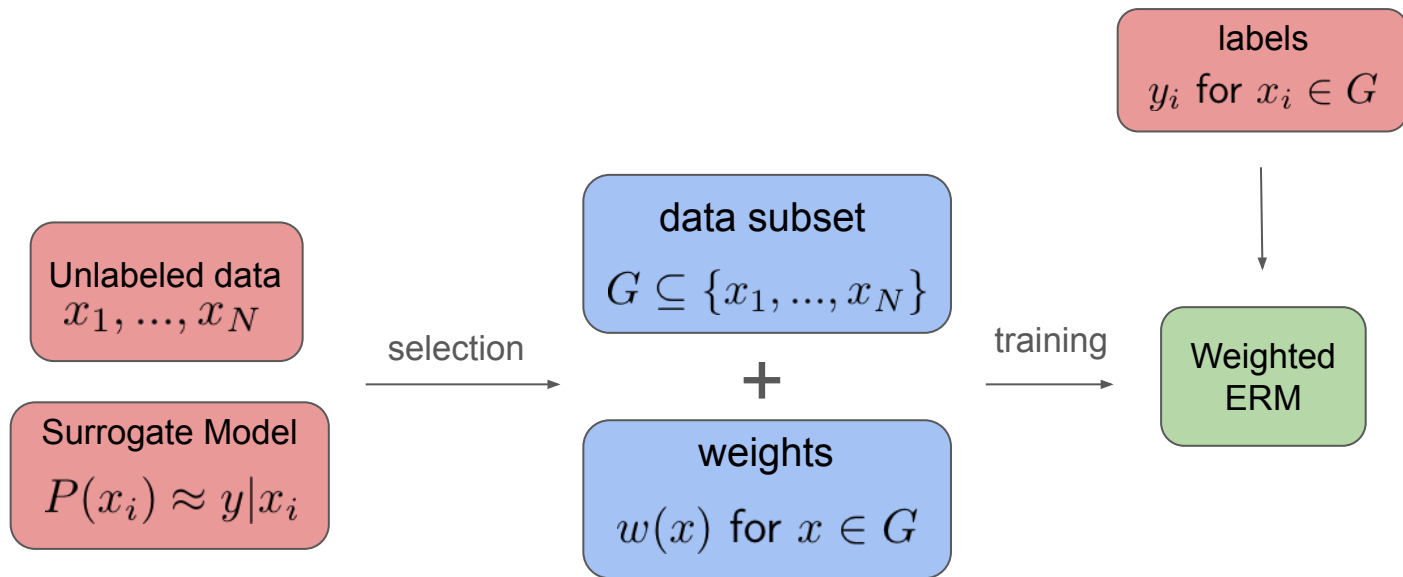
Surrogate Model

$$P(x_i) \approx y|x_i$$

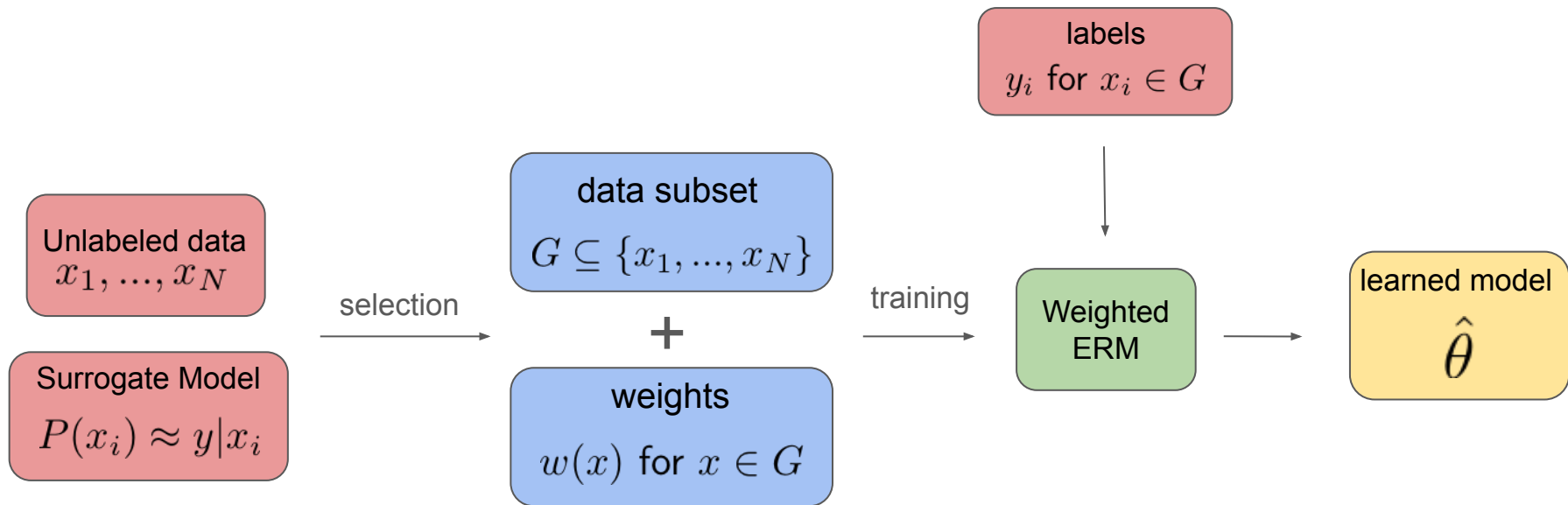
Data Selection Setup



Data Selection Setup



Data Selection Setup



Selection Algorithm Details

Selection scheme = (π, w)

Selection Algorithm Details

Selection scheme = (π, w)

Each x included in dataset
independently with probability

$$\pi(x, P(x))$$

If included, appears with weight

$$w(x, P(x)) > 0$$

Selection Algorithm Details

Selection scheme = (π, w)

Each x included in dataset
independently with probability

$$\pi(x, P(x))$$

If included, appears with weight

$$w(x, P(x)) > 0$$

Target size achieved
in expectation:

$$\sum_{i=1}^N \mathbb{E}[\pi(x_i)] = n$$

Selection Algorithm Details

Selection scheme = (π, w)

Each x included in dataset
independently with probability

$$\pi(x, P(x))$$

If included, appears with weight

$$w(x, P(x)) > 0$$

Target size achieved
in expectation:

$$\sum_{i=1}^N \mathbb{E}[\pi(x_i)] = n$$

Unbiased selection scheme: $w(x) = \frac{1}{\pi(x)}$

Selection Algorithm Details

Selection scheme = (π, w)

Each x included in dataset
independently with probability

$$\pi(x, P(x))$$

If included, appears with weight

$$w(x, P(x)) > 0$$

Target size achieved
in expectation:

$$\sum_{i=1}^N \mathbb{E}[\pi(x_i)] = n$$

Unbiased selection scheme: $w(x) = \frac{1}{\pi(x)}$

Non-reweighting selection scheme: $w(x) = 1$

Method of Analysis: Asymptotics

- Compare selections schemes based on performance as the original dataset size (N) grows to infinity.
 - Number of datapoints selected (n) also grows with N
 - Approaches some fixed fraction $\gamma \in (0, 1)$

Method of Analysis: Asymptotics

- Compare selections schemes based on performance as the original dataset size (N) grows to infinity.
