

Semi-supervised Learning : Inference

Gyumin Lee

Yonsei University

August 29, 2023

Introduction

Brief summary

The purpose of the presentation is to introduce semi-supervised inference and related studies.

Papers

1. Doubly Robust Self-Training

- Banghua Zhu, Mingyu Ding, Philip Jacobson, Ming Wu, Wei Zhan, Michael I. Jordan, and Jiantao Jiao(June 2023)
- <https://arxiv.org/abs/2306.00265>

2. Prediction-Powered Inference

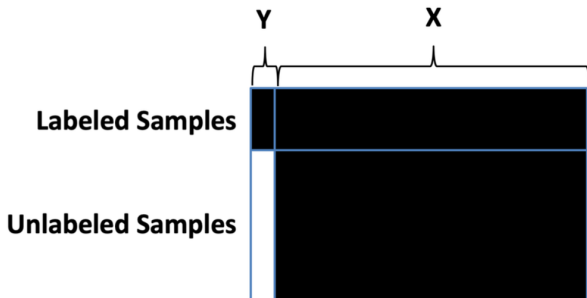
- Anastasios N. Angelopoulos, Stephen Bates, Clara Fannjiang, Michael I. Jordan, and Tijana Zrnic.(February 2023)
- <https://arxiv.org/abs/2301.09633>

Contents

1. Preliminaries: Semi-Supervised Inference
2. Related Works
3. Doubly Robust Self-Training
4. Prediction-Powered Inference

Preliminaries: Semi-Supervised Inference

- Basic assumption:
labeled data are more difficult/expensive to acquire than unlabeled data.
→ utilize unlabeled data!



- Examples
 - Survey sampling
 - Electronic health record
 - homeless consensus
 - ⋮

Preliminaries: Semi-Supervised Inference

Basic Setting

- Labeled data : $(X, Y) \in (\mathcal{X}, \mathcal{Y})^n$ where $X = (X_1, \dots, X_n)$ and $Y = (Y_1, \dots, Y_n)$.
- Unlabeled data : $(\tilde{X}, \tilde{Y}) \in (\mathcal{X}, \mathcal{Y})^m$ where \tilde{Y} is not observed.
- Assume that (X, Y) and (\tilde{X}, \tilde{Y}) are i.i.d. samples from a common distribution \mathbb{P} .
- Prediction rule : $f : \mathcal{X} \rightarrow \mathcal{Y}$, independent of the observed data.(e.g. pretrained model)

Goal

Our estimand of interest is denoted as θ^* which could be $E[Y]$ (mean estimation), $\min\{\theta : P(Y \leq \theta) \geq q\}$ (quantile estimation), or $\arg \min_{\theta \in \mathbb{R}^d} E[l_\theta(X, Y)]$ for some loss function l_θ .

Preliminaries: Semi-Supervised Inference

Example

Let our goal of estimation be $\theta^* = E[Y]$.

We have several candidates of estimators as follows

- $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$
- $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i)) + \frac{1}{n+m} \sum_{i=1}^{n+m} f(X_i)$

Note that both are unbiased estimators. How about variance?

- $Var(\bar{Y}) = Var(\frac{1}{n} \sum_{i=1}^n Y_i) = \frac{1}{n} Var(Y)$
- $Var(\hat{\theta}) = Var(\frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i)) + \frac{1}{n+m} \sum_{i=1}^{n+m} f(X_i))$
 $= \frac{1}{n} (Var(Y) + \frac{m-n}{m} (Var(f(X)) - 2Cov(Y, f(X))))$

→ If $Var(f(X)) < 2Cov(Y, f(X))$, then $\hat{\theta}$ improves the original one!

Preliminaries: Semi-Supervised Inference

Example(cont.)

Suppose we take an ideal prediction function $f(X) = E(Y|X)$

Since

$$\begin{aligned}\text{Cov}(Y, E(Y|X)) &= E[(Y - E[Y])(E[Y|X] - E[Y])] \\ &= E[E[(Y - E[Y])(E[Y|X] - E[Y])|X]] \\ &= E[(E[Y|X] - E[Y])^2] \\ &= \text{Var}(E(Y|X))\end{aligned}$$

, we obtain the result assuming $n \ll m$ as

$$\begin{aligned}\text{Var}(\hat{\theta}) &= \frac{1}{n}E(\text{Var}(Y|X)) + \frac{1}{m}\text{Var}(E(Y|X)) \\ &\leq \frac{1}{n}E(\text{Var}(Y|X)) + \frac{1}{n}\text{Var}(E(Y|X)) = \frac{1}{n}\text{Var}(Y) = \text{Var}(\bar{Y})\end{aligned}$$

Hence, $\hat{\theta}$ has smaller variance than the sample mean \bar{Y} .

Related Works

Semi-supervised inference: general theory and estimation of means (Anru Zhang, Lawrence D. Brown and T. Tony Cai)

For the ideal semi-supervised inference, where $m = \infty$, the proposed estimator is 'least square estimator', which is defined as

$$\hat{\theta}_{LS} = \boldsymbol{\mu}' \hat{\beta} = \hat{\beta}_1 + \boldsymbol{\mu}' \hat{\beta}_{(2)} = \bar{Y} - \hat{\beta}'_{(2)}(\bar{X} - \boldsymbol{\mu}).$$

Here, $\boldsymbol{\mu} = (1, \boldsymbol{\mu}')' = \mathbb{E}X$ is known and $\hat{\beta} = [\hat{\beta}_1 \hat{\beta}'_{(2)}]' = (X'X)^{-1}X'Y$, where X is $n \times (p+1)$ matrix including intercept column of ones.

