Blessing from Human-Al Interaction: Super Reinforcement Learning in Confounded Environments

Jiayi Wang

Joint work with Zhengling Qi and Chengchun Shi

02/03/2023

1

• An urgent care example



• An urgent care example



Basic Setup

- A random tuple $(S, U, A, \{R(a)\}_{a \in A})$
 - * $\{R(a)\}_{a \in A}$ denotes a set of the potential/counterfactual rewards under A = a.
 - * $R = \sum_{a \in \mathcal{A}} R(a) \mathbb{1}(A = a).$

• A random tuple $(S, U, A, \{R(a)\}_{a \in A})$

* $\{R(a)\}_{a \in A}$ denotes a set of the potential/counterfactual rewards under A = a.

*
$$R = \sum_{a \in \mathcal{A}} R(a) \mathbb{1}(A = a).$$

- Offline dataset
 - * i.i.d copies of (S, A, R): $\{S_i, A_i, R_i\}_{i=1}^n$.
 - * A is generated by some behavior policy $\pi^b : S \times U \to \mathcal{P}(A)$ that depends on both observed and unobserved features.

• A random tuple $(S, U, A, \{R(a)\}_{a \in \mathcal{A}})$

* $\{R(a)\}_{a \in A}$ denotes a set of the potential/counterfactual rewards under A = a.

*
$$R = \sum_{a \in \mathcal{A}} R(a) \mathbb{1}(A = a).$$

- Offline dataset
 - * i.i.d copies of (S, A, R): $\{S_i, A_i, R_i\}_{i=1}^n$.
 - * A is generated by some behavior policy $\pi^b : S \times U \to \mathcal{P}(A)$ that depends on both observed and unobserved features.
- Target: find the optimal policy π^* such that the value function

$$\mathbb{E}\left\{\sum_{\boldsymbol{a}\in\mathcal{A}}R(\boldsymbol{a})\pi^*(\boldsymbol{a}\mid\boldsymbol{S},\boldsymbol{U})\right\}$$

is maximized, using the offline dataset.

Common Policy

• Since U is unobserved, most existing solutions focus on finding an optimal policy $\pi^* \in \Pi = \{\pi : S \to \mathcal{P}(\mathcal{A})\}$ given by

$$\pi^*(a^*|s) = 1 ext{ if } a^* = rgmax_{a \in \mathcal{A}} \mathbb{E}[R(a)|S=s], orall s \in S.$$



Common Policy

• Since U is unobserved, most existing solutions focus on finding an optimal policy $\pi^* \in \Pi = \{\pi : S \to \mathcal{P}(\mathcal{A})\}$ given by

$$\pi^*(a^*|s) = 1 ext{ if } a^* = rgmax_{a \in \mathcal{A}} \mathbb{E}[R(a)|S=s], orall s \in S.$$



• Can we improve upon this?

Super Policy

• Leverages the input of human expertise, since actions generated by the behavior policy depend on the latent information $(\pi_b : S \times U \to \mathcal{P}(\mathcal{A}))$.

Super Policy

- Leverages the input of human expertise, since actions generated by the behavior policy depend on the latent information $(\pi_b : S \times U \to \mathcal{P}(\mathcal{A}))$.
- The urgent care example





Super Policy

- Leverages the input of human expertise, since actions generated by the behavior policy depend on the latent information $(\pi_b : S \times U \to \mathcal{P}(\mathcal{A}))$.
- The urgent care example





Super Policy

- Leverages the input of human expertise, since actions generated by the behavior policy depend on the latent information $(\pi_b : S \times U \to \mathcal{P}(\mathcal{A}))$.
- The urgent care example



Super Policy

• Aim to find a super-policy ν^* in a larger policy class $\Omega = \{\nu : S \times A \to \mathcal{P}(A)\}$ such that

$$u^*(a^*|s,a') = 1 ext{ if } a^* = rg\max_{a \in \mathcal{A}} \mathbb{E}\{R(a)|S=s,A=a'\}, orall(s,a') \in S imes A.$$



- $S \sim Bernoulli(p = 0.5)$ and $U \sim Bernoulli(p = 0.5)$ are independent.
- Action is binary and the behavior policy satisfies $P(A = 1|S, U = 1) = P(A = 0|S, U = 0) = 1 \epsilon$ for some $0 \le \epsilon \le 1$.
- R = 8(A 0.5)(S 0.2)(U 0.3).

