

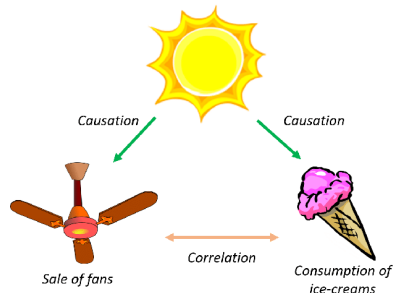
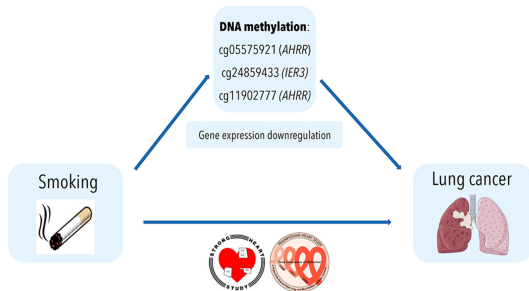
Towards a Unified Analysis of Neural Networks in Nonparametric Instrumental Variable Regression: Optimization and Generalization

Zonghao Chen Atsushi Nitanda Arthur Gretton Taiji Suzuki

December 14, 2025

Background

- Nonparametric regression is used everywhere in statistics.
 - Kernel based estimators, (deep) neural networks, nearest neighbours, etc.
- Regression fails when
 - there exist **confounding** that affects both the input and the output.
 - A classical example: does smoking cause lung cancer?



How to prove a direct relation under (unobserved) confounding?

Background: Causal Inference and Instrumental Variable

- The causal effect of smoking X on the risk of lung cancer Y .
- Unobserved confounding ϵ : gene, occupation, childhood.

$$Y = h_o(X) + \epsilon, \quad \text{where } \epsilon \not\perp X.$$

- h_o is the target of interest.
- In this talk, we focus on **additive** confounder.
- Regression always outputs a **biased** estimate $\mathbb{E}[Y | X] = h_o(X) + \mathbb{E}[\epsilon | X] \neq h_o(X)$.
- Instrumental variable Z that affects Y only through X : price of the cigarette.

$$\epsilon \perp Z$$

- Conditioning both sides on Z

$$\mathbb{E}[Y | Z] = \mathbb{E}[h_o(X) | Z]. \quad (\text{NPIV})$$

- **Nonparametric** way of estimating h_o .
- Given i.i.d. samples $\{\mathbf{x}_i, \mathbf{z}_i, y_i\}_{i=1}^n$.

Background: NPIV versus NPR

- Nonparametric regression (NPR): $Y = h_o(X) + \epsilon$ with $\epsilon \perp X$.
 - Conditioning both sides on X :

$$\mathbb{E}[Y | X] = h_o(X). \quad (\text{NPR})$$

- Target $h_o = \arg \min_h \mathbb{E}_{YX}[(Y - h(X))^2]$.
- Least squares estimator

$$\theta^* = \arg \min_{\theta} \frac{1}{n} \sum_{i=1}^n (h_{\theta}(\mathbf{x}_i) - y_i)^2, \quad \{\mathbf{x}_i, y_i\}_{i=1}^n \sim P_{XY}.$$

- h_{θ} is a neural network parameterized by θ .
- Nonparametric instrumental variable regression (NPIV): $Y = h_o(X) + \epsilon$ with $\epsilon \not\perp X$ but $\epsilon \perp Z$.
 - Conditioning both sides on Z :

$$\mathbb{E}[Y | Z] = \mathbb{E}[h_o(X) | Z]. \quad (\text{NPIV})$$

Problem: How to estimate h_o given $\{\mathbf{x}_i, y_i, \mathbf{z}_i\}_{i=1}^n \sim P_{XYZ}$?

Background: Offline reinforcement learning

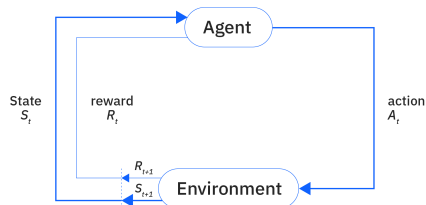
- s : state, a : action, r : reward, γ : discount factor.
- $Q(s, a) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a]$ denotes the expected long-term return when taking action a in state s .
- Bellman equation:

$$\mathbb{E}[r \mid s, a] = Q(s, a) - \gamma \mathbb{E}[Q(s', a') \mid s, a].$$

- The Bellman equation shares the same structure as NPIV:

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h_o(X) \mid Z].$$

- Correspondence: $Y \rightarrow r$, $Z \rightarrow (s, a)$, $X \rightarrow (s', a')$.



Background: Proximal causal inference

- X treatment, Y outcome, W outcome proxy, Z treatment proxy
- Conditional moment equation:

$$\mathbb{E}[Y \mid Z, X] = \mathbb{E}[h_o(W, X) \mid Z, X].$$

- Bridge function h_o
- Correspondence with NPIV

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h_o(X) \mid Z]$$

- $Y \rightarrow Y, Z \rightarrow (Z, X), X \rightarrow (W, X).$

How to solve NPIV

Question: How to find h_o in the NPIV equation $\mathbb{E}[Y | Z] = \mathbb{E}[h_o(X) | Z]$?

Background: NPIV and 2SLS

- Define $T : L^2(P_X) \rightarrow L^2(P_Z)$ as the **unknown** conditional expectation operator defined by $(Tf)(Z) = \mathbb{E}[f(X) | Z]$.

$$\mathbb{E}[Y | Z] = \mathbb{E}[h_o(X) | Z] =: (Th_o)(Z). \quad (\text{NPIV})$$

- Target $h_o = \arg \min_h \mathbb{E}_{YZ}[(Y - (Th)(Z))^2]$.
 - T is **unknown** yet it is a **conditional expectation** and hence can be learned via regression.
- Two-stage least squares (2SLS)
 - h_{θ_x} and h_{θ_z} are two neural networks.

$$\theta_z^*(\theta_x) = \arg \min_{\theta_z} \frac{1}{m} \sum_{i=1}^m (h_{\theta_z}(\mathbf{z}_i) - h_{\theta_x}(\mathbf{x}_i))^2, \quad \{\mathbf{x}_i, \mathbf{z}_i\}_{i=1}^m \sim P_{XZ} \quad (\text{2SLS})$$

$$\theta_x^* = \arg \min_{\theta_x} \frac{1}{n} \sum_{i=1}^n (h_{\theta_z^*(\theta_x)}(\mathbf{z}_i) - y_i)^2, \quad \{\mathbf{z}_i, y_i\}_{i=1}^n \sim P_{ZY}.$$

- $h_{\theta_z^*(\theta_x)}(\mathbf{z}_i) \approx \mathbb{E}[h_{\theta_x}(X) | Z = \mathbf{z}_i] = (Th_{\theta_x})(\mathbf{z}_i).$

Target: Optimization and Generalization

- Two-stage least squares (2SLS)
 - Both h_{θ_x} and h_{θ_z} are neural networks.

