

16/3/19

Policy Iteration Algorithm

- Policy Improvement

$$\text{find } \gamma_{k+1} \text{ s.t. } T\gamma_{k+1} \cdot \gamma_k = T\gamma_k$$

- Policy Evaluation

(given Policy γ_k)

(Iterative PE can take a long time for convergence)

Value Iteration Algorithm

$$V_{k+1}(s) = \max_a E \left[R_{t+1} + \gamma V_k(s_{t+1}) \right] \quad \begin{array}{l} s_t = s, \\ A_t = a \end{array}$$
$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$

~~if $s \in S$~~

VI can be thought of as spl. Case of PI when the IPE step is run exactly once before the improvement step is conducted.

Value Iteration for estimating $\pi = \pi^*$

Parameter: small threshold $\theta > 0$ determining the accuracy of estimation.

loop:

$$\Delta \leftarrow 0$$

loop for each $s \in S$:

$$v \leftarrow V(s)$$

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

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

until $\Delta < \theta$

output $V(s)$ for s and a deterministic policy $\pi \approx \pi^*$

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

Modified Policy Iteration

Fix same integers $m_0, m_1, m_2, \dots \geq 1$

Start with some initial estimate of policy π_0

At k^{th} step: suppose current policy is π_k

$$\pi_{k+1}(s)$$

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

Policy Evaluation

$$V_{k+1} = T_{\pi_k}^{\pi_k} V_k$$

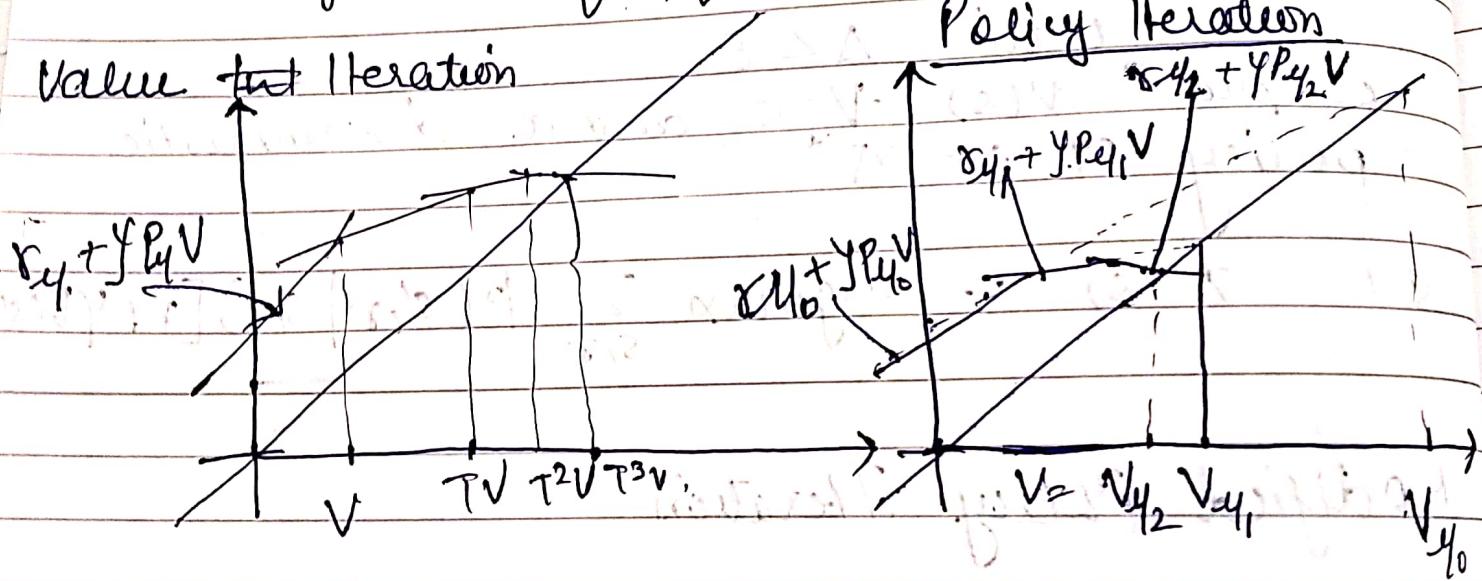
Policy Improvement (Gives π_{k+1})

