

Policy and Off-Policy Evaluation

Mohammad Sadegh Talebi
m.shahi@di.ku.dk

Department of Computer Science



Motivation

"OPPOSITE OF REINFORCEMENT LEARNING, IT MEANS
THE MODEL IS GIVEN TO YOU"

We've studied planning in a known discounted MDP:

- Using VI, PI, and their variants
- Planning is a slang word for 'solving MDP'

WE JUST HAVE T_C AND
WANT TO KNOW V^π
(POLICY EVALUATION)

What if the MDP is unknown but accessible only through collected data?

- RL deals with (near-)optimally solving an unknown MDP using offline/online data (experience).
- The first step is policy evaluation using offline/online data.



PE vs. OPE vs. OPO

Policy Evaluation (PE) from data: Estimate V^π using data sampled from π .

Two related problems:

- **Off-Policy Evaluation (OPE):** Estimate V^π using data collected according to some **fixed policy** $\pi_b \neq \pi$
 - π_b is called the **behavior** (or **logging**) policy — an exploratory policy.
 - $\pi \neq \pi_b$ is called the **target** policy (a.k.a. **estimation** policy).
- **Off-Policy Optimization (OPO):** Find an optimal policy using data collected according to some **behavior policy** π_b



OPE/OPO

Consider a company selling products according to some policy A .

- Interactions with the world can be modeled as an MDP.
- The transition function (determined by, e.g., customer arrivals, market dynamics) is unknown, but the company has a rich dataset logged via A .
- The expected revenue under A can be found by computing V^A (Policy Evaluation, *this lecture!*).

Shall the company switch to a new policy B or not?

- Yes, if B yields a higher revenue, i.e., $V^B > V^A$
- One can find the **unknown** V^B via the dataset of A (via OPE methods).
- **Also OPE gives confidence sets on $V^B \implies$ Better to switch to B only if**

$$V^B \geq V^A + \text{margin, with high probability}$$



Part 1: Policy Evaluation



Policy Evaluation

Policy Evaluation

Given: A dataset \mathcal{D} collected under some *fixed* policy π .

Mathematically, $\mathcal{D} = \{(s_t, a_t, r_t), 1 \leq t \leq n\}$ where

$$a_t \sim \pi(\cdot|s_t), \quad r_t \sim R(s_t, a_t), \quad s_{t+1} \sim P(\cdot|s_t, a_t)$$

Goal: Derive (point) estimate, and possibly confidence intervals, for V^π .

We study two algorithms:

- A model-based method, which we call MB-PE.
- A model-free method called Temporal Difference (TD) Learning.

→ MODEL-BASED: FIRST ESTIMATE MDP AND THEN ESTIMATE THE VALUE FUNCTION

→ MODEL-FREE: DIRECTLY TRY TO ESTIMATE THE VALUE FUNCTION,
WITHOUT ESTIMATING THE MODEL (=MDP)



MB-PE: A Model-Based Method



Known Model

Recall the definition of V^π , $\pi \in \Pi^{\text{SR}}$:

$$V^\pi(s) = \mathbb{E}^\pi \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \middle| s_1 = s \right]$$

and the Bellman equation:

$$V^\pi(s) = r_t + \gamma \mathbb{E}^\pi \left[\sum_{t=2}^{\infty} \gamma^{t-1} r_t \middle| s_1 = s \right] = r_t + \gamma \mathbb{E}^\pi \left[V^\pi(s_{t+1}) \middle| s_1 = s \right]$$

π induces an MRP (P^π, r^π) with:

$$P_{s,s'}^\pi = \sum_a \pi(a|s) P(s'|s, a), \quad r^\pi(s) = \sum_a \pi(a|s) r(s, a)$$

$$\text{Then, } V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$$



MB-PE: Idea

Idea: Define estimators for P^π and r^π and apply the certainty equivalence principle.