Stage I: $\theta_z^*(\theta_x) = \arg \min_{\theta_z} \frac{1}{m} \sum_{i=1}^m (h_{\theta_z}(\mathbf{z}_i) - h_{\theta_x}(\mathbf{x}_i))^2, \quad \{\mathbf{x}_i, \mathbf{z}_i\}_{i=1}^m \sim P_{XZ}$

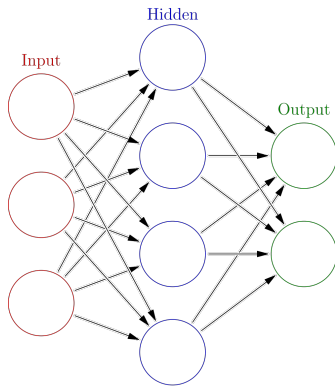
(2SLS)

Stage II: $\theta_x^* = \arg \min_{\theta_x} \frac{1}{n} \sum_{i=1}^n (h_{\theta_z^*(\theta_x)}(\mathbf{z}_i) - y_i)^2, \quad \{\mathbf{z}_i, y_i\}_{i=1}^n \sim P_{ZY}.$

Bilevel optimization theory: Does gradient based algorithm can actually find the **global optimum** θ_z^*, θ_x^* ? If it does, what is the **iteration complexity**?

Statistical theory: Given the global optimal θ_z^*, θ_x^* , is $h_{\theta_x^*}$ a **consistent** estimator of h_o ? If it does, what is the **sample complexity**?

Mean-field neural networks



Background: Mean-field two-layer neural networks

- Consider neural networks with a single hidden layer of size N :
 - $\mathcal{X} = [x^{(1)}, \dots, x^{(N)}] \in (\mathbb{R}^{d_x})^N$ are the network parameters and \mathbf{x} is the network input

$$h(\mathbf{x}, \mathcal{X}) = \frac{1}{N} \sum_{i=1}^N \Psi(\mathbf{x}, x^{(i)})$$

- Here, $\Psi(\mathbf{x}, x) = w_2 a(w_1^\top \mathbf{x} + b)$ with parameters $x = (w_1, w_2, b)$ and a being an activation function.
- As the empirical distribution $\frac{1}{N} \sum_{i=1}^N \delta_{x^{(i)}} \rightarrow \mu$ as $N \rightarrow \infty$:

$$h_\mu(\mathbf{x}) = \int \Psi(\mathbf{x}, x) d\mu(x) = \mathbb{E}_{X \sim \mu}[\Psi(\mathbf{x}, X)].$$

- So called "mean-field neural network".

Question: What is the purpose of considering the mean-field limit?

Background: Mean-field perspective of two-layer neural networks

- Consider the squared loss with ℓ_2 -norm regularization

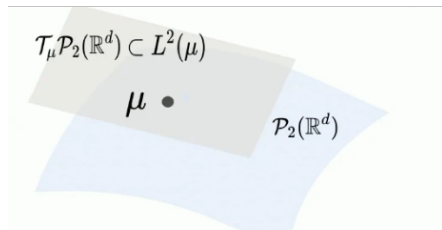
$$F(\mu) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}, y) \sim \rho} [(\mathbb{E}_{\mathbf{X} \sim \mu} [\Psi(\mathbf{x}, \mathbf{X})] - y)^2] + \frac{\zeta}{2} \mathbb{E}_{\mathbf{X} \sim \mu} [\|\mathbf{X}\|^2],$$

- ρ is a data distribution, e.g. $\rho = \frac{1}{n} \sum_{i=1}^n \delta_{(\mathbf{x}_i, y_i)}$.
- F is **linear convex** in μ : for any probability measures $\mu, \nu \in \mathcal{P}$,

$$F(\vartheta\mu + (1 - \vartheta)\nu) \leq \vartheta F(\mu) + (1 - \vartheta)F(\nu), \quad \forall \vartheta \in (0, 1).$$

- Optimization problem in \mathcal{P} : $\min F(\mu)$.
- Wasserstein gradient flow!

Background: Mean-field perspective of two-layer neural networks



Definition (Wasserstein gradient)

Let $\mathcal{G} : \mathcal{P}_2(\mathbb{R}^d) \rightarrow \mathbb{R}$ be a regular functional. The **Wasserstein gradient of \mathcal{G}** evaluated at $\mu \in \mathcal{P}_2(\mathbb{R}^d)$ is the unique function $\nabla \mathcal{G}(\mu) : \mathbb{R}^d \rightarrow \mathbb{R}^d$, s.t. for any $T \in \mathcal{T}_\mu \mathcal{P}_2(\mathbb{R}^d)$,

$$\lim_{\epsilon \rightarrow 0} \frac{1}{\epsilon} [\mathcal{G}((\text{Id} + \epsilon T)_\# \mu) - \mathcal{G}(\mu)] = \int [\nabla \mathcal{G}(\mu)](x)^\top T(x) \, d\mu(x) = \langle \nabla \mathcal{G}, T \rangle_{L^2(\mu)}.$$

- Wasserstein gradient of F at $\mu \in \mathcal{P}$ evaluated at $\mathbf{x} \in \mathbb{R}^d$.

