


ASSIGNMENT - 2

Que - 1)

Arms : $\{1, 2, \dots, 10\}$

a_e : even arm

$$P[a_e] = 2 P[a_o]$$

a_o : odd Arm

You pick arms 10 times with the above policy.

Reward of an arm $\sim N(i, 1)$

↳ arm number.

We have 5 odd arms & 5 even arms

$$\therefore (p + p + p + p) + (2p + 2p + 2p + 2p + 2p) = 1$$

$$15p = 1$$

$$p = 1/15$$

$$P[a_o] = \frac{1}{15} \quad P[a_e] = \frac{2}{15}$$

At any annual time step, let's say you use the policy to pick an arm ' A_t ', & receive R_t .

The Reward R_t is a R.V as policy is stochastic. \therefore we need $E[R_t]$ at every t .

This we can calculate by summing our ' f ' for all 10 arms. where

$$f_x^a = \left(\text{true mean of the arm } 'a' \right) \left(P[\text{choosing the arm } 'a'] \right)$$

at $x = 1$ $f_x = E[R_x]$ basically

$$f_1 = (\alpha^*(1)) P(1) + (\alpha^*(2)) P(2) + \dots (\alpha^*(10)) P(10)$$

$$f_1 = 1p + 2(2p) + 3(p) + \dots 10(2p)$$

$$f_1 = p(1+3+5+7+9) + 2p(2+4+6+8+10)$$

$$f_1 = \frac{25}{15} + \frac{60}{15}$$

$$f_1 = \frac{85}{15} = \frac{15}{3} \approx 5.66$$

Since, our policy & true mean is the same for every time step. \therefore

$$f_1 = f_2 = f_3 = \dots f_{10}$$

$$E[R_1 + R_2 + \dots R_{10}] = E[R_1] + E[R_2] + \dots E[R_{10}]$$

$$E \left[\underbrace{R_1 + R_2 + R_3 + \dots R_{10}}_{G_x} \right] = \underbrace{5.66 + 5.66 + \dots}_{10 \text{ times}}$$

$E[G_x] = 56.66$

Aw

2)

$$A = \{1, 2, 3, \dots, 10\}$$

Dividing set of
all arms into 2
subsets.

$$B = \{1, 2, 4, 5, 7, 9, 10\}$$

$$C = \{3, 6, 8\}$$

$$R_{\pi}(B) = \begin{cases} 0, & P[R_{\pi}(B) = 0] = 0.5 \\ 1, & P[R_{\pi}(B) = 1] = 0.5 \end{cases}$$

$$R_{\pi}(C) = \begin{cases} 0, & P[R_{\pi}(C) = 0] = 0.3 \\ 0.2, & P[R_{\pi}(C) = 0.2] = 0.3 \\ 1, & P[R_{\pi}(C) = 1] = 0.4 \end{cases}$$

Final Policies to minimize expected Reward.

$$E[R_{\pi}(B)] = 0(0.5) + 1(0.5) = 0.5$$

$$E[R_{\pi}(C)] = 0(0.3) + 0.2(0.3) + 1(0.4) = 0.46$$

$$E[R_{\pi}(B)] > E[R_{\pi}(C)]$$

∴ As long as our policy ensures that we pick an arm from the subset 'B'. We are good !

Our policies can be :-

1) only pick arms (1, 2, 4)

$$P(1) = \frac{1}{3} \quad P(2) = \frac{1}{3} \quad P(4) = \frac{1}{3}$$

rest zeros.

2) only pick arm 1.

$$P(1) = 1, \text{ rest zeros.}$$

3) $P[\text{picking } 1^{\text{st}} \text{ +ve multiple of 3}] = 0$

$P[\text{picking } 2^{\text{nd}} \text{ +ve multiple of 3}] = 0$

$P[\text{picking } 4^{\text{th}} \text{ +ve multiple of 2}] = 0$

Rest equally distributed $\frac{1}{7}$ each.

4) $P[9] = \frac{1}{2} \quad P[7] = \frac{1}{4} \quad P[1] = \frac{1}{4}$

rest zeros

5) $P[9] = 1, \text{ rest zeros}$

6) $P[9] = 0.8 \quad P[1] = 0.2, \text{ rest zeros}$

3)

Explore only non greedy choices.