A Toy Example

- $S \sim Bernoulli(p = 0.5)$ and $U \sim Bernoulli(p = 0.5)$ are independent.
- Action is binary and the behavior policy satisfies $P(A = 1|S, U = 1) = P(A = 0|S, U = 0) = 1 \epsilon$ for some $0 \le \epsilon \le 1$.

•
$$R = 8(A - 0.5)(S - 0.2)(U - 0.3).$$

Policy Value	$\mathcal{V}(\pi_b)$	$\mathcal{V}(\pi^*)$	$\mathcal{V}(u^*)$
$\epsilon = 0.5$	0.0	0.4	0.4
$\epsilon = 0$	0.6	0.4	1.0
$\epsilon = 1$	-0.6	0.4	1.0

- $S \sim Bernoulli(p = 0.5)$ and $U \sim Bernoulli(p = 0.5)$ are independent.
- Action is binary and the behavior policy satisfies $P(A = 1|S, U = 1) = P(A = 0|S, U = 0) = 1 \epsilon$ for some $0 \le \epsilon \le 1$.

•
$$R = 8(A - 0.5)(S - 0.2)(U - 0.3).$$

Policy Value	$\mathcal{V}(\pi_b)$	$\mathcal{V}(\pi^*)$	$\mathcal{V}(u^*)$
$\epsilon = 0.5$	0.0	0.4	0.4
$\epsilon = 0$	0.6	0.4	1.0
$\epsilon = 1$	-0.6	0.4	1.0

Theorem (Super-Optimality)

 $\mathcal{V}(\nu^*) \geq \max{\{\mathcal{V}(\pi^*), \mathcal{V}(\pi_b)\}}.$

Identification

• Adopt the proximal causal inference framework developed by Tchetgen et al. [2020].

Specifically, we assume the existence of certain action and reward proxies Z and W in additional to (S, A, R).

• Adopt the proximal causal inference framework developed by Tchetgen et al. [2020].

Specifically, we assume the existence of certain action and reward proxies Z and W in additional to (S, A, R).

Assumption

- (a) $R \perp Z \mid (S, U, A);$
- (b) $W \perp (Z, A) \mid (S, U), W \not\perp U \mid S;$
- (c) $R(a) \perp A \mid (S, U)$ for $a \in A$;

(d) There exists a bridge function $q: \mathcal{W} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ such that

$$\mathbb{E}\left[q(W,a,S) \mid U,S,A=a\right] = \mathbb{E}\left[R \mid U,S,A=a\right].$$
(1)

• Adopt the proximal causal inference framework developed by Tchetgen et al. [2020].

Specifically, we assume the existence of certain action and reward proxies Z and W in additional to (S, A, R).

Assumption

- (a) $R \perp Z \mid (S, U, A);$
- (b) $W \perp\!\!\!\perp (Z, A) \mid (S, U), W \not\!\!\perp U \mid S;$
- (c) $R(a) \perp A \mid (S, U)$ for $a \in A$;

(d) There exists a bridge function $q: \mathcal{W} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ such that

$$\mathbb{E}\left[q(W,a,S) \mid U,S,A=a\right] = \mathbb{E}\left[R \mid U,S,A=a\right].$$
(1)

•
$$\mathcal{V}(\nu) = \mathbb{E}\left[\sum_{a \in \mathcal{A}} q(W, a, S) \nu(a \mid S, A)\right]$$
 for any $\nu \in \Omega$.

• In practice, one may want to include as many confounders in the policy as possible to achieve the largest super-optimality.

We further extend the policy class to $\Omega = \{\nu : S \times Z \times A \to \mathcal{P}(A)\}$ and consider the corresponding super-policy ν^* .

• In practice, one may want to include as many confounders in the policy as possible to achieve the largest super-optimality.

We further extend the policy class to $\Omega = \{\nu : S \times Z \times A \to \mathcal{P}(A)\}$ and consider the corresponding super-policy ν^* .

