## Monte Carlo Methods:

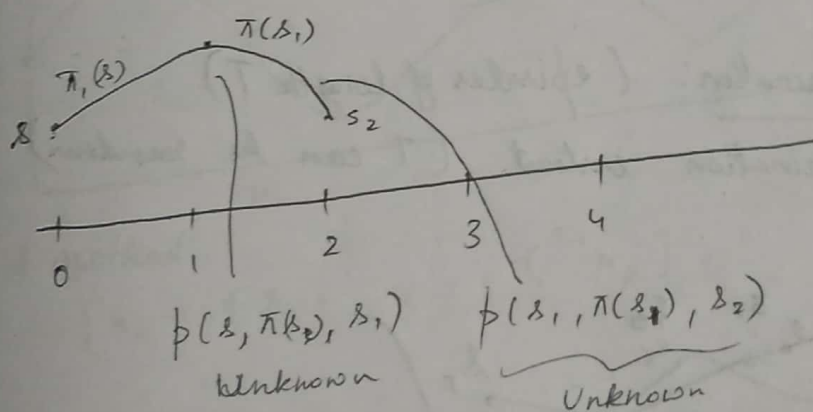— Prediction (first problem)
— Control (second problem)

_Idea:-_ Model of system is unknown. One has access to samples / transitions either through simulation / real data.

_Prediction Problem:_ Two approaches
- first visit method
- Every visit method.

~~Rec~~ Recall, value $f^n$ under a given policy $\pi$

$$V_\pi(s) = E\left[\sum_{k=0}^{\infty} \gamma^k \mathcal{H}_{R+1} \mid S_0 = s\right]$$



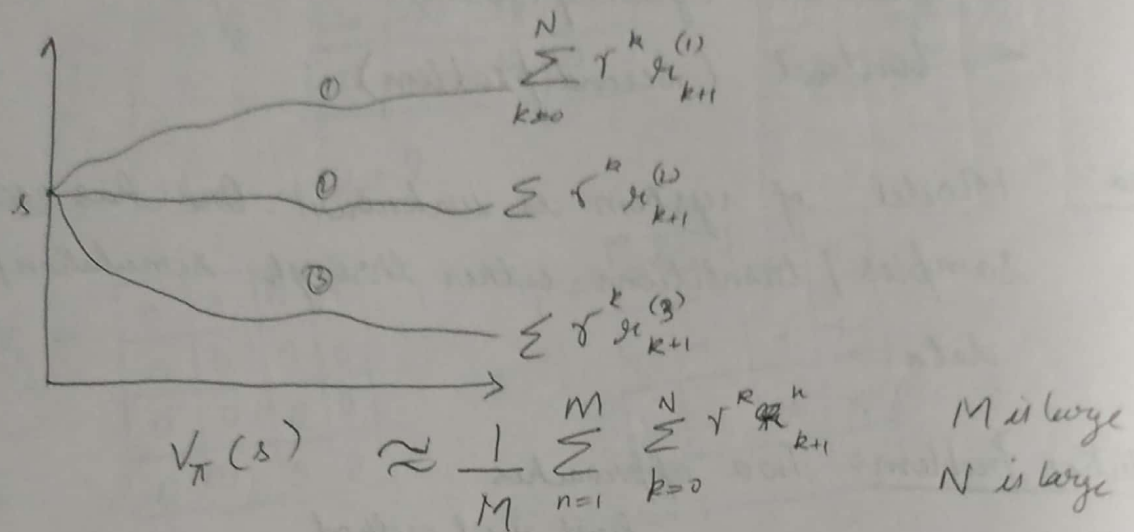$p(s, \pi(s_0), s_1)$      $p(s_1, \pi(s_1), s_2)$

Unknown        Unknown

In order to find $V_\pi(s)$, we need to calculate or approx $E[\cdot]$ but we cannot calculate $E[\cdot]$ since we don't know $p(s, \pi(s), \cdot)$

