

13<sup>th</sup> Oct

Mercedes-Benz R & D India

## Policy Iteration

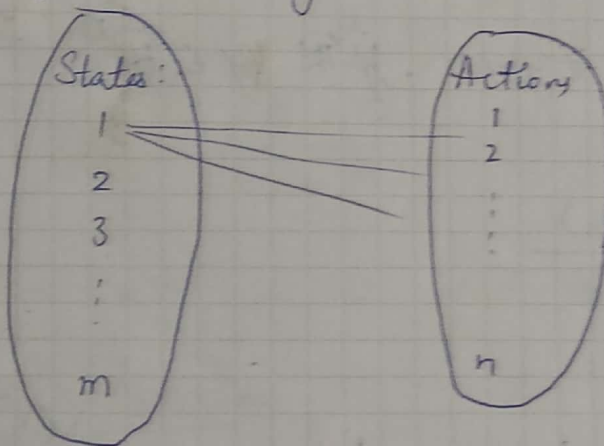
- Start with initial policy  $\pi_0$
- Evaluate policy  $\pi_0$  [find  $V_{\pi_0}$ ]
- Improvement of policy  $\pi_0$  to another policy  $\pi_1$
- Evaluate policy  $\pi_1$ , find  $V_{\pi_1}$

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \dots \xrightarrow{E} V_{\pi_k}$$

$\xrightarrow{E}$ : policy evaluation

$\xrightarrow{I}$ : policy improvement

2) After how many iterations, do we get  $V_{\pi^*}$



$$\underbrace{m \times m \times m \times \dots \times m}_n \Rightarrow m^n$$

Max no of ~~func~~ functions that map state space to action space =  $m^n$

Functions denote policies

$\therefore$  Max no of policies =  $m^n < \infty$

$$\pi : S \longrightarrow A$$

policy improvement:  $\exists$  at least one state  $s \in S$

$$(\pi_i \rightarrow \pi_{i+1}) \text{ s.t. } V_{\pi_{i+1}}(s) > V_{\pi_i}(s) \quad \left( \begin{array}{l} \text{does not mean} \\ \text{there is degradation} \end{array} \right)$$

As the procedure has converged

but not imply policy have converged. Since there can be multiple policy that are optimal.

$$q_{\pi_0}(s, b) = E_{\pi_0} [r + \gamma V_{\pi_0}(s') \mid x_0 = s, a_0 = b]$$

$\Rightarrow$  implies, pick initial reward due to 'b', next pick actions according to  $\pi_0$  in future.

$$\pi_1(s) = \underset{b \in A(s)}{\operatorname{argmax}} q_{\pi_0}(s, b) \quad \forall s \in S$$

$$V_{\pi_1}(s) \geq V_{\pi_0}(s) \quad \forall s \in S$$

$\Downarrow$  Either  $V_{\pi_1}(s) = V_{\pi_0}(s) \quad \forall s \in S$  or  $\exists \bar{s} \in S$  s.t.  $V_{\pi_1}(\bar{s}) > V_{\pi_0}(\bar{s})$

If  $V_{\pi_1} = V_{\pi_0} \quad \forall s$   
means  $\pi_1, \pi_0$  are converged.

In the above method, policies are

$$\pi_1, \pi_2, \dots, \pi_m$$

the search space is finite.

$$V_{\pi_1} \leq V_{\pi_2} \leq \dots \leq V_{\pi_m}$$

$\Rightarrow$  The procedure will converge in a finite number of policy improvement steps.

Unlike randomize policy

$$\pi(a|s) = \begin{matrix} a_1 & w.p. p_1 \\ a_2 & w.p. p_2 \\ \vdots & \vdots \\ a_m & w.p. p_m \end{matrix}$$

$$\text{st } p_1 + p_2 + p_3 + \dots + p_m = 1.$$

this method has a distribution of policies.

One might require to use gradient search. This method can have local max, but in previous method you get global max due to finite search space.