Theorem

Under above Assumptions, if we further suppose some completeness and regularity conditions, solving the following linear integral equation

$$\mathbb{E}\left[q(W,a,S) \mid Z,S,A=a\right] = \mathbb{E}\left[R \mid Z,S,A=a\right],\tag{2}$$

for every $a \in A$ with respect to q gives a valid bridge function that satisfies (1). And the optimal policy ν^* in class Ω is given by

$$\nu^*(a^* \mid s, z, a') = 1 \quad if \quad a^* = \arg\max_{a \in \mathcal{A}} \mathbb{E}\left[q(W, a, S) \mid S = s, Z = z, A = a'\right].$$
(3)

Algorithm 1: Learning Algorithm for the contextual bandits under unmeasured confounding

- 1 Input: Data $\mathcal{D} = (S_i, Z_i, A_i, R_i, W_i)_{i=1}^n$.
- 2 Obtain the estimation of the bridge function \hat{q} by solving the estimation equation (2) using data \mathcal{D}
- 3 Implement any supervised learning method for estimating $\mathbb{E}[\hat{q}(W, S, a) \mid S, Z, A]$. 4 Compute

 $a^* = \arg \max_{a \in \mathcal{A}} \hat{\mathbb{E}} \left[\hat{q}(W, S, a) \mid S = s, Z = z, A = a' \right] \quad \forall (s, z, a') \in \mathcal{S} \times \mathcal{Z} \times \mathcal{A}.$ 5 **Output:** $\hat{\nu}^*$ with $\hat{\nu}(a^* \mid s, z, a') = 1$ and $\hat{\nu}(\tilde{a} \mid s, z, a') = 0$ for $\tilde{a} \neq a^*$.

Identification

Algorithm

Algorithm 2: Learning Algorithm for the contextual bandits under unmeasured confounding

- 1 Input: Data $\mathcal{D} = (S_i, Z_i, A_i, R_i, W_i)_{i=1}^n$.
- 2 Obtain the estimation of the bridge function \hat{q} by solving the estimation equation (2) using data \mathcal{D}
- 3 Implement any supervised learning method for estimating $\mathbb{E}[\hat{q}(W, S, a) | S, Z, A]$.
- 4 Compute

 $a^* = \arg \max_{a \in \mathcal{A}} \hat{\mathbb{E}} \left[\hat{q}(W, S, a) \mid S = s, Z = z, A = a' \right] \quad \forall (s, z, a') \in \mathcal{S} \times \mathcal{Z} \times \mathcal{A}.$ 5 **Output:** $\hat{\nu}^*$ with $\hat{\nu}(a^* \mid s, z, a') = 1$ and $\hat{\nu}(\tilde{a} \mid s, z, a') = 0$ for $\tilde{a} \neq a^*$.

• Step 2: minimax estmation of conditional moment models.

$$q = \underset{q \in \mathcal{Q}}{\operatorname{argmin}} \sup_{g \in \mathcal{G}} \mathbb{E} \left\{ \left[(q(W, A, S) - R) \right] g(S, Z, A) - \frac{1}{2} g^2(S, Z, A) \right\}$$

Algorithm

Algorithm 3: Learning Algorithm for the contextual bandits under unmeasured confounding

- 1 Input: Data $\mathcal{D} = (S_i, Z_i, A_i, R_i, W_i)_{i=1}^n$.
- 2 Obtain the estimation of the bridge function \hat{q} by solving the estimation equation (2) using data \mathcal{D}
- 3 Implement any supervised learning method for estimating $\mathbb{E}[\hat{q}(W, S, a) \mid S, Z, A]$. 4 Compute

 $\begin{aligned} \mathbf{a}^* &= \arg \max_{\mathbf{a} \in \mathcal{A}} \hat{\mathbb{E}} \left[\hat{q}(W, S, \mathbf{a}) \mid S = s, Z = z, A = \mathbf{a}' \right] \quad \forall (s, z, \mathbf{a}') \in \mathcal{S} \times \mathcal{Z} \times \mathcal{A}. \\ \mathbf{5} \text{ Output: } \hat{\nu}^* \text{ with } \hat{\nu}(\mathbf{a}^* \mid s, z, \mathbf{a}') = 1 \text{ and } \hat{\nu}(\tilde{\mathbf{a}} \mid s, z, \mathbf{a}') = 0 \text{ for } \tilde{\mathbf{a}} \neq \mathbf{a}^*. \end{aligned}$