Smoothed Estimator for P^π : $\widehat{P}_{s,s'}^\pi = \frac{N(s, s') + \alpha}{N(s) + \alpha S}$, with We basically just counts how many time we were in s and ended up in s'

$$\widehat{P}_{s,s'}^\pi = \frac{N(s, s') + \alpha}{N(s) + \underbrace{\alpha S}_{\text{ALSO AVOID DIVISION BY 0}}}, \quad \text{with}$$

$$N(s, s') = \sum_{t=1}^{n-1} \mathbb{I}\{s_t = s, s_{t+1} = s'\} \quad \text{and} \quad N(s) = \sum_{s' \in \mathcal{S}} N(s, s')$$

- $\alpha \geq 0$ is an arbitrary choice controlling the level of smoothing.
- $\alpha = 0$ corresponds to Maximum Likelihood Estimator (unbiased).
- $\alpha = 1/S$ corresponds to Laplace Smoothed Estimator (biased, but the bias vanishes as $N(s)$ increases).
- **Consistency:** $\widehat{P}_{s,s'}^\pi$ converges to $P_{s,s'}$ as $N(s) \rightarrow \infty$ almost surely.



MB-PE: Idea

Idea: Define estimators for P^π and r^π and apply the certainty equivalence principle.

Smoothed Estimator for r^π :

$$\hat{r}^\pi(s) = \frac{\alpha + \sum_{t=1}^{n-1} r_t \mathbb{I}\{s_t = s\}}{\alpha + N(s)}$$

NO NEED TO EXPLICITLY COUNT THE ACTIONS, BECAUSE IF WE KNOW THE STATE AND THE ACTION IS FIXED WE ALREADY HAVE ALL THE INFORMATION TO KNOW THE ACTION TAKEN

- **Consistency:** $\hat{r}^\pi(s)$ converges to $r^\pi(s)$ as $N(s) \rightarrow \infty$ almost surely.
- Unbiased for $\alpha = 0$.

Then, the following is an estimate for V^π :

$$\hat{V}^\pi = (I - \gamma \hat{P}^\pi)^{-1} \hat{r}^\pi$$

↑

HERE THE POLICY IS
"INCORPORATED" IN \hat{P}^π
(I THINK? NOT SURE)

← ESTIMATED \hat{P}^π AND \hat{r}^π AND NOW WE CAN ESTIMATE \hat{V}^π
(IT'S A POINT ESTIMATE)

We have trajectory like:

s_{t+1}	r_t	s_{t+1}
1	0	1
2	0	3
3	0	3
3	0	3
...



MB-PE: Convergence

Theorem

If π visits all states infinitely often, then \hat{V}^π converges to V^π almost surely:

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} \hat{V}^\pi = V^\pi\right) = 1$$

- I.e., if π is exploratory enough, \hat{V}^π converges to V^π in the following sense:

$$\mathbb{P}\left(\exists \mathcal{D}, \exists s \in \mathcal{S} : \lim_{t \rightarrow \infty} \hat{V}^\pi(s; \mathcal{D}) \neq V^\pi(s)\right) = 0$$

I.e., datasets for which $\hat{V}^\pi \neq V^\pi$ will occur with probability 0.

- It follows from the a.s. convergence of \hat{P}^π to P^π and of \hat{r}^π and r^π .
- We can use concentration inequalities (e.g., Hoeffding's) to derive confidence interval(s) for V^π .
 - E.g., they could tell us how much data is needed to have

$$\forall s : |\hat{V}^\pi(s) - V^\pi(s)| \leq \varepsilon \quad \text{w.p. at least } 1 - \delta$$

for input (ε, δ) .

POINT ESTIMATE \hat{V}^π
 $V^\pi \in [a, b] \text{ w.p. } \geq 1 - \delta$



MB-PE: Pros and Cons

This is a **model-based** approach since it maintains an approximate model of MDP (or MRP) and then computes V^π for that.

Disadvantages of the model-based solution:

- It results in **value estimates with a large variance in practice**, which is undesirable.
- It **maintains estimates of $S^2 + S$ elements of MRP**, though we need to maintain S estimates to find V^π .
- **Computational complexity is $O(S^3)$** , and **space complexity is $O(S^2)$** .
- May not be easily converted into an incremental procedure.



Temporal Difference Learning

(SKIPPED FOR NOW)



Temporal Difference Learning

- Temporal Difference Learning was popularized and extended by Richard Sutton in 1988.
- However, the earliest reported use dates back to Arthur Samuel (1959).

