

Lecture - 04

Dynamic Programming :-

access to the dynamics of the world

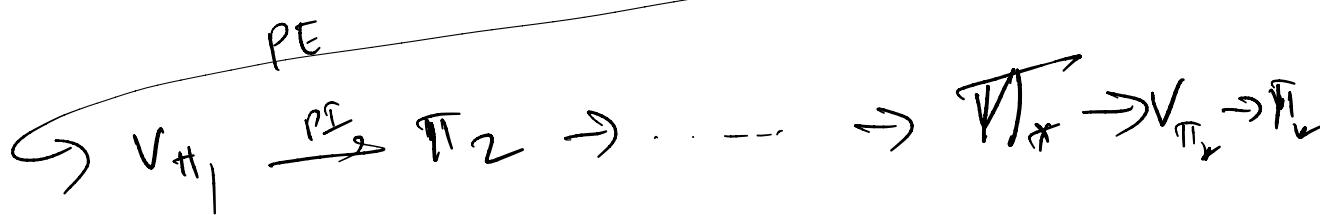
$$P(s', r | s, a)$$

Policy Iteration

Value Iteration

Policy evaluation: Given a fixed policy, find v_{π}

Policy improvement: Given v_{π} , find π'
 $\pi' \geq \pi$



Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

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

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

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$$old\text{-action} \leftarrow \pi(s)$$

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

If $old\text{-action} \neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation
Initialize $V(s)$, for all $s \in S^+$, arbitrarily except that $V(\text{terminal}) = 0$

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 $\Delta < \theta$

PE

Output a deterministic policy, $\pi \approx \pi_*$, such that

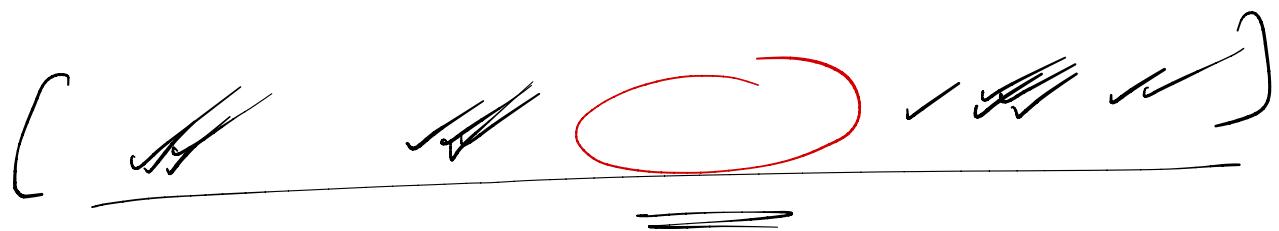
$$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$$

DT

Backgammon

10^{20} states.

Asynchronous DP :-



Using agent's experience in the MDP
to update values

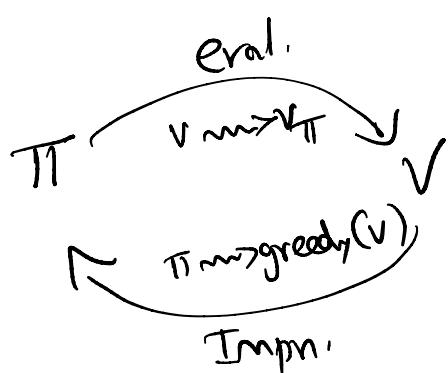
Agent's experience \rightarrow used to decide what state to update.

Latest value & Policy info from DP \rightarrow can guide the agent's decision making.

