

# Estimating Long-term Rewards by Off-policy Reinforcement Learning

---

Lihong Li

Amazon

December 15<sup>th</sup>, 2020

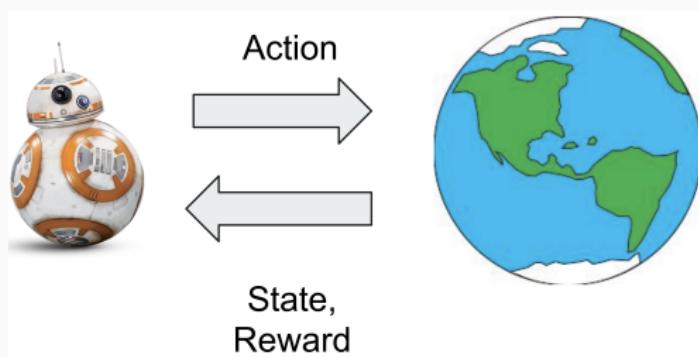


AI for Economics Seminar

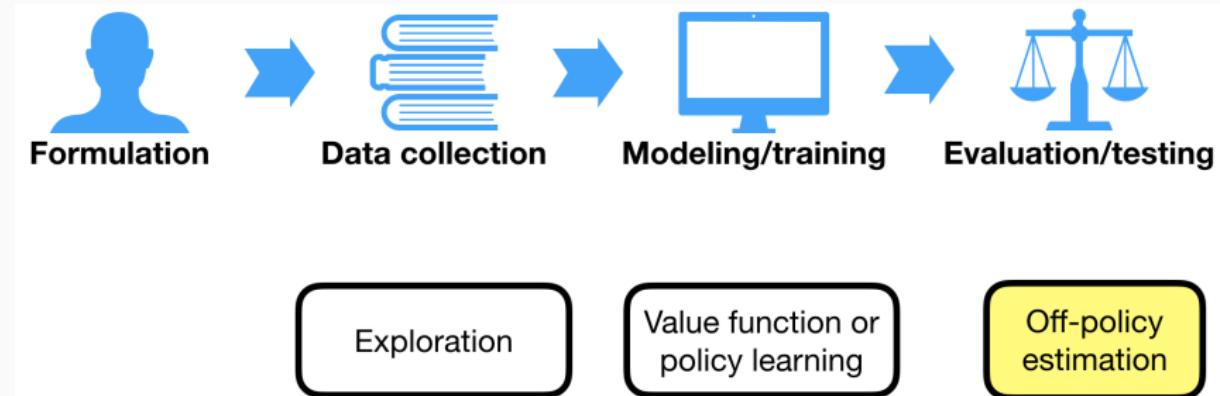
Thanks to Andrew Bennett, Yinlam Chow, Bo Dai, Nathan Kallus, Qiang Liu, Ali Mousavi, Ofir Nachum, Dale Schuurmans, Csaba Szepesvári

# RL: A General Learning Paradigm

- Games
- Robotics and control
- Recommendation
- Inventory management
- Medical treatment
- Education systems
- Economic decisions (tax policy, job training, etc.)
- ...



# A Simplified View of RL



Informally, off-policy estimation (OPE) asks the question:  
*How good is a given policy **without** actually running it?*

# OPE Example: Personalized News Recommendation



([www.yahoo.com](http://www.yahoo.com) on 2010/10/18)

Naturally an RL problem:

- state: user information
- action: news article
- reward:  $\{0, 1\}$  (click or not)
- policy: state  $\rightarrow$  action

A typical evaluation dilemma in practice:

- policy deployment requires performance evaluation
- performance evaluation requires policy deployment

**OPE:** Can we evaluate performance without policy deployment?

Fundamentally a counterfactual estimation problem.

# OPE with Simulators

- Games: (almost) perfect simulation
- Robotics: physical dynamics
- Recommendation: user click model
- Auto. driving: all traffic conditions
- Dialogue: human conversation model
- Healthcare: patient's medical model
- Economy: consumer behavior model
- ...



Reliable simulators are often hard to build.

This talk: a complementary approach without building a simulator.