Application to Backgammon game by Gerald Tesauro (TD-Gammon), read more [here](#).



source: Wikipedia



Temporal Difference Learning

Assume \hat{V} is some estimate for V^π — Hence, $\hat{V}(s_t)$ is an estimate for $V^\pi(s_t)$.

Now consider $r_t + \gamma \hat{V}(s_{t+1})$:

$$\mathbb{E}\left[r_t + \gamma \hat{V}(s_{t+1}) \middle| s_t, \hat{V}\right] = \mathbb{E}_{a \sim \pi(s_t)} \left[R(s_t, a) + \gamma \sum_{s'} P(s' | s_t, a) \hat{V}(s') \middle| s_t, \hat{V} \right]$$

Hence, $r_t + \gamma \hat{V}(s_{t+1})$ gives *another estimate* for $V^\pi(s_t)$.



Temporal Difference Learning

Ideally we would like to have an estimate \hat{V} so that:

$$\hat{V}(s_t) \approx r_t + \gamma \hat{V}(s_{t+1})$$

- Given $\hat{V}(s_t)$, in view of Bellman's equation $r_t + \gamma \hat{V}(s_{t+1})$ serves as a target estimate for $V^\pi(s_t)$.
- The **temporal difference** error is $\delta_t = r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t)$.

Hence, we may update $\hat{V}(s_t)$ to reduce the error δ_t :

$$\underbrace{\hat{V}(s_t)}_{\text{new value}} \leftarrow \underbrace{\hat{V}(s_t)}_{\text{old value}} + \alpha_t \underbrace{\left(r_t + \gamma \hat{V}(s_{t+1}) - \hat{V}(s_t) \right)}_{\text{estimation error}}$$

This method is called **Temporal Difference (TD)** learning — this is a form of **bootstrapping**, since we refined $\hat{V}(s_t)$ using another estimate.



TD: Learning Rate

To guarantee convergence, learning rates $(\alpha_t)_{t \geq 1}$ must satisfy the *Robbins-Monro conditions*:

$$\alpha_t > 0, \quad \sum_{t=1}^{\infty} \alpha_t = \infty, \quad \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$

(I.e., a positive sequence that is *square-summable-but-not-summable*.)

Examples:

- $\alpha_t = \frac{1}{t+1}$
- $\alpha_t = \frac{2}{\sqrt{t} \log(t+1)}$
- $\alpha_t = \frac{c}{t^a}$ for $a \in (\frac{1}{2}, 1]$ and $c > 0$



TD

- **input:** $\mathcal{D} = \{(s_t, r_t)\}_{1 \leq t \leq n}, (\alpha_t)_{t \geq 1}$
- **initialization:** Select V_1 arbitrarily
- **for** $t = 1, \dots, n - 1$ **Update:**

$$V_{t+1}(s) = \begin{cases} V_t(s) + \alpha_t(r_t + \gamma V_t(s_{t+1}) - V_t(s)) & s = s_t \\ V_t(s) & \text{else.} \end{cases}$$

- **output:** V_n



TD: Advantages

- TD is **model-free**: It does not require a model of the MDP, only relies on collected experience.
- TD can be incremental (unlike the model-based methods).
- Computational complexity (per-step) is $O(1)$. Space complexity is S . Much cheaper than the model-based method.
- TD results in estimates for V^π with low variance.



Is TD Gradient?

- TD update resembles Stochastic Gradient Descent (SGD).
- However, it can be shown that TD is not an SGD for *any objective function* (see Philip Thomas' Notes, p. 69).
- In fact, TD is a **Stochastic Approximation (SA)** algorithm and it inherits convergence guarantee from SA — we briefly overview SA in next lecture.



