

AML, 29 Oct 2019

Monte Carlo Methods

$\pi \rightarrow$ Calculate V_{π} (Policy Evaluation)

$\pi_0 \xrightarrow{e} V_{\pi_0} \xrightarrow{\text{greedy}} \pi_1 \xrightarrow{e} V_{\pi_1} \dots$ Policy Iteration

May be done incrementally - Generalized PI

Given π - generate random episodes

- observe rewards

- work backwards & update V_{π}

Instead of v_{π} , estimate $q_{\pi}(s,a)$

- Ensure all (s,a) pairs are visited often enough
- Exploring Starts : choose random initial (s,a)
- Or, used ϵ -soft / ϵ -greedy strategies

$$Q_{\pi}(s,a) \geq \frac{\epsilon}{|N(s)|}$$

always

of actions at s

\Downarrow

Choose non greedy with $\frac{1}{|N(s)|} \cdot \epsilon$

Choose greedy with $(1 - \epsilon + \frac{1}{|N(s)|})$

Want to converge to an optimum policy

- Typically policy is greedy / deterministic

- Explore vs exploit

- Need to consciously deviate from current optimum

Instead

- Use a different policy to generate samples

Off Policy strategy

Use b to sample
& update π

vs

On policy strategy

↓

Same π to sample & update

Sample according to b , but update π

- Estimates via b samples generate expected values wrt b , not π !

Basic constraint

If $\pi(a|s) > 0$ then $b(a|s) > 0$



Typically greedy,
deterministic



Must be stochastic,
in general

Suppose we are at state S_t

Remaining trajectory $A_t, S_{t+1}, A_{t+1}, \dots, A_{T-1}, S_T$

According to π

$$\begin{aligned} & \Pr(A_t, S_{t+1}, \dots, A_{T-1}, S_T \mid S_t, A_{t:T-1} \sim \pi) \\ &= \pi(A_t \mid S_t) P(S_{t+1} \mid S_t, A_t) \pi(A_{t+1} \mid S_{t+1}) \dots \\ &= \prod_{k=t}^{T-1} \pi(A_k \mid S_k) P(S_{k+1} \mid S_k, A_k) \end{aligned}$$

Instead, if we use b rather than π

$$\prod_{k=t}^{T-1} b(A_k | S_k) P(S_{k+1} | S_k, A_k)$$

$$\begin{aligned} S_{t:T-1} &= \frac{\prod_{k=t}^{T-1} \pi(A_k | S_k) P(S_{k+1} | S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k | S_k) P(S_{k+1} | S_k, A_k)} \\ &= \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)} \end{aligned}$$

$$E[G_t | S_t = s] = v_b(s)$$

Instead

$$E[S_{t:T-1} \cdot G_t | S_t = s] = v_\pi(s)$$

Importance Sampling

All visits to s contribute to $V_{\pi}(s)$

$$V(s) = \sum_{t \in \tau(s)} \underbrace{\sum_{t: T(t)-1} G_t}_{\text{Adjusted observation}}$$

$|\tau(s)|$

time points
where s is
visited

end point of the
episode
in which t occurs

Uniform count across episodes

$$e_1, e_2, \dots, e_T \quad \left| \quad e_{T+1}, e_{T+2}, \dots, e_{T+T'} \quad \right|$$

Episode 1 episode 2

Previous ratio is called ordinary importance sampling

Instead Weighted Importance Sampling

$$V(s) = \frac{\sum_{t \in \tau(s)} \gamma_{t:T(t)-1} \cdot G_t}{\sum_{t \in \tau(s)} \gamma_{t:T(t)-1}}$$