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

# Markov Decision Process

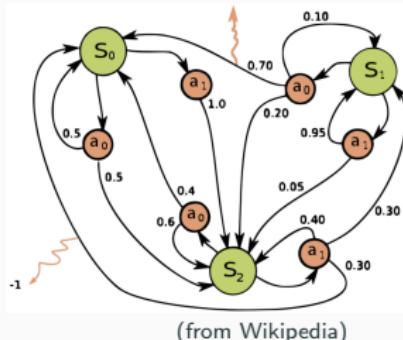
$$\text{MDP } M = \langle \mathcal{S}, \mathcal{A}, T, R, \mu_0, \gamma \rangle$$

- $\mathcal{S}$ : finite set of states
- $\mathcal{A}$ : finite set of actions
- $T : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ : transition kernel
- $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ : reward function
- $\mu_0 \in \Delta(\mathcal{S})$ : initial state distribution
- $\gamma \in [0, 1)$ : discount factor

See papers for treatments of infinite state-actions, and of the  $\gamma = 1$  case.

## Definitions

- Trajectory  $\tau = (s_0, a_0, r_0, s_1, a_1, r_1, s_2, \dots)$ , with  $s_0 \sim \mu_0$  controlled by policy  $\pi$ :  $\forall t, a_t \sim \pi(\cdot | s_t), s_{t+1} \sim T(\cdot | s_t, a_t)$
- Return:  $U(\tau) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t r_t$
- Policy value:  $v(\pi) := \mathbb{E}_{\tau \sim \pi}[U(\tau)]$  ← subject of study in this talk



## Estimating $v(\pi)$

### On-policy (Monte Carlo) estimation

- Run  $\pi$  directly to generate trajectories
- Estimate  $v(\pi)$  (by averaging rewards over trajectories)

E.g., randomized clinical trials, A/B testing, ...

Correct, and straightforward

Often costly, risky, or even impossible to run  $\pi$

Examples: healthcare, autonomous driving, ...

---

### Off-policy estimation

- Assume historical data  $\bar{\mathcal{D}}$  generated by behavior policy  $\pi_0$
- Estimate  $v(\pi)$  from  $\bar{\mathcal{D}}$

Cheap and safe

Biased: simple average over  $\bar{\mathcal{D}}$  estimates  $v(\pi_0)$ , not  $v(\pi)$

# Connection to Causal Effect Estimation

## A Special Case

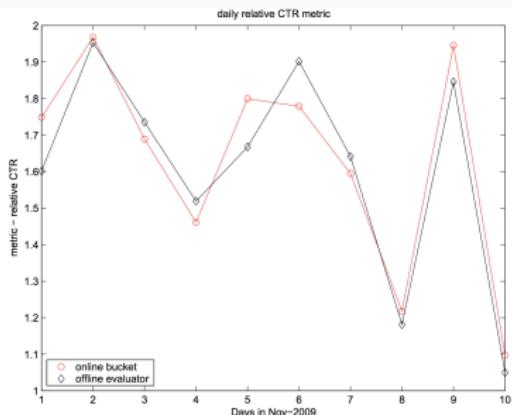
- $\mathcal{S}$  is set of patients
- $\mathcal{A} = \{\text{NewDrug}, \text{Placebo}\}$
- Reward is  $\{0, 1\}$  (cured or not)
- Discount  $\gamma = 0$
- $\pi_{\text{treatment}}(s) \equiv \text{NewDrug}$ ,  
 $\pi_{\text{control}}(s) \equiv \text{Placebo}$
- Want to estimate  $v(\pi_{\text{treatment}}) - v(\pi_{\text{control}})$

Same as the **average treatment effect**.

