

Date:

Subject:

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

discount $\rightarrow \gamma \in [0, 1]$

myopic

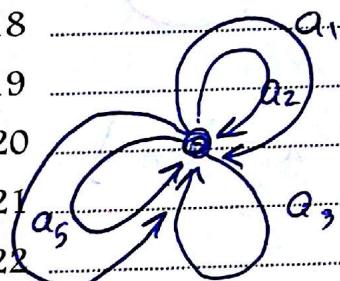
far-sighted

optimal policy

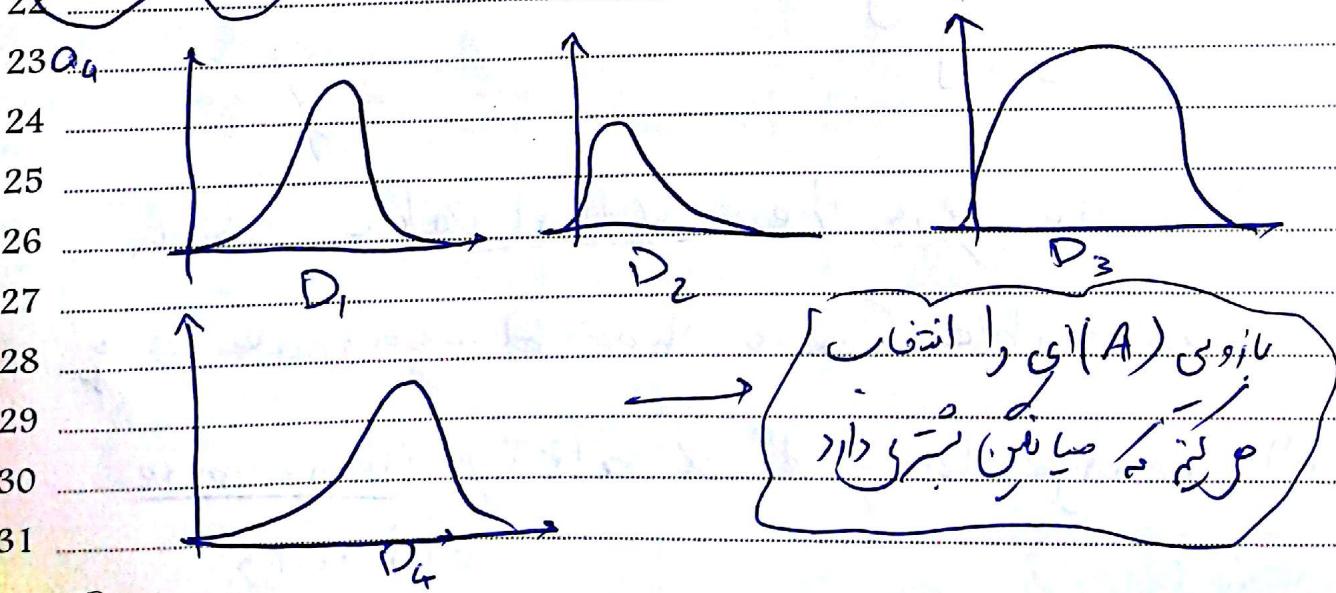
$$\pi_0 > \pi'_0 \text{ if } V_{\pi_0}(s) \geq V_{\pi'_0}(s), \forall s$$

$V_{\pi_0}(s)$ \rightarrow Value of a policy state (s) under π_0 Policy

Single state problem



\rightarrow with every action the state doesn't change



Date: Action Value

Subject:

$$q(a) = R_s^a = E \{ R_{t+1} | S_t = s, a_t = a \}$$

Action value estimation at time t

$$Q = \frac{\text{sum of all rewards}}{\text{sum of all actions}}$$

Batch mode

$$Q_{n+1} = Q_n + \frac{1}{n} [R_n - Q_n] \Rightarrow \text{online mode}$$

Current time step

Let's say Q_n is the true value
and R_n is the reward

$$\text{The core of RL} \Rightarrow Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

Learning Rate

$[0, 1]$

Note that the q is the actual value which

is not accessible. Q is the st estimate of q

Q will converge to q if α satisfy Robbins monro conditions

Kiran

Date:

Subject:

$$1. \sum_{n=1}^{\infty} \alpha_n = \alpha$$

$$2. \sum_{n=1}^{\infty} \alpha_n^2 < \infty$$

Robbins Monro
condition

$$\sum_{n=1}^{\infty} \frac{\alpha_n}{n} = \alpha$$

$$\sum_{n=1}^{\infty} \frac{1}{n^2} < \infty$$

So $\frac{1}{n}$ can be Learning
rate

Learning

Explorate

Learning

Explorate

10

12

Explorate 18

Weight

Don't Forget Exploring

while Exploiting after being Learned

Regret minimization as a goal

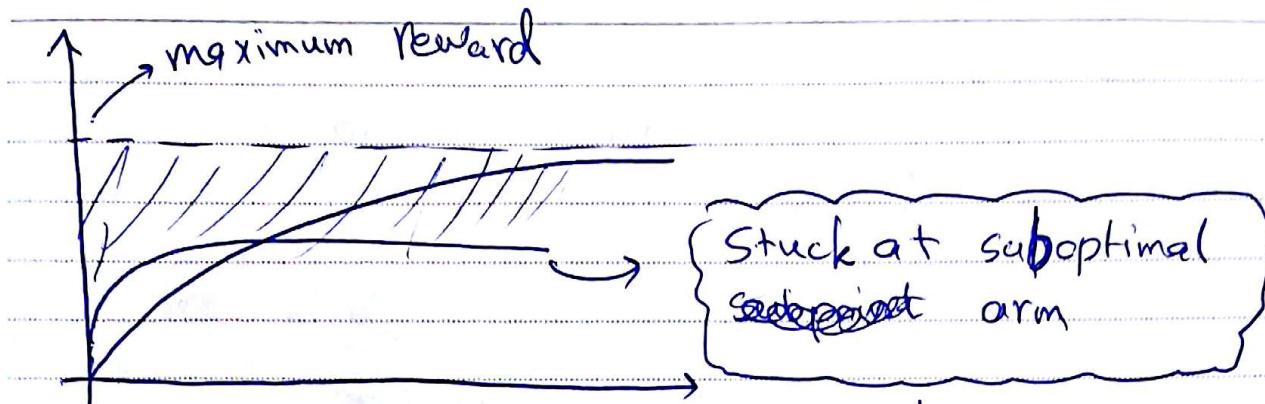
Cumulative Regret

Learning process

Kian

Date:

Subject:



in this case the ~~regard~~ will be regret

Exploration - Exploitation rate $\Rightarrow \beta$

Action 1 → $\frac{E}{m}$

E-Greedy

Action 2 → $\frac{e}{m}$

Action 3 → E/m

الesson لغة

$$\text{Action } i \rightarrow \frac{\Sigma}{m} + (\lambda - \varepsilon)$$

action m $\rightarrow \frac{S}{m}$

لتحت الماء

ابن ابي ابي داود) ابى سعيد البدھری روى عن عبد الرحمن بن معاذ قال: (لَا يَرْجِعُ الْمَاءُ إِلَى مَوْرِدِهِ)

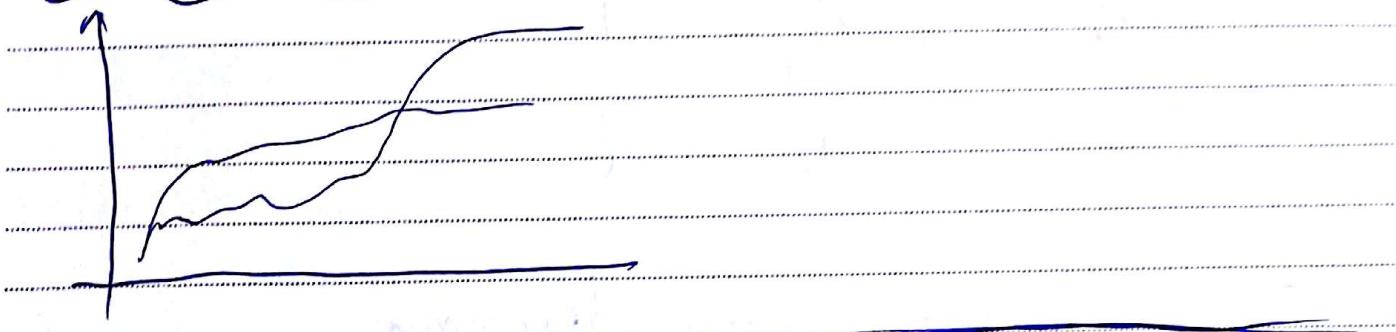
$$\varepsilon \in [0, 1]$$

Kian Greedy $\xrightarrow{\epsilon} \mathcal{E} \in [0, \frac{1}{2}]$ Randomized

Date:

Subject:

1 with initialization of $Q_0(a)$'s ~~value~~ ^{on} ~~of~~
2
3 gain more exploration power
4



11 Boltzman Soft max 8

$$12 \quad P(a_i) = \frac{e^{Q(a_i)/\beta}}{\sum_{j=1}^n e^{Q(a_j)/\beta}}$$

18 Reinforcement Comparison 8

$$21 \quad P_{C,\beta}(a_t) = P_t(a_t) + \beta(r_t - r'_t)$$

$$24 \quad \text{Preference} \quad r'_{t+1} = r'_t + \alpha(r_t - r'_t)$$

$$26 \quad P(a_t) = \frac{P_t(a_t)}{\sum e^{P_t(a_t)}}$$

Kian

Date: _____

Subject: _____

Actor Critic

$$a_t^* = \arg \max Q_{t-1}(a)$$

$$\pi_t(a_t^*) = \pi_{t-1}(a_{t-1}^*) - \beta(1 -$$

$$\pi_t(a_t) = \pi_{t-1}(a_{t-1}) - \beta($$

Exploration-Exploitation Dilemma

More on the Basics

→ Markov Decision Process

→ DQNs

→ Policy Gradients

→ Bellman Equations

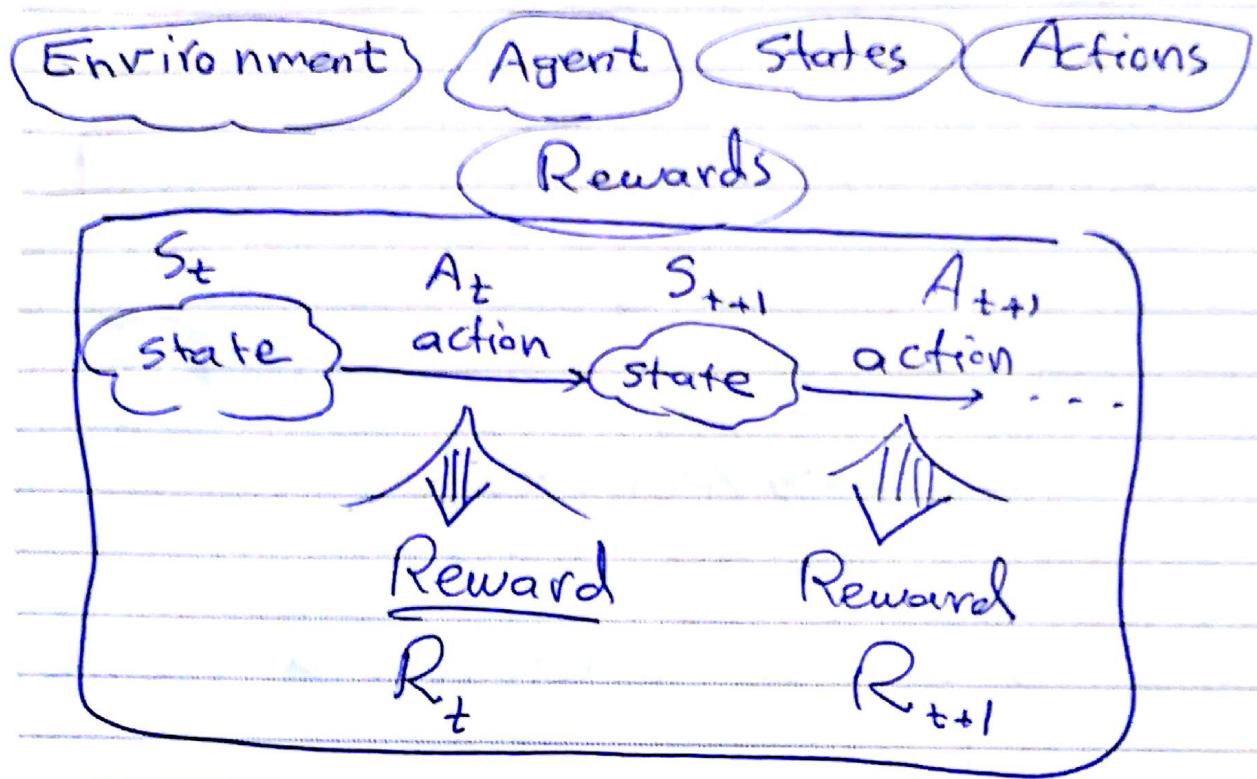
→ Q-Learning

Kian

Date:

Subject:

Markov Decision Making Processes (MDP)



Return

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T \rightarrow \text{Final Time Step}$$

Types of Tasks

→ Episodic Tasks

→ Continuing Tasks

It is the agent's goal to maximize the expected discounted return of rewards

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

Kian

Date:

Subject:

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

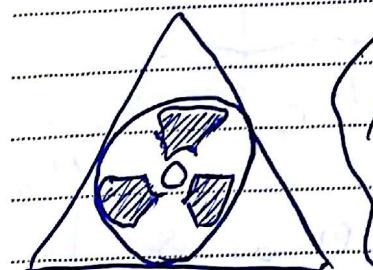
What is the probability that an agent will select a specific action from a specific state?!

