

STATE-ACTION BALANCING IN CAUSAL INFERENCE

Qiming Du

joint work with Gérard Biau, François Petit and Raphaël Porcher

Laboratoire de Probabilités, Statistique et Modélisation - UMR 8001
Sorbonne Université

(Please visit <https://mgimm.github.io/presentation> for additional information.)

January 13, 2022 Online

What we will talk about?

- Basic concepts of Causal Inference
 - Policy evaluation.
 - Unconfoundedness and Markov structure.
- Static Causal Inference
 - Inverse Probability Weighting (IPW) estimator.
 - G-computation estimator.
 - Doubly robust estimator.
 - Static state-action balancing (our methods).
- Dynamic Causal Inference (our new framework)
 - Measure flows in the change of policy.
 - Dynamical recursive balancing.
 - Error analysis of the policy evaluation estimators.

Notation and conventions

- $\xi(dx)$ (measure), f (test function), $\xi(f) = \int f(x)\xi(dx)$.
- $M(x, dy)$ (transition kernel), $M(f)(\cdot) = \int M(x, dy)f(y)$, and $\mu M(dy) = \int \mu(dx)M(x, dy)$.
- A transition kernel is a family of state-indexed (probability) measures.
- For any random variables X and Y , there exists a Markov transition kernel M that connects their distributions.
- For a probability measure $\xi(dx)$, the empirical measure ξ_n is given by $\frac{1}{n} \sum_{i=1}^n \delta_{X_i}$, where X_i is the i.i.d. sample according to $\xi(dx)$.

Causal Inference (naive):

Statistics

+

Ability to intervene the data generation procedure (action)

Policy: Assignment of actions.

Policy evaluation: Predict the average causal effects of a given policy.

Step 1: Static Causal Inference

without personalized treatment

Conceptual example of policy evaluation: Covid and Vaccination

Question: To vaccinate or not?

Model: Potential Outcomes framework (Rubin, 1974):

- X : Population covariate (age, sex, blood-type, etc).
- $A \in \{0, 1\}$: Action assignment indicator (1:vaccinated vs. 0:not-vaccinated)
- $(Y(0), Y(1))$: Effects associated to vaccination (e.g., infection, side effects, etc).

Goal: Estimate the average potential outcomes $\mathbb{E}[Y(1)]$ (and/or $\mathbb{E}[Y(0)]$).

\iff Policy evaluation (Everyone vaccinates (or not)).

Example: What is the average infection rate (on the whole population) after the vaccination?

\iff Estimation of $\mathbb{E}[Y(1)]$.

Why difficult?

Each individual has only one observation of two causal effects \implies missing value always exists.

Three Classical methods: IPW, G-computation, DRE

Assumptions on the data collection:

Dataset: $\mathcal{D}_n = \{(X^{(i)}, A^{(i)}, Y^{(i)}(A^{(i)})) : 1 \leq i \leq n\}$.

- Randomized study: Assume $A \perp\!\!\!\perp (Y(0), Y(1))$.
- **Observational study (unconfoundedness)**: Assume $A \perp\!\!\!\perp (Y(0), Y(1))$ given X .

1. Inverse Probability Weighting (IPW):

$$\mathbb{E} \left[Y(A) \frac{\mathbf{1}_{\{A=1\}}}{\mathbb{P}(A=1 | X)} \right] = \mathbb{E} \left[\mathbb{E} \left[Y(A) \frac{\mathbf{1}_{\{A=1\}}}{\mathbb{P}(A=1 | X)} \mid X \right] \right] = \mathbb{E} [Y(1)] .$$

Inverse Probability Weighting (IPW) estimator:

$$\text{IPWE} = \frac{1}{n} \sum_{i=1}^n Y^{(i)}(A^{(i)}) \frac{\mathbf{1}_{\{A^{(i)}=1\}}}{\hat{\mathbb{P}}(A=1 | X = X^{(i)})} .$$

The nuisance estimator $\hat{\mathbb{P}}(A=1 | X = X^{(i)})$ is estimated through a separated supervised learning (classification) problem.