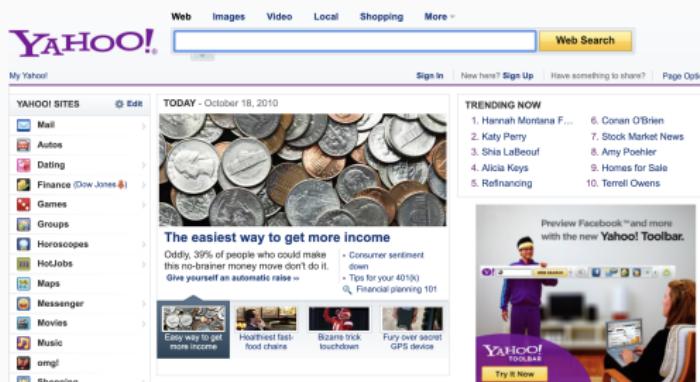
Also closely related to **covariate shift** in machine learning.

# Personalized News Recommendation

- $\mathcal{S}$  is all users
- $\mathcal{A}$  is news articles
- Reward is  $\{0, 1\}$  (click or not)
- $\gamma = 0$  (i.e., contextual bandits)
- $v(\pi)$  is click-through rate



[LCLW'11]



Survey on more applications on the Web [HLR'16]

# Job Training Programs

Impact of job training on long-term employment rates [Athey+'19]

- $\mathcal{S}$ : certain workforce population
- $\mathcal{A}$ : participation in training or not
- Reward: employment rate 9 years from now

Like a contextual bandit with long-delayed reward.

Policy value  $v(\pi)$  measures “long-term rewards” :

- accumulated over a long period of time
- allows repeated exposure to treatment
- more general

## Family Functioning

Fast Track prevention program [Murphy'03]

- Goal: prevent/reduce conduct disorder of children
- $S$ : child's behavioral condition at the end of semester
- $A$ : number of home visits scheduled in the next semester
- Reward: academic achievement in 10 semesters
- $\gamma = 0$

Each step  $t \in \{0, 1, \dots, 9\}$  corresponds to 1 semester.

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

# Precise Formulation

## Generation of data

$$\begin{aligned}\bar{\mathcal{D}} &= \left\{ \tau^{(i)} \right\}_{1 \leq i \leq m} \\ \tau^{(i)} &= \left( s_0^{(i)}, a_0^{(i)}, r_0^{(i)}, \dots, s_{H-1}^{(i)}, a_{H-1}^{(i)}, r_{H-1}^{(i)} \right),\end{aligned}$$

for some large enough **horizon  $H$**  and **behavior policy  $\pi_0$** ,  
where  $a_t \sim \pi_0(\cdot|s_t)$ ,  $r_t = R(s_t, a_t)$ ,  $s_{t+1} \sim T(\cdot|s_t, a_t)$ .

## OPE: Off-policy estimation

Estimate  $\hat{v}(\bar{\mathcal{D}}, \pi) \approx v(\pi)$  **without** knowing  $\{T, R\}$ .

When  $\pi_0$  is unknown, “**behavior-agnostic OPE**”.

# Bellman Equation

## Value functions

Generalizing policy value:

$$Q_\pi(s, a) = \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, a_0 = a \right]$$

$$\text{So } v(\pi) = (1 - \gamma) \mathbb{E}_{s_0 \sim \mu} [Q_\pi(s_0, \pi(a_0))].$$

## Bellman equation

$$Q_\pi(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim T(\cdot|s, a), a' \sim \pi(\cdot|s)} [Q_\pi(s', a')]$$

Compute  $Q_\pi$  from known  $(T, R)$ : dynamic programming, ...

# Regression Estimator

## A plug-in approach

- Use  $\bar{\mathcal{D}}$  to estimate  $(T, R)$  of MDP by  $(\hat{T}, \hat{R})$
- Compute  $\hat{Q}_\pi$  from  $(\hat{T}, \hat{R})$ , and estimate  $v(\pi)$  by

$$\begin{aligned}\hat{v}_{\text{Reg}} &= (1 - \gamma) \hat{\mathbb{E}}_{s_0} [\hat{Q}_\pi(s_0, \pi(s_0))] \\ &= \frac{1 - \gamma}{m} \sum_{\tau \in \bar{\mathcal{D}}} \hat{Q}_\pi(s_0, \pi(s_0))\end{aligned}$$