On the other hand, for the ordinary semi-supervised inference, where $m < \infty$, the proposed estimator is 'semi-supervised least squared estimator', which is defined as

$$\hat{\theta}_{SSLS} = \hat{\boldsymbol{\mu}}' \hat{\beta} = \bar{Y} - \hat{\beta}'_{(2)}(\bar{X} - \hat{\boldsymbol{\mu}}).$$

Here, $\hat{\boldsymbol{\mu}} = (1, \hat{\boldsymbol{\mu}}')' = \frac{1}{n+m} \sum_{k=1}^{n+m} X_k$.

Related Works

Semi-supervised inference: general theory and estimation of means (Anru Zhang, Lawrence D. Brown and T. Tony Cai)

- Results about l_2 risks

$$nE(\bar{Y} - \theta)^2 = \tau^2 + \beta_{(2)}^\top \Sigma_n \beta_{(2)}$$

$$nE(\hat{\theta}_{LS}^1 - \theta)^2 = \tau^2 + O\left(\frac{p^2}{n}\right)$$

$$\begin{aligned} nE\left(\hat{\theta}_{SSLS}^1 - \theta\right)^2 &= \tau_n^2 + \frac{n}{n+m} \beta_{(2)}^\top \Sigma_n \beta_{(2)} + O\left(\frac{p^2}{n}\right) \\ &\approx \frac{n}{n+m} E(\bar{Y} - \theta)^2 + \frac{m}{n+m} E(\hat{\theta}_{LS}^1 - \theta)^2 \end{aligned}$$

Related Works

Semi-supervised inference: general theory and estimation of means (Anru Zhang, Lawrence D. Brown and T. Tony Cai)

Let $(Y_1, X_1), \dots, (Y_n, X_n)$ be i.i.d. copies from P , and assume that $[Y, X]$ has finite second moments, Ξ is non-singular and $\tau^2 > 0$. Then, under the setting that P is fixed and $n \rightarrow \infty$

$$\frac{\hat{\theta}_{\text{LS}} - \theta}{\tau/\sqrt{n}} \xrightarrow{d} N(0, 1),$$

and

$$MSE/\tau^2 \xrightarrow{d} 1, \text{ where } MSE := \frac{\sum_{i=1}^n (Y_i - X_i^\top \hat{\beta})^2}{n - p - 1}, \tau^2 = E(Y_i - \vec{X}_i^\top \beta)^2$$

Related Works

Semi-supervised inference: general theory and estimation of means (Anru Zhang, Lawrence D. Brown and T. Tony Cai)

On the other hand, Let $(Y_1, X_1), \dots, (Y_n, X_n)$ be i.i.d. labeled samples from P , and let X_{n+1}, \dots, X_{n+m} be m additional unlabeled independent samples from P_X . Suppose Ξ is nonsingular and $\tau^2 > 0$. If P is fixed and $n \rightarrow \infty$, then

$$\frac{\sqrt{n}(\hat{\theta}_{\text{SSLS}} - \theta)}{\nu} \xrightarrow{d} N(0, 1),$$

and

$$\frac{\hat{\nu}^2}{\nu^2} \xrightarrow{d} 1$$

where $\hat{\nu}^2 = \frac{m}{m+n} \text{MSE} + \frac{n}{m+n} \hat{\sigma}_Y^2$ with $\text{MSE} = \frac{1}{n-p-1} \sum_{k=1}^n (Y_i - \vec{X}_k^\top \hat{\beta})^2$, $\nu^2 = \tau^2 + \frac{n}{n+m} \beta_{(2)}^\top \Sigma_n \beta_{(2)}$ and $\hat{\sigma}_Y^2 = \frac{1}{n-1} \sum_{k=1}^n (Y_i - \bar{\mathbf{Y}})^2$.

Related Works

- Semi-supervised inference: general theory and estimation of means(Anru Zhang, Lawrence D. Brown and T. Tony Cai, The Annals of Statistics) - mean estimation on regression setting
- "Semisupervised inference for explained variance in high dimensional linear regression and its applications(T. Tony Cai and Zijian Guo,JRSS)" - Variance estimation
- "High-dimensional Semi-supervised Learning: in Search for Optimal Inference of the Mean(Yuqian Zhang and Jelena Bradic, Biometrika)" - high-dimensional setting
- "The correlation-assisted missing data estimator(Timothy I. Cannings and Yingying Fan, JMLR)" - missing data approach, extension to U-statistics
- "Methods for correcting inference based on outcomes predicted by machine learning(Siruo Wang, Tyler H. McCormick, and Jeffrey T. Leek, Proceedings of the National Academy of Sciences, 2020)" - different approach using bootstrap
- "Valid inference after prediction(Keshav Motwani and Daniela Witten)" - showing the above method is a bad approach

Doubly Robust Self-Training

Motivation

- Given a teacher model, a large unlabeled dataset and a small labeled dataset, how can we design a principled learning process that ensures consistent and sample-efficient learning of the true model?