 - Number of datapoints selected (n) also grows with N
 - Approaches some fixed fraction $\gamma \in (0, 1)$

“Low-Dimensional Regime”: Keep dimension fixed as N grows

“High-Dimensional Regime”: Grow dimension with N, converging to fixed ratio

Setting 1: Low-Dimension, Perfect Surrogate

Quantity of Interest: Asymptotic Error Coefficient

$$\begin{array}{c}
 (\pi, w) \\
 \downarrow \\
 \rho(S, Q) = \lim_{M \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbb{E}[\min\{N \underbrace{\|\hat{\theta} - \theta_*\|_Q^2}_{\langle \hat{\theta} - \theta_*, Q(\hat{\theta} - \theta_*) \rangle}, M\}]
 \end{array}$$

Population minimizer \downarrow
 θ_*

Quantity of Interest: Asymptotic Error Coefficient

$$\begin{array}{c}
 (\pi, w) \\
 \downarrow \\
 \rho(S, Q) = \lim_{M \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbb{E}[\min\{N \|\hat{\theta} - \theta_*\|_Q^2, M\}]
 \end{array}$$

Population minimizer
↓

Under some assumptions:

- Limit exists
- Has a closed form in terms of
 - S
 - conditional gradient covariance of loss at θ_*
 - conditional hessian of the loss at θ_*

$$\langle \hat{\theta} - \theta_*, Q(\hat{\theta} - \theta_*) \rangle$$

↑

Quantity of Interest: Asymptotic Error Coefficient

$$\begin{array}{c}
 (\pi, w) \\
 \downarrow \\
 \rho(S, Q) = \lim_{M \rightarrow \infty} \lim_{N \rightarrow \infty} \mathbb{E}[\min\{N \|\hat{\theta} - \theta_*\|_Q^2, M\}]
 \end{array}$$

Population minimizer
↓
↑
 $\langle \hat{\theta} - \theta_*, Q(\hat{\theta} - \theta_*) \rangle$

Under some assumptions:

- Limit exists
- Has a closed form in terms of
 - S
 - conditional gradient covariance of loss at θ_*
 - conditional hessian of the loss at θ_*

Best selection scheme in a class C:

$$\arg \min_{S \in C} \rho(S, Q)$$

Up next: used closed-form to solve for optimal strategy in certain classes.