$$T_{\pi_{k+1}} V_{k+1} = T V_{k+1}$$

$$V_{n+1} = T V_n$$

π converges at $m \rightarrow \infty$
 $V \rightarrow V^*$

Value Iteration



Geometric Approach PI:

suppose $\pi \in \mathbb{R}^n$ is some test function

A. policy π is greedy w.r.t $T\pi$

$$Ty J = TJ$$

Thus

$$\begin{aligned} & \sum_j p_{ij}^y (y_{C(i)}) (r(i, y_{C(i)}, j) + \alpha J(j)) \\ &= \max_{y \in A(i)} \sum_{j \in S} p_{ij}^y (y) (r(i, y, j) + \alpha J(j)) \end{aligned}$$

Note :- There can be several greedy policies for a given J and there can be several J , for which a policy is greedy.

Let $R_y = \{J \mid y \text{ is greedy w.r.t } J\}$

Note $J \in R_y \iff \text{and only if } Ty J = TJ$.

In other words,

$$\begin{aligned} & \sum_{j \in S} p_{ij}^y (y_{C(i)}) (r(i, y_{C(i)}, j) + \alpha J(i)) \\ & \geq \sum_{j \in S} p_{ij}^y (y) (r(i, y, j) + \alpha J(i)) \\ & \quad \forall y \in A(i) \end{aligned}$$

Q Is R_y is a closed set?

Consider a sequence $\{J_n\}$ of function s.t $J_n \in R_y \forall n \geq 1$. Let $J_n \rightarrow J$ as $n \rightarrow \infty$.

$$J_n(x) \rightarrow J(x) \text{ as } n \rightarrow \infty$$

$$\max_{s \in S} \sum_{i \in I(s)} p_{is}^y (y_{C(i)})$$

since $J_n \in R_y \forall n$, we have

$$\sum_{j \in S} p_{ij}(y(c)) (r(i, y(c), j) + \alpha J_n(j))$$
$$\geq \sum_{j \in S} p_{ij}(y) (r(i, y(c), j) + \alpha J_n(j))$$

$\forall y \in A_{ij}$
 $\forall n \geq 1$

limits on both sides as $n \rightarrow \infty$, we obtain,

$$\sum_{j \in S} p_{ij}(y(c)) (r(i, y(c), j) + \alpha J(c)) \geq \sum_{j \in S} p_{ij}(y) (r(i, y, j) + \alpha J(j))$$

Property

Step 2: $f_n \geq g_n \forall n$.

Let some $f_n \leq g_n$.

point for the
limit to exist

Suppose. J_1 & J_2 are 2 functions in R_y

Let $\beta \in [0, 1]$ & let

$$J = \beta J_1 + (1-\beta) J_2$$

Is $J \in R_y$?

since $J \in R_y$, we have.

$$\begin{aligned}
 & \sum_{j \in S} p_{ij} (\gamma(i)) (\alpha(i, \gamma(i), j) + \alpha_{J_1}(j)) \\
 & \geq \sum_{j \in S} p_{ij} (\gamma) (\alpha(i, \gamma, j) + \alpha_{J_1}(j)) \quad \rightarrow \textcircled{1}
 \end{aligned}$$

Since $J_2 \in R_\gamma$, we have.

$$\begin{aligned}
 & \sum_{j \in S} p_{ij} (\gamma(i)) (\alpha(i, \gamma(i), j) + \alpha_{J_2}(j)) \\
 & \geq \sum_{j \in S} p_{ij} (\gamma) (\alpha(i, \gamma, j) + \alpha_{J_2}(j)) \quad \rightarrow \textcircled{2}
 \end{aligned}$$

$\beta \textcircled{1} + (1-\beta) \textcircled{2}$ gives

$$\begin{aligned}
 & \sum_{j \in S} p_{ij} (\gamma(i)) [\beta \alpha(i, \gamma(i), j) + \alpha_{J_1}(j)] \\
 & \quad + (1-\beta) \alpha(i, \gamma(i), j) + \alpha_{J_2}(j)
 \end{aligned}$$

$$\begin{aligned}
 & \geq \sum_{j \in S} p_{ij} (\gamma) [\beta \alpha(i, \gamma, j) + \alpha_{J_1}(j) \\
 & \quad + (1-\beta) \alpha(i, \gamma, j) + \alpha_{J_2}(j)]
 \end{aligned}$$

$$\begin{aligned}
 \text{or } & \sum_{j \in S} p_{ij} (\gamma(i)) (\alpha(i, \gamma(i), j) + \alpha_{J_1}(j)) \\
 & \geq \sum_{j \in S} p_{ij} (\gamma) (\alpha(i, \gamma, j) + \alpha_{J_1}(j)) \\
 & \quad + \alpha_{J_2}(j)
 \end{aligned}$$