Policies

How good is a specific action or a specific state for the agent?!

Value Function

$\pi(a|s)$ = the probability that $A_t = a$ if $S_t = s$



Important

Note that for each state $s \in S$,

π is a probability distribution

over ~~over~~ $a \in A(s)$

Kian

Date:

Subject:

State Value Function $\Rightarrow V_{\pi}(s)$

$$V_{\pi} = E\{G_t | S_t = s\}$$

$$= E\left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right\}$$

Action Value Function $\Rightarrow Q_{\pi}(s, a)$

$$Q_{\pi}(s, a) = E\{G_t | S_t = s, A_t = a\}$$

$$= E\left\{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right\}$$

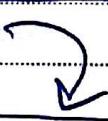
$$Q_{\pi}(s, a) = E\{G_t | S_t = s, A_t = a\}$$



Q -Function

Q -Value

Optimal Policy



$\pi^* > \pi$ if and only if $V_{\pi^*}(s) > V_{\pi}(s)$ for all $s \in S$

Optimal State Value Function

$$V^*(s) = \max_{\pi} V_{\pi}(s)$$

Kian

Date: Subject:

Optimal action value function

$$Q^*(s, a) = \max_{\pi} Q_{\pi}(s, a)$$

Bellman optimality equation for Q^*

$$Q^*(s, a) = E \left\{ R_{t+1} + \gamma \max_{a'} Q^*(s', a') \right\}$$

Q-Learning

(the state which the best

Possible next action can be taken at time t

exploration-exploitation dilemma

How To Solve?!

→ epsilon greedy

→ Actor Critic

→ Boltzmann Softmax

→ Reinforcement Comparison

| S | m1 | m2 | D | B | F | |
|-----|----|----|---|---|---|--|
| s1 | | | | | | |
| s2 | | | | | | |
| s3 | | | | | | |
| s4 | | | | | | |
| s5 | | | | | | |
| s6 | | | | | | |
| s7 | | | | | | |
| s8 | | | | | | |
| s9 | | | | | | |
| s10 | | | | | | |
| s11 | | | | | | |
| s12 | | | | | | |
| s13 | | | | | | |
| s14 | | | | | | |
| s15 | | | | | | |
| s16 | | | | | | |
| s17 | | | | | | |
| s18 | | | | | | |
| s19 | | | | | | |
| s20 | | | | | | |
| s21 | | | | | | |
| s22 | | | | | | |
| s23 | | | | | | |
| s24 | | | | | | |
| s25 | | | | | | |
| s26 | | | | | | |
| s27 | | | | | | |
| s28 | | | | | | |
| s29 | | | | | | |
| s30 | | | | | | |
| s31 | | | | | | |

Q-Value (e.g. $Q(s, a)$)

States

Reward

Date: Subject:

Learning rate

new

$$q(s, a) = \underbrace{(1 - \alpha)q(s, a)}_{\text{old value}} + \alpha(R_{t+1} + \gamma \max_a q(s', a))$$

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Rian

Date:

Subject:

Second Session 8

$$\boxed{G_t = R_{t+1} + \gamma G_{t+1}} \rightarrow \text{Recursive Return}$$

\rightarrow Bellman Equation For $V(s)$

$$\begin{aligned} V(s) &= E \{ G_t | S_{(t)} = s \} \\ &= E \{ R_{t+1} + \gamma V(S_{t+1}) | S_t = s \} \end{aligned}$$

$$V(s) = R_s + \gamma \sum P_{ss'} V(s')$$

$$\underbrace{s \sim S \text{ if } s' \in S}$$

$$\begin{pmatrix} V^{(1)} \\ V^{(2)} \\ \vdots \\ V^{(n)} \end{pmatrix} = \begin{pmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{pmatrix} + \gamma \begin{pmatrix} P_{11} & \cdots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \cdots & P_{nn} \end{pmatrix} \begin{pmatrix} V^{(1)} \\ V^{(2)} \\ \vdots \\ V^{(n)} \end{pmatrix}$$