$$\nabla F(\mu)(\mathbf{x}) = \mathbb{E}_{(\mathbf{x}, y) \sim \rho} [(\mathbb{E}_{X \sim \mu} [\Psi(\mathbf{x}, X)] - y) \nabla \Psi(\mathbf{x}, \mathbf{x})] + \zeta \mathbf{x}.$$

Background: Mean-field perspective of two-layer neural networks

- Wasserstein gradient of F at $\mu \in \mathcal{P}$ evaluated at $x \in \mathbb{R}^d$.

$$\nabla F(\mu)(\mathbf{x}) = \mathbb{E}_{(\mathbf{x}, y) \sim \rho} [(\mathbb{E}_{X \sim \mu} [\Psi(\mathbf{x}, X)] - y) \nabla \Psi(\mathbf{x}, \mathbf{x})] + \zeta \mathbf{x}.$$

- Wasserstein gradient of F at $\mu_{\mathcal{X}} = \frac{1}{N} \sum_{i=1}^N \delta_{x^{(i)}}$ evaluated at $x^{(i)} \in \mathbb{R}^d$.

$$\nabla F(\mu_{\mathcal{X}})(\mathbf{x}^{(i)}) = \mathbb{E}_{(\mathbf{x}, y) \sim \rho} \left[\left(\frac{1}{N} \sum_{i=1}^N \Psi(\mathbf{x}, \mathbf{x}^{(i)}) - y \right) \nabla \Psi(\mathbf{x}, \mathbf{x}^{(i)}) \right] + \zeta \mathbf{x}^{(i)}.$$

- Euclidean gradient of the loss

$$L(\mathcal{X}) := \frac{1}{2} \mathbb{E}_{(\mathbf{x}, y) \sim \rho} \left[\left(\frac{1}{N} \sum_{i=1}^N \Psi(\mathbf{x}, x^{(i)}) - y \right)^2 \right] + \frac{\zeta}{2} \frac{1}{N} \sum_{i=1}^N [\|x^{(i)}\|^2],$$

$$\nabla_{\mathbf{x}^{(i)}} L(\mathcal{X}) = \mathbb{E}_{(\mathbf{x}, y) \sim \rho} \left[\left(\frac{1}{N} \sum_{i=1}^N \Psi(\mathbf{x}, \mathbf{x}^{(i)}) - y \right) \frac{\nabla \Psi(\mathbf{x}, \mathbf{x}^{(i)})}{N} \right] + \frac{\zeta}{N} \mathbf{x}^{(i)}$$

- The Wasserstein gradient descent $x_{s+1}^{(i)} = x_s^{(i)} - \gamma \nabla F(\mu_{\mathcal{X}_s})(x_s^{(i)})$ coincides with the Euclidean gradient descent $x_{s+1}^{(i)} = x_s^{(i)} - \gamma \nabla_{x^{(i)}} L(\mathcal{X}_s)$ with a rescaled learning rate!

Background: Mean-field perspective of two-layer neural networks

- At iteration $s \in \{0, \dots, S\}$ and for any $i \in \{1, \dots, N\}$:

$$x_{s+1}^{(i)} = x_s^{(i)} - \gamma \nabla F(\mu_{\mathcal{X}})(x_s^{(i)}) + \sqrt{2\sigma\gamma} \xi_s^{(i)}.$$

- $\{\xi_s^{(i)}\}_{i=1}^N$ are N i.i.d samples from d dimensional unit Gaussian.
- Define an **entropic regularized** objective $\mathcal{F}(\mu) = F(\mu) + \sigma \text{Ent}(\mu)$.
 - $\text{Ent}(\mu) = \int \log \mu(x) \mu(x) dx$.
 - Noisy gradient descent is Wasserstein gradient descent of \mathcal{F} .
- The global optimum $\mu^* := \arg \min_{\mu} \mathcal{F}(\mu)$.

Question: Does noisy gradient descent can actually find the **global optimum** μ^* ? If it does, what is the **iteration complexity**?

Assumption (Bounded and smooth neural networks)

There exists a universal positive constant R such that $\sup_{x \in \mathbb{R}^{d_x}, x \in \mathcal{X}} |\Psi_x(x)| \leq R$ and $\sup_{x \in \mathbb{R}^{d_x}, x \in \mathcal{X}} |\nabla_x \Psi_x(x)| \leq R$.

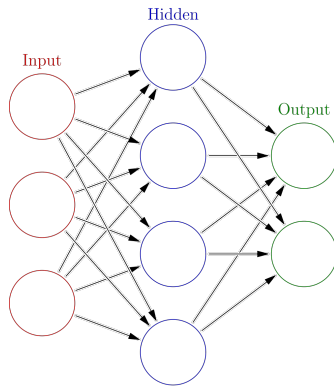
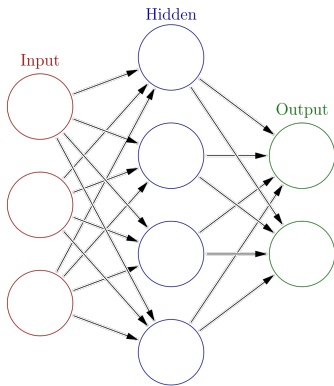
Background: Mean-field perspective of two-layer neural networks

- Define $h_*(\mathbf{x}) = \int \Psi(\mathbf{x}, x) d\mu_*(x)$ the output of the **optimal** mean-field neural network.
- Define $\hat{h}_S(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \Psi(\mathbf{x}, x_S^{(i)})$ the output of a **trained** neural network at time S .
- For any input $\mathbf{x} \in \mathcal{X}$,

$$\mathbb{E} \left[\left(\hat{h}_S(\mathbf{x}) - h_*(\mathbf{x}) \right)^2 \right] \leq \underbrace{\mathcal{O}(N^{-1})}_{\text{finite particle error}} + \underbrace{\mathcal{O} \left(\frac{\gamma^2 + \gamma \sigma d}{\mathcal{C}_{\text{LSI}} \sigma} \right)}_{\text{time discretization error}} + \underbrace{\mathcal{O}(\exp(-\gamma \mathcal{C}_{\text{LSI}} \sigma S))}_{\text{optimization error}}.$$

- Expectation is taken over the randomness in initialization and noise at each iteration.
- $\mathcal{C}_{\text{LSI}} = \Theta(\sigma^{-1} \exp(-\zeta^{-1} \sigma^{-1} \sqrt{d}))$ describes the 'difficulty' of learning μ_* .
 - It reflects the **curse of dimensionality**.