• Step 3: take $\hat{\mathbb{E}}\left[\hat{q}(W,S,a) \mid S=\cdot, Z=\cdot, A=\cdot\right]$ as the solution of

$$\underset{g \in \mathcal{G}}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^{n} \left[g(S_i, Z_i, A_i) - \hat{q}(\cdot, \cdot, a) \right]^2 + \mu \|g\|_{\mathcal{G}}^2,$$

Regret Bound

Theorem

Suppose q belongs to certain function class $Q \subset W \times S \times A$. Define the projection error as $\xi_n := \sup_{q \in Q, a \in A} \|g[\cdot, \cdot, \cdot; q(\cdot, \cdot, a)] - \hat{g}[\cdot, \cdot, \cdot; q(\cdot, \cdot, a)]\|_2$, and the bridge function estimation error as $\zeta_n := \|q - \hat{q}\|_2$. Then we obtain the following regret decomposition

$$\mathcal{V}(\nu^*) - \mathcal{V}(\hat{\nu}^*) \leq 2(\xi_n + p_{\max}\zeta_n),$$

where p_{max} is some overlap constant.

Regret Bound

Theorem

Suppose q belongs to certain function class $Q \subset W \times S \times A$. Define the projection error as $\xi_n := \sup_{q \in Q, a \in A} \|g[\cdot, \cdot, \cdot; q(\cdot, \cdot, a)] - \hat{g}[\cdot, \cdot, \cdot; q(\cdot, \cdot, a)]\|_2$, and the bridge function estimation error as $\zeta_n := \|q - \hat{q}\|_2$. Then we obtain the following regret decomposition

$$\mathcal{V}(\nu^*) - \mathcal{V}(\hat{\nu}^*) \leq 2(\xi_n + p_{\max}\zeta_n),$$

where p_{max} is some overlap constant.

Corollary

If the star-shaped spaces \mathcal{G} and \mathcal{Q} are VC-subgraph classes with VC dimensions $\mathbb{V}(\mathcal{G})$, and $\mathbb{V}(\mathcal{Q})$ respectively. Under certain technical assumptions, with probability at least $1 - \delta$,

$$\mathcal{V}(\hat{\nu}^*) - \mathcal{V}(\nu^*) \lesssim n^{-1/2} p_{\mathsf{max}} \sqrt{\log(1/\delta) + \max{\{\mathbb{V}(\mathcal{G}), \mathbb{V}(\mathcal{Q})\}}}.$$

Application

- A dataset from Study to Understand Prognoses and Preferences for Outcomes and Risks of Treatments [SUPPORT Connors et al., 1996].
 SUPPORT examined the effectiveness and safety of direct measurement of cardiac function by Right Heart Catheterization (RHC) for certain critically ill patients in intensive care units (ICU).
- Action: $A = 1 \rightarrow$ measured by RHC in the first 24 hours, otherwise A = 0. Response: $Y = 1 \rightarrow$ survived or censored at day 30; otherwise Y = -1.

• Covariates: demographics, estimated probability of survival, comorbidity, vital signs, physiological status, and functional status.

- Covariates: demographics, estimated probability of survival, comorbidity, vital signs, physiological status, and functional status.
- Ten variables measuring the patient's overall physiological status:
 - * subject to substantial measurement error
 - * single snapshot of underlying physiological state over time.

- Covariates: demographics, estimated probability of survival, comorbidity, vital signs, physiological status, and functional status.
- Ten variables measuring the patient's overall physiological status:
 - * subject to substantial measurement error
 - * single snapshot of underlying physiological state over time.
- Following Tchetgen et al. [2020]: W = (ph1, hema1), Z = (pafi1, paco21).

- Three types of policy classes are considered.
 - 1. Sonly: $S \to \mathcal{P}(\mathcal{A})$. The policy only depends on the observed state S.
 - 2. SZONLY: $S \times Z \to \mathcal{P}(A)$. The policy depends on on the observed state S and the action proxy Z.
 - 3. SUPER: $S \times Z \times A \to \mathcal{P}(A)$. The super-policy class where the policy depends on the observed state S, the action proxy Z, and observed action A.

• Three types of policy classes are considered.

1. Sonly: $S \to \mathcal{P}(\mathcal{A})$. The policy only depends on the observed state S.

2. SZONLY: $S \times Z \to \mathcal{P}(A)$. The policy depends on on the observed state S and the action proxy Z.

3. SUPER: $S \times Z \times A \to \mathcal{P}(A)$. The super-policy class where the policy depends on the observed state *S*, the action proxy *Z*, and observed action *A*.