TD: Convergence

Theorem

If all states are visited *infinitely often under π* and $(\alpha_t)_{t \geq 1}$ satisfies the Robbins-Monro conditions, then V_t converges to the true value function V^π almost surely:

$$\mathbb{P} \left(\forall s \in \mathcal{S}, \lim_{t \rightarrow \infty} V_t(s) = V^\pi(s) \right) = 1$$

In other words, if π is exploratory enough, V_t converges to V^π , in the following sense:

$$\mathbb{P} \left(\exists \mathcal{D}, \exists s \in \mathcal{S} : \lim_{t \rightarrow \infty} V_t(s; \mathcal{D}) \neq V^\pi(s) \right) = 0$$

i.e., datasets for which $V_\infty \neq V^\pi$ will occur with probability 0.



TD(λ)

TD only uses only r_t and $\hat{V}(s_{t+1})$ to refine $\hat{V}(s_t)$ — i.e., it looks *one-step into future*.

Why not looking into *ℓ -step into future*? using the target

$$\sum_{n=0}^{\ell} \gamma^n r_{t+n} + \gamma^{\ell+1} \hat{V}(s_{t+\ell+1})$$

The temporal difference error when using ℓ -step lookahead is:

$$\begin{aligned}\delta_t^\ell &= \sum_{n=0}^{\ell} \gamma^n r_{t+n} + \gamma^{\ell+1} \hat{V}(s_{t+\ell+1}) - \hat{V}(s_t) \\ &= \sum_{n=0}^{\ell} \gamma^n \left(r_{t+n} + \gamma \hat{V}(s_{t+n+1}) - \hat{V}(s_{t+n}) \right)\end{aligned}$$



TD(λ)

Looking into ℓ -step into future:

Now let's update $\widehat{V}(s_t)$ using a mixture of ℓ -steps information each weighted with $(1 - \lambda)\lambda^\ell$ for some $\lambda \in [0, 1]$:

$$\begin{aligned}\widehat{V}(s_t) &\leftarrow \widehat{V}(s_t) + \alpha_t \sum_{\ell=0}^{\infty} (1 - \lambda)\lambda^\ell \delta_t^\ell \\ &= \widehat{V}(s_t) + \alpha_t \sum_{n=0}^{\infty} \lambda^n \gamma^n (r_{t+n} + \gamma \widehat{V}(s_{t+n+1}) - \widehat{V}(s_{t+n}))\end{aligned}$$

This rule is called TD(λ) learning

- $\lambda = 0$ recovers TD (or TD(0)).
- $\lambda \rightarrow 1$ recovers the Monte-Carlo method.



Part 2: Off-Policy Evaluation



EVALUATING A POLICY π STARTING FROM
OBSERVATIONS ON A DIFFERENT POLICY π_b



OPE

OFF- Policy Evaluation

Given: A dataset \mathcal{D} of trajectories τ_1, \dots, τ_n , sampled from behavior policy π_b :

$$\tau_1 = (s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \dots, s_{T_1}^{(1)}, a_{T_1}^{(1)}, r_{T_1}^{(1)})$$

$$\vdots \qquad \vdots$$

$$\tau_n = (s_1^{(n)}, a_1^{(n)}, r_1^{(n)}, \dots, s_{T_n}^{(n)}, a_{T_n}^{(n)}, r_{T_n}^{(n)})$$

where

$$a_t^{(i)} \sim \pi_b(\cdot | s_t^{(i)}), \quad r_t^{(i)} \sim R(s_t^{(i)}, a_t^{(i)}), \quad s_{t+1}^{(i)} \sim P(\cdot | s_t^{(i)}, a_t^{(i)})$$

Goal: Derive (point) estimate, and possibly confidence intervals, for value of target policy π ($\neq \pi_b$).

Each trajectory could be even sampled from a different behavior policy.



OPE Assumptions

The main challenge of OPE is mismatch of distributions π_b and π

Coverage Assumption $\leftarrow \pi \text{ AND } \pi_b \text{ CAN'T BE TOO DIFFERENT!}$

For all $s \in \mathcal{S}$, if $\pi(a|s) > 0$ then $\pi_b(a|s) > 0$ \leftarrow OTHERWISE WE WON'T HAVE ENOUGH INFORMATION!

Implication: π is *absolutely continuous* with respect to π_b (thus a.k.a. **Absolute Continuity Assumption**).