Mean-field neural networks in 2SLS



Mean-field perspective of 2SLS

- Two-stage least squares (2SLS)

$$\begin{aligned}\text{Stage I:} \quad \mathcal{Z}^*(\mathcal{X}) &= \arg \min_{\mathcal{Z} \in (\mathbb{R}^{d_z})^{N_z}} \frac{1}{2} \mathbb{E}_\rho \left[(h(\mathbf{z}, \mathcal{Z}) - h(\mathbf{x}, \mathcal{X}))^2 \right], \\ \text{Stage II:} \quad \mathcal{X}^* &= \arg \min_{\mathcal{X} \in (\mathbb{R}^{d_x})^{N_x}} \frac{1}{2} \mathbb{E}_\rho \left[(h(\mathbf{z}, \mathcal{Z}^*(\mathcal{X})) - y)^2 \right].\end{aligned}\tag{1}$$

- $h(\mathbf{x}, \mathcal{X}) = \frac{1}{N_x} \sum_{i=1}^{N_x} \psi_{\mathbf{x}}(x^{(i)})$ where $\mathcal{X} = [x^{(1)}, \dots, x^{(N_x)}] \in (\mathbb{R}^{d_x})^{N_x}$ are the network parameters and \mathbf{x} is the network input.
- $h(\mathbf{z}, \mathcal{Z}) = \frac{1}{N_z} \sum_{i=1}^{N_z} \psi_{\mathbf{z}}(z^{(i)})$ where $\mathcal{Z} = [z^{(1)}, \dots, z^{(N_z)}] \in (\mathbb{R}^{d_z})^{N_z}$ are the network parameters and \mathbf{z} is the network input.
- ρ is the data distribution over $(\mathbf{x}, \mathbf{z}, y)$.
- A shorthand notation: $\psi_{\mathbf{x}}(x^{(i)}) = \Psi(\mathbf{x}, x^{(i)})$ and $\psi_{\mathbf{z}}(z^{(i)}) = \Psi(\mathbf{z}, z^{(i)})$.

Mean-field perspective of 2SLS

- **Mean field** neural networks $\int_{\mathbb{R}^{d_x}} \Psi_x(x) d\mu_x(x)$ and $\int_{\mathbb{R}^{d_z}} \Psi_z(z) d\mu_z(z)$ where μ_x, μ_z are the mean-field limit of the hidden layer.
- ℓ_2 and entropic regularizations for both stages:

$$\text{Stage I: } \mu_z^*(\mu_x) = \arg \min_{\mu_z \in \mathcal{P}(\mathbb{R}^{d_z})} \frac{1}{2} \mathbb{E}_\rho[(\int \Psi_z d\mu_z - \int \Psi_x d\mu_x)^2] + \frac{\zeta_1}{2} \mathbb{E}_{\mu_z}[\|z\|^2] + \sigma_1 \text{Ent}(\mu_z),$$

$$\text{Stage II: } \mu_x^* = \arg \min_{\mu_x \in \mathcal{P}(\mathbb{R}^{d_x})} \frac{1}{2} \mathbb{E}_\rho[(\int \Psi_z d\mu_z^*(\mu_x) - y)^2] + \frac{\zeta_2}{2} \mathbb{E}_{\mu_x}[\|x\|^2] + \sigma_2 \text{Ent}(\mu_x).$$

(Bi-MFLD)

- A **bilevel** optimization problem over $\mathcal{P}(\mathbb{R}^{d_x})$ and $\mathcal{P}(\mathbb{R}^{d_z})$.
 - Popular methods like explicit gradient (autodiff) and implicit gradient (high-order gradient) do not work.
 - For **fixed** μ_x , Stage I $\mu_z^*(\mu_x)$ can be solved via standard mean field Langevin dynamics.

Mean-field perspective of 2SLS

- Some notation:

$$F_1(\mu_x, \mu_z) = \frac{1}{2} \mathbb{E}_\rho[(\int \Psi_z d\mu_z - \int \Psi_x d\mu_x)^2] + \frac{\zeta_1}{2} \mathbb{E}_{\mu_z}[\|z\|^2]$$

$$F_2(\mu_x, \mu_z) = \frac{1}{2} \mathbb{E}_\rho[(\int \Psi_z d\mu_z - y)^2] + \frac{\zeta_2}{2} \mathbb{E}_{\mu_x}[\|x\|^2].$$

- $\mathcal{F}_1(\mu_x, \mu_z) = F_1(\mu_x, \mu_z) + \sigma_1 \text{Ent}(\mu_z)$ and $\mathcal{F}_2(\mu_x, \mu_z) = F_2(\mu_x, \mu_z) + \sigma_2 \text{Ent}(\mu_x)$.
- Bilevel optimization problem is

Stage I: $\mu_z^*(\mu_x) = \arg \min_{\mu_z \in \mathcal{P}(\mathbb{R}^{d_z})} \mathcal{F}_1(\mu_x, \mu_z)$, Stage II: $\mu_x^* = \arg \min_{\mu_x \in \mathcal{P}(\mathbb{R}^{d_x})} \mathcal{F}_2(\mu_x, \mu_z^*(\mu_x))$.

Observations:

- The **partial** Wasserstein gradients $\mu_x \mapsto F_1(\mu_x, \mu_z)$ and $\mu_x \mapsto F_2(\mu_x, \mu_z)$; $\mu_z \mapsto F_1(\mu_x, \mu_z)$ and $\mu_z \mapsto F_1(\mu_x, \mu_z)$ are simple.
- The **nested** Wasserstein gradient of $\mu_x \mapsto F_2(\mu_x, \mu_z^*(\mu_x))$ is nasty.