Consider a single episode

Ordinary

Single observation

Weighted

$$V(\hat{\beta}) =$$

$$S_{t:T-1} \cdot G_t$$

|

Lower bias

Higher variance

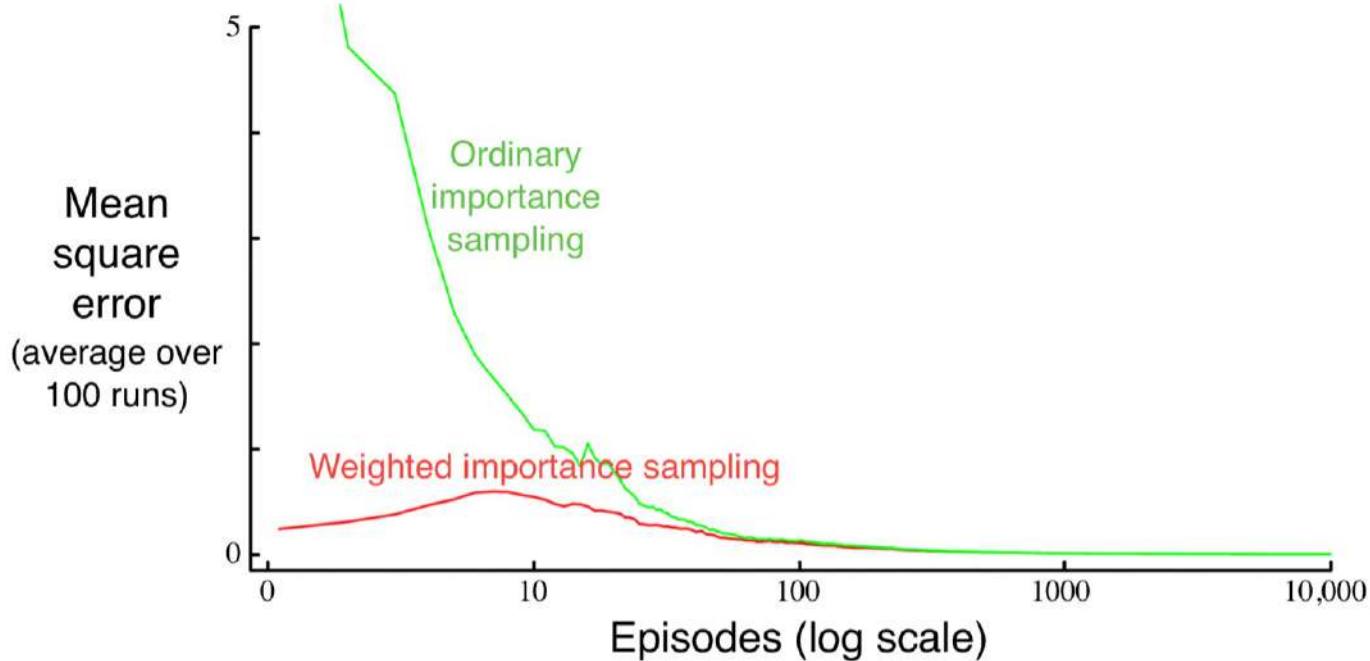
$$V(\hat{\beta}) = \frac{S_{t:T-1} \cdot G_t}{S_{t:T-1}}$$

$$S_{t:T-1}$$

|

Higher bias

lower variance



Exercise: Weighted importance sampling can be computed incrementally

Off-policy MC prediction (policy evaluation) for estimating $Q \approx q_\pi$

Input: an arbitrary target policy π

Initialize, for all $s \in \mathcal{S}$, $a \in \mathcal{A}(s)$:

$$Q(s, a) \in \mathbb{R} \text{ (arbitrarily)}$$

$$C(s, a) \leftarrow 0$$

Loop forever (for each episode):

✓ $b \leftarrow$ any policy with coverage of π

Generate an episode following b : $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$

$$G \leftarrow 0$$

$$W \leftarrow 1$$

Loop for each step of episode, $t = T-1, T-2, \dots, 0$, while $W \neq 0$:

$$G \leftarrow \gamma G + R_{t+1}$$

$$C(S_t, A_t) \leftarrow C(S_t, A_t) + W \leftarrow$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{W}{C(S_t, A_t)} [G - Q(S_t, A_t)]$$

$$W \leftarrow W \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$$

Generalize this to do policy iteration

Dynamic Programming



Update V_L to V_{L+1}

Bootstrapping

Monte Carlo



V_{L+1} replaces V_L

Not bootstrapped



MDP
theory



Temporal Difference Learning

a "true" contribution of

RL

Sampling + Bootstrapping

Monte Carlo

$$V(s_t) \leftarrow V(s_t) + \alpha [G_t - V(s_t)]$$

TD

$$V(s_t) \leftarrow V(s_t) + \alpha \left[R_{t+1} + \gamma V(s_{t+1}) - V(s_t) \right]$$

estimate

TD(0) - single step lookahead

$$V_{\pi}(s) \stackrel{\text{def}}{=} E[G_t \mid S_t = s] \quad - \text{Monte Carlo}$$
$$= E[R_{t+1} + \lambda G_{t+1} \mid S_t = s] \quad - \text{Defn of rewards \& values}$$
$$= E[R_{t+1} + \gamma V_{\pi}(S_{t+1}) \mid S_t = s] \quad - \text{TD}(0)$$

Example Driving time

Office \rightarrow Highway \rightarrow Side road \rightarrow Home

<i>State</i>	<i>Elapsed Time (minutes)</i>	<i>Predicted Time to Go</i>	<i>Predicted Total Time</i>
leaving office, friday at 6	0	30	30
reach car, raining	5	35	40
exiting highway	20	15	35
2ndary road, behind truck	30	10	40
entering home street	40	3	43
arrive home	43	0	43