**Pros:** simple and low variance.

**Cons:** uncontrolled bias due to model misspecification of  $(\hat{T}, \hat{R})$ .

# Inverse Propensity Score (IPS)

Key idea: change of measure

$$\begin{aligned} v(\pi) &= \mathbb{E}_{\tau \sim \pi}[U(\tau)] \\ &= \mathbb{E}_{\tau \sim \pi_0} \left[ \frac{\Pr_{\pi}(\tau)}{\Pr_{\pi_0}(\tau)} U(\tau) \right] \\ &= \mathbb{E}_{\tau \sim \pi_0} \left[ \frac{\prod_{t=0}^{H-1} T(s_t | s_{t-1}, a_{t-1}) \pi(a_t | s_t)}{\prod_{t=0}^{H-1} T(s_t | s_{t-1}, a_{t-1}) \pi_0(a_t | s_t)} U(\tau) \right] \\ &\approx \frac{1}{|\bar{\mathcal{D}}|} \sum_{\tau \in \bar{\mathcal{D}}} \left[ \underbrace{\prod_{t=0}^{H-1} \frac{\pi(a_t | s_t)}{\pi_0(a_t | s_t)}}_{\text{"propensity score"}} U(\tau) \right] =: \hat{v}_{\text{IPS}}. \end{aligned}$$

**Pros:** unbiased.

**Cons:** high variance (more discussion soon).

## Many Variants of IPS/DR

- Per-decision IPS [Precup+'01]
- Self-normalized IPS [Precup+'01, Swaminathan+'15]
- Doubly robust estimators [JL'16, Thomas+'16]
- More robust doubly robust estimators [Farajtabar+'18]
- Use of temporal abstraction to reduce variance [Guo+'17]
- ...

Unfortunately, not widely used in practice unless  $\gamma \ll 1$ .

## Curse of Horizon

- For  $\tau = (s_0, a_0, r_0, \dots, s_{H-1}, a_{H-1}, r_{H-1})$ ,

$$\frac{\Pr_\pi(\tau)}{\Pr_{\pi_0}(\tau)} = \prod_{t=0}^{H-1} \frac{\pi(a_t|s_t)}{\pi_0(a_t|s_t)}$$

whose variance easily blows up (**exponentially** in horizon  $H$ ).

- Exponential blow-up is **unavoidable** in *minimax* rate [LMS'15]
- Similar for doubly robust estimator [JL'16]
- Can truncate ratios when they grow too large, at the cost of introducing large bias (e.g. ReTrace [Munos+'16])

Can we break the **Curse of Horizon** (at least in “benign” cases)?

Many recent developments [LLTZ'18, NCDL'19, DNCLSS'20, ...].

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

## Behavior-agnostic OPE

Now consider a more general setting where data consists of  $n$  transition tuples:  $\mathcal{D} = \{(s_i, a_i, r_i, s'_i)\}_{1 \leq i \leq n}$

Assumptions:

- $(s_i, a_i) \sim d_{\mathcal{D}}$  (unknown behavior distribution)
- $\mathbb{E}[r_i] = R(s_i, a_i)$
- $s'_i \sim T(s_i, a_i)$

Otherwise, no assumption on  $\pi_0$  (other than regularity conditions).  
Data do not even have to be in the form of trajectories.

# Linear Program Characterizations of Policy Value

Two equivalent LPs for the policy value [Puterman'94]:

## Primal LP

$$\begin{aligned} v(\pi) = \min_{Q: S \times A \rightarrow \mathbb{R}} \quad & (1 - \gamma) \mathbb{E}_{(s_0, a_0) \sim \mu_0 \pi} [Q(s_0, a_0)] \\ \text{s.t.} \quad & Q \geq R + \gamma \mathcal{P}_\pi Q \end{aligned}$$

## Dual LP

$$\begin{aligned} v(\pi) = \max_{d \in \Delta(S \times A)} \quad & \mathbb{E}_{(s, a) \sim d} [R(s, a)] \\ \text{s.t.} \quad & d = (1 - \gamma) \mu_0 \pi + \gamma \mathcal{P}_\pi^* d \end{aligned}$$

where

$$\mathcal{P}_\pi Q(s, a) := \sum_{s', a'} T(s'|s, a) \pi(a'|s') Q(s', a')$$

$$\mathcal{P}_\pi^* d(s, a) := \sum_{\bar{s}, \bar{a}} \pi(a|s) T(s|\bar{s}, \bar{a}) d(\bar{s}, \bar{a})$$