Self-Training

- involves using a teacher model to generate pseudo-labels for all unlabeled data, and then training a new model on a mixture of both pseudo-labeled and labeled data.
 - can lead to overreliance on the teacher model and can miss important information provided by the labeled data.
- highly sensitive to the accuracy of the teacher model

Doubly Robust Self-Training

Problem setting

- Given

Unlabeled samples $\mathcal{D}_1 = \{X_1, \dots, X_m\}$ drawn from \mathbb{P}_X ,

Labeled samples $\mathcal{D}_2 = \{(X_{m+1}, Y_{m+1}), \dots, (X_{m+n}, Y_{m+n})\}$ drawn from $\mathbb{P}_X \times \mathbb{P}_{Y|X}$

Pre-trained model $\hat{f} : \mathcal{X} \rightarrow \mathcal{Y}$

- Goal: find θ^* such that $\theta^* = \operatorname{argmin} E[l_\theta(X, Y)]$ for some loss function l_θ

Main result

$$\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{DR} = \frac{1}{m+n} \sum_{i=1}^{m+n} l_\theta(X_i, \hat{f}(X_i)) - \frac{1}{n} \sum_{i=m+1}^{m+n} l_\theta(X_i, \hat{f}(X_i)) + \frac{1}{n} \sum_{i=m+1}^{m+n} l_\theta(X_i, Y_i)$$

is doubly robust!

Doubly Robust Self-Training

Main result

■ Traditional:

$$\begin{aligned}\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{SL}}(\theta) &= \frac{1}{m+n} \left(\sum_{i=1}^m \ell_{\theta}(X_i, \hat{f}(X_i)) + \sum_{i=m+1}^{m+n} \ell_{\theta}(X_i, Y_i) \right) \\ &= \frac{1}{m+n} \sum_{i=1}^{m+n} \ell_{\theta}(X_i, \hat{f}(X_i)) - \frac{1}{m+n} \sum_{i=m+1}^{m+n} \ell_{\theta}(X_i, \hat{f}(X_i)) + \frac{1}{m+n} \sum_{i=m+1}^{m+n} \ell_{\theta}(X_i, Y_i)\end{aligned}$$

■ Doubly Robust:

$$\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR}} = \frac{1}{m+n} \sum_{i=1}^{m+n} l_{\theta}(X_i, \hat{f}(X_i)) - \frac{1}{n} \sum_{i=m+1}^{m+n} l_{\theta}(X_i, \hat{f}(X_i)) + \frac{1}{n} \sum_{i=m+1}^{m+n} l_{\theta}(X_i, Y_i)$$

Doubly Robust Self-Training

Mean estimation

- Goal: find θ^* such that $\theta^* = \operatorname{argmin} E[(\theta - Y)^2]$

- Loss only from labeled data

$$: \mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{TL} = \frac{1}{n} \sum_{i=m+1}^{m+n} (\theta - Y_i)^2$$

$$\rightarrow \hat{\theta}_{TL} = \frac{1}{n} \sum_{i=m+1}^{m+n} Y_i$$

- Loss from self-training

$$: \mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{SL} = \frac{1}{m+n} \left(\sum_{i=1}^m (\theta - \hat{f}(X_i))^2 + \sum_{i=m+1}^{m+n} (\theta - Y_i)^2 \right)$$

$$\rightarrow \hat{\theta}_{SL} = \frac{1}{m+n} \left(\sum_{i=1}^m \hat{f}(X_i) + \sum_{i=m+1}^{m+n} Y_i \right)$$

- Doubly robust loss: $\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{DR} =$

$$\frac{1}{m+n} \sum_{i=1}^{m+n} l_{\theta}(X_i, \hat{f}(X_i)) - \frac{1}{n} \sum_{i=m+1}^{m+n} l_{\theta}(X_i, \hat{f}(X_i)) + \frac{1}{n} \sum_{i=m+1}^{m+n} l_{\theta}(X_i, Y_i)$$

$$\rightarrow \hat{\theta}_{DR} = \frac{1}{m+n} \sum_{i=1}^{m+n} \hat{f}(X_i) - \frac{1}{n} \sum_{i=m+1}^{m+n} \hat{f}(X_i) - Y_i$$

Doubly Robust Self-Training

Mean estimation(cont.)

- $E[(\theta^* - \hat{\theta}_{TL})^2] = \frac{1}{n} \text{Var}(Y)$
 - $E[(\theta^* - \hat{\theta}_{SL})^2] \leq \frac{2m^2}{(m+n)^2} E[(\hat{f}(X) - Y)^2] + \frac{2m}{(m+n)^2} \text{Var}(\hat{f}(X) - Y) + \frac{2n}{(m+n)^2} \text{Var}(Y)$
 - $E[(\theta^* - \hat{\theta}_{DR})^2] \leq 2 \min \left(\frac{1}{n} \text{Var}(Y) + \frac{m+2n}{(m+n)n} \text{Var}(\hat{f}(X)), \frac{m+2n}{(m+n)n} \text{Var}(\hat{f}(X) - Y) + \frac{1}{m+n} \text{Var}(Y) \right)$
- $E[(\theta^* - \hat{\theta}_{DR})^2] \leq \frac{4}{n} (\text{Var}[Y] + \text{Var}[\hat{f}(X)])$ no matter how poor the estimator $\hat{f}(X)$ is.
- when $\text{Var}[\hat{f}(X) - Y]$ is small, $E[(\theta^* - \hat{\theta}_{DR})^2] \leq \frac{2}{m+n} \text{Var}[Y]$.
- $E[(\theta^* - \hat{\theta}_{SL})^2]$ always has a non-vanishing term, $\frac{2m^2}{(m+n)^2} \mathbb{E}[(\hat{f}(X) - Y)]^2$ unless the predictor is accurate.