Table: Evaluation results of the optimal policies learned from three different policy classes using the RHC data. The averages of evaluation values over 20 random splits are presented. Larger values indicate better performances. Values in the parentheses are standard deviations.

Sonly	SZonly	Super
-------	--------	-------

0.55 (5.80e-02) 0.55 (5.78e-02) 0.69 (1.10e-02)

Confounded Sequantial Decision Making

Confounded Sequantial Decision Making

Basic Setup

- Confounded POMDP $\mathcal{M} = (\mathcal{S}, \mathcal{U}, \mathcal{A}, T, \mathcal{P}, r).$
 - * $\mathcal{P} = \{\mathcal{P}_t\}_{t=1}^T$, \mathcal{P}_t denotes transition kernel from $\mathcal{S} \times \mathcal{U} \times \mathcal{A}$ to $\mathcal{S} \times \mathcal{U}$ at time t

Assumption

(Markovianity) The process $\{S_t, U_t, A_t, R_t\}_{t=1}^T$ satisfies the Markov property, i.e., for any t, (R_t, S_{t+1}, U_{t+1}) depends on the past history only through (S_t, U_t, A_t) .

Confounded Sequantial Decision Making

Basic Setup

- Confounded POMDP $\mathcal{M} = (\mathcal{S}, \mathcal{U}, \mathcal{A}, T, \mathcal{P}, r).$
 - * $\mathcal{P} = \{\mathcal{P}_t\}_{t=1}^T$, \mathcal{P}_t denotes transition kernel from $\mathcal{S} \times \mathcal{U} \times \mathcal{A}$ to $\mathcal{S} \times \mathcal{U}$ at time t

Assumption

(Markovianity) The process $\{S_t, U_t, A_t, R_t\}_{t=1}^T$ satisfies the Markov property, i.e., for any t, (R_t, S_{t+1}, U_{t+1}) depends on the past history only through (S_t, U_t, A_t) .



- Offline dataset:
 - * i.i.d. episodes $\{S_{i,t}, A_{i,t}, R_{i,t}\}_{t=1}^{T}$, i = 1, ..., n.
 - * A_t is generated by some behavior policy $\pi_t^b : S \times U \to \mathcal{P}(A)$.

- Offline dataset:
 - * i.i.d. episodes $\{S_{i,t}, A_{i,t}, R_{i,t}\}_{t=1}^{T}$, i = 1, ..., n.
 - * A_t is generated by some behavior policy $\pi_t^b: S \times U \to \mathcal{P}(A)$.
- Value function given a generic policy $\{\pi_t\}_{t=1}^T$:

$$V^{\pi}_t(s,u) = \mathbb{E}^{\pi}[\sum_{t'=t}^T R_{t'} \mid S_t = s, U_t = u]$$

- Offline dataset:
 - * i.i.d. episodes $\{S_{i,t}, A_{i,t}, R_{i,t}\}_{t=1}^{T}$, i = 1, ..., n.
 - * A_t is generated by some behavior policy $\pi_t^b: S \times U \to \mathcal{P}(A)$.
- Value function given a generic policy $\{\pi_t\}_{t=1}^T$:

$$V^{\pi}_t(s,u) = \mathbb{E}^{\pi}[\sum_{t'=t}^T R_{t'} \mid S_t = s, U_t = u]$$

• Target: Estimate an optimal policy that maximizes

$$\mathcal{V}(\pi) = \mathbb{E}[V_1^{\pi}(S_1, U_1)]$$

using the batch data.

- Offline dataset:
 - * i.i.d. episodes $\{S_{i,t}, A_{i,t}, R_{i,t}\}_{t=1}^{T}$, i = 1, ..., n.
 - * A_t is generated by some behavior policy $\pi_t^b: S \times U \to \mathcal{P}(A)$.
- Value function given a generic policy $\{\pi_t\}_{t=1}^T$:

$$V^{\pi}_t(s,u) = \mathbb{E}^{\pi}[\sum_{t'=t}^T R_{t'} \mid S_t = s, U_t = u]$$

• Target: Estimate an optimal policy that maximizes

$$\mathcal{V}(\pi) = \mathbb{E}[V_1^{\pi}(S_1, U_1)]$$

using the batch data.