A Model-Based Method



Known Model

If MDP M known, for any $\pi \in \Pi^{\text{SR}}$: $V^\pi = (I - \gamma P^\pi)^{-1} r^\pi$

Idea: Estimate P and R via \mathcal{D} and apply the certainty equivalence principle.

For simplicity, for now assume that \mathcal{D} contains only one trajectory:

$$\mathcal{D} = \{(s_t, a_t, r_t), t = 1, \dots, n\}$$

where:

$$a_t \sim \pi_b(\cdot | s_t), \quad r_t \sim R(s_t, a_t), \quad s_{t+1} \sim P(\cdot | s_t, a_t)$$



A Model-Based Solution (I)

Idea: Estimate P and R via \mathcal{D} and apply the certainty equivalence principle.

Introduce counts: For all (s, a, s')

$$N(s, a, s') = \sum_{t=1}^{n-1} \mathbb{I}\{s_t = s, a_t = a, s_{t+1} = s'\} \quad \text{and} \quad N(s, a) = \sum_{s' \in \mathcal{S}} N(s, a, s')$$

Smoothed Estimator for P and R :

$$\hat{P}(s'|s, a) = \frac{N(s, a, s') + \alpha}{N(s, a) + \alpha S}, \quad \hat{R}(s, a) = \frac{\alpha + \sum_{t=1}^{n-1} r_t \mathbb{I}\{s_t = s, a_t = a\}}{\alpha + N(s, a)}$$

THE ACTION a IS EXPLICITLY USED, BECAUSE WE NEED TO ESTIMATE P AND R , NOT P^π AND R^π

with $\alpha > 0$ an arbitrary smoothing parameter.

For any (s, a) , if $\pi_b(a|s) > 0$, then



$$\hat{P}(\cdot|s, a) \rightarrow_{N(s) \rightarrow \infty} P(\cdot|s, a) \quad \text{and} \quad \hat{R}(s, a) \rightarrow_{N(s) \rightarrow \infty} R(s, a), \quad \text{almost surely.}$$

A Model-Based Solution (II)

Smoothed Estimator for P and R :

$$\hat{P}(s'|s, a) = \frac{N(s, a, s') + \alpha}{N(s, a) + \alpha S}, \quad \hat{R}(s, a) = \frac{\alpha + \sum_{t=1}^{n-1} r_t \mathbb{I}\{s_t = s, a_t = a\}}{\alpha + N(s, a)}$$

⇒ Build the empirical MDP $\widehat{M} = (\mathcal{S}, \mathcal{A}, \widehat{P}, \widehat{R}, \gamma)$.

Then, the following is an estimate for V^π :

$$\widehat{V}^\pi = (I - \gamma \widehat{P}^\pi)^{-1} \widehat{r}^\pi$$

with $\widehat{P}_{s,s'}^\pi = \sum_{a \in \mathcal{A}} \pi(a|s) \widehat{P}(s'|s, a)$

W_E ARE INCORPORATING π IN
 P AND r

$$\widehat{r}^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \widehat{R}(s, a)$$

Theorem

Under the coverage assumption and that all states are visited infinitely often under π_b , \widehat{V}^π converges to V^π almost surely:

$$\mathbb{P}\left(\lim_{n \rightarrow \infty} \widehat{V}^\pi = V^\pi\right) = 1$$



Model-Free Methods



Importance Sampling: Basic Facts

Consider two distributions P and Q defined on \mathcal{X} , with $P \ll Q$.

$$\mathbb{E}_{x \sim P}[f(x)] = \int_x f(x)P(x)dx = \int_x f(x)Q(x) \underbrace{\frac{P(x)}{Q(x)}}_{\substack{\text{LIKE A MISMATCH CORRECTION} \\ \text{WEIGHT USED FOR ESTIMATIONS}}} dx = \mathbb{E}_{x \sim Q}\left[\frac{P(x)}{Q(x)}f(x)\right]$$

Note that importance weight $\frac{P(x)}{Q(x)}$ is well-defined due to $P \ll Q$.