Doubly Robust Self-Training

General loss

With probability at least $1 - \delta$,

$$\begin{aligned} \|\nabla_{\theta} \mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR}}(\theta^*)\|_2 \leq C \min & \left(\|\Sigma_{\theta^*}^{\hat{f}}\|_2 \sqrt{\frac{d}{(m+n)\delta}} + \|\Sigma_{\theta^*}^{Y-\hat{f}}\|_2 \sqrt{\frac{d}{n\delta}}, \right. \\ & \left. \|\Sigma_{\theta^*}^{\hat{f}}\|_2 \left(\sqrt{\frac{d}{(m+n)\delta}} + \sqrt{\frac{d}{n\delta}} \right) + \|\Sigma_{\theta^*}^Y\|_2 \sqrt{\frac{d}{n\delta}} \right), \end{aligned}$$

where C is a universal constant, and we denote

$\Sigma_{\theta}^{Y-\hat{f}} = \text{Cov} [\nabla_{\theta} \ell_{\theta}(X, \hat{f}(X)) - \nabla_{\theta} \ell_{\theta}(X, Y)]$ and let $\Sigma_{\theta}^{\hat{f}} = \text{Cov} [\nabla_{\theta} \ell_{\theta}(X, \hat{f}(X))]$,
 $\Sigma_{\theta}^Y = \text{Cov} [\nabla_{\theta} \ell_{\theta}(X, Y)]$.

- $\|\nabla_{\theta} \mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{SL}}(\theta^*)\|_2 \geq C$ for some positive constant C
- When \hat{f} is a perfect predictor, one has $\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR}}(\theta^*) = \frac{1}{m+n} \sum_{i=1}^{m+n} \ell_{\theta}(X_i, Y_i)$.

Doubly Robust Self-Training

The case of Distribution mismatch

$$\begin{aligned}\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR2}}(\theta) &= \frac{1}{m} \sum_{i=1}^m \ell_{\theta}(X_i, \hat{f}(X_i)) - \frac{1}{n} \sum_{i=m+1}^{m+n} \frac{1}{\pi(X_i)} \ell_{\theta}(X_i, \hat{f}(X_i)) \\ &\quad + \frac{1}{n} \sum_{i=m+1}^{m+n} \frac{1}{\pi(X_i)} \ell_{\theta}(X_i, Y_i)\end{aligned}$$

$\mathbb{E}[\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR2}}(\theta)] = \mathbb{E}_{\mathbb{P}_{X,Y}}[\ell_{\theta}(X, Y)]$ as long as one of the two assumptions hold:

- For any x , $\pi(x) = \frac{\mathbb{P}_X(x)}{\mathbb{Q}_X(x)}$.
- For any x , $\ell_{\theta}(x, \hat{f}(x)) = \mathbb{E}_{Y \sim \mathbb{P}_{Y|X=x}}[\ell_{\theta}(x, Y)]$.

Doubly Robust Self-Training

Experiment: ImageNet

Minimize the curriculum-based loss in epoch

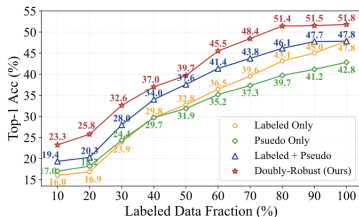
$$\mathcal{L}_{\mathcal{D}_1, \mathcal{D}_2}^{\text{DR}, t}(\theta) = \frac{1}{m+n} \sum_{i=1}^{m+n} \ell_{\theta}(X_i, \hat{f}(X_i)) \\ - \alpha_t \cdot \left(\frac{1}{n} \sum_{i=m+1}^{m+n} \ell_{\theta}(X_i, \hat{f}(X_i)) - \frac{1}{n} \sum_{i=m+1}^{m+n} \ell_{\theta}(X_i, Y_i) \right).$$

Table 1: Comparisons on mini-ImageNet100, all models trained for 100 epochs.

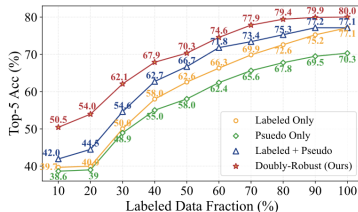
Labeled Data Percent	Labeled Only		Pseudo Only		Labeled + Pseudo		Doubly robust Loss	
	top1	top5	top1	top5	top1	top5	top1	top5
1	2.72	9.18	2.81	9.57	2.73	9.55	2.75	9.73
5	3.92	13.34	4.27	13.66	4.27	14.4	4.89	16.38
10	6.76	20.84	7.27	21.64	7.65	22.48	8.01	21.90
20	12.3	31.3	13.46	30.79	13.94	32.63	13.50	32.17
50	20.69	46.86	20.92	45.2	24.9	50.77	25.31	51.61
80	27.37	55.57	25.57	50.85	30.63	58.85	30.75	59.41
100	31.07	60.62	28.95	55.35	34.33	62.78	34.01	63.04

Doubly Robust Self-Training

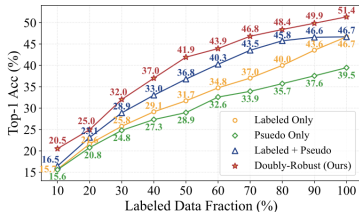
Experiment: ImageNet



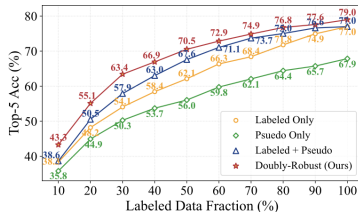
(a) Top-1 on DaViT



(b) Top-5 on DaViT



(c) Top-1 on ResNet50



(d) Top-5 on ResNet50

Figure 1: Comparisons on ImageNet100 using two different network architectures. Both Top-1 and Top-5 accuracies are reported. All models are trained for 20 epochs.