$\Rightarrow J \in R_\gamma$. Thus R_γ is a convex set.
Thus, R_γ is both closed and convex.

$$\Delta \leftarrow \max_{\mathbf{u}_{\text{opt}}} \sum_{i \in S} \lambda_i (y_{iS}) \sum_{j \in S} V(j)$$

They is obtained their 'greedy' policy
 \bar{J}_Y w.r.t J

$$J = (J(1), J(2), \dots, J(n)) \quad J: \mathbb{R}^n \rightarrow \mathbb{R}$$

Q Do R_Y form a partition of \mathbb{R}^n

$$R_Y = \{J \mid \bar{J}_Y \text{ is greedy w.r.t } J\}$$

we showed R_Y is closed and convex

$$J \in R_Y \Leftrightarrow T_Y J = T J$$

Suppose $J_Y \in R_Y$

$$T_Y J_Y = T J_Y$$

$$J_Y = T J_Y$$

$$J_Y = T_Y J_Y$$

$$\text{or } J_Y = J_Y$$

Example

Consider a discounted reward problem with 2 states 1 & 2. Actions available in each state

1st action - Stay there at a reward of -1

2nd action - Move to other state at a reward of 0.

There are 4 policies

2 states.

M_L



Possible

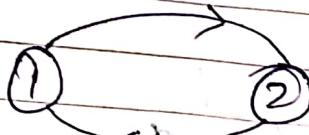
②

4_L

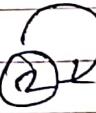
①

②

M_g



4_g



consider an arbitrary reward vector $J \in \mathbb{R}^4$. In state 1, a greedy policy can choose to stay if and only if:

$$-1 + \alpha J(1) \geq 0 + \alpha J(2)$$

$$\text{or } \alpha J(1) \geq \alpha J(2) + 1$$

$$\text{or } J(1) \geq J(2) + \frac{1}{\alpha} \rightarrow (1)$$

similarly in state 2 a greedy policy chooses to stay if and only if:

$$-1 + \alpha J(2) \geq 0 + \alpha J(1)$$

$$\text{or } J(2) \geq J(1) + \frac{1}{\alpha} \rightarrow (2)$$

We now identify sets R_{4_L} , R_{4_g} , R_{4_g} and R_{4_w} associated with four policies.

$V(k)$

for each
 k

For worst policy γ_w

$$J^{\gamma_w}(1) = -1 - \alpha - \alpha^2 - \alpha^3 - \dots \\ \Rightarrow -\frac{1}{1-\alpha}$$

$$J^{\gamma_w}(2) = -\frac{1}{(1-\alpha)}$$

For γ_w to be greedy w.r.t J_J we require both (1) & (2) to hold. So we require

$$J(1) \geq J(1) + \frac{2}{\alpha}$$

Similarly

$$J(2) \geq J(2) + \frac{2}{\alpha}$$

This is impossible

$$\Rightarrow \gamma_w = \emptyset$$

Consider now the best policy γ_L

$$J^{\gamma_L}(1) = -\frac{1}{(1-\alpha)}$$

$$J^{\gamma_L}(2) = -\frac{\alpha}{(1-\alpha)}$$

π_L is greedy w.r.t J if and only if

$$J(1) \geq J(2) + \frac{1}{\alpha}.$$

$$0 + \alpha J(1) \geq -1 + \alpha J(2)$$

Thus

$$\begin{cases} J(1) \geq J(2) + \frac{1}{\alpha} \\ J(0) \geq J(2) - \frac{1}{\alpha} \end{cases} \Rightarrow \boxed{J(1) \geq J(2) + \frac{1}{\alpha}}$$

If true then
 π is the optimal policy.