$A = 3 \text{ arms} : \{a_1, a_2, a_3\}$

$$R_t(A) = \begin{cases} 0, & p = 1/2 \\ 1, & p = 1/2 \end{cases}$$

$Q_t(a)$ is the initial estimate (non zero)

$$Q_t(a) = \{q_1, q_2, q_3\}$$

\uparrow \uparrow \uparrow
 $Q_t(a_1)$ $Q_t(a_2)$ $Q_t(a_3)$

$q_1 + q_2 + q_3 \neq 0$

Assuming

$q_3 > q_2 > q_1$

Explain a sequence of $Q_t(a)$, A_t , R_t for $t = 1, 2, 3, 4, 5, 6$.

Explore at odd times.

Exploit at even times.

Ans

The time sequence is following.

We start with initial Q : $Q_t(A)$.

On the basis of this and our policy we make our Action A_t & get a reward R_t at $t=1$. We will then calculate Q_2 for the arms & so on

$t = 1$, we explore (As per our policy)

As per our assumption $q_3 > q_2 > q_1$

∴ We can't choose a_3 as we are exploring & we can't pick greedy Action while exploring in this question.

$$\therefore P[A_1 = a_1] = P[A_1 = a_2] = \frac{1}{2}$$

$$E[R_1] = \sum_r r P[r]$$

$$E[R_1] = 0(P[R_1 = 0]) + 1(P[R_1 = 1])$$

$$E[R_1] = 0(P[A_1 = a_1] P[R_1 = 0 | A_1 = a_1] + P[A_1 = a_2] P[R_1 = 0 | A_1 = a_2])$$

$$+ 1(P[A_1 = a_1] P[R_1 = 1 | A_1 = a_1] + P[A_1 = a_2] P[R_1 = 1 | A_1 = a_2])$$

$$E[R_1] = 0 + 1\left(\left(\frac{1}{2} \times \frac{1}{2}\right) + \left(\frac{1}{2} \times \frac{1}{2}\right)\right)$$

$$E[R_1] = \frac{1}{2}$$

This is highly obvious as only rewards possible are 0 or 1 & both have equal p.

However, to continue we need to assume some stuff.

Let's say we exploited and picked arm 2 ' a_2 '. Reward we got was τ .

$$A_1 = a_2 \quad R_1 = 1$$

(0, 1, 0)

$$\therefore Q_2(a_1) = \varphi_1$$

$$Q_2(a_2) = \varphi_2 + 1$$

$$Q_2(a_3) = \varphi_3$$

$t = 2$

Now let's say $\varphi_2 + 1$ exceeded φ_3 .

$\therefore a_2$ is now the greedy choice.

We exploit at even times.

$$A_2 = a_2$$

Say, we get $R_2 = 0$ this time.

$E[R_2] = \frac{1}{2}$ still BTW. Even if we don't the if $Q_2(a_2) > Q_2(a_3)$, still $E[R_2] = 1/2$ as there are just 2

rewards $\{0, 1\}$ & both have equal chance of being picked.

But sticking to our assumption

$$Q_2(a_2) > Q_2(a_3)$$

$(0, 2, 0)$

$$\begin{aligned} \therefore A_2 &= a_2 \\ R_2 &= 0 \end{aligned} \quad (\text{let's say})$$

$$Q_3(a_1) = q_1$$

$$Q_3(a_2) = \frac{q_2 + 1 + 0}{2}$$

$$Q_3(a_3) = q_3$$

$$\underline{t = 3}$$

Let's say

$Q_3(a_2)$ is the largest now

No more considering
 $E[R_t]$ or
 $E[R_t] = \frac{1}{2}$
 always

We have to explore now. Say we pick arm a_1 by $\phi = 1/2$ and get the reward 1.