Doubly Robust Self-Training

Experiment: nuScenes

$$\mathcal{L}_{obj}^{DR}(\theta) = \frac{1}{M + N_{ps}} \sum_{i=1}^{M+N_{ps}} \ell_{\theta}(X_i, f(X_i)) \\ - \frac{1}{N_{ps}} \sum_{i=M+1}^{M+N_{ps}} \ell_{\theta}(X'_i, f(X'_i)) + \frac{1}{N} \sum_{i=M+1}^{M+N} \ell_{\theta}(X_i, Y_i),$$

Table 2: Performance comparison on nuScenes *val* set.

Labeled Data Fraction	Labeled Only		Labeled + Pseudo		Doubly robust Loss	
	mAP↑	NDS↑	mAP↑	NDS↑	mAP↑	NDS↑
1/24	7.56	18.01	7.60	17.32	8.18	18.33
1/16	11.15	20.55	11.60	21.03	12.30	22.10
1/4	25.66	41.41	28.36	43.88	27.48	43.18

Table 3: Per-class mAP (%) comparison on nuScenes *val* set using 1/16 of total labels in training.

	Car	Ped	Truck	Bus	Trailer	Barrier	Traffic Cone
Labeled Only	48.6	30.6	8.5	6.2	4.0	6.8	4.4
Labeled + Pseudo	48.8	30.9	8.8	7.5	5.7	6.7	4.0
Improvement	+0.2	+0.3	+0.3	+1.3	+1.7	-0.1	-0.4
Doubly robust Loss	51.5	32.9	9.6	8.2	5.2	7.2	4.5
Improvement	+2.9	+2.3	+1.1	+2.0	+1.2	+0.4	+0.1

Prediction-Powered Inference

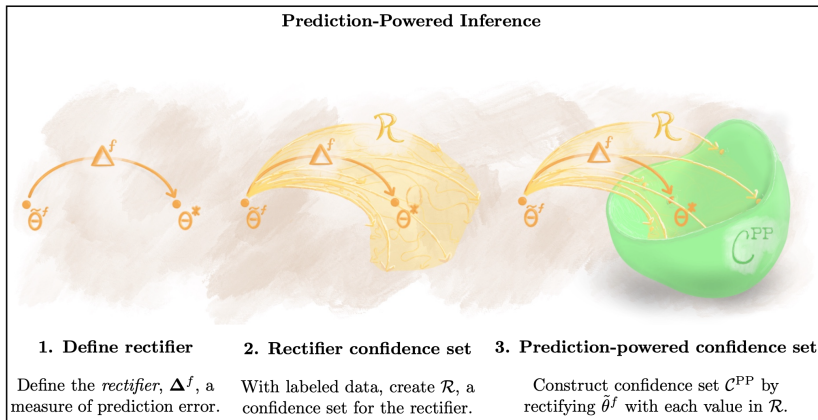
Motivation

How can we assess the role of prediction in terms of basic principles of statistical inference? Is it possible to exploit predictions from a machine-learning system while still providing guarantees of statistical validity?

- suggest a framework for performing valid statistical inference when an experimental data set is supplemented with predictions.
- use predictions from the model to perform inference, leverage the immense number of predictions to improve their confidence in a scientific conclusion.

Prediction-Powered Inference

General Framework



Prediction-Powered Inference

Mean estimation

Our estimates

$$\hat{\theta}^{class} = \frac{1}{n} \sum_{i=1}^n Y_i \quad \text{vs} \quad \hat{\theta}^{PP} = \underbrace{\frac{1}{N} \sum_{i=1}^N \tilde{f}_i}_{\tilde{\theta}^f} - \underbrace{\frac{1}{n} \sum_{i=1}^n (f_i - Y_i)}_{\hat{\Delta}^f}.$$

and their corresponding confidence intervals

$$\underbrace{\hat{\theta}^{class} \pm 1.96 \sqrt{\frac{\hat{\sigma}_Y^2}{n}}}_{\text{classical interval}} \quad \text{or} \quad \underbrace{\hat{\theta}^{PP} \pm 1.96 \sqrt{\frac{\hat{\sigma}_{f-Y}^2}{n} + \frac{\hat{\sigma}_{\tilde{f}}^2}{N}}}_{\text{prediction-powered interval}}$$

Prediction-Powered Inference

Convex estimation

Consider estimands of the form

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \mathbb{E} [\ell_{\theta} (X_1, Y_1)],$$

where a loss function ℓ_{θ} is convex.

From convexity, it holds that

$$\mathbb{E} [g_{\theta^*} (X_1, Y_1)] = 0$$

where $g_{\theta} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^p$ is a subgradient of ℓ_{θ} with respect to θ .