Mean-field perspective of 2SLS: Lagrangian formulation

- The bilevel optimization problem

$$\mu_z^*(\mu_x) = \arg \min_{\mu_z \in \mathcal{P}(\mathbb{R}^{d_z})} \mathcal{F}_1(\mu_x, \mu_z), \quad \mu_x^* = \arg \min_{\mu_x \in \mathcal{P}(\mathbb{R}^{d_x})} \mathcal{F}_2(\mu_x, \mu_z^*(\mu_x)). \quad (\text{Bilevel})$$

- A **constrained** optimization problem
 - Stage I problem re-casted as a constraint.

$$\min_{\mu_x, \mu_z} \mathcal{F}_2(\mu_x, \mu_z), \quad \mathcal{F}_1(\mu_x, \mu_z) - \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x)) \leq \varepsilon. \quad (\varepsilon\text{-constrained})$$

- A **Lagrangian** optimization problem

$$(\mu_{x,\lambda}^*, \mu_{z,\lambda}^*) = \arg \min_{\mu_x, \mu_z} \mathcal{L}_\lambda(\mu_x, \mu_z) := \mathcal{F}_2(\mu_x, \mu_z) + \lambda (\mathcal{F}_1(\mu_x, \mu_z) - \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x))). \quad (\lambda\text{-penalty})$$

- When $\lambda = +\infty$, it recovers the bilevel optimization problem.
- When $\lambda < +\infty$, one needs to take into account an additional approximation error.

Mean-field perspective of 2SLS: Lagrangian formulation

- Main challenge:

$$\begin{aligned}(\mu_{x,\lambda}^*, \mu_{z,\lambda}^*) &= \arg \min_{\mu_x, \mu_z} \mathcal{L}_\lambda(\mu_x, \mu_z) \\ &= \arg \min_{\mu_x, \mu_z} \mathcal{F}_2(\mu_x, \mu_z) + \lambda (\mathcal{F}_1(\mu_x, \mu_z) - \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x)))\end{aligned}$$

Proposition 1 (Wasserstein gradient of \mathcal{L}_λ)

Let $\mu_z^*(\mu_x) = \arg \min_{\mu_z} \mathcal{F}_1(\mu_x, \mu_z)$ be the solution to the stage I optimization. Then,

$$\begin{aligned}\nabla_1 \mathcal{L}_\lambda(\mu_x, \mu_z) &= \nabla_1 \mathcal{F}_2(\mu_x, \mu_z) + \lambda \nabla_1 \mathcal{F}_1(\mu_x, \mu_z) - \lambda \nabla_1 \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x)), \\ \nabla_2 \mathcal{L}_\lambda(\mu_x, \mu_z) &= \nabla_2 \mathcal{F}_2(\mu_x, \mu_z) + \lambda \nabla_2 \mathcal{F}_1(\mu_x, \mu_z).\end{aligned}$$

∇_1 (resp. ∇_2) denotes the Wasserstein gradient with the first (resp. second) argument.

- The Wasserstein gradient of the mapping $\mu_x \mapsto \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x))$ only involves the partial derivative with the **first argument** (envelope theorem).
- We avoid the nasty Wasserstein gradient of $\mu_x \mapsto \mathcal{F}_2(\mu_x, \mu_z^*(\mu_x))$.

Mean-field perspective of 2SLS: Lagrangian formulation

- Convexity of $\mathcal{L}_\lambda(\mu_x, \mu_z) = \mathcal{F}_2(\mu_x, \mu_z) + \lambda(\mathcal{F}_1(\mu_x, \mu_z) - \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x)))$.

Observations:

1. The partial mapping $\mu_z \mapsto \mathcal{L}_\lambda(\mu_x, \mu_z)$ is **convex**, for any fixed $\mu_x \in \mathcal{P}_2(\mathbb{R}^{d_x})$.
2. The partial mapping $\mu_x \mapsto \mathcal{L}_\lambda(\mu_x, \mu_z)$ is **not convex**, for any fixed $\mu_z \in \mathcal{P}_2(\mathbb{R}^{d_x})$.
3. The joint mapping $(\mu_x, \mu_z) \mapsto \mathcal{L}_\lambda(\mu_x, \mu_z)$ is **not convex**.

Question: How to exploit this **partial** convexity $\mu_z \mapsto \mathcal{L}_\lambda(\mu_x, \mu_z)$?

- Innerloop: $\mu_z^*(\mu_x) = \arg \min_{\mu_z} \mathcal{F}_1(\mu_x, \mu_z)$, $\tilde{\mu}_z^*(\mu_x) = \arg \min_{\mu_z} \mathcal{L}_\lambda(\mu_x, \mu_z)$.
- Outerloop: $\mu_{x,\lambda}^* = \arg \min_{\mu_x} \mathcal{L}_\lambda(\mu_x, \tilde{\mu}_z^*(\mu_x), \mu_z^*(\mu_x))$.
- Noisy gradient descent!

Mean-field perspective of 2SLS: Lagrangian formulation

- Inner-loop algorithm

$$\mu_z^*(\mu_x) = \arg \min_{\mu_z} \mathcal{F}_1(\mu_x, \mu_z) = \arg \min_{\mu_z} F_1(\mu_x, \mu_z) + \sigma_1 \text{Ent}(\mu_z),$$

$$\tilde{\mu}_z^*(\mu_x) = \arg \min_{\mu_z} \mathcal{L}_\lambda(\mu_x, \mu_z) = \arg \min_{\mu_z} F_2(\mu_x, \mu_z) + \lambda F_1(\mu_x, \mu_z) + \lambda \sigma_1 \text{Ent}(\mu_z).$$

- Fast convergence due to **partial** convexity of $\mu_z \mapsto F_1(\mu_x, \mu_z)$ and $\mu_z \mapsto F_2(\mu_x, \mu_z) + \lambda F_1(\mu_x, \mu_z)$ for fixed μ_x .