2. G-computation:

Denote by $\mu_1(\cdot) = \mathbb{E}[Y(1) \mid X = \cdot]$, then

$$\mathbb{E}[\mu_1(X)] = \mathbb{E}[Y(1)].$$

G-computation estimator

$$\text{GE} = \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X^{(i)}),$$

where the nuisance estimator $\hat{\mu}_1$ can also be obtained by a separate regression:

$$\{X^{(i)} \text{ to } Y^{(i)}(1) : 1 \leq i \leq n \text{ and } A^{(i)} = 1\}$$

Can we combine these two ideas? yes!

3. Doubly robust estimator:

$$\text{DRE} = \frac{1}{n} \sum_{i=1}^n (Y^{(i)}(A^{(i)}) - \hat{\mu}_1(X^{(i)})) \frac{\mathbf{1}_{\{A^{(i)}=1\}}}{\hat{\mathbb{P}}(A=1 \mid X=X^{(i)})} + \frac{1}{n} \sum_{i=1}^n \hat{\mu}_1(X^{(i)}).$$

Why interesting?

- Faster convergence rate (product of the two nuisance estimators).
- Semiparametric efficiency (optimal asymptotic variance over all parametric models) at optimal rate ($\mathcal{O}_{\mathbb{P}}(1/\sqrt{n})$).

Q: What can we improve (i.e., **our contributions**) in this static setting?

A: 1. The IPW part, and 2. more complicated policy evaluation.

Q: Why the IPW part should be improved?

A: Numerical instability when the “inversed” probability is close to 0— a single poorly estimated probability in the whole data set may completely destroy the whole estimator! (This will be illustrated later!)

A closer look at IPW

Q: Why it works?

A: It balances two subpopulations, i.e.,

X (the whole population) and $X(1)$ (people vaccinated),

through re-weighting. The associated weight function is

$$\dot{\eta}(\cdot) = \frac{\mathbb{E}[A]}{\mathbb{P}(A = 1 \mid X = \cdot)}.$$

Re-weighting transformation:

$$\Psi_{\eta}(\mu) : \mu \mapsto \mu(\eta \times \cdot)$$

We have

$$\Psi_{\dot{\eta}}(\xi_1) = \xi \quad \text{and} \quad \frac{d\xi}{d\xi_1} = \dot{\eta}$$

where ξ_1 (resp. ξ) is the probability measure of $X(1)$ (resp. X).

Reformulation of IPWE

We have

$$\begin{aligned}\text{IPWE} &= \frac{1}{n} \sum_{i=1}^n Y^{(i)}(A^{(i)}) \frac{\mathbf{1}_{\{A^{(i)}=1\}}}{\hat{\mathbb{P}}(A=1 \mid X=X^{(i)})} \\ &= \frac{1}{N_1} \sum_{i=1}^{N_1} \hat{\eta}(X(1)^{(i)}) Y^{(i)}(1),\end{aligned}$$

where $\hat{\eta}(X(1)^{(i)}) = \frac{N_1/n}{\hat{\mathbb{P}}(A=1 \mid X=\cdot)}$, which is an empirical version of η .

Natural idea to improve IPW/get rid of the inverse manipulation:

Directly estimate the weight function that corrects the difference between two measures!

Idea (source measure: ξ_1 ; target measure: ξ):

$$\hat{\eta} = \arg \min_{\eta \in H} \text{some-loss-between-measures}(\xi, \Psi_{\eta}(\xi_1)).$$

Are we able to generalize this simple idea to more general cases (e.g., with more complex policy/action assignment)? Yes! And this is our first contribution in the static setting!

Step 2: Static Causal Inference

with personalized treatment

General policy evaluation

Policy is now modeled by a transition kernel $\pi(x, da)$ from state space to action space.

- Sampling policy $\pi(x, da)$: policy that generates the data set \mathcal{D}_n .
- Target policy $\dot{\pi}(x, da)$: The policy to be evaluated.

Case 1: Finite-valued action (e.g., whether to vaccinate, Moderna or Pfizer?).

Example: What is the average infection rate of the whole population if people with age > 60 get vaccinated?

No need to change the framework. For example, one may replace all the $A = 1$ by $A = \dot{A}^{(i)}$ where $\dot{A}^{(i)} \sim \dot{\pi}(X^{(i)}, \cdot)$.