$$\begin{aligned} \therefore A_3 &= a_1 \\ R_3 &= 1 \end{aligned}$$

$(1, 2, 0)$

$$Q_4(a_1) = q_1 + 1$$

$$Q_4(a_2) = \frac{q_2 + 1}{2}$$

$$Q_4(a_3) = q_3$$

$t=4$ (exploit)

previous Assumption

initially

$$q_3 > q_2 > q_1$$

but at $t=2$ we said

$$q_2 + 1 > q_3 > q_1$$

Now let's say

$Q_4(a_3)$ is the largest
(doesn't contradict previous assumption)

So, we pick $a_3 \leftarrow$ get $R_4 = 1$.

$$A_4 = a_3 \quad R_4 = 1$$

(1,2,1)

This is the
first time
we picked
 a_3

$$Q_5(a_1) = (q_1 + 1)/1$$

$$Q_5(a_2) = \frac{q_2 + 1}{2}$$

$$Q_5(a_3) = \frac{q_3 + 1}{1}$$

$t = 5$

Obviously if φ_3 was the largest before. $\varphi_3 + 1$ is the largest now.

We implement & get :-

$$\boxed{\begin{array}{l} A_5 = \alpha_1 \\ R_5 = 1 \end{array}} \quad (\text{second time } \alpha_1 \text{ is picked}) \quad (2, 2, 1)$$

$$Q_6(\alpha_1) = \frac{\varphi_1 + 1 + 0}{2} = \frac{\varphi_1 + 1}{2}$$

$$Q_6(\alpha_1) = \frac{\varphi_1 + 1}{2}$$

$$Q_6(\alpha_2) = \frac{\varphi_2 + 1}{2}$$

$$Q_6(\alpha_3) = \varphi_3 + 1$$

$t = 6$

Obviously $\varphi_3 + 1$ is the longest & we implement.

$$R_6 = 1 \quad A_6 = \alpha_3$$

$$Q_7(\alpha_1) = \frac{\varphi_1 + 1}{2}$$

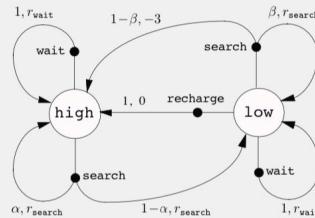
$$Q_7(\alpha_2) = \frac{\varphi_2 + 1}{2}$$

$$Q_7(\alpha_3) = \frac{\varphi_3 + 2}{2}$$

A mobile robot has the job of collecting empty soda cans in an office environment. It has sensors for detecting cans, and an arm and gripper that can pick them up and place them in an onboard bin; it runs on a rechargeable battery. The robot's control system has components for interpreting sensory information, for navigating, and for controlling the arm and gripper. High-level decisions about how to search for cans are made by a reinforcement learning agent based on the current charge level of the battery. To make a simple example, we assume that only two charge levels can be distinguished, comprising a small state set $\mathcal{S} = \{\text{high}, \text{low}\}$. In each state, the agent can decide whether to (1) actively **search** for a can for a certain period of time, (2) remain stationary and **wait** for someone to bring it a can, or (3) head back to its home base to **recharge** its battery. When the energy level is **high**, recharging would always be foolish, so we do not include it in the action set for this state. The action sets are then $\mathcal{A}(\text{high}) = \{\text{search}, \text{wait}\}$ and $\mathcal{A}(\text{low}) = \{\text{search}, \text{wait}, \text{recharge}\}$.

The rewards are zero most of the time, but become positive when the robot secures an empty can, or large and negative if the battery runs all the way down. The best way to find cans is to actively search for them, but this runs down the robot's battery, whereas waiting does not. Whenever the robot is searching, the possibility exists that its battery will become depleted. In this case the robot must shut down and wait to be rescued (producing a low reward). If the energy level is **high**, then a period of active search can always be completed without risk of depleting the battery. A period of searching that begins with a **high** energy level leaves the energy level **high** with probability α and reduces it to **low** with probability $1 - \alpha$. On the other hand, a period of searching undertaken when the energy level is **low** leaves it **low** with probability β and depletes the battery with probability $1 - \beta$. In the latter case, the robot must be rescued, and the battery is then recharged back to **high**. Each can collected by the robot counts as a unit reward, whereas a reward of -3 results whenever the robot has to be rescued. Let r_{search} and r_{wait} , with $r_{\text{search}} > r_{\text{wait}}$, denote the expected number of cans the robot will collect (and hence the expected reward) while searching and while waiting respectively. Finally, suppose that no cans can be collected during a run home for recharging, and that no cans can be collected on a step in which the battery is depleted. This system is then a finite MDP, and we can write down the transition probabilities and the expected rewards, with dynamics as indicated in the table on the left:

s	a	s'	$p(s' s, a)$	$r(s, a, s')$
high	search	high	α	r_{search}
high	search	low	$1 - \alpha$	-3
low	search	high	$1 - \beta$	r_{search}
low	search	low	β	-3
high	wait	high	1	r_{wait}
high	wait	low	0	-
low	wait	high	0	-
low	wait	low	1	r_{wait}
low	recharge	high	1	0
low	recharge	low	0	-



Note that there is a row in the table for each possible combination of current state, s , action, $a \in \mathcal{A}(s)$, and next state, s' . Some transitions have zero probability of occurring, so no expected reward is specified for them. Shown on the right is another useful way of

summarizing the dynamics of a finite MDP, as a *transition graph*. There are two kinds of nodes: *state nodes* and *action nodes*. There is a state node for each possible state (a large open circle labeled by the name of the state), and an action node for each state-action pair (a small solid circle labeled by the name of the action and connected by a line to the state node). Starting in state s and taking action a moves you along the line from state node s to action node (s, a) . Then the environment responds with a transition to the next state's node via one of the arrows leaving action node (s, a) . Each arrow corresponds to a triple (s, s', a) , where s' is the next state, and we label the arrow with the transition probability, $p(s'|s, a)$, and the expected reward for that transition, $r(s, a, s')$. Note that the transition probabilities labeling the arrows leaving an action node always sum to 1.

4)

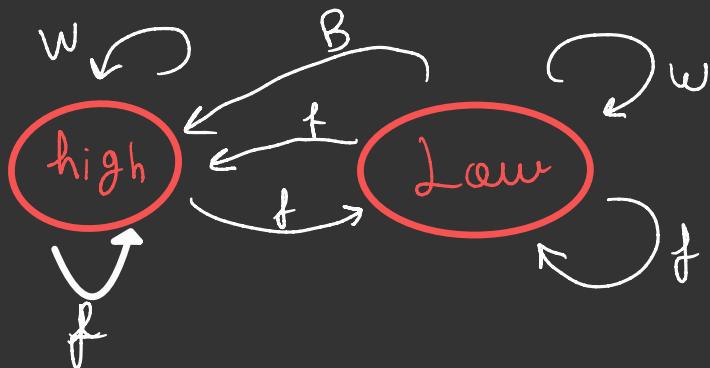
$$S = \{ \text{high}, \text{low} \}$$

$$A = \{ \text{search}, \text{wait}, \text{recharge} \}$$

$$A(\text{high}) = \{ \text{search}, \text{wait} \}$$

$$A(\text{low}) = \{ \text{search}, \text{wait}, \text{recharge} \}$$

If battery depleted then rescue.



Search $\rightarrow f$
Wait $\rightarrow w$
Recharge $\rightarrow B$

- * $P[S_{t+1} = h | h, w] = 1$ * $P[S_{t+1} = h | l, B] = 1$
- * $P[S_{t+1} = h | h, f] = \alpha$
- * $P[S_{t+1} = l | h, f] = 1 - \alpha$
- * $P[S_{t+1} = h | l, f] = 1 - \beta$ (Battery depletes, so it gets rescued
 \therefore high energy next)
- * $P[S_{t+1} = l | l, f] = \beta$
- * $P[S_{t+1} = l | l, w] = 1$

	$S_t = s$	$S_{t+1} = s'$	$A_t = a$ $a \in \{f, w, B\}$	$P[s' s, a]$	$R_t = \sigma_r$ read the prob to unknown its set.
1)	H	H	f	α	σ_f
2)	H	L	f	$1 - \alpha$	σ_f
3)	H	H	w	1	σ_w
4)	H	L	w	0	-
5)	L	L	f	β	σ_f
6)	L	L	w	1	σ_w
7)	L	H	f	$1 - \beta$	-3
8)	L	H	w	0	-
9)	L	H	B	1	0
10)	L	L	B	0	-

reward for entry 7 is met by -3,
because it gives no coins can be collected
if a step taken results in battery depletion