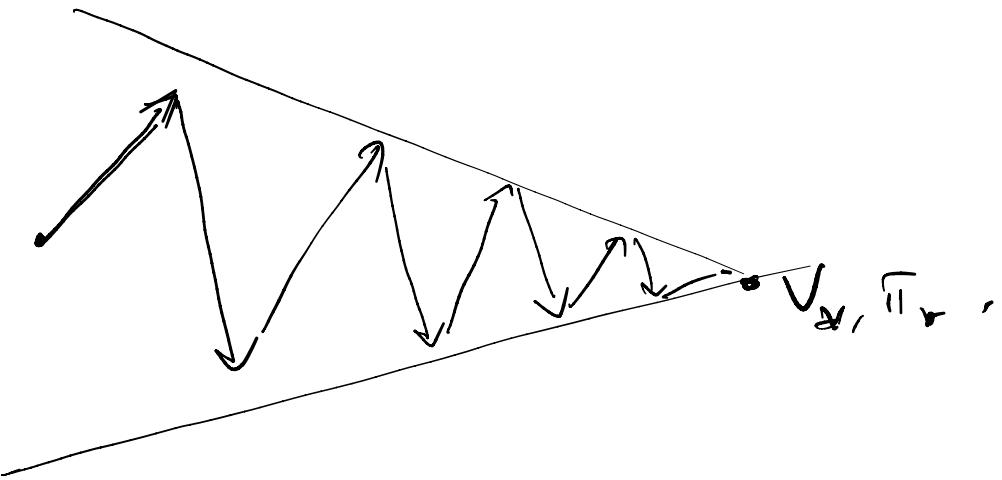
Generalized Policy Iteration (GPI):

\rightarrow P.E (make the value fn. consistent with current policy)

\rightarrow P.I (make policy greedy wrt curr. val. fn.)



$$\pi_{\text{old}} \rightarrow v_{\text{new}}$$



Efficiency of DP:-

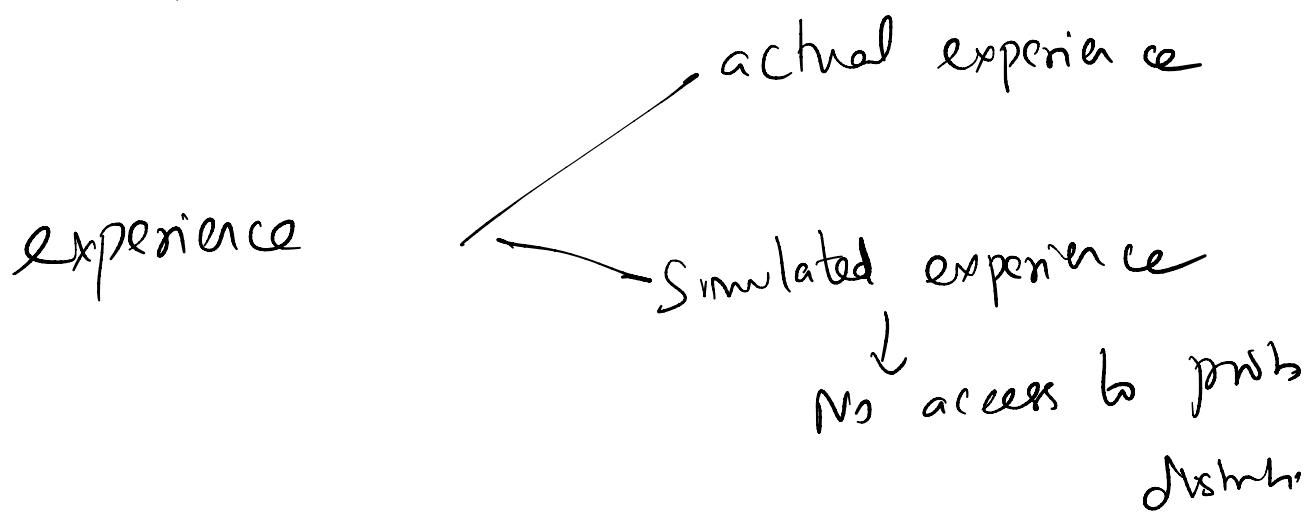
$n = \# \text{ of states}$

$k = \# \text{ of actions}$

k^n det. pol. -

Monte Carlo Methods -

Can we directly learn the value fn.
from experience ?