## OPE based on Dual LP

$$\begin{aligned} v(\pi) = & \max_{d \in \Delta(\mathcal{S} \times \mathcal{A})} \mathbb{E}_{(s,a) \sim d}[R(s,a)] \\ \text{s.t.} & \quad d = (1 - \gamma)\mu_0\pi + \gamma\mathcal{P}_\pi^*d \quad (*) \end{aligned}$$

### Observations

- There is a unique solution  $d = d_\pi$  that satisfies (\*).  
Known as stationary distribution, steady-state distribution, ...
- A dual approach to obtaining policy value:

$$v(\pi) = \mathbb{E}_{d_\pi}[R(s,a)] = \mathbb{E}_{d_D}[\zeta_\pi(s,a)R(s,a)]$$

where  $\zeta_\pi(s,a) = \frac{d_\pi(s,a)}{d_D(s,a)}$  is Radon-Nikodym derivative (RND)  
(under absolute continuity condition  $d_\pi \ll d_D$ )

- OPE is now reduced to estimating  $\zeta(s,a)$

# DualDICE Approach to OPE

## Observation

$\zeta_\pi$  is the solution to a linear system:

$$\zeta d_{\mathcal{D}} = (1 - \gamma)\mu_0\pi + \gamma\mathcal{P}_\pi^*\zeta d_{\mathcal{D}}$$

But solving it is nontrivial:

- The dimension ( $|\mathcal{S} \times \mathcal{A}|$ ) is large or even  $\infty$
- We need to approximate  $d_{\mathcal{D}}$  and  $\mathcal{P}_\pi^*$  with empirical data (and proper regularization is critical)

## Solution

DualDICE (Dual DIstribution Correction Estimation) formulates it into an optimization problem to solve for  $\zeta_\pi$  [NCDL'19].

## More on DualDICE

### Convergence [NCDL'19]

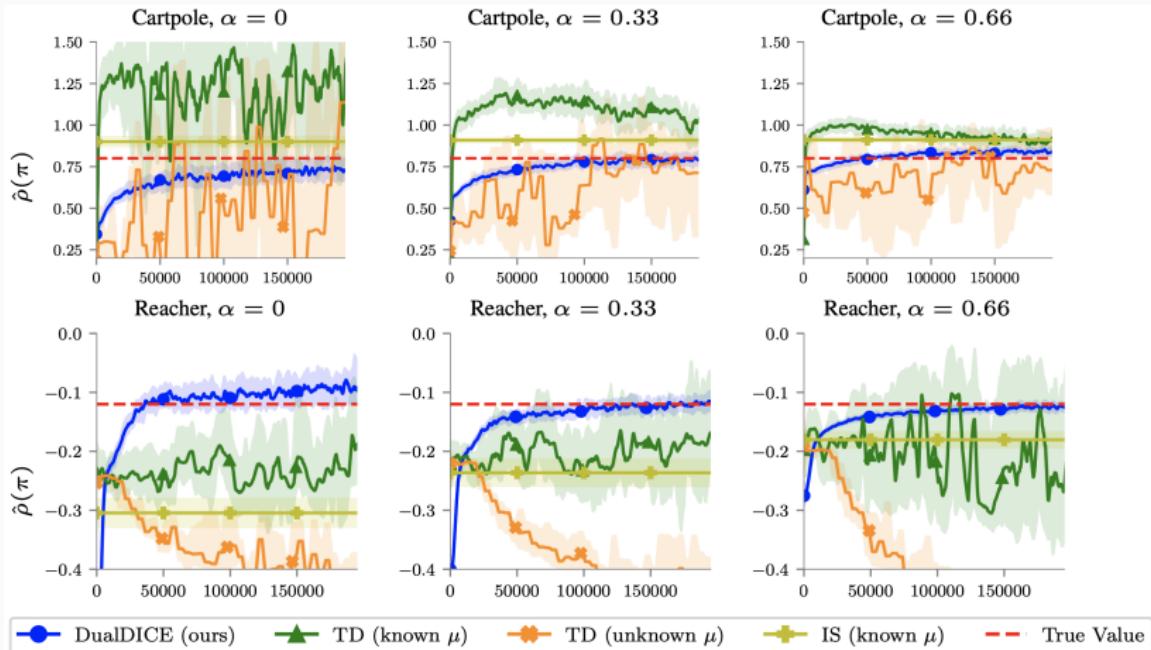
Under proper assumptions, the mean squared error