Optimal Unbiased Strategy

$$w(x) = \frac{1}{\pi(x)}$$

- Can simplify and optimize the asymptotic error coefficient to solve for the optimal unbiased strategy:

Optimal Unbiased Strategy

$$w(x) = \frac{1}{\pi(x)}$$

- Can simplify and optimize the asymptotic error coefficient to solve for the optimal unbiased strategy:
 - Recovers selection based on influence function

Influence function

$$\psi(\mathbf{x}, y) = -\mathbf{H}^{-1} \nabla_{\theta} L(\theta_*; y, \mathbf{x})$$

$$\longrightarrow \pi(\mathbf{x}_i) \propto \mathbb{E} \left\{ \left\| \psi(\mathbf{x}_i, y_i) \right\|_{\mathbf{Q}}^2 \mid \mathbf{x}_i \right\}^{1/2}$$

Optimal Unbiased Strategy

$$w(x) = \frac{1}{\pi(x)}$$

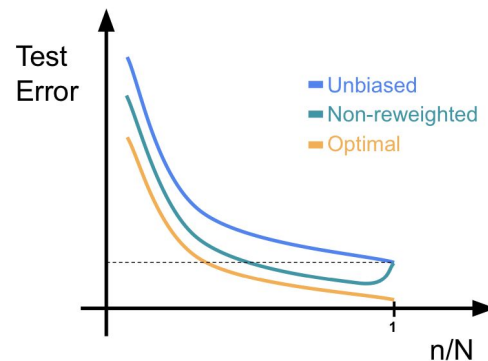
- Can simplify and optimize the asymptotic error coefficient to solve for the optimal unbiased strategy:
 - Recovers selection based on influence function

Influence function

$$\psi(\mathbf{x}, y) = -\mathbf{H}^{-1} \nabla_{\theta} L(\theta_*; y, \mathbf{x})$$

$$\longrightarrow \pi(\mathbf{x}_i) \propto \mathbb{E} \left\{ \left\| \psi(\mathbf{x}_i, y_i) \right\|_Q^2 \mid \mathbf{x}_i \right\}^{1/2}$$

- Monotone non-increasing in selected proportion
- Better than just random unbiased sampling



Optimal Non-reweighting Strategy

$$w(x) = 1$$

- Again, simplify and optimize the asymptotic error coefficient to characterize the optimal non-reweighting strategy.
- Show that optimal strategy must be a fixed point of the following process:

Optimal Non-reweighting Strategy

$$w(x) = 1$$

- Again, simplify and optimize the asymptotic error coefficient to characterize the optimal non-reweighting strategy.
- Show that optimal strategy must be a fixed point of the following process:
 - Strategy induces a score function for each point $Z(x; \pi)$
 - Optimal strategy must be a threshold decision based on score
 - Threshold parameters chosen to meet selection proportion

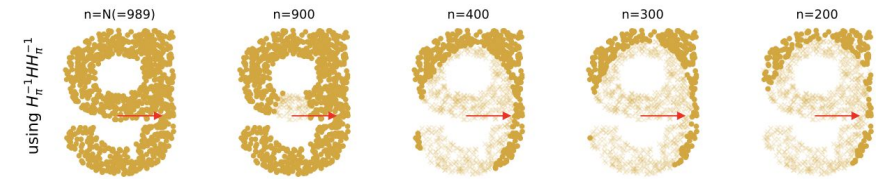
$$\pi_{\text{nr}}(\mathbf{x}) = \begin{cases} 1 & \text{if } Z(\mathbf{x}; \pi_{\text{nr}}) > \lambda, \\ 0 & \text{if } Z(\mathbf{x}; \pi_{\text{nr}}) < \lambda, \\ b(\mathbf{x}) \in [0, 1] & \text{if } Z(\mathbf{x}; \pi_{\text{nr}}) = \lambda. \end{cases}$$

Optimal Non-reweighting Strategy

$$w(x) = 1$$

- Again, simplify and optimize the asymptotic error coefficient to characterize the optimal non-reweighting strategy.
- Show that optimal strategy must be a fixed point of the following process:
 - Strategy induces a score function for each point $Z(x; \pi)$
 - Optimal strategy must be a threshold decision based on score
 - Threshold parameters chosen to meet selection proportion

$$\pi_{nr}(\mathbf{x}) = \begin{cases} 1 & \text{if } Z(\mathbf{x}; \pi_{nr}) > \lambda, \\ 0 & \text{if } Z(\mathbf{x}; \pi_{nr}) < \lambda, \\ b(\mathbf{x}) \in [0, 1] & \text{if } Z(\mathbf{x}; \pi_{nr}) = \lambda. \end{cases}$$