Monte Carlo methods:-

— Solving RL based on averaging
Sample returns

Episodic setting :-

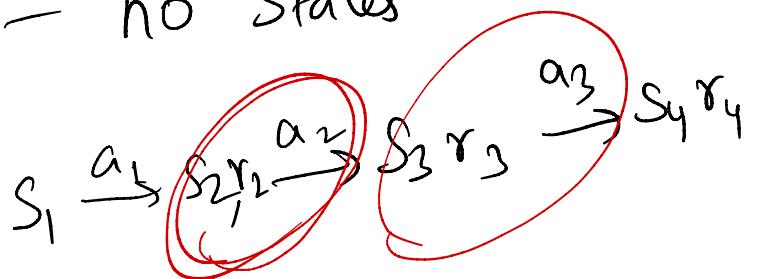
Similar to reward estimation in bandits.

bandits

MC - for RL

- rewards.

- no States



- returns.

- Many States

intw related bandit prob.

$$q(S, a)$$

non-stationary target

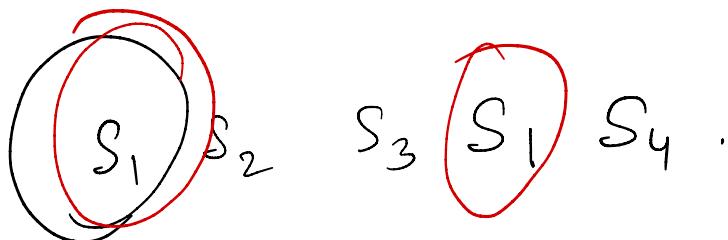
Monte Carlo Prediction:-

Given π , estimate $V_\pi(s)$

$V_\pi(s)$ — Value of a state 's'. under policy π , given a set of episodes obtained by following π and passing through s .

Each occurrence of a state 's' in an episode \rightarrow "visit" to's'

First visit MC \rightarrow estimates $V_\pi(s)$ as the avg of returns following first visit to's'.



Every visit MC \rightarrow average return following all visits to s .

First-visit MC prediction, for estimating $V \approx v_\pi$

Input: a policy π to be evaluated

Initialize:

$V(s) \in \mathbb{R}$, arbitrarily, for all $s \in \mathcal{S}$

$Returns(s) \leftarrow$ an empty list, for all $s \in \mathcal{S}$

Loop forever (for each episode):

Generate an episode following π : $S_0, A_0, R_1, S_1, A_1, R_2, \dots, S_{T-1}, A_{T-1}, R_T$

$G \leftarrow 0$

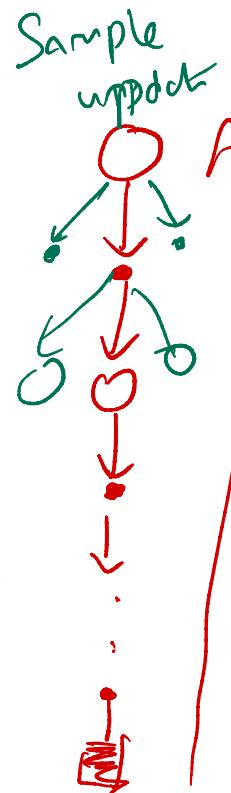
Loop for each step of episode, $t = T-1, T-2, \dots, 0$:

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

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :

Append G to $Returns(S_t)$

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



$$G_t = R_t + \gamma G_{t+1}$$

Rewards

1 2 3 4

return

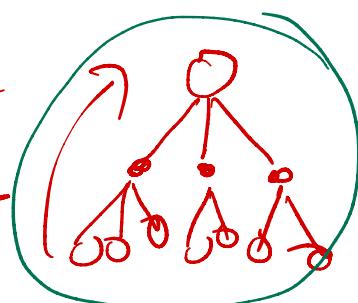
18 14 12 9

=

=

5

5



$\gamma = 1$

DP - expected update

MC - Sample update

Note 1: DP \rightarrow bootstrapping.

MC \rightarrow estimate of each state
are independent.

Note 2! If you want to estimate value
of one state

MC estimation of action values:-

$$\pi(a|s) = \arg \max_a q_{\pi}(s, a)$$

$q_{\pi}(s, a)$ — expected return when
starting in state 's' and
taking action 'a' and thereafter
following policy π .

"Exploring Starts" assumption?

the episodes start in a state-action pair and that every (s, a) has non-zero prob. of selected as start.