$$\mathbb{E} [\hat{v}_{\text{DualDICE}}(\pi) - v(\pi)]^2 = \tilde{O} \left( \varepsilon_{\text{approx}} + \varepsilon_{\text{opt}} + \frac{1}{\sqrt{n}} \right)$$

### Variations and Extensions

- Extension to  $\gamma = 1$  by GenDICE [ZDL'S'20]
- Doubly robustness [ZFLZL'20, Uehara+'20]
- A unification [YNDL'S'20] (DualDICE, GenDICE, LSTDQ [Lagoudakis&Parr'03], MQL/MWL [Uehara+'20], ...)
- Variational power iteration for easier optimization [PDL'S'20]
- Consistent policy gradient estimation (for **policy optimization**) from off-policy data [NDKCL'S'19]

# DualDICE Numerical Experiments



- Mujoco: simulated control
- $\alpha$  controls the degree of “off-policy-ness” ( $d_\pi = d_D$  if  $\alpha = 1$ )

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

# Uncertainty Estimation

Point estimate alone (like DualDICE) has limited use.  
Variance matters in most applications!

## Popular methods

- Normal approximation based on CLT [Bottou+'13, LCKG'15]
- Concentration inequalities [Thomas+'15]
- Bootstrap [Thomas+'15; Hanna+'17]

## Our approach: CoinDICE [DNCLSS'20]

- Based on empirical likelihood method
- Computationally tractable
- Avoids the curse of horizon

## Basic ideas

- Perturb the linear system with noise  $w$  within some range:

$$\zeta d_{\mathcal{D}} = (1 - \gamma)\mu_0\pi + \gamma \mathcal{P}_{\pi}^* \zeta d_{\mathcal{D}}$$

- Solve for  $\zeta$  while allowing  $w$  to vary
- The range of solved  $\zeta$  induces an (approximate) confidence interval for  $v(\pi)$

# CoinDICE: Statistical Coverage

## Asymptotic coverage

Under proper assumptions:

$$\lim_{n \rightarrow \infty} \Pr(\tilde{v}(\pi) \in C_{n,\xi}) = \Pr\left(\chi^2_{(1)} \leq \xi\right)$$