Algorithm INNERLOOP($\mu_x, T, \alpha, \beta, \lambda, \sigma_1$)

- 1: Initialize $\mu_{\mathcal{Z},0} = \frac{1}{N_z} \sum_{j=1}^{N_z} \delta_{z_0^{(j)}}$ and $\tilde{\mu}_{\mathcal{Z},0} = \frac{1}{N_z} \sum_{j=1}^{N_z} \delta_{\tilde{z}_0^{(j)}}$.
- 2: **for** $t = 0, \dots, T$ **do**
- 3: **for** $i = 1, \dots, N_z$ **do**
- 4: $z_{t+1}^{(i)} = z_t^{(i)} - \alpha \nabla_2 F_1(\mu_x, \mu_{\mathcal{Z},t})(z_t^{(i)}) + \sqrt{2\alpha\sigma_1} \xi_{z,t}^{(i)}$.
- 5: $\tilde{z}_{t+1}^{(i)} = \tilde{z}_t^{(i)} - \beta \nabla_2 F_2(\mu_x, \tilde{\mu}_{\mathcal{Z},t})(\tilde{z}_t^{(i)}) - \beta \lambda \nabla_2 F_1(\mu_x, \tilde{\mu}_{\mathcal{Z},t})(\tilde{z}_t^{(i)}) + \sqrt{2\beta\lambda\sigma_1} \tilde{\xi}_{z,t}^{(i)}$.
- 6: **end for**
- 7: **end for**

Mean-field perspective of 2SLS: Lagrangian formulation

Assumption 1 (Bounded and smooth neural networks)

There exists a universal positive constant R such that $\sup_{x \in \mathbb{R}^{d_x}, x \in \mathcal{X}} |\Psi_x(x)| \leq R$ and $\sup_{z \in \mathbb{R}^{d_z}, z \in \mathcal{Z}} |\Psi_z(z)| \leq R$. Also, $\sup_{x \in \mathbb{R}^{d_x}, x \in \mathcal{X}} |\nabla_x \Psi_x(x)| \leq R$ and $\sup_{z \in \mathbb{R}^{d_z}, z \in \mathcal{Z}} |\nabla_z \Psi_z(z)| \leq R$.

- It works for two-layer neural networks with tanh/ReLU plus smooth output clipping.

Assumption 2 (Bounded target)

There exists a universal constant M such that the target random variable $|Y| \leq M$ and $|h_o(X)| \leq M$ almost surely.

- Boundedness of $|Y|$ can be relaxed to sub-Gaussian residual $Y - (Th_o)(Z)$.

Mean-field perspective of 2SLS: Lagrangian formulation

Proposition 2 (Inner-loop convergence towards $\mu_z^*(\mu_x)$ and $\tilde{\mu}_z^*(\mu_x)$)

Suppose Assumption 1 and 2 hold. Given a fixed $\mu_x \in \mathcal{P}_2(\mathbb{R}^{d_x})$. Let $\mathcal{Z} = \{z^{(i)}\}_{i=1}^{N_z}$ and $\tilde{\mathcal{Z}} = \{\tilde{z}^{(i)}\}_{i=1}^{N_z}$ be the output of the inner-loop algorithm $\text{INNERLOOP}(\mu_x, T, \alpha, \beta, \lambda, \sigma_1)$. Denote $\mu_z^{(N_z)}$ and $\tilde{\mu}_z^{(N_z)}$ as the joint distribution of these N_z particles \mathcal{Z} . Suppose the step size satisfy $\alpha \leq \frac{1}{\zeta_1}$ and $\beta \leq \frac{1}{\lambda \zeta_2}$. For any $T > 0$,

$$\begin{aligned} \frac{\sigma_1}{N_z} \text{KL} \left(\mu_z^{(N_z)}, (\mu_z^*(\mu_x))^{\otimes N_z} \right) &\leq \frac{R^2}{N_z} + \frac{\alpha^2 + \alpha \sigma_1 d_z}{C_{\text{LSI},z} \sigma_1} + \mathcal{O}(\exp(-C_{\text{LSI},z} \sigma_1 \alpha T)) \\ \frac{\sigma_2}{N_z} \text{KL} \left(\tilde{\mu}_z^{(N_z)}, (\tilde{\mu}_z^*(\mu_x))^{\otimes N_z} \right) &\leq \frac{R^2}{N_z} + \frac{\beta^2 + \beta \sigma_1 d_z}{C_{\text{LSI},z} \sigma_1} + \mathcal{O}(\exp(-C_{\text{LSI},z} \sigma_1 \beta T)). \end{aligned}$$

- $C_{\text{LSI},z} = \Theta(\frac{\zeta_1}{\sigma_1} \exp(-\frac{R^2}{\zeta_1 \sigma_1} \sqrt{d_z/\pi}))$.
- Direct application of mean-field results.

Mean-field perspective of 2SLS: Lagrangian formulation

- Outer-loop algorithm

$$\mu_{x,\lambda}^* = \arg \min_{\mu_x} \mathcal{L}_\lambda(\mu_x, \tilde{\mu}_z^*(\mu_x), \mu_z^*(\mu_x)) = \arg \min_{\mu_x} L_\lambda(\mu_x, \tilde{\mu}_z^*(\mu_x), \mu_z^*(\mu_x)) + \sigma_2 \text{Ent}(\mu_x)$$

- Its Wasserstein gradient $\nabla L_\lambda(\mu_x, \tilde{\mu}_z^*(\mu_x))(x)$ equals (via envelope theorem)

$$\nabla_1 F_2(\mu_x, \tilde{\mu}_z^*(\mu_x))(x) + \lambda(\nabla_1 \mathcal{F}_1(\mu_x, \tilde{\mu}_z^*(\mu_x))(x) - \nabla_1 \mathcal{F}_1(\mu_x, \mu_z^*(\mu_x))(x))$$

Algorithm OUTERLOOP: $F^2\text{BMLD}$ (Fully-first order Bilevel MFLD)

- 1: Initialize $\mu_{\mathcal{X},0} = \frac{1}{N_x} \sum_{j=1}^{N_x} \delta_{x_0^{(j)}}$.
- 2: **for** $s = 0, \dots, S$ **do**
- 3: $\tilde{\mu}_{\mathcal{Z},s}, \mu_{\mathcal{Z},s} \leftarrow \text{INNERLOOP}(\mu_{\mathcal{X},s})$.
- 4: **for** $i = 1, \dots, N_x$ **do**
 $x_{s+1}^{(i)} = x_s^{(i)} - \gamma \left(\nabla_1 F_2(\mu_{\mathcal{X},s}, \tilde{\mu}_{\mathcal{Z},s})(x_s^{(i)}) + \lambda(\nabla_1 \mathcal{F}_1(\mu_{\mathcal{X},s}, \tilde{\mu}_{\mathcal{Z},s})(x_s^{(i)}) \right.$
5: $\left. - \nabla_1 \mathcal{F}_1(\mu_{\mathcal{X},s}, \mu_{\mathcal{Z},s})(x_s^{(i)}) \right) + \sqrt{2\gamma\sigma_2} \xi_{x,s}^{(i)}$.
- 6: **end for**
- 7: **end for**

Algorithm ✓ Theory ?