Monte Carlo Control :-

$$\pi_0 \xrightarrow{E} q_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} q_{\pi_1}$$
$$\dots \xrightarrow{E} \pi_r \xrightarrow{E} q_{\pi_r}$$

Policy improvement :

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

$$q_{\pi_K}(s, \pi_{k+1}(s)) = q_{\pi_K}(s, \underset{a}{\operatorname{argmax}} q_{\pi_k}(s, a))$$
$$= \max_a q_{\pi_k}(s, a)$$
$$\geq q_{\pi_K}(s, \pi_k(s)) \geq v_{\pi_k}(s)$$

- 2 strong assumptions.

- ① Episodes have exploring starts.
- ② Policy ev. Could be done with infinite # of episodes.

First
visit

Monte Carlo ES (Exploring Starts), for estimating $\pi \approx \pi_*$

Initialize:

$\pi(s) \in \mathcal{A}(s)$ (arbitrarily), for all $s \in \mathcal{S}$

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Loop forever (for each episode):

Choose $S_0 \in \mathcal{S}, A_0 \in \mathcal{A}(S_0)$ randomly such that all pairs have probability > 0

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

$G \leftarrow 0$

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

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

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

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

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

ES.

Monte-Carlo Control without ES :-

Stochastic / soft policies.

$$\pi(a|s) > 0 \quad \forall s \in \mathcal{S} \quad a \in \mathcal{A}(s)$$

~~G-greedy Policy.~~

non-greedy action $A := \frac{\epsilon}{|\mathcal{A}(s)|}$

greedy action : $1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$

Esoft Policies :-

$$\pi(a|s) \geq \frac{\epsilon}{|A(s)|} \quad \text{if } s, a, \\ \epsilon > 0.$$

first-visit MC control (for ε -soft policies), estimates $\pi \approx \pi_*$

Algorithm parameter: small $\varepsilon > 0$

Initialize:

$\pi \leftarrow$ an arbitrary ε -soft policy

$Q(s, a) \in \mathbb{R}$ (arbitrarily), for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

$Returns(s, a) \leftarrow$ empty list, for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$

Repeat forever (for each episode):

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

$G \leftarrow 0$

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

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

Unless the pair S_t, A_t appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$:

Append G to $Returns(S_t, A_t)$

$Q(S_t, A_t) \leftarrow$ average($Returns(S_t, A_t)$)

$A^* \leftarrow \arg \max_a Q(S_t, a)$

(with ties broken arbitrarily)

For all $a \in \mathcal{A}(S_t)$:

$$\pi(a|S_t) \leftarrow \begin{cases} 1 - \varepsilon + \varepsilon/|\mathcal{A}(S_t)| & \text{if } a = A^* \\ \varepsilon/|\mathcal{A}(S_t)| & \text{if } a \neq A^* \end{cases}$$

Any ε -greedy policy wrt' π is ε -optimal

is an improvement over any

ε -soft policy π' .

$\pi' = \varepsilon$ -greedy Policy.

$$q_{\pi'}(s, \pi'(s)) = \sum_a \pi'(a|s) \cdot q_{\pi}(s, a)$$

$$= \frac{\varepsilon}{|\mathcal{A}(s)|} \sum_a q_{\pi}(s, a) + (1-\varepsilon) \max_a q_{\pi}(s, a)$$

$$q_{\pi'}(s, \pi'(s)) = \sum_a \pi'(a|s) \cdot q_{\pi}(s, a)$$

$$= \frac{\epsilon}{|A(s)|} \sum_a q_{\pi}(s, a) + (1-\epsilon) \max_a q_{\pi}(s, a)$$

$$\geq \frac{\epsilon}{|A(s)|} \sum_a q_{\pi}(s, a) + (1-\epsilon) \sum_a \underbrace{\frac{\pi(a|s) - \epsilon}{|A(s)|}}_{1-\epsilon} q_{\pi}(s, a)$$

$$\geq \frac{\epsilon}{|A(s)|} \sum_a q_{\pi}(s, a) - \frac{\epsilon}{|A(s)|} q_{\pi}(s, a) + \sum_a \pi(a|s) q_{\pi}(s, a)$$

$$\geq \sum_a \overline{\pi(a|s)} q_{\pi}(s, a) = V_{\pi}(s)$$

$\pi' \geq \pi$

On-policy / Off-policy

"behavior policy" to explore

"forget policy" being learned.

data generated

from "off" the
target policy.