• Common optimal policy $\pi^* \in \Pi \equiv \{\pi = \{\pi_t\}_{t=1}^T \mid \pi_t : S \times \mathcal{Z}_t \to \mathbb{P}(\mathcal{A})\}$

Super policy $\nu^* \in \Omega \equiv \{\nu = \{\nu_t\}_{t=1}^T \mid \nu_t : S \times \mathcal{Z}_t \times \mathcal{A} \to \mathcal{P}(\mathcal{A})\}$

Identification

• We assume the existence of certain action and reward proxies $\{Z_t\}_{t=1}^T$ and $\{W_t\}_{t=1}^T$ that can help identify policy values.

Identification

Proxy Variables

• We assume the existence of certain action and reward proxies $\{Z_t\}_{t=1}^T$ and $\{W_t\}_{t=1}^T$ that can help identify policy values.

Assumption

- (a) (Reward proxy) $W_t \perp (A_t, U_{t-1}, S_{t-1}) \mid (U_t, S_t)$, $W_t \perp U_t \mid S_t$, for $1 \leq t \leq T$.
- (b) (Action proxy) $Z_t \perp (R_t, W_t, S_{t+1}, U_{t+1}, W_{t+1}) \mid (U_t, S_t, A_t)$ for $1 \le t \le T$.

Identification

Proxy Variables

• We assume the existence of certain action and reward proxies $\{Z_t\}_{t=1}^T$ and $\{W_t\}_{t=1}^T$ that can help identify policy values.

Assumption

(a) (Reward proxy) $W_t \perp (A_t, U_{t-1}, S_{t-1}) \mid (U_t, S_t)$, $W_t \not\perp U_t \mid S_t$, for $1 \leq t \leq T$.

(b) (Action proxy) $Z_t \perp (R_t, W_t, S_{t+1}, U_{t+1}, W_{t+1}) \mid (U_t, S_t, A_t)$ for $1 \le t \le T$.



Theorem

Under appropriate assumptions, there always exist Q-bridge functions $\{q^{\nu}\}_{t=1}^{T}$ satisfying

$$\mathbb{E}^{\nu}\left[\sum_{t'=t}^{T} R_{t'} \mid U_t, S_t, A_t\right] = \mathbb{E}\left[\sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a)\nu_t(a \mid S_t, Z_t, A_t) \mid U_t, S_t, A_t\right],$$
(4)

for t = 1, ..., T. In particular, set $q_{T+1}^{\nu} = 0$, q_t^{ν} can be obtained by solving the following linear integral equations for t = T, ..., 1,

$$\mathbb{E}\{q_t^{\nu}(W_t, S_t, A_t) - R_t - V_{t+1}^{\nu}(W_{t+1}, S_{t+1}, Z_{t+1}, A_{t+1}) \mid Z_t, S_t, A_t\} = 0, \quad (5)$$

where $V_t^{\nu}(W_t, S_t, Z_t, A_t) = \sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a) \nu_t(a \mid S_t, Z_t, A_t).$

Identification

Linear Integral Equations

Theorem

Under appropriate assumptions, there always exist Q-bridge functions $\{q^{\nu}\}_{t=1}^{T}$ satisfying

$$\mathbb{E}^{\nu}\left[\sum_{t'=t}^{T} R_{t'} \mid U_t, S_t, A_t\right] = \mathbb{E}\left[\sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a)\nu_t(a \mid S_t, Z_t, A_t) \mid U_t, S_t, A_t\right],\tag{4}$$

for t = 1, ..., T. In particular, set $q_{T+1}^{\nu} = 0$, q_t^{ν} can be obtained by solving the following linear integral equations for t = T, ..., 1,

$$\mathbb{E}\{q_t^{\nu}(W_t, S_t, A_t) - R_t - V_{t+1}^{\nu}(W_{t+1}, S_{t+1}, Z_{t+1}, A_{t+1}) \mid Z_t, S_t, A_t\} = 0, \quad (5)$$

where $V_t^{\nu}(W_t, S_t, Z_t, A_t) = \sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a) \nu_t(a \mid S_t, Z_t, A_t).$

$$\mathcal{V}(
u) = \mathbb{E}\left[\sum_{\pmb{a}\in\mathcal{A}} q_1^
u(W_1,S_1,\pmb{a})
u_1(\pmb{a}\mid S_1,Z_1,A_1)
ight], \quad orall
u\in\Omega.$$