Given are samples $X_i \sim Q, i = 1, \dots, n$:

- Importance Sampling (IS) estimator of $\mathbb{E}_{x \sim P}[f(x)]$:

$$\hat{f}_{\text{IS}} = \frac{1}{n} \sum_{i=1}^n f(X_i) \frac{P(X_i)}{Q(X_i)}$$

- Weighted Importance Sampling (wIS) estimator of $\mathbb{E}_{x \sim P}[f(x)]$:

$$\hat{f}_{\text{wIS}} = \frac{1}{\sum_{i=1}^n \frac{P(X_i)}{Q(X_i)}} \sum_{i=1}^n f(X_i) \frac{P(X_i)}{Q(X_i)}$$



Contrast these to $\hat{f} = \frac{1}{n} \sum_{i=1}^n f(X_i)$ built using $X_i \sim P, i = 1, \dots, n$.

Importance Weight Estimators: Properties

Lemma

\hat{f}_{IS} is consistent and unbiased.

Proof. Consistency follows from the SLLN. Unbiased since

$$\mathbb{E}[\hat{f}_{\text{IS}}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_Q[f(X_i) \frac{P(X_i)}{Q(X_i)}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_P[f(X_i)] = \mathbb{E}_P[f(X)]$$

Lemma

\hat{f}_{wIS} is consistent and biased.

Proof. To prove consistency, observe that by the SLLN, $\frac{1}{n} \sum_{i=1}^n \frac{P(X_i)}{Q(X_i)}$ converges to 1 w.p. 1 (since $X_i \sim Q$) and $\frac{1}{n} \sum_{i=1}^n f(X_i) \frac{P(X_i)}{Q(X_i)}$ converges to $\mathbb{E}_P[f(X)]$ w.p. 1. Showing biased via counter example: taking $X_1 = \dots = X_n$,

$$\hat{f}_{\text{wIS}} = \frac{1}{\sum_{i=1}^n \frac{P(X_i)}{Q(X_i)}} \sum_{i=1}^n f(X_i) \frac{P(X_i)}{Q(X_i)} = f(X_1)$$

Hence, $\mathbb{E}[\hat{f}_{\text{wIS}}] = \mathbb{E}[f(X_1)] \neq \mathbb{E}[f(X_1) \frac{P(X_1)}{Q(X_1)}]$.



Importance Sampling Estimator for OPE

Consider a trajectory $\tau = (s_1, a_1, r_1, \dots, s_T, a_T, r_T) \sim \pi_b$ (with $s_1 = s$).

- $\sum_{t=1}^T \gamma^{t-1} r_t$ is an estimator for $V^{\pi_b}(s)$. \leftarrow A NAIVE ESTIMATION
- To estimate V^π , we apply Importance Sampling \Rightarrow the entire τ corresponds to a sample:

$$\Rightarrow \frac{\mathbb{P}(\tau|\pi)}{\mathbb{P}(\tau|\pi_b)} \sum_{t=1}^T \gamma^{t-1} r_t \text{ is IS estimator of } V^\pi(s). \quad X \equiv \tau$$

- Note that $\begin{cases} \mathbb{P}(\tau|\pi) &= \prod_{t=1}^T \pi(a_t|s_t) P(s_{t+1}|s_t, a_t) \mathbb{P}(r_t|s_t, a_t) \\ \mathbb{P}(\tau|\pi_b) &= \prod_{t=1}^T \pi_b(a_t|s_t) P(s_{t+1}|s_t, a_t) \mathbb{P}(r_t|s_t, a_t) \end{cases}$

$$\begin{array}{ccc} X \sim \pi & X \not\sim \pi_b & \downarrow \\ \pi \neq \pi_b \end{array}$$

\Rightarrow Importance sampling estimator of $V^\pi(s)$ built using τ :

$$\frac{P(\tau)}{Q(x)} X \equiv \frac{\mathbb{P}(\tau|\pi)}{\mathbb{P}(\tau|\pi_b)} X$$

$$\hat{V}_{IS}^\pi(s; \tau) = \prod_{t=1}^T \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \sum_{t=1}^T \gamma^{t-1} r_t := \rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t$$

In general, we define $\rho_{1:t} = \prod_{t'=1}^t \frac{\pi(a_{t'}|s_{t'})}{\pi_b(a_{t'}|s_{t'})}$ for any t .