Linear regression setting: Z measures how different x is from selected data (related to *leverage scores*)

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

$$\rho(\pi_t; \mathbf{Q}) = \rho(\pi_{nr}; \mathbf{Q}) + t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

↑
Comes from closed-form
for asymptotic error

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

$$\rho(\pi_t; \mathbf{Q}) = \rho(\pi_{nr}; \mathbf{Q}) + t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

$$\rho(\pi; \mathbf{Q}) = \text{Tr}\left(\mathbb{E}\{\pi(\mathbf{x})\mathbf{G}(\mathbf{x})\}\mathbb{E}\{\pi(\mathbf{x})\mathbf{H}(\mathbf{x})\}^{-1}\mathbf{Q}\mathbb{E}\{\pi(\mathbf{x})\mathbf{H}(\mathbf{x})\}^{-1}\right)$$

↑
Comes from closed-form
for asymptotic error

Conditional hessian

$$\mathbf{H}(\mathbf{x}) := \mathbb{E}\{\nabla_{\boldsymbol{\theta}}^2 L(\boldsymbol{\theta}_*; y, \mathbf{x}) | \mathbf{x}\}$$

Conditional gradient covariance

$$\mathbf{G}(\mathbf{x}) := \mathbb{E}\{\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_*; y, \mathbf{x}) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_*; y, \mathbf{x})^\top | \mathbf{x}\}$$

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

$$\rho(\pi_t; \mathbf{Q}) = \rho(\pi_{nr}; \mathbf{Q}) - t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

$$\rho(\pi; \mathbf{Q}) = \text{Tr}\left(\mathbb{E}\{\pi(\mathbf{x})\mathbf{G}(\mathbf{x})\}\mathbb{E}\{\pi(\mathbf{x})\mathbf{H}(\mathbf{x})\}^{-1}\mathbf{Q}\mathbb{E}\{\pi(\mathbf{x})\mathbf{H}(\mathbf{x})\}^{-1}\right)$$

↑
Comes from closed-form
for asymptotic error

$$Z(\mathbf{x}; \pi) := -\text{Tr}\{\mathbf{G}(\mathbf{x})\mathbf{H}_\pi^{-1}\mathbf{Q}\mathbf{H}_\pi^{-1}\} + 2\text{Tr}\{\mathbf{H}(\mathbf{x})\mathbf{H}_\pi^{-1}\mathbf{Q}\mathbf{H}_\pi^{-1}\mathbf{G}_\pi\mathbf{H}_\pi^{-1}\}$$

Conditional hessian

$$\mathbf{H}(\mathbf{x}) := \mathbb{E}\{\nabla_{\boldsymbol{\theta}}^2 L(\boldsymbol{\theta}_*; y, \mathbf{x}) | \mathbf{x}\}$$

Conditional gradient covariance

$$\mathbf{G}_\pi := \mathbb{E}_\pi \mathbf{G}(\mathbf{x}), \quad \mathbf{H}_\pi := \mathbb{E}_\pi \mathbf{H}(\mathbf{x}), \quad \text{where} \quad \mathbb{E}_\pi f(\mathbf{x}) := \frac{\mathbb{E}\{f(\mathbf{x})\pi(\mathbf{x})\}}{\mathbb{E}\{\pi(\mathbf{x})\}}$$

$$\mathbf{G}(\mathbf{x}) := \mathbb{E}\{\nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_*; y, \mathbf{x}) \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}_*; y, \mathbf{x})^\top | \mathbf{x}\}$$

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

Must be non-negative

$$\rho(\pi_t; \mathcal{Q}) = \rho(\pi_{nr}; \mathcal{Q}) - t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

Must be non-negative

$$\rho(\pi_t; \mathbf{Q}) = \rho(\pi_{nr}; \mathbf{Q}) - t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

- Claim: $Z(\mathbf{x}; \pi_{nr})$ constant on all \mathbf{x} with $\pi_{nr}(\mathbf{x}) \in (0, 1)$
 - If not, can construct feasible strategy that breaks non-negativity

Optimal Non-reweighting Strategy: Proof Idea

- Suppose we have a minimizer π_{nr}
- Imagine perturbing π_{nr} to some $\pi_t := (1 - t)\pi_{nr} + t\pi$.
- Can express perturbed error value as

Must be non-positive

$$\rho(\pi_t; \mathbf{Q}) = \rho(\pi_{nr}; \mathbf{Q}) - t \int (\pi(\mathbf{x}) - \pi_{nr}(\mathbf{x})) Z(\mathbf{x}; \pi_{nr}) \mathbb{P}(d\mathbf{x}) + o(t)$$