PI procedure

Policy Improvement procedure:

1. Initialization :  $V(s) \in \mathbb{R}$  and  $\pi(s)$  arbitrary  $\forall s \in S$

2. Policy Eval : Loop

$$\Delta \leftarrow 0$$

Loop for each  $s \in S$

$$v \leftarrow V(s)$$

$$V(s) = \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until  $\Delta < \theta$  (small +ve no)



## 3) Policy Improve:

for each  $s \in S$ old action  $\leftarrow \pi(s)$ 

$$\pi(s) \leftarrow \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

- If old action =  $\pi(s) \forall s$ then stop & return  $V \equiv V^* \quad \& \quad \pi \equiv \pi^*$ - If old action  $\neq \pi(s)$  for some  $(s)$ ,  
go back to step 2 & perform P.E for policy  $\pi$ - If  $V_{\pi}(s) = V_{old}(s) \forall s$ , then $V_{\pi} = V^*$  &  $\pi$  is old-policy are both optimal.Ex: Jack's Car Rental. - Implement policy iteration.Value Iteration:

$$V_{k+1}(s) = \max_a E[R_{t+1} + \gamma V_k(s_{t+1}) | S_t = s, A_t = a]$$

$$= \max_a \sum p(s', r | s, a) [r + \gamma V_k(s')] \quad k \geq 0$$

$$V_0 = (V_0(s), s \in S)$$

Note:

combines both evaluation &amp; iteration

eg: at  $t=0$ . $V_1(s)$  is calculated, later  $\max_a$  is taken~~this~~ first step is eval & second is improvement

Similarity between PI & VI:

Suppose in PI, we run iterative PE step exactly for one iteration  $\Delta$  then do improvement.

$\Rightarrow$  same as VI

### Modified Policy Iteration (MPI)

Suppose  $m_1, m_2, m_3, \dots$  are a sequence of integers (all  $a \geq 1$ ) s.t. iterative P.E in PI procedure is run exactly  $m_1$  steps, the first time,  $m_2$  steps the second time etc...

MPI is guaranteed to converge to optimal policy-value function tuple.

Policy evaluated

$$P(T_{\pi} V)(s) = \sum_{s', r} p(s', r | s, \pi(s)) [r + \gamma V_{\pi}(s')] \quad \forall s \in S$$

Define  $T(V)(s) = \max_{a \in A(s)} \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$

$V = TV$  — Bellman equation.

$T$  is a contraction.

$V, TV, T^2V, \dots \rightarrow V^*$  as  $n \rightarrow \infty$  — value iteration

$$V_{n+1} = TV_n \quad \forall n$$

$V^*$  is the unique sol<sup>n</sup> to bellman eqn.  $V = TV$



Value Iteration:

Initialize:  $V(s) \forall s \in S$ 

Loop:

$$\Delta \leftarrow 0$$

loop  $\forall s \in S$ 

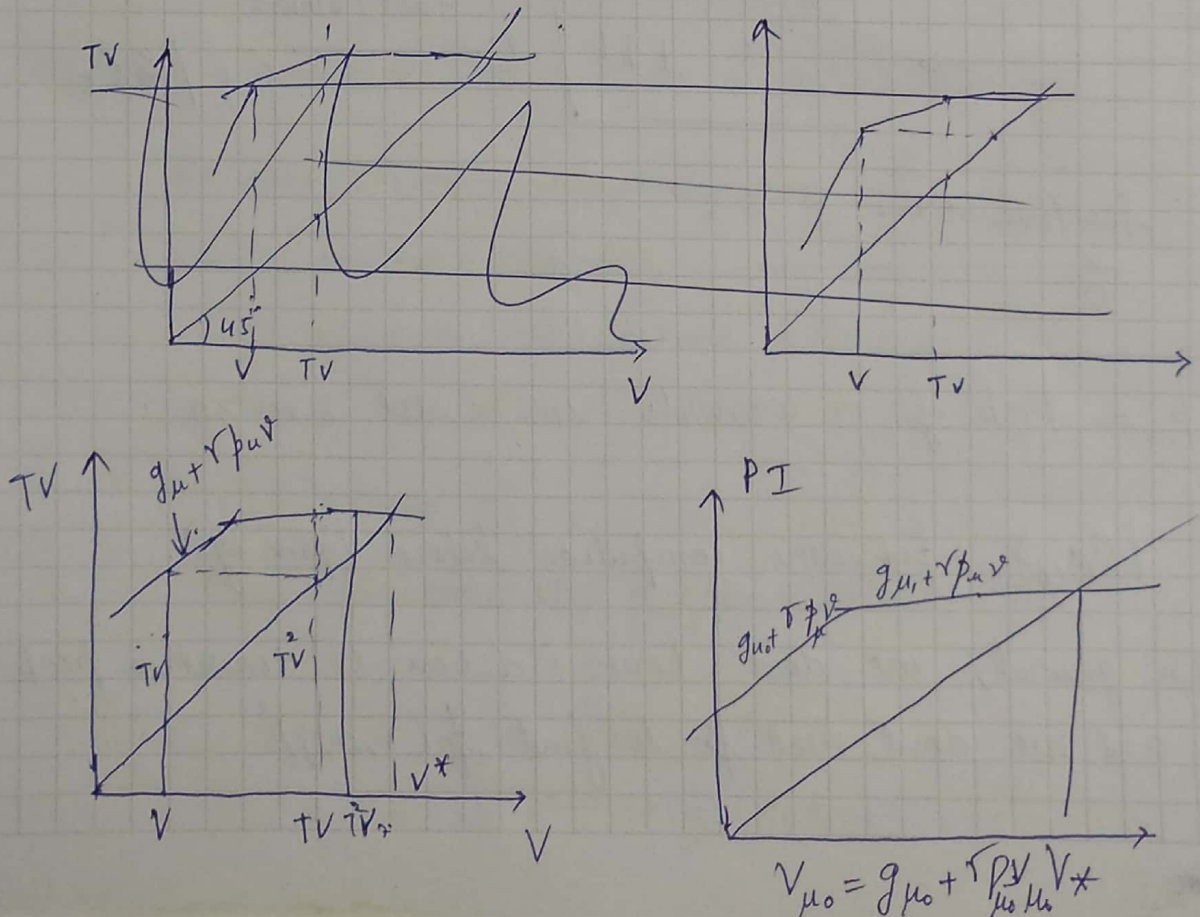
$$v \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum_{s', r} p(s', r | s, a) [r + V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until  $\Delta < \theta$ Output  $V(s) \forall s$  and deterministic policy  $\pi$ 

$$s.t. \pi(s) = \argmax_a \sum_{s', r} p(s', r | s, a) [r + \gamma V(s')]$$

VI vs PI

5

## Monte Carlo Methods.

Learning Theory: "Model of System is unknown."

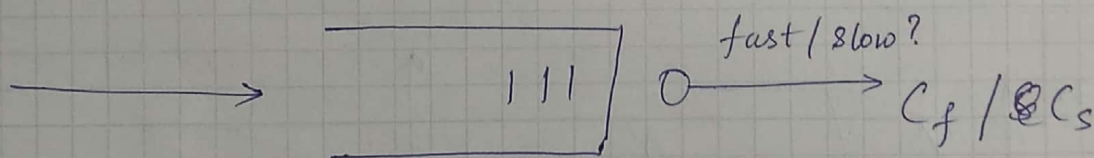
i.e.,  $p(i, a, j)$  is unknown

$r(i, a, j)$  is unknown.

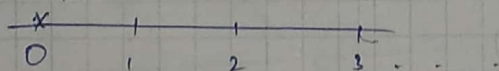
but it is possible to observe next state given current state and action.

Observe immediate random rewards,  $R_n$

We have access to either real data or data from a simulator



fast/slow at  $t=0$



$p_m$  = prob of  $m$  arrivals in a slot,  $m \geq 0$

$p(i, a, j)$  prob were computed before (prev eg)

In general, we don't have access to transition prob and we don't need  $p$  to find  $p(i, a, j)^s$



## Monte - Carlo Prediction.

Aim is to learn value function for a given policy from sample rewards.