Define the rectifier

$$\Delta^f(\theta) = \mathbb{E} [g_{\theta} (X_1, Y_1) - g_{\theta} (X_1, f_1)].$$

Prediction-Powered Inference

Convex estimation

Create the confidence set for the rectifier, $\mathcal{R}_\delta(\theta)$,

$$P\left(\Delta^f(\theta) \in \mathcal{R}_\delta(\theta)\right) \geq 1 - \delta.$$

Also, for every θ , we want a confidence set $\mathcal{T}_{\alpha-\delta}(\theta)$ for $\mathbb{E}[g_\theta(X_1, f_1)]$, satisfying

$$P(\mathbb{E}[g_\theta(X_1, f_1)] \in \mathcal{T}_{\alpha-\delta}(\theta)) \geq 1 - (\alpha - \delta).$$

Using the results, finally form the confidence set for θ^* as follows

$$P(\theta^* \in \mathcal{C}_\alpha^{\text{PP}}) \geq 1 - \alpha$$

where $\mathcal{C}_\alpha^{\text{PP}} = \{\theta : 0 \in \mathcal{R}_\delta(\theta) + \mathcal{T}_{\alpha-\delta}(\theta)\}$ where $+$ denotes the Minkowski sum.

Prediction-Powered Inference

Various algorithms

Algorithm 1 Prediction-powered mean estimation

Input: labeled data (X, Y) , unlabeled features \tilde{X} , predictor f , error level $\alpha \in (0, 1)$

- 1: $\hat{\theta}^{\text{PP}} \leftarrow \hat{\theta}^f - \hat{\Delta}^f := \frac{1}{N} \sum_{i=1}^N \tilde{f}_i - \frac{1}{n} \sum_{i=1}^n (f_i - Y_i)$ ▷ prediction-powered estimator
- 2: $\hat{\sigma}_f^2 \leftarrow \frac{1}{N} \sum_{i=1}^N (\tilde{f}_i - \hat{\theta}^f)^2$ ▷ empirical variance of imputed estimate
- 3: $\hat{\sigma}_{f-Y}^2 \leftarrow \frac{1}{n} \sum_{i=1}^n (f_i - Y_i - \hat{\Delta}^f)^2$ ▷ empirical variance of empirical rectifier
- 4: $w_\alpha \leftarrow z_{1-\alpha/2} \sqrt{\frac{\hat{\sigma}_{f-Y}^2}{n} + \frac{\hat{\sigma}_f^2}{N}}$ ▷ normal approximation

Output: prediction-powered confidence set $\mathcal{C}_\alpha^{\text{PP}} = (\hat{\theta}^{\text{PP}} \pm w_\alpha)$

Algorithm 2 Prediction-powered quantile estimation

Input: labeled data (X, Y) , unlabeled features \tilde{X} , predictor f , quantile $q \in (0, 1)$, error level $\alpha \in (0, 1)$

- 1: Construct fine grid Θ_{grid} between $\min_{i \in [N]} \tilde{f}_i$ and $\max_{i \in [N]} \tilde{f}_i$
- 2: **for** $\theta \in \Theta_{\text{grid}}$ **do**
- 3: $\hat{\Delta}^f(\theta) \leftarrow \frac{1}{n} \sum_{i=1}^n (\mathbb{1}\{Y_i \leq \theta\} - \mathbb{1}\{f_i \leq \theta\})$ ▷ empirical rectifier
- 4: $\hat{F}(\theta) \leftarrow \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{\tilde{f}_i \leq \theta\}$ ▷ imputed CDF
- 5: $\hat{\sigma}_{\hat{\Delta}}^2(\theta) \leftarrow \frac{1}{n} \sum_{i=1}^n (\mathbb{1}\{Y_i \leq \theta\} - \mathbb{1}\{f_i \leq \theta\} - \hat{\Delta}^f(\theta))^2$ ▷ empirical variance of empirical rectifier
- 6: $\hat{\sigma}_{\hat{F}}^2(\theta) \leftarrow \frac{1}{N} \sum_{i=1}^N (\mathbb{1}\{\tilde{f}_i \leq \theta\} - \hat{F}(\theta))^2$ ▷ empirical variance of imputed CDF
- 7: $w_\alpha(\theta) \leftarrow z_{1-\alpha/2} \sqrt{\frac{\hat{\sigma}_{\hat{\Delta}}^2(\theta)}{n} + \frac{\hat{\sigma}_{\hat{F}}^2(\theta)}{N}}$ ▷ normal approximation

Output: prediction-powered confidence set $\mathcal{C}_\alpha^{\text{PP}} = \{\theta \in \Theta_{\text{grid}} : |\hat{F}(\theta) + \hat{\Delta}^f(\theta) - q| \leq w_\alpha(\theta)\}$

Prediction-Powered Inference

Various algorithms(cont.)