Consider J_{π_L}

π is $J_{\pi_L} \in R_{\pi_L}$ if Yes then it is optimal.

$$\text{i.e. } 1 \leq \frac{-1}{1-\alpha} \geq \frac{-\alpha}{1-\alpha} + \frac{1}{\alpha}$$

$$\frac{-\alpha^2 + 1 - \alpha}{\alpha(1-\alpha)} = \frac{(1-\alpha)(-\alpha^2 + 1 - \alpha)}{\alpha(1-\alpha)} = \frac{(1-\alpha)(1+\alpha)}{\alpha(1-\alpha)}$$

$$\frac{-1}{1-\alpha} \neq 1 + \alpha \quad \text{Not true.}$$

No optimal policy in this case.

$$\frac{-1}{1-\alpha} \geq \frac{-\alpha^2 + 1 - \alpha}{\alpha(1-\alpha)}$$

$$-\alpha(1-\alpha) \geq (1-\alpha)(-\alpha^2 + 1 - \alpha)$$

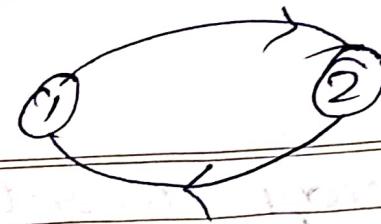
$$\text{or } -\alpha^2 + \alpha^2 \geq -\alpha^2 + \alpha^3 + (1-\alpha) - \alpha + \alpha^2$$
$$(1-\alpha)(1+\alpha) \leq \alpha(1-\alpha)(1+\alpha)$$

If $\alpha=1$ then only it is true.

V_{π_L}

for car

consider π_g



$$J^{\pi_g(1)} = 0; J^{\pi_g(2)} = 0.$$

for π_g to be greedy w.r.t J

$$J(1) \geq J(2) - \frac{1}{2}$$

$$\& J(2) \geq J(1) - \frac{1}{2}$$

$$\begin{aligned} 0 &\geq -\frac{1}{2} \\ 0 &\geq -\frac{1}{2} \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{Both true.}$$

∴ π_g is the optimal policy.

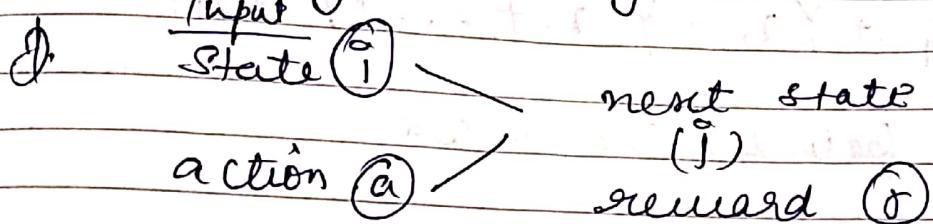
$$J^{\pi_g} \in R_{\pi_g}$$

$$J^{\pi_g} = J^*$$

$$\textcircled{2} \quad \pi_g = \pi^*$$

Chapter 5 : Monte Carlo Methods

- we do not know system model.
- we don't know transition probabilities $P_{ij}(a)$
- we don't know reward function $\delta(i, a)$
- we have a generative or simulation model of the system.



Goal : learn value function from the sample rewards.

Monte Carlo Predictions

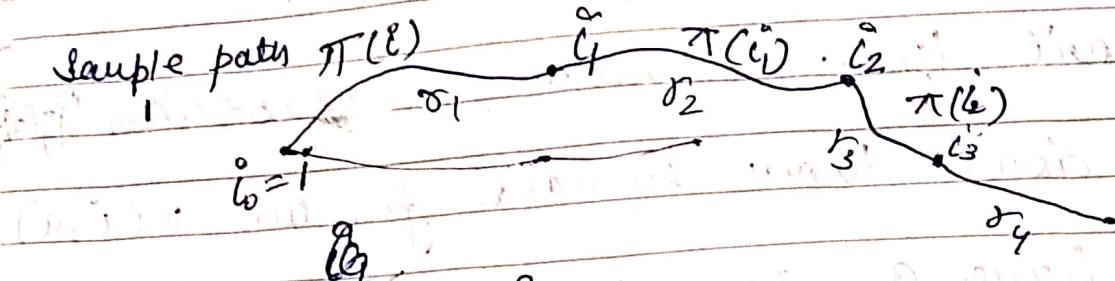
This is estimating value of a given policy π

$$V^\pi(i) = E \left[\sum_{k=0}^{\infty} \gamma^k \delta(i_n, a_n, i_{n+1}) \right]_{i_0=i, a_0=a}$$

$$0 < \gamma < 1$$

① The generative model is giving the new state j according to a distribution a .