Importance Sampling Estimators for OPE

Given a dataset \mathcal{D} of n trajectories τ_1, \dots, τ_n :

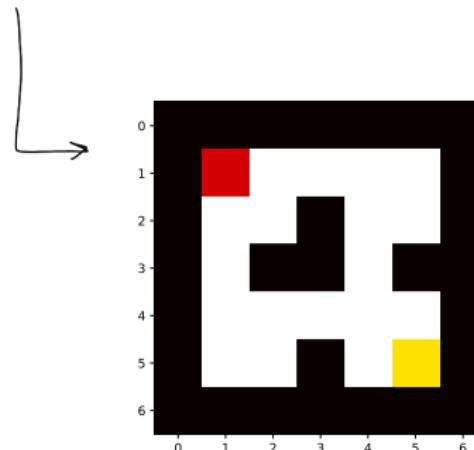
$$\tau_1 = (s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \dots, s_{T_1}^{(1)}, a_{T_1}^{(1)}, r_{T_1}^{(1)})$$

$$\vdots \qquad \vdots$$

$$\tau_n = (s_1^{(n)}, a_1^{(n)}, r_1^{(n)}, \dots, s_{T_n}^{(n)}, a_{T_n}^{(n)}, r_{T_n}^{(n)})$$

all starting in $s_{\text{init}} \in \mathcal{S}$ (i.e., $s_1^{(1)} = \dots = s_1^{(n)} = s_{\text{init}}$).

We are mostly interested in estimating $V^\pi(s_{\text{init}})$ for some π .



Example: 4-room Grid-World

- s_{init} is ■.
- Terminal state ■.
- Our interest is to estimate $V^\pi(\blacksquare)$



Importance Sampling Estimators for OPE

Given a dataset \mathcal{D} of n trajectories τ_1, \dots, τ_n :

$$\begin{aligned} \tau_1 &= (s_{\text{init}}, a_1^{(1)}, r_1^{(1)}, \dots, s_{T_1}^{(1)}, a_{T_1}^{(1)}, r_{T_1}^{(1)}) \\ &\quad \vdots \qquad \vdots \\ \tau_n &= (s_{\text{init}}, a_1^{(n)}, r_1^{(n)}, \dots, s_{T_n}^{(n)}, a_{T_n}^{(n)}, r_{T_n}^{(n)}) \end{aligned} \quad \left. \right\} \begin{matrix} \text{THEY CAN BE OF} \\ \text{DIFFERENT LENGTH} \end{matrix}$$

- IS estimator of $V^\pi(s_{\text{init}})$ built using \mathcal{D} :

$$\widehat{V}_{\text{IS}}^\pi(s_{\text{init}}; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \widehat{V}_{\text{IS}}^\pi(s_{\text{init}}; \tau_i) = \frac{1}{n} \sum_{i=1}^n \rho_{1:T_i}^{(i)} \sum_{t=1}^{T_i} \gamma^{t-1} r_t^{(i)}$$

(consistent and unbiased, but typically with high variance)

- wIS estimator of $V^\pi(s_{\text{init}})$ built using \mathcal{D} :

$$\widehat{V}_{\text{wIS}}^\pi(s_{\text{init}}; \mathcal{D}) = \frac{\sum_{i=1}^n \rho_{1:T_i}^{(i)} \sum_{t=1}^{T_i} \gamma^{t-1} r_t^{(i)}}{\sum_{i=1}^n \rho_{1:T_i}^{(i)}}$$

(consistent, slightly biased, but with lower variance)

WE LIK THIS MORE (HIGH VARIANCE IS MORE ANNOYING THAN A SIGHT BIAS)



Per-Decision IS Estimator

Consider IS: Note that

$$\rho_{1:T} \sum_{t=1}^T \gamma^{t-1} r_t = \rho_{1:T} r_1 + \rho_{1:T} \gamma r_2 + \dots + \rho_{1:T} \gamma^{T-1} r_T$$

Observe that for each t ,

$$\mathbb{E}[\rho_{\mathbf{1}:T} r_t] = \mathbb{E}[\rho_{\mathbf{1}:t} r_t]$$

Because r_t is independent of future actions and states (and hence $\rho_{t+1:T}$).