Algorithm 3 Prediction-powered logistic regression

Input: labeled data (X, Y) , unlabeled features \tilde{X} , predictor f , error level $\alpha \in (0, 1)$

- 1: Construct fine grid $\Theta_{\text{grid}} \subset \mathbb{R}^d$ of possible coefficients
- 2: $\hat{\Delta}_j^f \leftarrow \frac{1}{n} \sum_{i=1}^n X_{i,j} (f_i - Y_i)$, $j \in [d]$ ▷ empirical rectifier
- 3: $\hat{\sigma}_{\Delta,j}^2 \leftarrow \frac{1}{n} \sum_{i=1}^n \left(X_{i,j} (f_i - Y_i) - \hat{\Delta}_j^f \right)^2$, $j \in [d]$ ▷ empirical variance of empirical rectifier
- 4: **for** $\theta \in \Theta_{\text{grid}}$ **do**
- 5: $\hat{g}_j^f(\theta) \leftarrow \frac{1}{N} \sum_{i=1}^N \tilde{X}_{i,j} \left(\mu_\theta(\tilde{X}_i) - \tilde{f}_i \right)$, $j \in [d]$, where $\mu_\theta(x) = \frac{1}{1 + \exp(-x^\top \theta)}$ ▷ imputed gradient
- 6: $\hat{\sigma}_{g,j}^2(\theta) \leftarrow \frac{1}{N} \sum_{i=1}^N \left(\tilde{X}_{i,j} (\mu_\theta(\tilde{X}_i) - \tilde{f}_i) - \hat{g}_j^f(\theta) \right)^2$, $j \in [d]$ ▷ empirical variance of imputed gradient
- 7: $w_{\alpha,j}(\theta) \leftarrow z_{1-\alpha/(2d)} \sqrt{\frac{\hat{\sigma}_{\Delta,j}^2}{n} + \frac{\hat{\sigma}_{g,j}^2(\theta)}{N}}$, $j \in [d]$ ▷ normal approximation

Output: prediction-powered confidence set $\mathcal{C}_\alpha^{\text{PP}} = \left\{ \theta \in \Theta_{\text{grid}} : |\hat{g}_j^f(\theta) + \hat{\Delta}_j^f| \leq w_{\alpha,j}(\theta), \forall j \in [d] \right\}$

Algorithm 4 Prediction-powered linear regression

Input: labeled data (X, Y) , unlabeled features \tilde{X} , predictor f , coefficient $j^* \in [d]$, error level $\alpha \in (0, 1)$

- 1: $\hat{\theta}^{\text{PP}} \leftarrow \tilde{\theta}^f - \hat{\Delta}^f := \tilde{X}^\top \tilde{f} - X^\top (f - Y)$ ▷ prediction-powered estimator
- 2: $\tilde{\Sigma} \leftarrow \frac{1}{N} \tilde{X}^\top \tilde{X}$, $\tilde{M} \leftarrow \frac{1}{N} \sum_{i=1}^N (\tilde{f}_i - \tilde{X}_i^\top \tilde{\theta}^f)^2 \tilde{X}_i \tilde{X}_i^\top$
- 3: $\tilde{V} \leftarrow \tilde{\Sigma}^{-1} \tilde{M} \tilde{\Sigma}^{-1}$ ▷ “sandwich” variance estimator for imputed estimate
- 4: $\Sigma \leftarrow \frac{1}{n} X^\top X$, $M \leftarrow \frac{1}{n} \sum_{i=1}^n (f_i - Y_i - X_i^\top \hat{\Delta}^f)^2 X_i X_i^\top$
- 5: $V \leftarrow \Sigma^{-1} M \Sigma^{-1}$ ▷ “sandwich” variance estimator for empirical rectifier
- 6: $w_\alpha \leftarrow z_{1-\alpha/2} \sqrt{\frac{V_{j^*,j^*}}{n} + \frac{\tilde{V}_{j^*,j^*}}{N}}$ ▷ normal approximation

Output: prediction-powered confidence set $\mathcal{C}_\alpha^{\text{PP}} = \left(\hat{\theta}_{j^*}^{\text{PP}} \pm w_\alpha \right)$

Prediction-Powered Inference

Various algorithms(cont.)

- mean: $\theta^* = \arg \min_{\theta \in \mathbb{R}} \mathbb{E} [\ell_{\theta} (Y_1)] = \arg \min_{\theta \in \mathbb{R}} \mathbb{E} \left[\frac{1}{2} (Y_1 - \theta)^2 \right]$
- quantile: $\theta^* = \arg \min_{\theta \in \mathbb{R}} \mathbb{E} [\ell_{\theta} (Y_1)] =$
 $\arg \min_{\theta \in \mathbb{R}} \mathbb{E} [q (Y_1 - \theta) 1\{Y_1 > \theta\} + (1 - q) (\theta - Y_1) 1\{Y_1 \leq \theta\}]$
- logistic regression:
 $\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E} [\ell_{\theta} (X_1, Y_1)] = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E} [-Y_1 \theta^{\top} X + \log (1 + \exp (\theta^{\top} X_1))]$
- linear regression: $\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E} [\ell_{\theta} (X_1, Y_1)] = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E} [(Y_1 - X_1^{\top} \theta)^2]$

Then powered confidence set in algorithms have valid coverage:

$$\liminf_{n, N \rightarrow \infty} P(\theta^* \in \mathcal{C}_{\alpha}^{\text{PP}}) \geq 1 - \alpha$$

Prediction-Powered Inference

Beyond convex estimation

1. estimate $\mathbb{E}[\ell_{\theta^*}(X_1, Y_1)]$ by approximating θ^* with an imputed estimate on the first $N/2$ unlabeled data points

$$\tilde{\theta}^f = \arg \min_{\theta \in \Theta} \frac{2}{N} \sum_{i=1}^{N/2} \ell_{\theta}(\tilde{X}_i, \tilde{f}_i), \quad \tilde{L}^f(\theta) := \frac{2}{N} \sum_{i=N/2+1}^N \ell_{\theta}(\tilde{X}_i, \tilde{f}_i).$$