Case 2: General cases (e.g., continuous-valued action when considering dosage, expenses, etc.)

Example: What is the average infection rate of the whole population if people with age > 60 get vaccinated with various dosage (from 50ml to 500ml)?

State-Action Markov reformulation.

Markov structure of causal dynamics

One (tiny) step further from the Markov Decision Process.

- State-action variable: $X^{\mathfrak{h}} = (X, A)$.
- State-action space: $\mathcal{X}^{\mathfrak{h}} = \mathcal{X} \times \mathcal{A}$.
- $\pi^{\mathfrak{h}} = \text{id}_{\mathcal{X}} \times \pi$.
- Goal: estimate $\mathcal{V}^{\pi} = \mathbb{E}[\dot{Y}] = \mathbb{E}[r(\dot{Z})]$.

Causal effects Y (resp. \dot{Y}) is now modeled by $r(Z)$ (resp. $r(\dot{Z})$).

Why? Consistent with the dynamical setting.

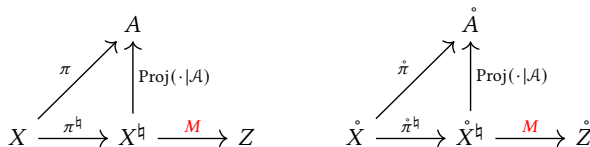


Figure: Markov structure of static causal model

When $\#\mathcal{A} < \infty$, the dynamic is equivalent to the potential outcomes framework with unconfoundedness assumption.

Covariate shifts?

Denote by ξ the population distribution (X) of sampling policy. Denote by $\overset{\circ}{\xi}$ the population distribution ($\overset{\circ}{X}$) of target policy. It is possible that $\xi \neq \overset{\circ}{\xi}$. (Lacking of external validation (Pearl et al., 2014).)

Example: The vaccination data set is collected in the US, how can we calibrate it so that it can be used to conduct causal inference in France?

What are we collecting?

- State-action variables under sampling policy (**US**): $\{(X^{(i)}, A^{(i)}) : 1 \leq i \leq n\}$
- State variables of target policy (**FR**): $\{\overset{\circ}{X}^{(i)} : 1 \leq i \leq m\}$.
- Causal effects under sampling policy (**US**): $\{r(Z^{(i)}) : 1 \leq i \leq n\}$.

What are we balancing (in order to get the weight function estimation)?

The (empirical) **state-action** distribution of X^{\natural} and $\overset{\circ}{X}^{\natural}$, denoted respectively by

re-weighted source measure: $\Psi_{\eta^{\natural}}(\xi_n^{\natural})$ and target measure: $\overset{\circ}{\xi}_m^{\natural}$.

Estimators: DE and DRE

Direct estimator:

$$\hat{V}_{\text{DE}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \hat{\eta}^{\natural}(X^{\natural(i)}) r(Z^{(i)}),$$

Doubly robust estimator:

$$\hat{V}_{\text{DRE}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \hat{\eta}^{\natural}(X^{\natural(i)}) (r(Z^{(i)}) - \hat{r}^{\natural(i)}(X^{\natural(i)})) + \frac{1}{m} \sum_{j=1}^m \pi^{\natural}(\hat{r}^{\natural(j)}) (\dot{X}^{(j)}),$$

where

- $\hat{\eta}^{\natural}$ is the estimated weight function $\dot{\eta}^{\natural}$ that corrects the difference between the two state-action distributions, i.e., $\Psi_{\hat{\eta}^{\natural}}(\xi^{\pi^{\natural}}) = \xi^{\pi^{\natural}}$;
- $\hat{r}^{\natural(j)}$ is the estimated conditional expectation function $\mathbb{E}[r(Z) \mid X^{\natural} = \cdot]$ (or simply $r^{\natural} = M(r)$) by a separated regression.

How to estimate the weight function?

Generalized IPW (through density ratio estimation, see, e.g., [Sugiyama et al. \(2012\)](#)).
or Balancing (our method)!