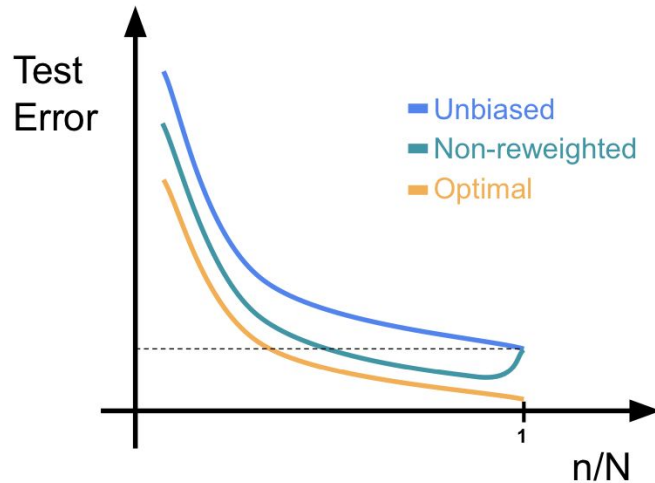
- Claim: $Z(\mathbf{x}; \pi_{nr})$ constant on all \mathbf{x} with $\pi_{nr}(\mathbf{x}) \in (0, 1)$
 - If not, can construct feasible strategy that breaks non-positivity
- Implies the threshold structure

$$\pi_{nr}(\mathbf{x}) = \begin{cases} 1 & \text{if } Z(\mathbf{x}; \pi_{nr}) > \lambda, \\ 0 & \text{if } Z(\mathbf{x}; \pi_{nr}) < \lambda, \\ b(\mathbf{x}) \in [0, 1] & \text{if } Z(\mathbf{x}; \pi_{nr}) = \lambda. \end{cases} \quad \text{the constant value}$$

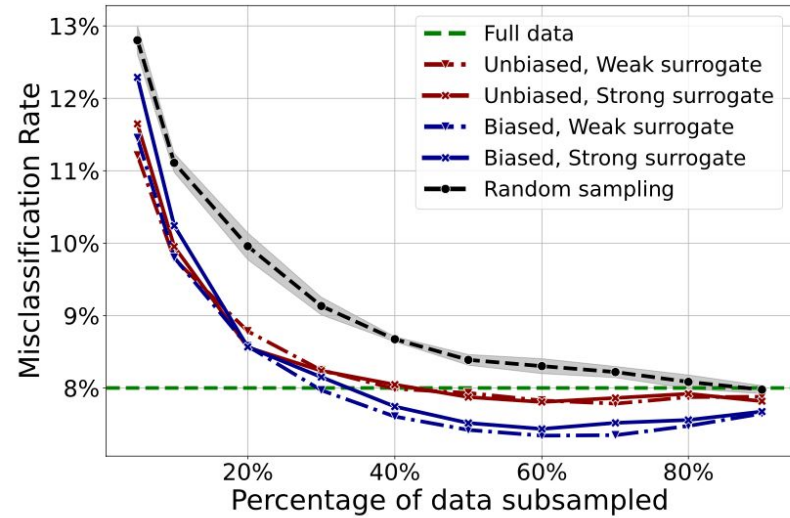
Non-reweighting vs Unbiased Strategies

Thm: Unbiased can be arbitrarily worse than non-reweighting strategies.

- In terms of ratio of asymptotic error coefficients



theorized



empirically

Part 2: Imperfect surrogates, High-dimensional asymptotics

Imperfect surrogates

So far, we assume access to (essentially) $P(y|x)$ via a perfect surrogate

What happens if the surrogate is imperfect, so $P_{\text{su}}(y|x) \neq P(y|x)$

First idea: “Plug in” the surrogate (treat it as if it were truly $P(y|x)$)

Main result from this paper: This is suboptimal, but something close to it is roughly (minimax) optimal

First approach: Plug-in estimation

Plugin unbiased data selection. We form

$$\mathbf{G}_{\text{su}}(\mathbf{x}) := \mathbb{E}_{\text{su}} \left\{ \nabla_{\boldsymbol{\theta}} L(\hat{\boldsymbol{\theta}}^{\text{su}}; y, \mathbf{x}) \nabla_{\boldsymbol{\theta}} L(\hat{\boldsymbol{\theta}}^{\text{su}}; y, \mathbf{x})^{\top} | \mathbf{x} \right\},$$

$$\mathbf{H}_{\text{su}}(\mathbf{x}) := \mathbb{E}_{\text{su}} \left\{ \nabla_{\boldsymbol{\theta}}^2 L(\hat{\boldsymbol{\theta}}^{\text{su}}; y, \mathbf{x}) | \mathbf{x} \right\},$$

and subsample according to

Either “read off” or

$$\hat{\boldsymbol{\theta}}^{\text{su}} := \arg \min \sum_{i=1}^N L(\boldsymbol{\theta}; y_i^{\text{su}}, \mathbf{x}_i)$$