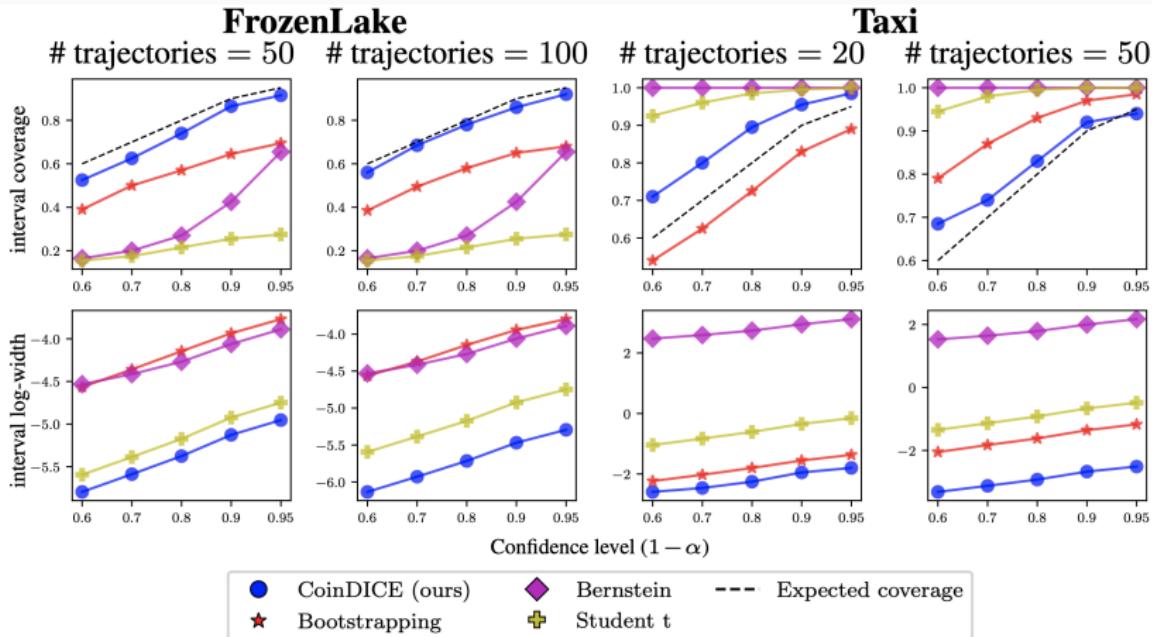
## Non-asymptotic coverage

$C_{n,\xi}$  is a valid confidence interval, up to error of order  $O(n^{-1})$ .

## Word of caution

The CI considers statistical error, **but not approximation error**.

# CoinDICE Numerical Experiments

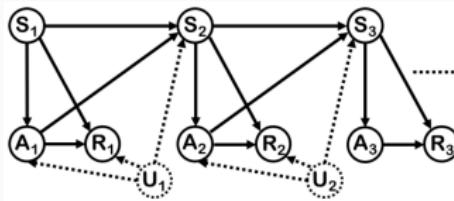


- Two standard benchmark problems
- $\alpha$  is the intended confidence interval coverage

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

# MDPs with Unmeasured Confounding



- Described by tuple  $\langle \mathcal{S}, \mathcal{A}, \dots, \mathcal{U} \rangle$
- $(u_1, u_2, u_3, \dots)$  are IID
- Same as MDP except for  $\Pr(s'|s, a, \mathbf{u})$ ,  $R(s, a, \mathbf{u})$ ,  $\pi_0(a|s, \mathbf{u})$  and  $\pi(a|s, \mathbf{u})$
- MDPs  $\subset$  MDPUCs  $\subset$  partially observable MDPs
- Simplified example: insulin injection for type-1 diabetes patients
  - $\mathcal{S}$ : blood glucose levels
  - $\mathcal{U}$ : exogenous unmeasured events (e.g., food intake, exercise)

## OPE in MDPUC

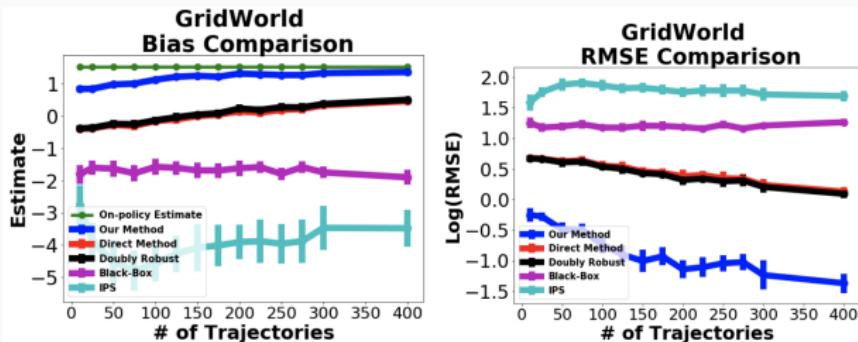
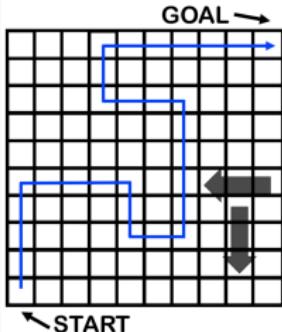
- Same form of estimator:

$$\hat{v}_{\text{MDPUC}} = \frac{1}{n} \sum_{i=1}^n \zeta_i r_i$$

with  $\zeta_i$  is the importance ratio as before ( $u$  gets canceled thanks to IID assumption)