Teaser of balancing

Direct estimator and worst-case-error interpretation

Idea: $r(Z^{(i)})$ can be regarded as a noisy version of $r^{\natural}(X^{\natural(i)})$. Hence, one considers to minimize the worst-case-error, i.e.,

$$\hat{H} = \arg \min_{\eta^{\natural} \in H} \sup_{\gamma \in \Gamma} \left| \Psi_{\eta^{\natural}}(\xi_n^{\pi^{\natural}})(\gamma) - \xi_{mN}^{\pi^{\natural}}(\gamma) \right|,$$

Characterization of balancing:

- H : family of candidates of the weight function.
- Γ : family of test functions (integrable).

It is well-known that that sup-term is called Integral Probability Metric (IPM).

$$\hat{H} = \arg \min_{\eta^{\natural} \in H} \left(\text{IPM}_{\Gamma} \left(\Psi_{\eta^{\natural}}(\xi_n^{\pi^{\natural}}), \xi_{mN}^{\pi^{\natural}} \right)^2 + \lambda \left\| \eta^{\natural} \right\|_{L^2(\xi_n^{\pi^{\natural}})}^2 \right),$$

(Bias-Variance decomposition)

In this case, H is a $L^2(\xi^{\pi^{\natural}})$ -ball.

Well-specifiedness: $r^{\natural} \in \Gamma$.

Balancing for the DE

Example: Optimal Transport (OT) balancing without L^2 penalty, see, e.g., [Reygner and Touboul \(2020\)](#).

Theorem (Informal)

When well-specified, the error of the DE is controlled by the sampling complexity of the chosen IPM.

Sampling complexity?

rate of convergence of $\text{IPM}_\Gamma(\xi_n, \xi)$ for some empirical measure ξ_n , which is determined by the richness of Γ .

In the OT case, the sampling complexity is in general $n^{-1/d}$ when d , the dimension of the state-action space, is large.

However, there is in general no reason that $\hat{\eta}^\natural$ will converge to the ideal weight function η^\natural in an L^2 sense, which is required by the DRE.

Why L^2 convergence of weight function estimation matters?

Denote by $\mathcal{V}_{\text{ODRE}}^{\pi^{\circ}}$ the oracle version of the DRE, namely, by replacing the nuisance estimators $\hat{\eta}^{\natural}$ and \hat{r}^{\natural} by their oracle/ideal counterparts η° and r° .

Theorem (Informal)

Under mild assumptions, we have

$$\left| \mathcal{V}_{\text{ODRE}}^{\pi^{\circ}} - \hat{\mathcal{V}}_{\text{DRE}}^{\pi^{\circ}} \right| \leq C \left\| \hat{\eta}^{\natural} - \eta^{\circ} \right\|_{L^2(\xi_n^{\pi^{\circ}})} \left\| \hat{r}^{\natural} - r^{\circ} \right\|_{L^2(\xi_n^{\pi^{\circ}})} + o_{\mathbb{P}} \left(\sqrt{\frac{n+m}{nm}} \right).$$

In addition, when no covariate shifts are involved, $\mathcal{V}_{\text{ODRE}}^{\pi^{\circ}}$ achieves semiparametric efficiency.

Error analysis of the DRE:

$$|\text{DRE} - \text{REF}| \leq \underbrace{|\text{DRE} - \text{ODRE}|}_{\text{Theorem above}} + \underbrace{|\text{ODRE} - \text{REF}|}_{\text{converges at parametric/optimal rate}}$$

Riesz representable measure space: weight \iff Riesz representer

- Source measure: ξ .
- Target measure: $\check{\xi}$.
- Riesz representable measure space: $\Xi(\xi) := \{\Psi_\eta(\xi) : \eta \in L^2(\xi)\}$.

Assume $\check{\xi} \in \Xi(\xi)$, then $\Xi(\xi)$ serves as the candidate measure space of all the re-weighted source measure with an L^2 weight.

Proposition

Let ξ be a positive finite measure on \mathcal{X} . Denote by U the unit ball in $L^2(\xi)$. We have the following isometric isomorphism between the metric spaces $L^2(\xi)$ and $\Xi(\xi)$:

$$(L^2(\xi), \|\cdot - \cdot\|_{L^2(\xi)}) \xrightleftharpoons[\frac{d\cdot}{d\xi}]{\Psi_\cdot(\xi)} (\Xi(\xi), \text{IPM}_U(\cdot, \cdot))$$

Heuristic:

On a compact state-action space, what if we use a RELU network to approximate the L^2 unit ball?