$\pi(\cdot | s)$

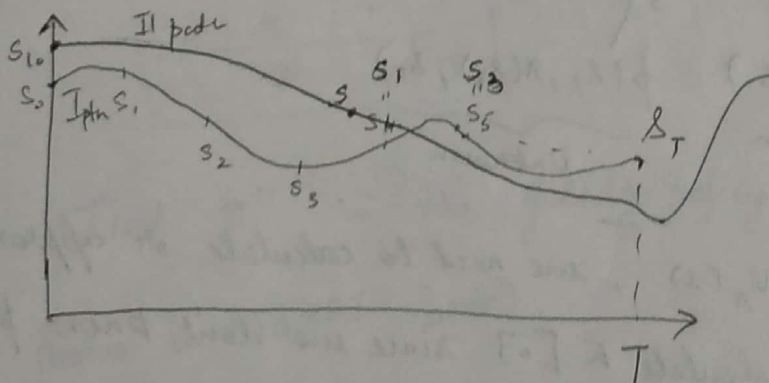[Most elementary procedure],

Run sample paths with intial state $s$.



$$\sum_{k=0}^{N} \gamma^{k} r_{k+1}^{(1)}$$

$$\sum \gamma^{k} r_{k+1}^{(2)}$$

$$\sum \gamma^{k} r_{k+1}^{(3)}$$

$$V_{\pi}(s) \approx \frac{1}{M} \sum_{n=1}^{M} \sum_{k=0}^{N} \gamma^{k} r_{k+1}^{n} \qquad \begin{array}{l} M \text{ is large} \\ N \text{ is large} \end{array}$$

$$V_{\pi}(s) = E\left[ \sum_{k=0}^{\infty} \gamma^{k} r_{k+1} \mid s_0 \right] \Big/ \!\!\! \text{(proor)}$$

this method is inefficient since everytime intial state $s$ is changed, the steps of $M$ & $N$ has to ~~been~~ repeated again

Problems with termination: (episodes of length $T$)

$T$ is the termination instant. ($T$ can be random)

## First visit method:

= state $s_0$ is visited first time at time 0

— state $s_1$ is visited — — — — 1

— state $s_2$ — — — — — 2

— state $s_3$ — — — — — 3

— state $s_1$ is visited second time at time 4

## I path:

state $s$ $\quad r_1 + \gamma \, r_2 + \ldots \gamma^{T-1} \cdot r_T$

## II path

state $s$ $\quad r_1^{(2)} + \gamma \, r_2^{(2)} \ldots \gamma^{T-1} r_T^{(2)}$

## Every visit method:



$s_0 = s, \; s_{11} = s, \; s_{18} = s, \; s_{2a} = s$

## I visit method:

$$[r_1 + \gamma \, r_2 + \ldots \gamma^{T-1} r_{k-}]$$

## Every visit method:

$$r_1 + \gamma \, r_2 + \ldots \gamma^{T-1} r_{T} \longrightarrow X_1$$

$$r_{11} + \gamma \, r_{13} + \ldots \gamma^{T-1} r_T \longrightarrow X_2$$

$$r_{19} + \gamma \, r_{20} + \ldots \gamma^{T-1} \cdot r_T \longrightarrow X_3$$

$$\frac{1}{N} \sum_{i=1}^{N} X_i$$

Kavitha

# Exercise: Grid World Example

state

| 0 | 1 | 2 | 3 |
|---|---|---|---|
| 4 | 5 | 6 | 7 |
| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 0 |

Action
(up)

optimal policy

$$p(1, \leftarrow, 0) = 0.7 \qquad p(1, \leftarrow, 2) = 0.1 \qquad p(i, a, i)$$
$$\qquad\qquad\qquad\qquad\quad p(1, \leftarrow, 1) = 0.1$$
$$p(1, \leftarrow, 5) = 0.1$$

$$p(5, \nwarrow, 4) = 0.4$$
$$p(5, \nwarrow, 1) = 0.4 \qquad \leftarrow \text{ similar rule for all other actions.}$$
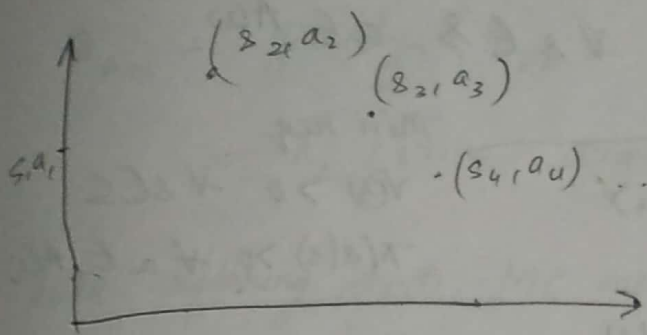$$p(5, \nwarrow, 6) = 0.1$$
$$p(5, \nwarrow, 9) = 0.1$$

Apply monte-Carlo transitions first visit method averaged over 28 independent sample paths to estimate $V_\pi(s)$ $\forall s \{1, 2, \ldots 14\}$

## Note:

- Monte carlo estimate can be used to estimate q-values $(q_*)$ [one can obtain optimal policy here]

- Makes sense to look at randomized policies instead of at deterministic $q(s, a)$ is being estimated for various $(s, a)$ tuples.