$$q(s, a) \rightarrow E \{ G_t | S_{(t)} = s, A_{(t)} = a \}$$

$$\star V_{\pi_0}(s) = \sum_{a \in A} \pi(a|s) q_{\pi_0}(s, a)$$

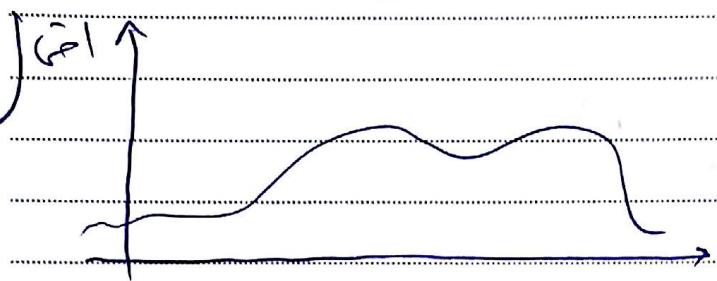
$$P q_{\pi_0}(s, a) = \mathbb{E}_R R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi_0}(s')$$

Kian

میز
ساعی
کنید

Date: Subject:

الآن ندخل في دروس الـ Q-learning



actions in states

$q_{\pi}(s, a)$

$$V_{\pi}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \right)$$

$$R_s^a = \sum_{a' \in A} \pi(a'|s) R_s^{a'}$$

$V_{\pi}(s)$ (جذب الـ π)

$$P_{ss'}^a = \sum_{a' \in A} \pi(a'|s) P_{s's'}^a$$

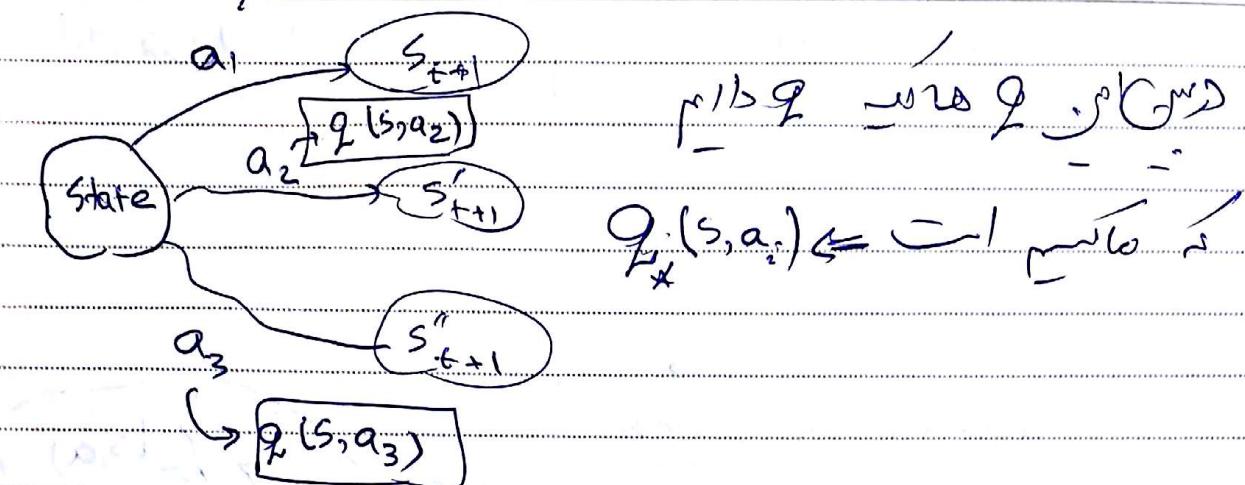
$$q_{\pi}(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s')$$

$$= R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \sum_{a' \in A} \pi(a'|s) q_{\pi}(s', a')$$

$q_{\pi}(s, a)$ (جذب الـ π)

Rian

Date: Subject: $q_*(s, a)$



$$V_*(s) = \max_a q_*(s, a)$$

$$q_*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_*(s')$$

Dynamic Programming 8

Solve complex problems by Breaking them to subproblems

I) iterative policy evaluation

assumption \Rightarrow we have fully defined environment

$$V_{k+1}(s) = \sum_{a \in A} \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$$

Kian

Date: _____
 Policy evaluation

Subject: _____
 Loop

loop:

~~closed~~

Policy evaluation

Policy improvement

$Q_p(s, a)$

$$V(b) = \sum_{a \in A(b)} Q_p(a(b)) \sum_{s', r} P(s', r | b, a) [r + \gamma V_0(s')]$$