$$\pi(\mathbf{x}) = \min \left(1; c(\gamma) Z_{\text{su}}(\mathbf{x})^{1/2} \right),$$

$$Z_{\text{su}}(\mathbf{x}) := \text{Tr}(\mathbf{G}_{\text{su}}(\mathbf{x}) \mathbf{H}_{1,\text{su}}^{-1} \mathbf{Q} \mathbf{H}_{1,\text{su}}^{-1}),$$

$$\mathbf{H}_{1,\text{su}} := \mathbb{E}\{\mathbf{H}_{\text{su}}(\mathbf{x})\}.$$

Approach for studying optimality: minimax framework

Assume that the surrogate predictor is close to (but not equal to) true model

$$\mathcal{K}_d(\mathbb{P}_{\text{su}}; r) := \left\{ \mathbb{P} : \mathbb{E}_{\mathbf{x}} \|\mathbb{P}(\cdot | \mathbf{x}) - \mathbb{P}_{\text{su}}(\cdot | \mathbf{x})\|_{\text{TV}} \leq r \right\}.$$


(Doesn't have to be TV, but we need the set to be convex)

This allows us to define minimax risk:

$$R_*(S; \mathcal{K}_d) := \sup_{\mathbb{P} \in \mathcal{K}_d} \mathbb{E}_{\mathbf{y}, \mathbf{X} \sim \mathbb{P}(\mathbb{P})} R_{\#}(S; \mathbf{y}, \mathbf{X}),$$

$$R_{\text{MM}}(\mathcal{K}_d) := \inf_{S \in \mathcal{A}} R_*(S; \mathcal{K}_d).$$

Test risk of estimator
under selection
scheme S



What is the optimal way to use the surrogate?

Theorem 5. Assume that any $P_N \in \mathcal{K}_{d,N}$ is supported on $\|\mathbf{y}\| \leq M$, and that $(\mathbf{y}, \mathbf{X}) \mapsto R(\hat{\boldsymbol{\theta}}_A(\mathbf{y}, \mathbf{X}))$ is continuous for any A . Define

$$\bar{R}_{\text{MM}}(\mathcal{K}_d) := \inf_{S \in \mathcal{A}} \bar{R}_*(S; \mathcal{K}_d) := \inf_{S \in \mathcal{A}} \sup_{P_N \in \mathcal{K}_{d,N}} \mathbb{E}_{\mathbf{y}, \mathbf{X} \sim \mathbb{P}(P_N)} R_{\#}(S; \mathbf{y}, \mathbf{X}). \quad (5.14)$$

Then we have

$$\bar{R}_{\text{MM}}(\mathcal{K}_d) = \sup_{P_N \in \mathcal{K}_{d,N}} \inf_{S \in \mathcal{A}} \mathbb{E}_{\mathbf{y}, \mathbf{X} \sim \mathbb{P}(P_N)} R_{\#}(S; \mathbf{y}, \mathbf{X}). \quad (5.15)$$

Sion's minimax theorem

Further, assume P_{MM} achieves the supremum over \mathcal{K}_d above. Then any

$$S_{\text{MM}} \in \arg \min_{S \in \mathcal{A}} \mathbb{E}_{\mathbf{y}, \mathbf{X} \sim \mathbb{P}(P_{\text{MM}})} R_{\#}(S; \mathbf{y}, \mathbf{X}) \quad (5.16)$$

achieves the minimax error.

Intuition: we should use the “worst $P(\mathbf{y}|\mathbf{x})$ ” that is near the given surrogate

High-dimensional asymptotics

$$\frac{n}{N} \rightarrow \gamma, \quad \frac{N}{p} \rightarrow \delta_0,$$

Specify to:

- Gaussian covariates (so \mathbf{x} is drawn from isotropic Gaussian with dim p)
- Response dependent only on linear function of \mathbf{x} : $\mathbb{P}(y_i \in A | \mathbf{x}_i) = \mathbb{P}(A | \langle \boldsymbol{\theta}_0, \mathbf{x}_i \rangle)$
- Generalized linear models + Ridge

$$\hat{R}_N(\boldsymbol{\theta}) := \frac{1}{N} \sum_{i=1}^N S_i(\langle \hat{\boldsymbol{\theta}}^{\text{su}}, \mathbf{x}_i \rangle) L(\langle \boldsymbol{\theta}, \mathbf{x}_i \rangle, y_i) + \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2.$$

In the next slide, we will just ~~try to understand~~ show the theorem statement :)