Hence, we consider the following Per-Decision IS estimator:

$$\begin{aligned}\widehat{V}_{\text{PDIS}}^{\pi}(s; \tau) &= \rho_{\mathbf{1}:1} r_1 + \rho_{\mathbf{1}:2} \gamma r_2 + \dots + \rho_{\mathbf{1}:t} \gamma^{t-1} r_t + \dots + \rho_{\mathbf{1}:T} \gamma^{T-1} r_T \\ &= \sum_{t=1}^T \rho_{\mathbf{1}:t} \gamma^{t-1} r_t\end{aligned}$$

Contrast it with $\widehat{V}_{\text{IS}}^{\pi}(s; \tau) = \rho_{\mathbf{1}:T} \sum_{t=1}^T \gamma^{t-1} r_t$



Importance Sampling Estimators for OPE

Given a dataset \mathcal{D} of n trajectories τ_1, \dots, τ_n :

$$\tau_1 = (\mathbf{s}_{\text{init}}, a_1^{(1)}, r_1^{(1)}, \dots, s_{T_1}^{(1)}, a_{T_1}^{(1)}, r_{T_1}^{(1)})$$

$$\vdots \qquad \vdots$$

$$\tau_n = (\mathbf{s}_{\text{init}}, a_1^{(n)}, r_1^{(n)}, \dots, s_{T_n}^{(n)}, a_{T_n}^{(n)}, r_{T_n}^{(n)})$$

- Per-Decision IS estimator of $V^\pi(s_{\text{init}})$ built using \mathcal{D} :

$$\widehat{V}_{\text{PDIS}}^\pi(s_{\text{init}}; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \widehat{V}_{\text{PDIS}}^\pi(s_{\text{init}}; \tau_i) = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{T_i} \boldsymbol{\rho}_{1:t}^{(\mathbf{i})} \gamma^{t-1} r_t^{(i)}$$

(consistent and unbiased; expected to yield lower variance than IS and wIS.)



Importance Sampling Estimators for OPE

Summary of importance sampling estimators built using \mathcal{D} comprising of τ_1, \dots, τ_n , starting in s_{init} :

$$\widehat{V}_{\text{IS}}^{\pi}(s_{\text{init}}; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \boldsymbol{\rho}_{\mathbf{1}: \mathbf{T}_i}^{(i)} \sum_{t=1}^{T_i} \gamma^{t-1} r_t^{(i)}$$

(consistent and unbiased, but typically with high variance)

$$\widehat{V}_{\text{wIS}}^{\pi}(s_{\text{init}}; \mathcal{D}) = \frac{\sum_{i=1}^n \boldsymbol{\rho}_{\mathbf{1}: \mathbf{T}_i}^{(i)} \sum_{t=1}^{T_i} \gamma^{t-1} r_t^{(i)}}{\sum_{i=1}^n \boldsymbol{\rho}_{\mathbf{1}: \mathbf{T}_i}^{(i)}}$$

(consistent, slightly biased, but with lower variance)

$$\widehat{V}_{\text{PDIS}}^{\pi}(s_{\text{init}}; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^{T_i} \boldsymbol{\rho}_{\mathbf{1}: \mathbf{t}}^{(i)} \gamma^{t-1} r_t^{(i)}$$

(consistent and unbiased;
expected to yield lower variance than IS and wIS.)



Summary

- PE vs. OPE vs. OPO
 - Data-driven approaches where data at hand is “off” the target policy.
 - PE and OPE are **prediction** (= estimation) problems, not **learning**.
 - In contrast, OPO is a **learning** problems.
- For PE, we studied
 - MB-PE: model-based, implementing certainty-equivalence
 - TD: model-free, implementing a form of bootstrapping, but a stochastic approximation method
- For OPE, we studied
 - MB-OPE: model-based, implementing certainty-equivalence
 - IS, wIS, and PDIS: estimators based on importance sampling and Monte-Carlo
- We mostly discussed asymptotic convergence guarantees and consistency of estimators. PAC-type bounds exist and are more relevant in practice.