~~other~~ $\frac{1}{4}$ $\forall a$

Random policy

$$= 0.25 \times \left[-P(2|b,u) - P(1|b,d) - P(5|b,1) - P(7|b,r) \right]$$

$$= \frac{1}{4} \times (-4) = -1$$

| | | | |
|----|----|----|----|
| 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 |
| 9 | 10 | 11 | 12 |
| 13 | 14 | 15 | 16 |

مثلاً $b = 5$ فـ $V(b) = -1$

فرحان: ريو (ارج) = 1

| | | | |
|-----|-----|-----|-----|
| 0 | -14 | -2 | -22 |
| -14 | -18 | -20 | -20 |
| -20 | -20 | -18 | -14 |
| -22 | -20 | -14 | 0 |

After
Convergence

سید علی بن ابی طالب
جعفر بن ابی طالب

Policy improvement 8

$$q_{\pi^*}(s, \pi^*(s)) = \max_a q_{\pi}(s, a) > q_{\pi}(s, \pi(s)), V(s)$$

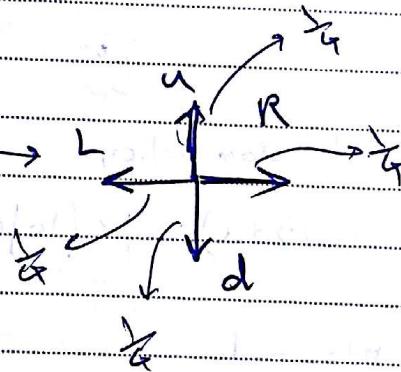
$$\rightarrow \pi' > \pi$$

in a greedy manner we choose the action that

has a better $q(s, a)$ $\pi'(s) = \arg \max_a (q(s, a))$

$$q(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s')$$

| | | | |
|-----|-----|-----|-----|
| 0 | -4 | -2 | -22 |
| -4 | -18 | -20 | -20 |
| -20 | -20 | -18 | -14 |



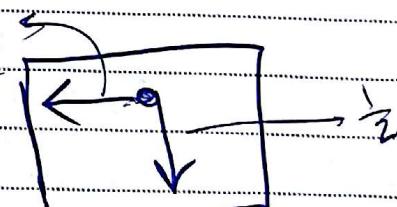
$$q(s, a) = \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$$

$$\rightarrow q(7, R) = 1 \times (-18 - 20) = -38$$

$$q(7, U) = -21$$

$$q(7, D) = -19$$

$$q(7, L) = -19$$



Date: Subject:

1 $q_h(s, a) \rightarrow \text{Optimal Policy } \pi^*$ این درستی که s را با a می‌کند.
2 $q_h(s, \pi(s)) \rightarrow \text{Policy } \pi$ را توصیه می‌کند.

3
4
5
6
7 $q_h(s, \pi(s)) \rightarrow \text{Optimal Policy } \pi^*$ این درستی که s را با a می‌کند.
8 $\pi(s)$ را توصیه می‌کند.

Value Iteration &

loop :

$$\Delta \leftarrow 0$$

for each $s \in S$

$$v \leftarrow V(s)$$

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

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

exit

until $\Delta < \theta$

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

$$\pi(s) = \arg \max_a (q_h(s, a))$$

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

Kian

Date:

Subject:

a note on how the Equations and Formulas of $Q(s, a)$ and $V(s)$ are related to each other

The main Formula
and definition

$$V(s) = E\{G_t \mid S_t = s\}$$

$$Q(s, a) = E\{G_t \mid S_t = s, A_t = a\}$$

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

$$= R_{t+1} + \gamma G_{t+1} \Rightarrow G_t = R_{t+1} + \gamma G_{t+1}$$

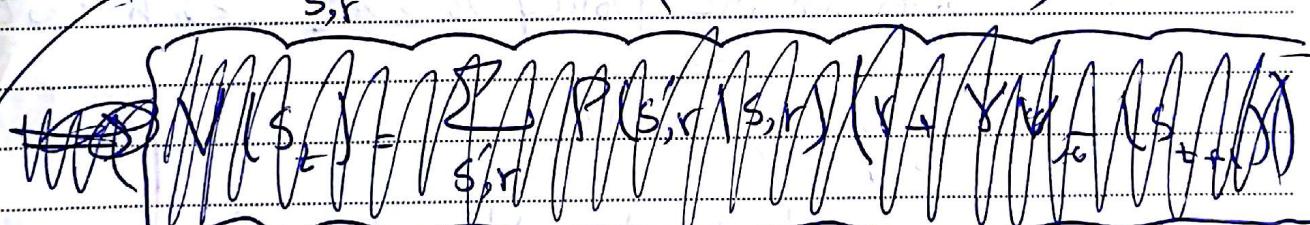
$$V(s_{t+1}) = E\{G_{t+1} \mid S_{t+1} = s\}$$