$s \{ s' \} \{ a \} \{ \sigma \} \{ P[s', \sigma | s, a] \}$ $H \quad H \quad w \quad o \quad 1$ $H \quad H \quad f \quad \sigma_f \quad \alpha$ $H \quad H \quad f \quad \sigma_r + \sigma_f \quad \text{○}$ $H \quad H \quad w \quad \begin{matrix} \sigma = \sigma_w \\ \sigma \neq \sigma_w \end{matrix} \quad \begin{matrix} ! \\ \text{○} \end{matrix}$ $H \quad L \quad f \quad \sigma_r + \sigma_f \quad 1 - \alpha$ $H \quad L \quad f \quad \sigma_r + \sigma_f \quad \text{○}$ $H \quad L \quad w \quad \text{any } \sigma_r \quad \text{○}$ $H \quad L \quad \beta \quad \text{any } \sigma_r \quad \text{○}$ $L \quad H \quad f \quad \sigma_r = -3 \quad 1 - \beta$ $L \quad H \quad f \quad \sigma_r \neq -3 \quad \text{○}$ $L \quad H \quad w \quad \text{any } \sigma_r \quad \text{○}$ $L \quad H \quad \beta \quad \begin{matrix} \sigma = 0 \\ \sigma \neq 0 \end{matrix} \quad \text{○}$ $L \quad L \quad \beta \quad \text{any } \sigma_r \quad \text{○}$ $L \quad L \quad w \quad \begin{matrix} \sigma_w \\ \sigma \neq \sigma_w \end{matrix} \quad \begin{matrix} ! \\ \text{○} \end{matrix}$ $L \quad L \quad f \quad \sigma_r = \sigma_f \quad \beta$ $L \quad L \quad f \quad \sigma_r + \sigma_f \quad \text{○}$

8) u_* in terms of q_*

$$u_*(s) = \max_{a \in A} \left(E \left[R_{t+1} + \gamma u_*(s_{t+1}) \mid s_t = s, A_t = a \right] \right)$$

$$u_*(s) = \max_{a \in A} q_{*a}(s)$$

Ans

11) R_{t+1} depends on s_t , a_t . But does R_{t+2} depend on s_t , a_t ?

$$P[R_{t+2} = r'' | s_t = s, a_t = a]$$

$$P[r'' | s, a] = \sum_{a'} \sum_{s'} P[r'', s', a' | s, a]$$

$$P[r'' | s, a] = \underbrace{\sum_{a'} \sum_{s'} \left(P[r'' | s, a, s', a'] \right) P[a', s' | s, a]}_{\downarrow}$$

Here if we know s' , a'
then we don't have to
condition upon
(s, a)

$$P[r'' | s, a] = \sum_{a'} \sum_{s'} P[r'' | s', a'] P[s', a' | s, a]$$

more generally :- \downarrow MDP form

$$P[r'' | s, a] = \sum_{s''} \sum_{a'} \sum_{s'} P[r'', s'' | s', a'] P[s', a' | s, a]$$

$$2) E[R_{t+2} \mid S_t = s, A_t = a] = \sum_{\pi'' \in R(s'')} \pi'' P[\pi'' \mid s, a]$$

$$= \sum_{\pi'' \in R(s'')} \sum_{s''} \sum_{a'} \sum_{s'} P[\pi'', s'' \mid s', a'] P[s', a' \mid s, a]$$