Identification

Linear Integral Equations

Theorem

Under appropriate assumptions, there always exist Q-bridge functions $\{q^{\nu}\}_{t=1}^{T}$ satisfying

$$\mathbb{E}^{\nu}\left[\sum_{t'=t}^{T} R_{t'} \mid U_t, S_t, A_t\right] = \mathbb{E}\left[\sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a)\nu_t(a \mid S_t, Z_t, A_t) \mid U_t, S_t, A_t\right],\tag{4}$$

for t = 1, ..., T. In particular, set $q_{T+1}^{\nu} = 0$, q_t^{ν} can be obtained by solving the following linear integral equations for t = T, ..., 1,

$$\mathbb{E}\{q_t^{\nu}(W_t, S_t, A_t) - R_t - V_{t+1}^{\nu}(W_{t+1}, S_{t+1}, Z_{t+1}, A_{t+1}) \mid Z_t, S_t, A_t\} = 0, \quad (5)$$

where $V_t^{\nu}(W_t, S_t, Z_t, A_t) = \sum_{a \in \mathcal{A}} q_t^{\nu}(W_t, S_t, a) \nu_t(a \mid S_t, Z_t, A_t).$

$$u^* = rgmax_{
u \in \Omega} \mathbb{E}\left[\sum_{a \in \mathcal{A}} q_1^{
u}(W_1, S_1, a)
u_1(a \mid S_1, Z_1, A_1)
ight]$$

.

A Practical Algorithm

Assumption (Memoryless Unmeasured Confounding)

For $2 \le t \le T$, U_t is independent of past data history (including latent factors in the past) up to time t - 1 given S_t .

Assumption (Memoryless Unmeasured Confounding)

For $2 \le t \le T$, U_t is independent of past data history (including latent factors in the past) up to time t - 1 given S_t .



Algorithm 4: Super RL for the confounded POMDP

- 1 Input: Data $\mathcal{D} = \{\mathcal{D}_t\}_{t=1}^T$ with
 - $\mathcal{D}_t = \{(S_{i,t}, Z_{i,t}, A_{i,t}, R_{i,t}, W_{i,t}, S_{i,t+1}, Z_{i,t+1}, W_{i,t+1})\}_{i=1}^n.$
- 2 Let $\hat{q}_{\mathcal{T}+1} = 0$ and $\hat{\nu}^*_{\mathcal{T}}$ be an arbitrary policy.

3 Repeat for
$$t = T, \ldots, 1$$
:

- 4 Obtain an estimator \hat{q}_t for q_t by solving (5) using data \mathcal{D}_t and \hat{q}_{t+1} obtained from the last iteration.

Theorem

Suppose $q_t \in Q^{(t)}$ for $1 \le t \le T$ and $\hat{\nu}^*$ is computed via Algorithm 4. Then under above Assumptions and appropriate technical conditions, we obtain the following regret decomposition,

$$\mathcal{V}(\nu^*) - \mathcal{V}(\hat{
u}^*) \lesssim \left(\sum_{t=1}^T 2p_{t,\max}\xi_{t,n}
ight) + \sqrt{T\sum_{t=1}^T (p_{t,\max}^\omega)^2 (\zeta_{t,n})^2}.$$

Theorem

Suppose $q_t \in Q^{(t)}$ for $1 \le t \le T$ and $\hat{\nu}^*$ is computed via Algorithm 4. Then under above Assumptions and appropriate technical conditions, we obtain the following regret decomposition,

$$\mathcal{V}(\nu^*) - \mathcal{V}(\hat{
u}^*) \lesssim \left(\sum_{t=1}^T 2p_{t,\max}\xi_{t,n}\right) + \sqrt{T\sum_{t=1}^T (p_{t,\max}^\omega)^2 (\zeta_{t,n})^2}.$$

- $p_{t,\max}$, $p_{t,\max}^{\omega}$: overlap constants.
- $\zeta_{t,n}$: *Q*-bridge function estimation error.
- $\xi_{t,n}$: projection estimation error.