Theorem statement: Setup

$$\beta_0 := \lim_{N, p \rightarrow \infty} \frac{\langle \hat{\boldsymbol{\theta}}^{\text{su}}, \boldsymbol{\theta}_0 \rangle}{\|\boldsymbol{\theta}_0\|}, \quad \beta_s := \lim_{N, p \rightarrow \infty} \left\| \mathbf{P}_0^\perp \hat{\boldsymbol{\theta}}^{\text{su}} \right\|_2$$

The high-dimensional asymptotics of the test error is determined by a saddle point of the following Lagrangian (here and below $\boldsymbol{\alpha} := (\alpha_0, \alpha_s, \alpha_\perp)$, $\boldsymbol{\beta} := (\beta_0, \beta_s, 0)$):

$$\mathcal{L}(\boldsymbol{\alpha}, \mu, \omega) := \frac{\lambda}{2} \|\boldsymbol{\alpha}\|^2 - \frac{1}{2\delta_0} \mu \alpha_\perp^2 + \mathbb{E} \left\{ \min_{u \in \mathbb{R}} \left[S(\langle \boldsymbol{\beta}, \mathbf{G} \rangle) L(\alpha_0 G_0 + \alpha_s G_s + u, Y) + \frac{1}{2} \mu (\alpha_\perp G_\perp - u)^2 \right] \right\} \quad (6.7)$$

Here expectation is with respect to

$$\mathbf{g} = (G_0, G_s, G_\perp) \sim \mathbf{N}(\mathbf{0}, \mathbf{I}_3), \quad Y \sim \mathbf{P}(\cdot \mid \|\boldsymbol{\theta}_0\|_2 G_0). \quad (6.8)$$

as well as the randomness in S .

Theorem statement: result

Theorem 6. Assume $u \mapsto L(u, y)$ is convex, continuous, with at most quadratic growth, and $\lambda > 0$. Further denote by α^*, μ^* the solution of the following minimax problem (α^* is uniquely defined by this condition.)

**Given a selection strategy, this tells us the (asymptotic) error!
(Importantly, does not identify optimal strategy)**

(a) Extremely vague idea: Decompose into θ_0 direction, θ_s direction, and the rest (which is all indistinguishable bc of Gaussianity)

The rest is an application of Gordon's Gaussian comparison inequality (generalization of Slepian's inequality)

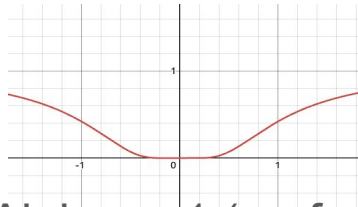
Special case: Gaussian \mathbf{x} , binary y

Posit data selection mechanism:

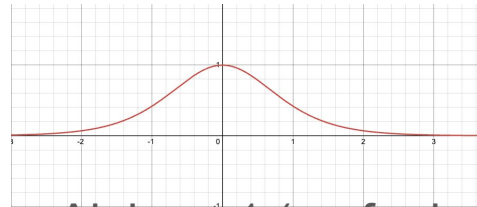
$$\pi(\mathbf{x}_i) = \min \left(\overset{\text{Normalizing constant}}{c(\gamma)} \overset{\text{Second derivative of log-MGF (for logistic regression, } 1 - \tanh(t)^2)}{\phi''(\langle \hat{\boldsymbol{\theta}}^{\text{su}}, \mathbf{x}_i \rangle)}^\alpha; 1 \right)$$

For $\alpha = 1/2$, roughly equivalent to influence function-based sampling (only because the data is Gaussian and so the Hessian has a simple closed form)

Alpha controls whether we prefer hard examples or easy:



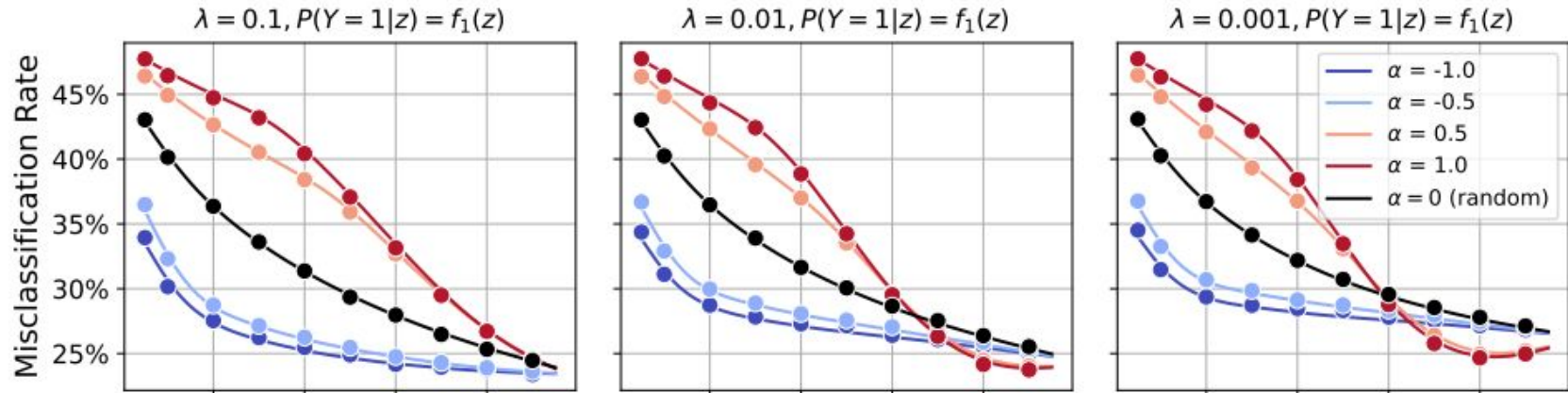
Alpha = -1 (prefer easy)