13) $U_D(s) = E[G_t \mid S_t = s]$

$$U_D(s) = E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s]$$

$$U_D(s) = E[R_{t+1} + \gamma(R_{t+2} + \gamma R_{t+3} + \dots) \mid S_t = s]$$

$$U_D(s) = E[R_{t+1} + \gamma G_{t+1} \mid S_t = s]$$

$$U_D(s) = E[R_{t+1} \mid S_t = s] + \gamma E[G_{t+1} \mid S_t = s]$$

double iterated expectations

$$U_D(s) = E[R_{t+1} \mid S_t = s] + \gamma E[U_D(s_{t+1}) \mid S_t = s]$$

$$U_D(s) = E[R_{t+1} + \gamma U_D(s_{t+1}) \mid S_t = s]$$

$$U_D(s) = \sum_{\pi} \sum_{s'} (\pi + \gamma U_D(s')) P[R_{t+1} = \pi, S_{t+1} = s' \mid S_t = s]$$

$$U_{\pi}(s') = \sum_{a} \sum_{r} \sum_{s'} (r + \gamma U_{\pi}(s')) P[r, s' | s, a]$$

$$U_{\pi}(s') = \sum_{a} \sum_{r} \sum_{s'} \left[(\pi(a|s)) (r + \gamma U_{\pi}(s')) \left(P[s', a | s, a] \right) \right]$$

$$U_{\pi}(s') = \sum_{a} \pi(a|s) \sum_{r} \sum_{s'} (r + \gamma U_{\pi}(s')) P[r, s' | s, a]$$

$$14) \quad R_1 = 2 \quad R_2 = -1 \quad R_3 = 10$$

$$R_4 = -3 \quad \gamma = 0.5$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots$$

$$\therefore G_4 = G_5 = G_6 = G_7 = \dots = 0$$

* $G_3 = R_4 + \gamma R_5 + \dots \rightarrow 0$

$$G_3 = -3$$

* $G_2 = R_3 + \gamma R_4 = 10 + \frac{1}{2}(-3)$

$$G_2 = 8.5$$

Alternatively

$$G_2 = R_3 + \gamma G_3$$

$$G_2 = 10 + \frac{1}{2}(-3) = 8.5$$

* $G_1 = R_2 + \gamma R_3 + \gamma^2 R_4$

$$G_1 = -1 + \frac{10}{2} + \left(-\frac{3}{4}\right) = -\frac{4+20-3}{4}$$

$$G_1 = 3.25$$

Alternatively

$$G_1 = R_2 + \gamma G_3$$

$$G_1 = -1 + \frac{8.5}{2} = 3.25$$

$$* G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4$$

$$\text{or } G_0 = 2 + -\frac{1}{2} + \frac{10}{4} + -\frac{3}{8} = 3.625$$

$$G_0 = R_1 + \gamma G_1$$

$$G_0 = 2 + \frac{3.625}{2} = \frac{7.25}{2}$$

$$G_0 = 3.625$$

$$* R_x = C + \gamma C$$

$$\therefore G_x = C + \gamma C + \gamma^2 C + \gamma^3 C + \dots$$

$$G_x = C(1 + \gamma + \gamma^2 + \gamma^3 + \dots)$$

$\gamma < 1 \Rightarrow \infty$ gp which converges

$$G_x = C \left(\frac{1}{1-\gamma} \right) = \frac{C}{1-\gamma}$$

$$G_x = \frac{C}{1-\gamma}$$

Sum of infinite gp with $\gamma < 1 : \frac{a}{1-\gamma}$

15) OPTIMAL POLICY FROM OAVF

Given : $q_{\pi^*}(s, a)$ & (s, a) pairs
 $s \in S, a \in A(s)$

$$\pi^*(s) = \underset{a \in A(s)}{\operatorname{argmax}} (q_{\pi^*}(s, a))$$

From $\pi^* \rightarrow q_{\pi^*}$ or v_{π^*}

write the Bellman equation

Corresponding to π^* & Solve them for both
 v_{π^*} & q_{π^*} . Later in course

15) OPTIMAL POLICY FROM OSVF

Given : $v_{\pi^*}(s, a)$ & $s \in S$, Find $\pi^*(s)$

Argmax $q_{\pi^*}(s, a) = \underset{a \in A(s)}{\operatorname{argmax}} E \left[R_{t+1} + \gamma v_{\pi^*}(s_{t+1}) \mid s_t = s, A_t = a \right]$

Argmax $q_{\pi^*}(s, a) = \underset{a \in A(s)}{\operatorname{argmax}} \sum_{s'} \sum_{a'} (r + \gamma v_{\pi^*}(s')) P(s', a' | s, a)$
Calculate this sum for all $a \in A(s)$ & choose a which maximizes it.