Simulation and Application

- Existing setting from Miao et al. [2022].
- n = 1000 and T = 20.

Table: Simulation results for the sequential decision making problem. The simulation is performed over 50 simulated datasets. Mean regret values for estimated optimal policies under different policy classes are provided. The smaller regret values indicate better performances. Values in the parentheses are the standard deviations of the regret values.

Sonly	SZonly	Super
F = A = (1 - 0 - 01)		22 (4.0.01)

5.4 (1.9e-01) 5.3 (4.7e-01) **2.2** (4.9e-01)

 Multi-parameter Intelligent Monitoring in Intensive Care (MIMIC-III) dataset (https://physionet.org/content/mimiciii/1.4/) records the longitudinal information (including information of demographics, vitals, labs and scores) of patients who satisfied the sepsis criteria, and the goal is to learn an optimal personalized treatment strategy for sepsis.

- Multi-parameter Intelligent Monitoring in Intensive Care (MIMIC-III) dataset (https://physionet.org/content/mimiciii/1.4/) records the longitudinal information (including information of demographics, vitals, labs and scores) of patients who satisfied the sepsis criteria, and the goal is to learn an optimal personalized treatment strategy for sepsis.
- Despite the richness of data collected at the ICU, the mapping between true patient states and clinical observations is usually ambiguous [Nanayakkara et al., 2022], and therefore makes this dataset fit into the setting of a confounded POMDP.

- Multi-parameter Intelligent Monitoring in Intensive Care (MIMIC-III) dataset (https://physionet.org/content/mimiciii/1.4/) records the longitudinal information (including information of demographics, vitals, labs and scores) of patients who satisfied the sepsis criteria, and the goal is to learn an optimal personalized treatment strategy for sepsis.
- Despite the richness of data collected at the ICU, the mapping between true patient states and clinical observations is usually ambiguous [Nanayakkara et al., 2022], and therefore makes this dataset fit into the setting of a confounded POMDP.
- Simplify the action space into 4-dimensional. Fix the horizon at T = 2.

Table: Evaluation results of the optimal policies learned from three different policy classes using the MIMIC-III data. The averages of evaluation values over 20 random splits are presented. Larger values indicate better performances. Values in the parentheses are standard deviations.

Sonly	SZonly	Super	
-2.83 (5.30e-02)	-2.81 (5.03e-02)	- 1.75 (1.14e-02)	

Conclusion

- We introduce super reinforcement learning, which takes the observed action in the offline data as input to enhanced policy learning under endogeneity.
- We establish the identification results for the super-policy in various confounded environments.
- Practical algorithms are provided to perform the super-policy learning with finite-sample regret guarantees.
- Our super policy can be used in "human-in-the-loop" and "machine-in-the-loop".

Machine-in-the-loop Human



Machine-in-the-loop Human



Machine-in-the-loop Human



Thank You!

- Alfred F Connors, Theodore Speroff, Neal V Dawson, Charles Thomas, Frank E Harrell, Douglas Wagner, Norman Desbiens, Lee Goldman, Albert W Wu, Robert M Califf, et al. The effectiveness of right heart catheterization in the initial care of critically iii patients. *Jama*, 276(11):889–897, 1996.
- Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. Human decisions and machine predictions. *The quarterly journal of economics*, 133(1):237–293, 2018.
- Rui Miao, Zhengling Qi, and Xiaoke Zhang. Off-policy evaluation for episodic partially observable markov decision processes under non-parametric models, 2022. URL https://arxiv.org/abs/2209.10064.
- Thesath Nanayakkara, Gilles Clermont, Christopher James Langmead, and David Swigon. Unifying cardiovascular modelling with deep reinforcement learning for uncertainty aware control of sepsis treatment. *PLOS Digital Health*, 1(2):e0000012, 2022.

- Eric J Tchetgen Tchetgen, Andrew Ying, Yifan Cui, Xu Shi, and Wang Miao. An introduction to proximal causal learning. *arXiv preprint arXiv:2009.10982*, 2020.
- Guy Tennenholtz, Uri Shalit, and Shie Mannor. Off-policy evaluation in partially observable environments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10276–10283, 2020.
- Jiayi Wang, Zhengling Qi, and Chengchun Shi. Blessing from experts: Super reinforcement learning in confounded environments. *arXiv preprint arXiv:2209.15448*, 2022.