One step further to the L^2 convergence of weight estimation

Q: Do we really need an isometry in order to have an L^2 convergence of η when conducting IPM-based optimization?

A: Not really!

Notation: $\partial H = \{\eta - \eta' : \eta, \eta' \in H\}$.

Lemma (Informal)

If $\eta^{\circ} \in H$ and there exists $\alpha > 0$ and $\beta > 0$ such that $\alpha(\partial H) \cap \beta U \subset \Gamma$, then we have

$$\forall \eta^{\natural} \in L^2(\xi^{\pi^{\natural}}), \quad \left\| \eta^{\natural} - \eta^{\circ} \right\|_{L^2(\xi^{\pi^{\natural}})} \leq C_{\alpha, \beta} \text{IPM}_{\Gamma} \left(\Psi_{\eta^{\natural}}(\xi_n^{\pi^{\natural}}), \Psi_{\eta^{\circ}}(\xi_n^{\pi^{\natural}}) \right)$$

The construction is similar to a dual version (in a Fenchel sense) of [Chernozhukov et al. \(2020\)](#).

Take home message:

when Γ is rich enough, the L^2 -error of weight function estimation given by the IPM optimization is controlled by the sampling complexity of the chosen IPM.

Practical considerations

Pipeline:

- 1 Fix H to ensure that $\eta^{\natural} \in H$.
- 2 Construct Γ such that $\alpha(\partial H) \cap \beta U$
- 3 Solve the adversarial optimization, i.e., arg min max-optimization.

Example (of our new method):

- H : Intersection of RKHS ball and L^2 -ball.
- Γ : RKHS ball (which recovers MMD).
- Computation: Quadratic programming, explicitly solvable for small scale problem/
Gradient descend-based method for large scale problem.

DE vs. DRE with our balancing method

DE:

- Requires that $r^{\natural} \in \Gamma$.
- Allows to optimized in $H = L^2(\xi^{\pi^{\natural}})$, i.e., n -value optimization.
- The error of the DE is controlled by the sampling complexity of the chosen IPM.

DRE:

- Requires that $\dot{\eta}^{\natural} \in H$.
- Requires to model properly the candidate space H .
- The L^2 -error of the weight function estimation is controlled by the sampling complexity of the chosen IPM.

Q: When the optimal (parametric) rate is achieved, the DRE is always better than the DE?

A: No, their asymptotic variances are not comparable in general (see, e.g., [Kallus and Uehara \(2020\)](#)).

Step 3: Dynamical Causal Inference

Offline Reinforcement Learning/Dynamical Treatment Regimes

Dynamical Causal Inference under our new framework

Example: What is the average infection rate if we apply 3-vaccination policy with personalized treatment assignment (age<30 Pfizer, age>30 Moderna, with various dosage)?

- Infection may occur after each injection.
- Infection will influence the future vaccination (infected \implies no vaccination).

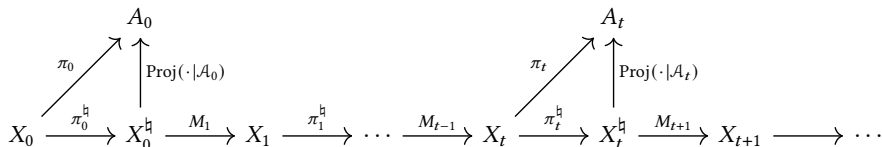


Figure: Markov structure of causal dynamics

Goal: Estimate $\mathcal{V}^{\hat{\pi}} = \mathbb{E} \left[\sum_{t=0}^T r_t(\dot{X}_{t+1}) \right]$.

Models

What is the data set?