We first consider episodic tasks —  $\mathcal{T}$  a transition state

two procedures

first visit  
MC method

every visit  
MC method.

~~MC~~ First visit MC prediction

Input : a policy  $\pi$  to be evaluated.

Initialize:  $V(s) \in \mathbb{R} \quad \forall s \in \mathcal{S}$

Returns( $s$ )  $\leftarrow$  an empty list  $\forall s \in \mathcal{S}$

loop (for each episode)

Generate an episode following  $\pi$

$s_0, a_0, r_1, s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T$

$G \leftarrow 0$

loop for each step of episode  $t = T-1, T-2, \dots, 0$

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

Append  $G$  to Returns( $s_t$ )

$V(s_t) \leftarrow \text{Avg}(\text{Returns}(s_t))$



Note:  $T$  is a random time.

In the first pass,  $G = R_T$

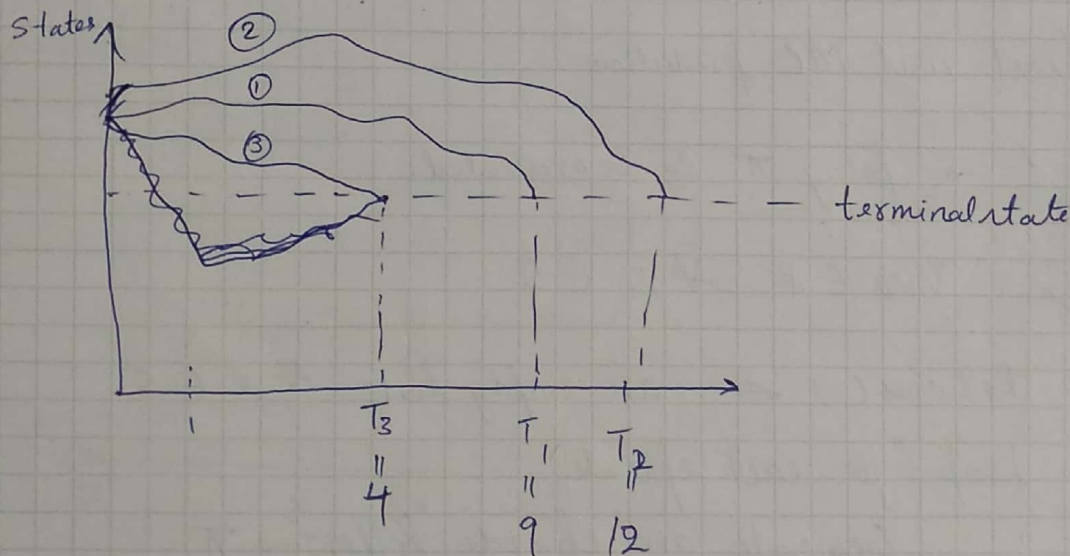
Next:  $G := \gamma R_T + R_{T-1}$

Next:  $\gamma(\gamma R_T + R_{T-1}) + R_{T-2}$   
 $= \gamma^2 R_T + \gamma R_{T-1} + R_{T-2}$

Eventually,

$$G := R_1 + \gamma R_2 + \gamma^2 R_3 \dots \gamma^{T-1} R_T$$

$$V_\pi(s) = E_\pi[G_T | s_0 = s]$$



$$G_1 = R_1 + \gamma R_2 + \dots + \gamma^{T_1-1} R_{T_1}$$

$$G_2 = R_1^2 + \gamma R_2^2 + \dots + \gamma^{T_2-1} R_{T_2}^2$$

$$G_m = R_1^m + \gamma R_2^m + \dots + \gamma^{T_m-1} R_{T_m}^m$$

$$\frac{1}{M} \sum_{i=1}^M h_i \approx V_{\pi}(s)$$

~~For~~ A primary goal: Estimate  $q_{*}$  (optimal action-values)

Consider policy evaluation for given policy  $\pi$ .