2. construct $(\mathcal{R}_{\delta/2}^l(\theta), \mathcal{R}_{\delta/2}^u(\theta))$ and $\mathcal{T}_{\alpha-\delta}(\theta)$ such that

$$P\left(\Delta^f(\theta) \leq \mathcal{R}_{\delta/2}^u(\theta)\right) \geq 1 - \delta/2; \quad P\left(\Delta^f(\theta) \geq \mathcal{R}_{\delta/2}^l(\theta)\right) \geq 1 - \delta/2;$$
$$P\left(\tilde{L}^f(\theta) - \mathbb{E}[\ell_{\theta}(X_1, f_1)] \leq \mathcal{T}_{\alpha-\delta}(\theta)\right) \geq 1 - (\alpha - \delta).$$

Prediction-Powered Inference

Beyond convex estimation

3. combining 1 2, obtain

$$\mathcal{C}_\alpha^{\text{PP}} = \left\{ \theta \in \Theta : \tilde{L}^f(\theta) \leq \tilde{L}^f(\tilde{\theta}^f) - \mathcal{R}_{\delta/2}^l(\theta) + \mathcal{R}_{\delta/2}^u(\tilde{\theta}^f) + \mathcal{T}_{\alpha-\delta}(\theta) \right\}$$

such that

$$P(\theta^* \in \mathcal{C}_\alpha^{\text{PP}}) \geq 1 - \alpha.$$

Prediction-Powered Inference

Experiments

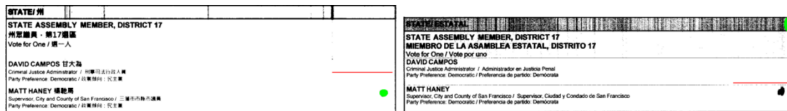


Figure 3: **Examples of ballots** correctly and incorrectly classified. The raw ballot is black and white, the voter's marking is automatically identified by a computer vision algorithm with a green annotation, and markings below the red line annotation will be considered votes for Matt Haney (and vice versa). The instructional portion of the ballots was cropped out.

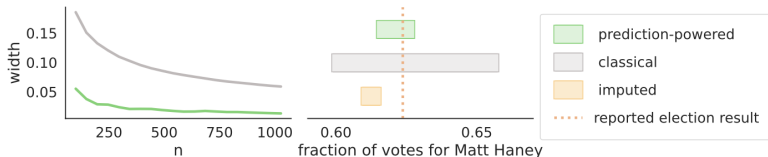


Figure 4: **Election results** produced by prediction-powered inference and the classical and imputed baselines at level 95%. Left: width of intervals as a function of n . Right: confidence intervals with $n = 1024$.

Prediction-Powered Inference

Experiments

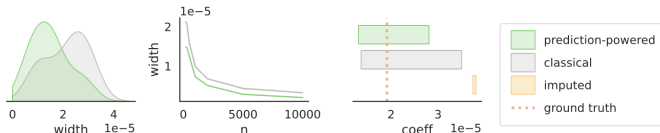


Figure 9: **Confidence intervals for the logistic regression coefficient** relating income and private health insurance coverage at the 95% level. Left: distribution of interval widths with $n = 200$. Middle: mean width as a function of n . Right: intervals with $n = 200$.

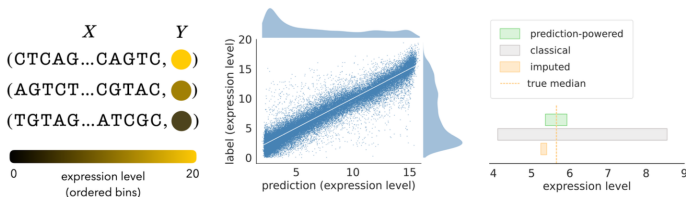


Figure 10: **Predicting gene expression levels from a promoter sequence** [34]. Left: each data point consists of a promoter sequence, X_i , and an expression level, Y_i . Middle: predictive performance of the transformer model on the native yeast promoters used in our experiments (RMSE 2.18, Pearson 0.963, Spearman 0.946). Right: confidence intervals for the median native yeast promoter expression level with $n = 75$ and $\alpha = 0.1$.

Further Studies

- High-dimension setting: what if $p \rightarrow \infty$?
- Using sample splitting?
- General framework for other estimators such as kernel mean embedding, MMD and other U-statistics
- Applying such for estimating mean difference or other test statistics
- Computing statistical validity for other semi-/self-supervised learning models
-

Reference

- Kim, I. (2023). Selective Topics in Mathematical Statistics. 33-35.
- Zhang, A., Brown, L. D., Cai, T. T. (2019). Semi-supervised inference: General theory and estimation of means.
- Zhang, Y., Bradic, J. (2022). High-dimensional semi-supervised learning: in search of optimal inference of the mean. *Biometrika*, 109(2), 387-403.
- Tony Cai, T., Guo, Z. (2020). Semisupervised inference for explained variance in high dimensional linear regression and its applications. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(2), 391-419.
- Angelopoulos, A. N., Bates, S., Fannjiang, C., Jordan, M. I., Zrnic, T. (2023). Prediction-powered inference. arXiv preprint arXiv:2301.09633.
- Zhu, B., Ding, M., Jacobson, P., Wu, M., Zhan, W., Jordan, M., Jiao, J. (2023). Doubly Robust Self-Training. arXiv preprint arXiv:2306.00265.