Mean-field perspective of 2SLS: Lagrangian formulation

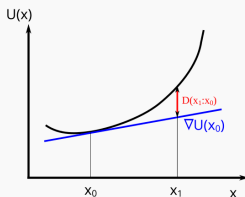
Question: How to prove convergence $\mu_{\mathcal{X},S} = \frac{1}{N_x} \sum_{i=1}^{N_x} \delta_{x_S^{(i)}} \rightarrow \mu_{x,\lambda}^*$?

- The key is convexity!
- $\mu_x \mapsto L_\lambda(\mu_x, \tilde{\mu}_z^*(\mu_x), \mu_z^*(\mu_x))$ is only **weakly** convex.

Lemma (Lower-bound on the Bregman divergence of L_λ)

Suppose Assumption 1 holds. Then, we have $B_{L_\lambda}(\mu_x, \mu'_x) \geq -\frac{R^3\lambda}{4\sigma_1} \text{TV}^2(\mu_x, \mu'_x)$.

- L_λ is more convex as σ_1 increases yet **less convex as λ increases**.



Mean-field perspective of 2SLS: Lagrangian formulation

Theorem 3 (Convergence bound)

Suppose Assumption 1 and 2 hold. Let $\mathfrak{c} > 0$ and assume that $\sigma_1 \sigma_2 \mathfrak{c} \geq 4R^3 \lambda$. Suppose the step size $\gamma \leq \zeta_2^{-1}$. Given a fixed $\lambda > 0$, for any number of iterations $S \in \mathbb{N}^+$, we have

$$\mathcal{H}(S) \lesssim \exp(-\sigma_2 C_{\text{LSI},x} S \gamma) + \frac{\lambda R^2}{\sigma_1 N_x} + \frac{\lambda^2 \left(\sqrt{\frac{\mathfrak{KL}}{N_z}} + \sqrt{\frac{\tilde{\mathfrak{KL}}}{N_z}} \right)}{\sigma_2 C_{\text{LSI},x}} + \frac{\lambda^2 (\gamma^2 + \gamma \sigma_2 d_x)}{\sigma_2 C_{\text{LSI},x}} + \mathfrak{c}^2 C_{\text{LSI},x}.$$

- \mathfrak{KL} and $\tilde{\mathfrak{KL}}$ represents the inner-loop optimization error.
- \mathfrak{c} is a **slack** parameter arising from the weak convexity of L_λ .
- Define $h_{*,\lambda}(\mathbf{x}) = \int \Psi_{\mathbf{x}}(x) d\mu_{\mathbf{x},\lambda}^*(x)$ the **global optimum** mean field network. Define $\hat{h}_S(\mathbf{x}) = \frac{1}{N_x} \sum_{i=1}^{N_x} \Psi_{\mathbf{x}}(x_S^{(i)})$ where $\{x_S^{(i)}\}_{i=1}^{N_x}$ are the output of $F^2\text{BMLD}$.

$$\forall \mathbf{x} \in \mathcal{X}, \quad \mathbb{E} \left[\left(\hat{h}_S(\mathbf{x}) - h_{*,\lambda}(\mathbf{x}) \right)^2 \right] \leq \sqrt{\sigma_2^{-1} \mathcal{H}(S) + \frac{\lambda \mathfrak{c}}{\sigma_1 \sigma_2} + \frac{\lambda}{N_x \sigma_1 \sigma_2}} + \frac{1}{N_x}.$$

- The optimization bound wants **small** λ .

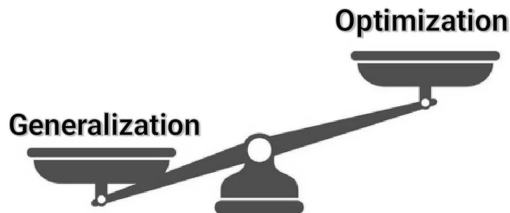
Mean-field perspective of 2SLS: Generalization

Optimization theory:

We have proved that F^2_{BMLD} can indeed find the global optimum solution $h_{*,\lambda}$.

Statistical theory:

How well does $h_{*,\lambda}$ generalize towards h_o when given finite i.i.d samples over $(\mathbf{x}, \mathbf{z}, y)$?



Mean-field perspective of 2SLS: Generalization

- Given m i.i.d samples $\{\mathbf{z}_i, \mathbf{x}_i\}_{i=1}^m \sim P_{ZX}$ in stage I and n i.i.d samples $\{\mathbf{z}_i, y_i\}_{i=1}^n \sim P_{ZY}$ in stage II:

$$\mathcal{F}_1(\mu_x, \mu_z) = \sum_{i=1}^m \frac{1}{2m} \left(\int \Psi(\mathbf{z}_i, z) d\mu_z - \int \Psi(\mathbf{x}_i, x) d\mu_x \right)^2 + \frac{\zeta_1}{2} \mathbb{E}_{\mu_z}[\|z\|^2] + \sigma_1 \text{Ent}(\mu_z),$$

$$\mathcal{F}_2(\mu_x, \mu_z) = \sum_{i=1}^n \frac{1}{2n} \left(\int \Psi(\mathbf{z}_i, z) d\mu_z^*(\mu_x) - y_i \right)^2 + \frac{\zeta_2}{2} \mathbb{E}_{\mu_x}[\|x\|^2] + \sigma_2 \text{Ent}(\mu_x).$$

- Recall that $T : L^2(P_X) \rightarrow L^2(P_Z)$ defined as $T : f \mapsto \mathbb{E}[f(X) \mid Z]$ and NPIV:

$$\mathbb{E}[Y \mid Z] = \mathbb{E}[h_o(X) \mid Z]. \quad (\text{NPIV})$$

Mean-field perspective of 2SLS: Generalization

Assumption 3 (Stage II well-specifiedness)

h_o belongs to a KL restricted Barron space $\mathcal{B}_{M_x} := \{\int \Psi(\cdot, x) d\mu_x(x) \mid \text{KL}(\mu_x, \nu_x) \leq M_x\}$, where $\nu_x = \mathcal{N}(0, \zeta_2 \sigma_2^{-1} \text{Id}_{d_x})$. That is, there exists a measure $\mu_x^\circ \in \mathcal{B}_{M_x}$ such that $h_o(\cdot) = \int \Psi(\cdot, x) d\mu_x^\circ$.