Estimate  $q_{\pi}(s, a)$  [expected return starting in state  $s$  & picking action  $a$  & subsequently picking actions according to  $\pi$ ]

M.C methods for

action-value functions are same as for state-value functions except that we now look at state-action tuples  $(s, a)$  that are visited.

Note: If  $\pi$  is deterministic, then in each state ' $s$ ', we pick only the action  $\pi(s) \triangleq a$

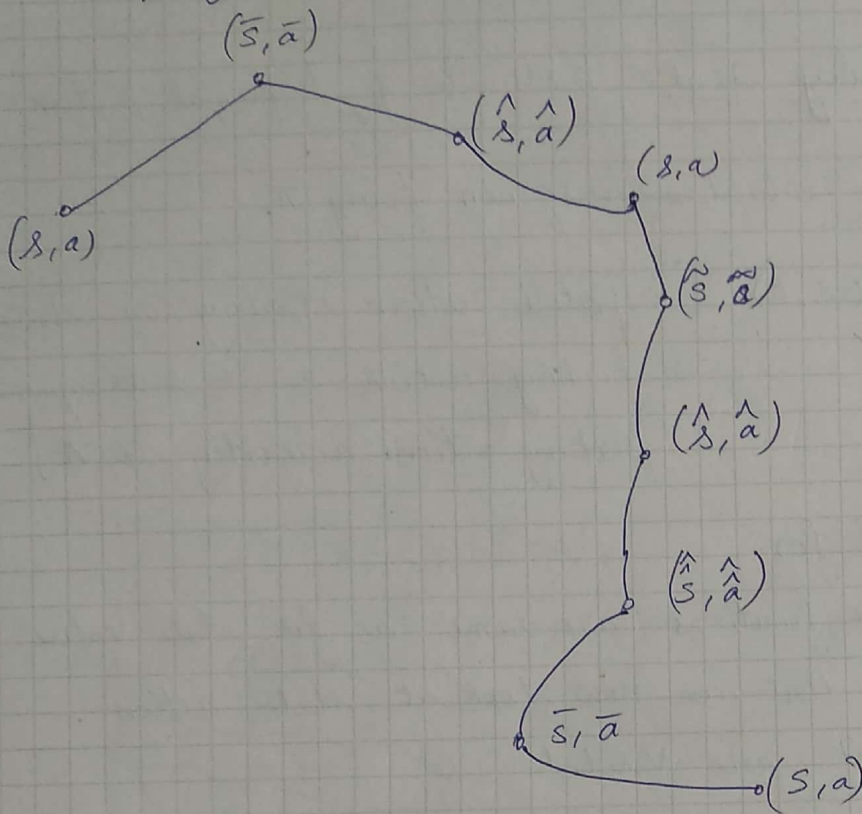
In such case, we will ~~be~~ not learn values for actions

$$b \in A(s) \setminus \{a\}$$

$$\hookrightarrow A(s) \text{ differs from } A(s) \cap \{a\}^c$$



Every visit MC Method for action-values. under a policy



Sequence of state-action tuples  $\{(x_n, A_n)\}$  where  $A_n \approx \pi(\cdot | s_n)$  forms a Markov chain.

Transition prob for  $(x_n, A_n)$ ,  $n \geq 0$

$$P(X_{n+1}=j, A_{n+1}=b \mid X_n=i, A_n=a)$$

$$= P(A_{n+1}=b \mid X_{n+1}=j, X_n=i, A_n=a) * P(X_{n+1}=j \mid X_n=i, A_n=a)$$

$$= \pi(b|j) p(i, a, j) \quad \text{--- } i, a \text{ doesn't matter for } \underline{b}$$

$$= \sum_j \underbrace{\sum_b \pi(b|j) p(i, a, j)}_1 = \sum_j 1 \cdot p(i, a, j) = 1$$