- (US) State-Action trajectories under sampling policy $\pi = (\pi_t; 0 \leq t \leq T)$ (a sequence of Markov kernels): $\{(X_t^{(i)}, A_t^{(i)}; 0 \leq t \leq T) : 1 \leq i \leq n\}$.
- (FR) Target initial state covariates: $\{\overset{\circ}{X}_0^{(i)} : 1 \leq i \leq m\}$.
- (US) Observed causal effects at each time step: $\{r_t(X_{t+1}^{(i)}) : 0 \leq t \leq T, 1 \leq i \leq n\}$.

Terminal measures:

- ξ_t^π : State terminal measure of X_t under sampling policy π .
- $\xi_t^{\pi^h}$: State-action terminal measure of X_t^h under sampling policy π .
- $\overset{\circ}{\xi}_t^\pi$: State terminal measure of $\overset{\circ}{X}_t$ under target policy $\overset{\circ}{\pi}$.
- $\overset{\circ}{\xi}_t^{\pi^h}$: State-action terminal measure of $\overset{\circ}{X}_t^h$ under target policy $\overset{\circ}{\pi}$.

The objective re-writes

$$\mathcal{V}^{\overset{\circ}{\pi}} = \sum_{t=0}^T \overset{\circ}{\xi}_{t+1}^{\overset{\circ}{\pi}}(r_t) = \sum_{t=0}^T \overset{\circ}{\xi}_t^{\overset{\circ}{\pi}^h}(r_t^h).$$

What happens if we change the policy?

Measure flows in the change of policy

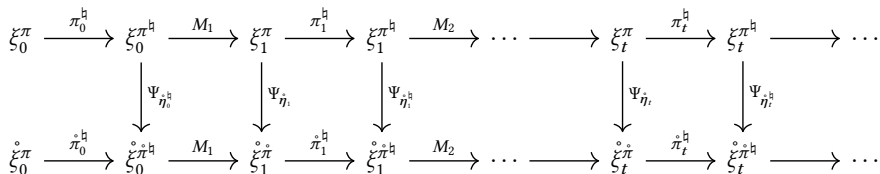


Figure: Measure flows in the change of policy

What are these weight functions (when exist)?

$$\forall 1 \leq t \leq T, \forall x_t^h = (x_t, a_t) \in \mathcal{X}_t^h, \quad \dot{e}_t^h(x_t^h) = \frac{d\overset{\circ}{\pi}_t(x_t, \cdot)}{d\pi_t(x_t, \cdot)}(a_t).$$

For $t = 0$, we let, taking into account the covariate shifts,

$$\forall x_0^h = (x_0, a_0) \in \mathcal{X}_0^h, \quad \dot{e}_0^h(x_0^h) = \frac{d\overset{\circ}{\xi}\pi_0}{d\xi\pi_0}(x_0) \frac{d\overset{\circ}{\pi}_0(x_0, \cdot)}{d\pi_0(x_0, \cdot)}(a_0).$$

Measure flows in the change of policy

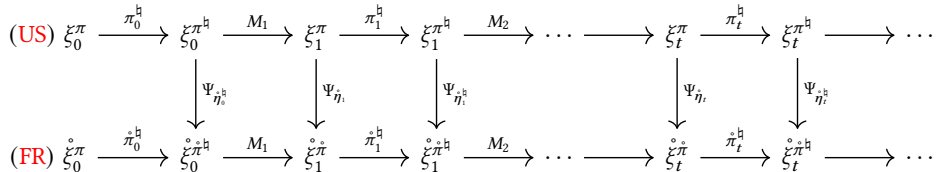


Figure: Measure flows in the change of policy

Proposition

Under mild assumptions, the weight functions $\dot{\eta}_t$ and $\dot{\eta}_t^h$ are well-defined respectively in $L^1(\mathcal{X}_t)$ and $L^1(\mathcal{X}_t^h)$. In addition, for any $0 \leq t \leq T$, we have

$$\dot{\eta}_t^h(\cdot) = \mathbb{E} \left[\prod_{s=0}^t \dot{e}_s^h(X_s^h) \mid X_t^h = \cdot \right] \quad \text{and} \quad \dot{\eta}_{t+1}(\cdot) = \mathbb{E} \left[\dot{\eta}_t^h(X_t^h) \mid X_{t+1} = \cdot \right].$$

So, possible to implement balancing?

In a smart way, yes!