Graph showing points $(s_2, a_2)$, $(s_2, a_3)$, $(s_4, a_4)$ with axes labeled $s_i a_i$

- Actions $a_i$ in states $s_i$ will be picked up from the given policy $\pi$
- If $\pi$ is deterministic, then in state $s_i$, action picked $= \pi(s_i) = a_i$
- This is a disadvantage since we only learn about how good action $a_i = \pi(s_i)$ is in state $s_i$ but not about other actions feasible in state $s_i$

- Thus, we should consider randomized policies $\pi(a|s)$, $\in A(s)$

$$s.t \quad \pi(a|s) = P(A_i = a | S_i = s) > 0 \quad \forall a \in A(s).$$

## Exploring Starts:

Episodes start with certain state-action pairs. Assume that all state action pairs have non-zero probability of the selection at the start of episode.

$\Rightarrow$ All state action pairs will be ~~time limited~~ visited a infinite number of times in the limit as number of episodes $\to \infty$

Suppose, $V(s) = p(s_o = s)$, $s \in S$ be the initial distribution on states.

Then, $\mu(s, a) = p(s_o = s, A_o = a) = P(s_o = s) P(A_o = a | s_o = s)$

$$= V(s) \pi(a|s)$$

Kavitha

We assume that,

$$\mu(s,a) = \nu(s)\,\pi(a|s) > 0 \quad \forall s \in S,\; a \in A(s)$$

min req

$$\nu(s) > 0 \quad \forall s \in S$$
$$\pi(a|s) > 0 \quad \forall a \in A(s)$$

$$\sum_{s,a} \mu(s,a) = \sum_{s}\sum_{a} \nu(s)\,\pi(a|s)$$

$$= \sum_{s} \nu(s) \sum_{a} \pi(a|s)$$

$$= 1$$

(N.G)

① Inverse R.L problem — construct optimal reward

→ Avg cost MDP — Q learning.    raghub @ iisc.ac.in

## Monte Carlo Control:

Evaluation
$Q \sim q_\pi$

$\pi$            $Q$

Improvement
$\pi \sim$ Greedy $Q$

Given an intial policy $\pi_0$,

$$\pi_0 \xrightarrow{\;E\;} q_0 \xrightarrow{\;I\;} \pi_1 \xrightarrow{\;E\;} q_{\pi_1} \xrightarrow{\;I\;} \dots \xrightarrow{I} \pi_* \xrightarrow{E} q_*$$

goal: Find $(\pi_*, q_{\pi_*})$

$$(\pi, q_\pi) \xrightarrow[\text{Improvement}]{\text{policy Eval}} (\pi^*, q_{\pi^*})$$

$$\pi_{k+1} = \text{greedy}(q_{\pi_k})$$

$$= \pi_{k+1}(s) = \underset{a}{\text{argmax}} \; q_{\pi_k}(s,a) \quad \forall \, s \in S$$

with exploring starts, MC methods will compute $q_{\pi_k}$ exactly for arbitrary $\pi_k$

$\pi_{k+1}$ is a better policy than $\pi_k$ because,

$$q_{\pi_k}(s, \pi_{k+1}(s)) = q_{\pi_k}(s, \underset{a}{\text{argmax}} \; q_{\pi_k}(s,a))$$

$$= \underset{a}{\text{max}} \; q_{\pi_k}(s,a)$$

$$\geqslant q_{\pi_k}(s, \pi_k(s))$$

$$= v_{\pi_k}(s)$$

If $q_{\pi_k}(s, \pi_{k+1}(s)) = q_{\pi_k}(s, \pi_k(s))$

then both $\pi_k$ & $\pi_{k+1}$ are optimal policies

Else, $\exists \, s_0 \in S$ s.t

$$q_{\pi_k}(s_0, \pi_{k+1}(s_0)) > q_{\pi_k}(s_0, \pi_k(s_0))$$

Note that $q_{\pi_k}(s,a)$ are estimated using Monte-Carlo based policy evaluation. We dont need an infinite no of iterates for PE to converge for given policy.