Assumption 4 (Stage I well-specifiedness)

The conditional expectation $T[\int \Psi(\cdot, x) d\mu_x(x)](\mathbf{z}) = \int \mathbb{E}[\Psi(X, x) \mid Z = \mathbf{z}] d\mu_x(x)$ belongs to a KL restricted Barron space $\mathcal{B}_{M_z} := \{\int \Psi(\cdot, z) d\mu_z(z) \mid \text{KL}(\mu_z, \nu_z) \leq M_z\}$, where $\nu_z = \mathcal{N}(0, \zeta_1 \sigma_1^{-1} \text{Id}_{d_z})$. That is, there exists a measure $\mu_z^\circ(\mu_x) \in \mathcal{B}_{M_z}$ such that $T[\int \Psi(\cdot, x) d\mu_x(x)](\mathbf{z}) = \int \Psi(\cdot, z) d\mu_z^\circ(\mu_x)$.

- M_x, M_z are universal constants that control the size of the Barron spaces.

Mean-field perspective of 2SLS: Generalization

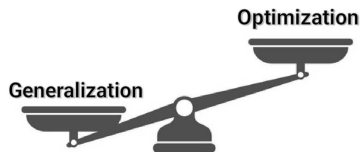
Theorem 4 (Generalization bound)

Suppose Assumption 1,2,3,4 hold. For $\lambda > 0$, let $\mu_{x,\lambda}^*$ be the optimal solution to the Lagrangian problem and $h_{*,\lambda}(\mathbf{x}) = \int \Psi(\mathbf{x}, x) d\mu_{x,\lambda}^*(x)$ be its associated mean field neural network. Then, with $P_{XZY}^{\otimes(m+n)}$ probability at least $1 - 8\delta$,

$$\mathbb{E}_{P_Z} \left[\left((Th_{*,\lambda})(Z) - (Th_o)(Z) \right)^2 \right] \lesssim \sigma_2 M_x + \sigma_1 M_z + \frac{R^2(R+M)^2}{\sigma_1 \lambda} \\ + \sqrt{\frac{M_z + \frac{1}{\sigma_1} + \log(\delta^{-1})}{m}} + \sqrt{\frac{M_x + \frac{1}{\sigma_2} + \log(\delta^{-1})}{n}}.$$

- The generalization bound wants **large λ** so the Lagrangian problem is more faithful to the original bilevel optimization problem.
- $\mathcal{O}(m^{-\frac{1}{2}})$ and $\mathcal{O}(n^{-\frac{1}{2}})$ arise from Rademacher complexity bound.
 - Two-stage regression so we need both $m, n \rightarrow \infty$.

Mean-field perspective of 2SLS: Optimization and Generalization



- **Trade-off** on $\lambda, \sigma_1, \sigma_2$ in terms of optimization and generalization.
- Optimization bound: $\mathbb{E} \left[\left(\hat{h}_S(\mathbf{x}) - h_{*,\lambda}(\mathbf{x}) \right)^2 \right] = \mathcal{O}(\lambda^2 + \sigma_1^{-1} + \sigma_2^{-1})$.
- Generalization bound: $\mathbb{E}_{P_Z} \left[\left((Th_{*,\lambda})(Z) - (Th_{\circ})(Z) \right)^2 \right] = \mathcal{O}(\lambda^{-1} + \sigma_1 + \sigma_2)$.
- Unfortunately, there does not exist a pair of $\lambda, \sigma_1, \sigma_2$ such that both errors vanish.

Experiments: Offline RL on Cartpole

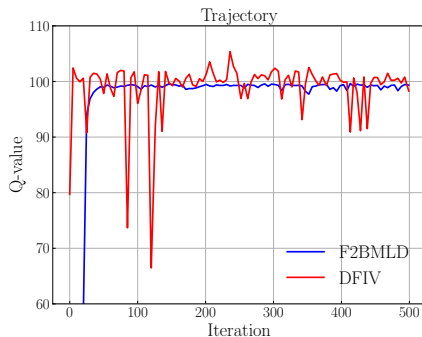
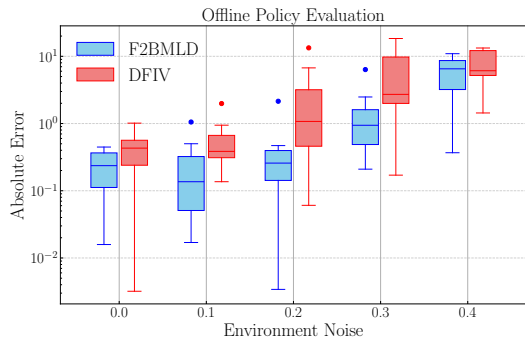


Figure: **Left:** Comparison of DFIV and F2BMLD in terms of target policy value. **Right:** Comparison of DFIV and F2BMLD training trajectories.

- λ is selected from a set $\{0.1, 1.0, 10.0\}$.
- More stable trajectory because of **fully-first order gradient** in optimization.

About Me



- Zonghao Chen
- 4th year PhD Student at University College London (UCL)
 - Foundational AI Centre
 - Gatsby Computational Neuroscience Unit
- Graduated from Tsinghua University in 2022
 - Department of EE
- Kernel (nonparametric) methods, causal inference, statistical learning theory