$$\rightarrow E\{V(s_{t+1})\} = E\{G_{t+1} \mid S_{t+1} = s\}$$

$$\rightarrow V(s_t) = E\{R_{t+1} + \gamma G_{t+1} \mid S_{t+1} = s\}$$

$$= E\{R_{t+1} + \gamma V(s_{t+1}) \mid S_{t+1} = s\}$$

$$= \sum_{s', r} P(s', r | s, a) (r + \gamma V(s_{t+1}))$$



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

Kian

Date:

Subject:

1 Therefore we will have

$$2 Q(\bar{\pi}, s, a) = \sum_{s', r} (s', r | s, a) (R + \gamma V_{\bar{\pi}}(s'))$$

Model Free Learning 8

→ unlike the previous method (dynamic programming)

in this case we don't have the model of env

→ the MDP is unknown

Monte Carlo Learning 8

$$V_{\bar{\pi}}(s) = E[G_{\bar{\pi}} | S_t = s]$$

input: policy $\bar{\pi}$ to be evaluated

initialize g

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

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

Loop forever (For each episode):

Generate an episode following $\bar{\pi}: S_0, A_0, S_1, A_1, S_2, A_2, \dots$

$G \leftarrow \emptyset$

Loop for each step of episode: $t = T_1, T_2, T_3, \dots$

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

append G to $\text{Returns}(S_t)$

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

Rian

Date:

Subject:

AI

$$V(S_t) = V(S_t) + \alpha (G_t - V(S_t))$$

Temporal Difference 8

$$Q(S_t) = E_{\pi}[R_{t+1} + \gamma Q(S_{t+1}) | S_t]$$

$$Q(S_t) = E_{\pi}[R_{t+1} + \gamma Q(S_{t+1}) | S_t]$$

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_t)]$$



$$R_{t+1} + \gamma V(S_{t+1})$$

$$\Rightarrow V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

② Monte Carlo \rightarrow Unbiased - Higher Variance

Temporal Difference \rightarrow Biased - Lower Variance

h

Kian

Date: Subject:

1 - $\hat{V}(s) = \mathbb{E}[R_t + \gamma V(s_{t+1})]$ $\forall s$

2 - $\hat{V}(s) = \mathbb{E}[R_t + \gamma V(s_{t+1})]$ $\forall s$

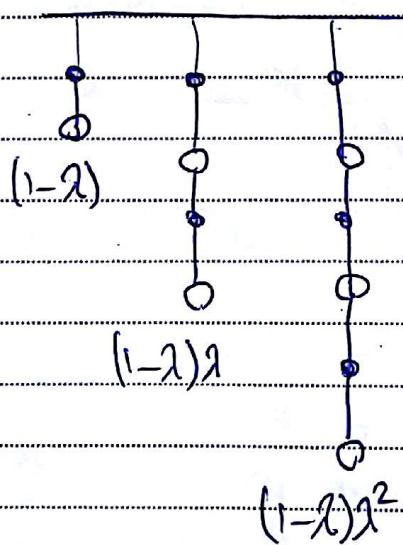
3 - $\boxed{\text{Wavy line}}$

8 n-Step Bootstrapping

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(s_{t+n})$$

13 n-Step return

16 → Averaging n-steps



29 Eligibility Returns

31 Frequency heuristic / Recency heuristic

Kian

Date:

Subject:

Egibility Trace

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Monte Carlo Control & (off policy)

Sudo Code 3

Kian

Date: Subject:

Because the policy is Greedy, so Monte Carlo Control

Will not have any Exploration



عندها لا يوجد أي اكتشاف

for ϵ -greedy

Off Policy Monte Carlo Control

π' → Behavior policy \rightarrow Soft

π' → Estimation policy \rightarrow Greedy

If $\pi(a|s) > 0 \Rightarrow \pi'(a|s) > 0$.