Kavitha

## Work around 1:

$$\text{Stop when } \left| V_{\pi_{k+1}}(s) - V_{\pi_R}(s) \right| < \delta$$

## Work around 2:

Use a priori defined integer $M_1, M_2, M_3$ etc steps of iterate P.E which is followed by policy improvement (Modified Policy Iteration)

This takes less no of steps of P.E before an improvement before an improvement step is conducted.

## Monte Carlo Exploring Starts (ES) for Estimating $\pi \approx \pi_*$

### Initialize:

$\pi(s) \in A(s) \quad \forall s \in S$

$Q(s,a) \in R \text{ (arbitrary)} , \forall s \in S, a \in A(s)$

$\text{Returns}(s,a) \leftarrow \text{Empty list} \quad \forall s \in S, a \in A(s)$

### Loop for each episode:     (similar to $Q$-Value iteration)

- choose $S_0 \in S$, $A_0 \in A(S_0)$ randomly such that all pairs of states & actions have prob $> 0$

- Generate an episode from $S_0, A_0$ following $\pi$:
$$\mathcal{D}: S_0, A_0, R_1, S_1, A_1, R_2 \cdots \cdots S_{T-1}, A_{T-1}, R_T$$

- $G \leftarrow 0$

  Loop ( for each step of of episode $t = T-1, T-2 \cdots 0$)
  - $G \leftarrow \gamma G + R_{t+1}$

  - Append $G$ to returns $(S_t, A_t)$

$$Q(S_t, A_t) \leftarrow \text{Average (Returns } (S_t, A_t))$$

$$\pi(S_t) \leftarrow \arg\max_a Q(S_t, a)$$

↰ for each episode

---

$$G \leftarrow 0$$

For $t = T-1, T-2 \ldots 0$

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

$$G_{T-1} = R_T$$

$$G_{T-2} = r G_{T-1} + R_{T-1} = r R_T + R_{T-1}$$

$$G = R_1 + r R_2 + r^2 R_3 + \ldots r^{T-1} R_T$$

$$\bar{A}_n = \frac{1}{n} \sum_{i=1}^{n} G_i \quad \text{avg of } G_t's \text{ over } n \text{ episodes}$$

$$\bar{A}_{n+1} = \frac{1}{n+1} \sum_{i=1}^{n+1} G_i$$

$$= \frac{1}{n+1} \left( \sum_{i=1}^{n} G_i + G_{n+1} \right) = \frac{n}{n+1} \bar{A}_n + \frac{1}{n+1} G_{n+1}$$

$$= \bar{A}_n + \frac{1}{n+1} (G_{n+1} - \bar{A}_n)$$

---

Q) Can Monte- Carlo control procedure (with ES) converge to a sub optimal policy.

No, MC control with ES gives us optimal policy since otherwise correspondig value $f^n$ will be sub optimal & policy improvement step happening on that value $f^n$ will give a better policy ⟹ convergence didn't happen.

(Inverted pendulum — MOC with E-s)

Kavitha

# Monte Carlo "without" Exploring Starts:

If we don't use exploring starts, alternatively:

- All actions need to be selected infinitely often in each action.

### R.L methods:

On-Policy Methods     Off policy methods.

### On policy method:

Generate an episodes using policy $\pi$

$$S_0, A_0, R_1, S_1, A_1, \ldots S_{T-1}, A_{T-1}, R_T$$

### Goal:

Estimate $V_\pi(S_0)$ [Value of state $S_0$ under policy $\pi$)

$$\boxed{\begin{array}{l} G_{T_1}^{(1)} = R_1^{(1)} + \gamma R_2^{(1)} \\ + \ldots \gamma^{T-1} R_{T_1}^{(1)} \end{array}}$$

If $G_{T_1}^{(1)}$, $G_{T_2}^{(2)}$, $\ldots \ldots G_{Tm}^{(m)}$ are estimate

of $V_\pi(S_0)$ from $M$ episodes, then

$$V_\pi(S_0) \simeq \frac{1}{M} \sum_{i=1}^{M} G_{T_i}^{(i)}$$

$T_i$ is termination for it$^{th}$ episode.

### Off policy Methods:

Generate episodes using policy $\pi$

$$S_0, A_0, R_1 \ldots \ldots S_{T-1}, A_{T-1}, R_T$$

Goal: Estimate $v_b(S_0)$ value of state $S_0$ under policy $b$

Note: → On-policy methods _evaluate_ or improve policy
used to make decisions.

→ Off-policy methods evaluate or improve policy
different from policy used to generate data.