$$i_t \sim P_C \cdot \{i, \pi(i)\}$$



Rewards are Random.

$$\hat{V}_1^{\pi(i)} = r_1 + \gamma r_2 + \gamma^2 r_3 + \dots + \gamma^{N-1} r_N$$

for sample path 2 : Sample cost :-

$$\hat{V}_2^{\pi}(i) = \dots$$

$$\hat{V}_M^{\pi}(i)$$

$$\frac{1}{M} \sum_{R=1}^M \hat{V}_R^{\pi}(i)$$

- taking avg.
of all sample cost

Two types of Methods in Monte Carlo

① First visit methods.

$$i = i_0 \rightarrow i_1 \rightarrow i_2 \rightarrow i_3 \dots$$

i_1, i_2 can be same states.

states



take: Single estimate \rightarrow time.

take: Single estimate from a trajectory.

(2) Every visit methods.

$\hat{C} = \hat{C}_1$ (multiple estimates from same trajectory)

First visit of MC Prediction for estimating V_{π}

input: a policy π to be evaluated.

initialize:

$V(s)$: EIR, availability, π , $s \in S$

Returns (s) \leftarrow an empty list, π , $s \in S$

loop forever (for each episode):

Generate an episode following π :

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

$G \leftarrow 0$

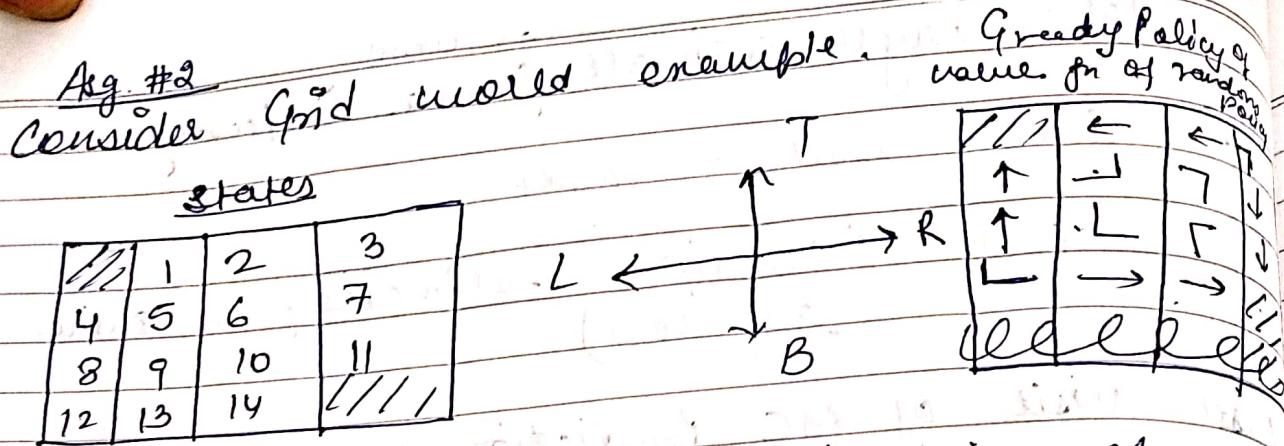
loop for each step of episode, $t = T-1, T-2, \dots, 0$

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

Unless s_t appears in s_0, s_1, \dots, s_{t-1}

Append G to return (s_t)

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



Recall the example from previous class where the random policy was evaluated.

Consider same reward structure as before.

Transition Probabilities

$$p(0|1, \leftarrow) = 0.7$$

$$p(5|1, \leftarrow) = 0.1$$

$$p(2|1, \leftarrow) = 0.1$$

$$p(1|1, \leftarrow) = 0.1$$

going to 1 from

Transition Probabilities

$$p(1|5, \uparrow) = 0.4$$

$$p(4|5, \uparrow) = 0.4$$

$$p(9|5, \uparrow) = 0.1$$

$$p(6|5, \uparrow) = 0.1$$

Apply Monte-Carlo first visit method averaged over 20 iterations to estimate.

$$V_T(s) \leftarrow V(s \in \{1, \dots, 14\})$$

Algorithm P Pick start state randomly on each trajectory.