New recursive balancing strategy

- Initial balancing: Compare $\xi_n^{\pi^h}$ with $\xi_{0,mN}^{\pi^h}$ to get the estimation of $\{\hat{\eta}_0^h(X_0^{(i)}) : 1 \leq i \leq n\}$.
- Weight smoothing: After getting $\{\hat{\eta}_t^h(X_t^{(i)}) : 1 \leq i \leq n\}$ for $t \geq 0$, we run a separate regression (or do nothing, i.e., let $\hat{\eta}_{t+1}(X_{t+1}^{(i)}) = \hat{\eta}_t^h(X_t^{(i)})$) to get $\{\hat{\eta}_{t+1}(X_{t+1}^{(i)}) : 1 \leq i \leq n\}$.
- Update balancing: Once we have $\{\hat{\eta}_{t+1}(X_{t+1}^{(i)}) : 1 \leq i \leq n\}$, we compare $\xi_{t+1,n}^{\pi^h}$ with $\xi_{t+1,nN}^{\pi^h}$ to get the estimation $\{\hat{\eta}_{t+1}^h(X_{t+1}^{(i)}) : 1 \leq i \leq n\}$.

PRO:

- Compatible with general Polish action space.
- No need to worry about covariate shifts (built-in solution).

CON:

- Can be computationally heavy in some setting.

Now that we have estimated the weight functions, what about the actual estimators?

Our results on the error analysis of the DE and the DRE

Direct estimator:

$$\hat{\mathcal{V}}_{\text{DE}}^{\pi^{\circ}} = \sum_{t=0}^T \Psi_{\hat{\eta}_{t+1}}(\xi_{t+1,n}^{\pi})(r_t) = \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^T \hat{\eta}_{t+1}^{(i)}(X_{t+1}^{(i)}) r_t(X_{t+1}^{(i)}).$$

Doubly robust estimator:

$$\hat{\mathcal{V}}_{\text{DRE}}^{\pi^{\circ}} = \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^T \left(\hat{\eta}_t^{\natural}(X_t^{\natural(i)}) \left(r_t(X_{t+1}^{(i)}) - \hat{r}_t^{\natural(i)}(X_t^{\natural(i)}) \right) + \hat{\eta}_{t-1}^{\natural}(X_{t-1}^{\natural(i)}) \pi_t^{\circ \natural}(\hat{r}_t^{\natural(i)})(X_t^{(i)}) \right),$$

Sampling complexity of IPM:

$$\forall t \geq 1, \quad \sigma_t^{\text{IPM}}(n) = \text{IPM}_{\Gamma_t} \left(\Psi_{\hat{\eta}_t^{\natural}}(\xi_{t,n}^{\pi^{\natural}}), \Psi_{\hat{\eta}_t^{\natural}}(\xi_t^{\pi^{\natural}}) \right).$$

With a slight abuse of notation, we omit m and N at time 0, i.e.,

$$\sigma_0^{\text{IPM}}(n) = \sigma_0^{\text{IPM}}(n, m, N) = \text{IPM}_{\Gamma_0} \left(\Psi_{\hat{\eta}_0^{\natural}}(\xi_{0,n}^{\pi^{\natural}}), \Psi_{\hat{\eta}_0^{\natural}}(\xi_0^{\pi^{\natural}}) \right) + \text{IPM}_{\Gamma_0} \left(\xi_0^{\pi^{\circ \natural}}, \xi_{0,mN}^{\pi^{\circ \natural}} \right).$$

Direct estimator

Well-specifiedness (dynamical version):

We say that the causal dynamics are well-specified by a sequence of collections of test functions $\Gamma = (\Gamma_t; 0 \leq t \leq T)$ if

$$\forall 0 \leq t \leq T, \quad r_t^{\mathfrak{h}} = M_{t+1}(r_t) \in \Gamma_t,$$

and

$$\forall 1 \leq t \leq T, \quad \forall \gamma_t \in \Gamma_t, \quad M_t^{\pi^{\mathfrak{h}}}(\gamma_t) \in \Gamma_{t-1}.$$

Theorem (Informal)

Under mild assumptions, if the causal dynamics are well-specified by $\Gamma = (\Gamma_t; 0 \leq t \leq T)$, then we have

$$\left| \hat{\mathcal{V}}_{\text{DE}}^{\pi} - \mathcal{V}^{\pi} \right| \leq C \left(\sum_{t=0}^T (T - t + 1) \left(\sigma_t^{\text{IPM}}(n) + \frac{1}{\sqrt{N}} \right) \right),$$

where $C > 0$ is a constant that is independent to n, m, N and T .