Off-Policy prediction via Importance Sampling

π - target pol. } fixed.

b - behavior policy }

"Coverage" assumption:

$$\pi(a|s) > 0 \Rightarrow b(a|s) > 0$$

π can be deterministic.

$$b \quad V_b(s)$$

$$\underline{V_\pi(s)}$$

Importance Sampling!

Given a starting state s_t ,
 the prob. of subseq. action traj.
 under any policy π is

$$\Pr \left\{ \underbrace{A_t, s_{t+1}, A_{t+1}, \dots, s_T}_{=}, A_{t:T} \sim \pi \right\}$$

$$= \pi(A_t | s_t) P(s_{t+1} | s_t, A_t) \pi(A_{t+1} | s_{t+1}) \dots P(s_T | s_{T-1}, A_{T-1})$$

$$= \prod_{k=t}^{T-1} \pi(A_k | s_k) P(s_{k+1} | s_k, A_k)$$

$$P_{t:T} = \frac{\prod_{k=t}^{T-1} \pi(A_k | s_k) P(s_{k+1} | s_k, A_k)}{\prod_{k=t}^{T-1} \cancel{\pi}(A_k | s_k) P(s_{k+1} | s_k, A_k)}$$

$$P_{t:T-1} = \overbrace{\prod_{k=t}^{T-1} \frac{\pi(A_k | s_k)}{b(A_k | s_k)}}^{\text{Red circled term}}$$

$$\mathbb{E}[G_t \mid S_t = s] = v_b(s)$$

$$\mathbb{E}\left[\underbrace{P_{t:T-1}}_{\text{Red double underline}} G_t \mid S_t = s\right] = v_{\pi}(s)$$

~~V(s)~~ $T(s)$ — set of all time steps
when state 's' is visited

$$V(s) = \frac{\sum_{t \in T(s)} P_{t:T-1} G_t}{|T(s)|}$$

Ordinary
Important
Sampling (OIS)

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

Weighted
Importance
Sampling (WIS)

first visit off-Policy MC with OIS

WIS.

OIS

WLS

~~Get~~

estimate $v_{\pi}^{(s)}$
all the time

estimate V_b
for first interact

No bias

biased

high variance

less variance

Incremental implementation of WIS :-

$$G_1, G_2, \dots, G_{n-1}$$

$$w_1, w_2, \dots, w_{n-1}$$

$$V_n = \frac{\sum_{k=1}^{n-1} w_k G_k}{\sum_{k=1}^{n-1} w_k} \quad n \geq 2.$$

$$V_{n+1} = V_n + \frac{w_n}{c_n} [G_n - V_n]$$

$$\text{and } c_{n+1} = c_n + w_{n+1}$$

$$c_0 = 0.$$

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

Input: an arbitrary target policy π

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

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

$$\underline{C(s, a) \leftarrow 0}$$

Loop forever (for each episode):

$$\underline{b \leftarrow \text{any policy with coverage of } \pi}$$

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 $W \neq 0$:

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

$$\underline{C(S_t, A_t) \leftarrow C(S_t, A_t) + W}$$

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

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

Off-policy MC control, for estimating $\pi \approx \pi_*$

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

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

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

$$\pi(s) \leftarrow \arg \max_a Q(s, a) \quad (\text{with ties broken consistently})$$

Loop forever (for each episode):

$$b \leftarrow \text{any soft policy}$$

Generate an episode using 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, $t = T-1, T-2, \dots, 0$:

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

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

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

$$\pi(S_t) \leftarrow \arg \max_a Q(S_t, a) \quad (\text{with ties broken consistently})$$

If $A_t \neq \pi(S_t)$ then exit inner Loop (proceed to next episode)

$$W \leftarrow W \frac{1}{b(A_t | S_t)}$$