Alpha = 1 (prefer hard)

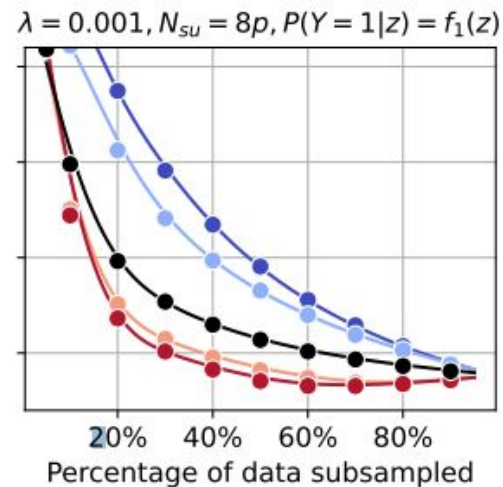
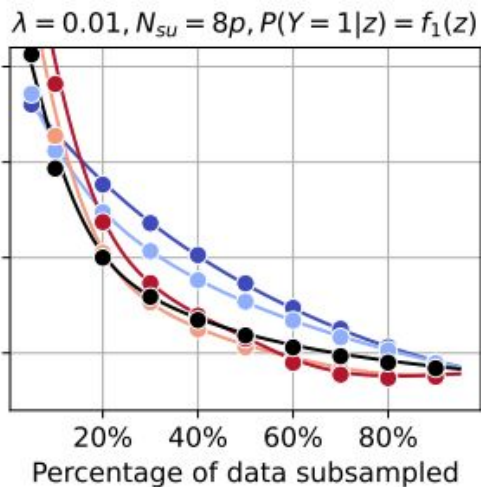
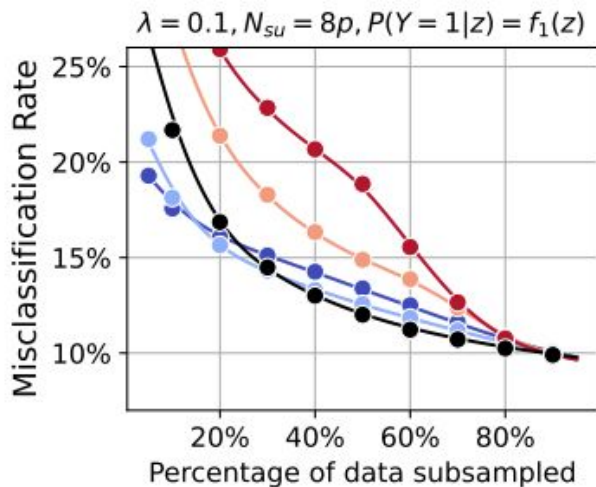
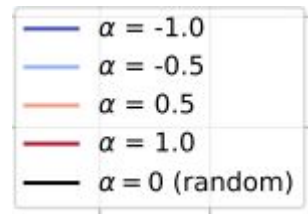
Results: perfect surrogate

Setting: “misspecified” linear model, so $P(y=1|x) = f_1(x^T \theta)$



Results: imperfect surrogate

Surrogate is trained with additional $N_{\{su\}}$ samples



Takeaways (imperfect surrogate)

1. Learning after data selection often outperforms learning on the full sample.
2. Upsampling ‘hard’ datapoints (i.e. using $\alpha > 0$) is often the optimal strategy. This appears to be more common than in the well-specified case.
3. As shown in Figure 8, the performance of data selection-based learning degrades gracefully with the quality of the surrogate.
4. In particular, we observe once more the striking phenomenon of Figure 1, cf. bottom row, rightmost plot of Figure 8. At subsampling fraction $n/N = 60\%$, learning on selected data outperforms learning on the full data, even if the surrogate model only used additional $N_{\text{su}}/N \approx 21.7\%$ samples. As shown in next section, this effect is even stronger with real data.

Experimental Verification