Doubly robust estimator

For the error term given by weight smoothing, we denote

$$\sigma_t^{\text{WS}}(n) = \|\hat{\eta}_{t+1} - \overset{\circ}{\eta}_{t+1}\|_{L^2(\xi_{t+1,n}^\pi)}.$$

we denote the error of the additional regression of the average reward function r^{\natural} by

$$\forall 0 \leq t \leq T, \quad \sigma_t^{\text{REG}}(n) = \|\hat{r}_t^{\natural} - r_t^{\natural}\|_{L^2(\xi_{t,n}^{\pi^{\natural}})}.$$

Theorem (Informal)

If the weight functions are well specified in H_t and if the implemented balancing satisfies

$$\forall 0 \leq t \leq T, \quad \exists \alpha_t, \beta_t > 0, \quad \alpha_t(\partial H_t) \cap \beta_t U_t \subset \Gamma_t,$$

then we have

$$\left| \hat{\mathcal{V}}_{\text{DRE}}^{\pi} - \mathcal{V}^{\pi} \right| \leq C \left(\sum_{t=0}^T (T-t+1) \left(\sigma_t^{\text{IPM}}(n) + \frac{1}{\sqrt{N}} + \sigma_t^{\text{WS}}(n) \right) \sigma_t^{\text{REG}}(n) \right)$$

where $C > 0$ is a constant that is independent from n, m, N , and T .

Conclusion

Our contributions:

- 1 New state-action Markov reformulation of causal dynamics, that is capable of dealing with general action spaces and covariate shifts.
- 2 New theoretical framework for balancing method through Riesz representable measure space arguments (connection between existing methods and inspiration of new methods).
- 3 Recursive balancing strategy, with transparent error analysis for the DE and the DRE.

More details:

State-Action Balancing in Multi-Stage Causal Inference. Q. Du, G. Biau, F. Petit, R. Porcher. (preprint, 2022).

Perspectives

- 1 Balancing extracts features?
- 2 Connection with the Feynman-Kac formalism/Sequential Monte Carlo
 - How to efficiently interact with the environment with the collected offline data?
 - *Variance Estimation in Adaptive Sequential Monte Carlo*. Q. Du, A. Guyader. Annals of Applied Probability, 2021.
 - *Asymmetric Sequential Monte Carlo*. Q. Du. Under review, 2021.
 - What if we estimate directly the Markov kernel M_t ?
 - *Wasserstein Random Forests and Applications in Heterogeneous Treatment Effects*. Q. Du, G. Biau, F. Petit, R. Porcher. AISTATS 2021.
- 3 Transductive Transfer learning?

$$\arg \min_{f \in F} \text{IPM}_{\Gamma} \left(\frac{1}{n} \sum_{i=1}^n \delta_{X_i} Y_i, \frac{1}{n} \sum_{i=1}^n \delta_{X_i} f(X_i) \right)$$

Thank you for your attention!

References

- Chernozhukov, V., Newey, W., Singh, R., and Syrgkanis, V. (2020). Adversarial estimation of riesz representers. *arXiv e-prints*.
- Kallus, N. and Uehara, M. (2020). Double reinforcement learning for efficient off-policy evaluation in markov decision processes. *Journal of Machine Learning Research*, 21(167):1–63.
- Pearl, J., Bareinboim, E., et al. (2014). External validity: From do-calculus to transportability across populations. *Statistical Science*, 29(4):579–595.
- Reygner, J. and Touboul, A. (2020). Reweighting samples under covariate shift using a wasserstein distance criterion.
- Rubin, D. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66:688–701.
- Sugiyama, M., Suzuki, T., and Kanamori, T. (2012). *Density ratio estimation in machine learning*. Cambridge University Press.