$\pi(a|s)$



action انتخاب

· تنسی داده $\pi(a|s)$

$P_{ss'}$

· احتمال دیدن از این انتخاب $\pi(a|s')$

$a_{\text{جديد}}$

π' → ϵ -greedy



$\pi' \rightarrow$ ϵ -greedy



Kian

Date:

Subject:

How to Compute $\bar{V}_\pi(s) = E_{\pi_t} [G_t | S_t = s]$

$$V_\pi(s) = E_{\pi_t} [G_t | S_t = s] = E_{\pi_t} [\text{?} | S_t = s]$$

$$\left\{ \begin{aligned} E_{x \sim P} [F(x)] &= \sum P(x) F(x) = \sum Q(x) \frac{P(x)}{Q(x)} F(x) \\ &= E_{x \sim Q} \left[\frac{P(x)}{Q(x)} F(x) \right] \end{aligned} \right.$$

$$V_\pi(s) = E_{\pi_t} [G_t | S_t = s] = E_{\pi_t} [W_t G_t | S_t = s]$$

$$= \frac{\sum w_t G_t}{\sum w_t}$$

$$P'(E_{st}) = \prod_{k=t}^{T-1} \pi'(s_k, a_k) P_{s_k, s_{k+1}}^{a_k q_k}$$

 P

$$w_t = \frac{P}{P'} = \frac{\prod_{k=t}^{T-1} \pi'(s_k, a_k)}{\prod_{k=t}^{T-1} \pi'(s_k, q_k)}$$

Kian

Date:

Subject:

$$V_{n+1}(s) = V_n(s) + \frac{\alpha}{C_n} [G_n - V_n(s)]$$

$$G_{n+1} = G_n + W_{n+1}$$

Updating V_n incrementally

Sudo Code :-

initialize For all $s \in S, a \in A(s)$

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

$$\pi(s) \leftarrow \underset{a}{\operatorname{argmax}} Q(s, a)$$

Loop For ever (For each episode) :

π' \leftarrow update using any soft policy update (like ϵ -greedy)

Generate an episode : $\pi' : S_0, A_0; S_1, A_1; S_2, A_2; \dots$

$$G \leftarrow 0$$

$$w \leftarrow 1$$

Loop For each step of episode = $T=1, 2, \dots$

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

$$C(s_t, A_t) \leftarrow C(s_t, A_t) + w$$

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

$$\pi(s_t) \leftarrow \underset{a}{\operatorname{argmax}} (Q(s_t, a))$$

If $A_t \neq \pi(s_t)$ then exit the loop

$$w \leftarrow w \times \frac{1}{\pi'(A_t, s_t)}$$

Date:

Subject:

How To use TD instead of MC in MC Control ??

Sarsa \Rightarrow on policy TD control

Q, V of
Terminal
states
are always
zero

Sudo Code 8

Algorithm parameters: step size $\alpha \in [0, 1]$, small ϵ .

initialize $Q(s, a)$, for all $s \in S^+$, $a \in A(s)$

$$Q(\text{terminal}, \cdot) = 0$$

Loop for each episode:

b) initialize s

choose A from s using policy derived from Q

loop for each step of episode:

i) Take action A , observe R, s' (e.g.: ϵ -greedy)

i) choose A' from s' using policy derived from Q

i) $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma Q(s', A') - Q(s, A)]$

until s is terminal

n-step Sarsa

$$Q(s, A_t) \leftarrow Q_{(n)}(s_t, A_t) + \alpha [Q^{(n)} - Q(s_t, A_t)]$$

$$Q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^n Q(s_{t+n}, A_{t+n})$$

Sarsa & 8

$$Q_t^A = (1-\gamma) \sum_{n=1}^{\infty} \gamma^{n-1} Q_t^{(n)}$$

Kian

Q-Learning & off policy TD Control

$\pi \rightarrow$ Behavioral (Soft.) (will be updated using)
 ϵ -Greedy

$\pi \rightarrow$ Estimation!

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_a Q(S_{t+1}, a) - Q(S_t, A_t)]$$

Sudo Code

Algorithm parameters: step size $\alpha \in (0, 1]$, small $\epsilon > 0$.

initialize $Q(s, a)$, for all $s \in S^+$, $a \in A(s)$,

$Q(\text{terminal}, \cdot) = 0$

Loop For each step of episode:

choose A from S using policy derived from Q (e.g. ϵ -Greedy)

Take action A , observe R, S'

1 $Q(S, A) \leftarrow Q(S, A) + \alpha [R + \gamma \max_a Q(S', a)]$

$- Q(S, A)]$

$S \leftarrow S'$

until S is terminal

π , Q (policy, value function)

$\pi \leftarrow \text{argmax}_a Q(s, a)$

Kian