- Previous method can be extended (under certain technical assumptions such as ergodicity), to achieve  $1/n^{-1/2}$  convergence rate.
- See paper for details [BKLM'20]

# Experiments: Windy Gridworld

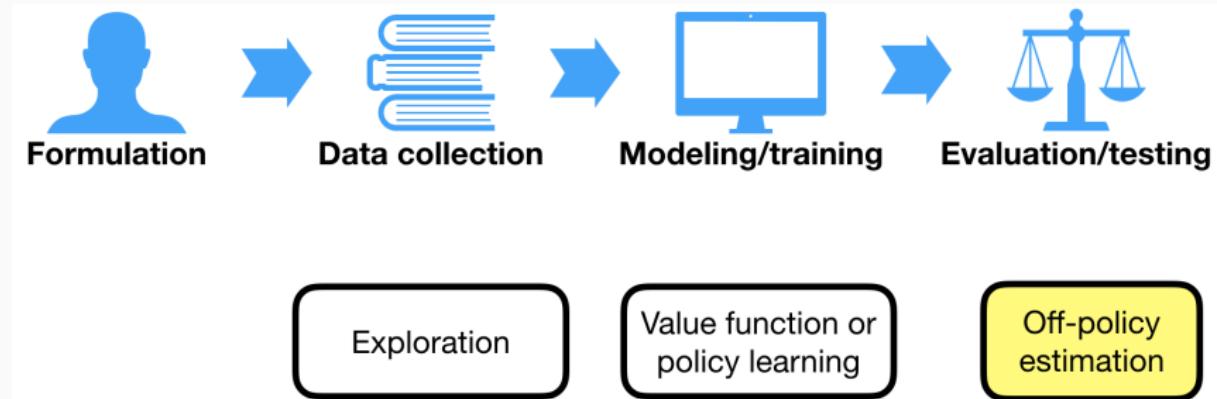


- Binary confounders that affect behavior policy, reward and transition probabilities
- Our method (blue) consistently estimates the true value, but not other baselines
- Sensitivity analysis and more results in the paper [BKLM'20]

# Outline

- The OPE problem
- Existing techniques
- OPE for long-term rewards
  - DualDICE: efficient point estimation
  - CoinDICE: confidence interval estimation
  - Dealing with unmeasured confounding
- Conclusions

# A Simplified View of RL



# Conclusions

## Recap

- Off-policy estimation is a critical bottleneck in “real-life” RL
- Many successes, mostly in short-horizon problems
- We break the curse of horizon based on duality structure in RL
- Extensions to confidence interval, and confounded cases

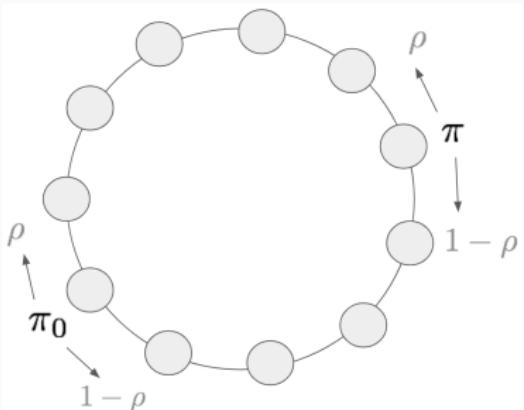
## Future work

- OPE in more general confounded cases
- OPE in multiagent systems [Zheng+ '20]
- Empirical tests in economic applications
- Connection with related work in econometrics

**THANK YOU!**

# Is Importance-Weighting Trajectories Necessary?

## “Circle” MDP



- $n = |\mathcal{S}|$  states arranged on circle
- Two actions:  
clockwise, counterclockwise
- Deterministic transitions

As  $H \rightarrow \infty$ :

- IPS/DR variance goes to  $\infty$
- But both policies visit every state equally often