

(1) Probability of picking an even arm =  $2x$   
u u u u odd arm =  $x$

No. of even arms = 5  
u u odd arms = 5

$\Rightarrow$

$$\textcircled{P} \quad 5(2x) + 5(x) = 1$$

$$\Rightarrow x = \frac{1}{15}$$

R<sub>a</sub>: ~~R~~ Reward drawn upon choosing arm a.  
(RV)

~~for each~~

$$E[R_a | a=i] = i$$

$$E[R_a] = \sum_i R_a=i \cdot P(a=i)$$

$$= \cancel{\frac{1}{15} \times 1} + \cancel{\frac{3}{15} \times 2} + \cancel{\frac{5}{15} \times 3} + \cancel{\frac{7}{15} \times 4} + \cancel{\frac{9}{15} \times 5}$$

$$= \frac{1}{15} [1 + 3 + 5 + 7 + 9] + \frac{2}{15} [2 + 4 + 6 + 8 + 10]$$

$$= \frac{1}{15} [1 + 3 + 5 + 7 + 9 + 4 + 8 + 12 + 16 + 20]$$

$$= \frac{1}{15} [85] = \frac{85}{15} = \frac{17}{3}$$

As each picking of arm is different from one another  
and the distribution from which reward is drawn  
is same.  $\therefore$  Rewards at different instants are  
iid Random Variables.

$$\Rightarrow \cancel{E[R_a]} \quad E\left[\sum_{t=1}^{10} R_a\right] = 10 \sum E[R_a] = \frac{170}{3}$$

②  $S_1 = \{1, 2, 4, 5, 7, 9, 10\}$   
 $S_2 = \{3, 6, 8\}$

$R_a$ : Reward drawn upon choosing arm  $a$ .

$$\begin{aligned} E_1 &= E[R_a | a \in S_1] = P(0) \times 0 + P(1) \times 1 \\ &= 0.5 \times 0 + 0.5 \times 1 \\ &= 0.5 \end{aligned}$$

$$\begin{aligned} E_2 &= E[R_a | a \in S_2] = P(0) \times 0 + P(0.2) \times 0.2 + P(1) \times 1 \\ &= 0.2 \times 0.3 + 1 \times 0.4 \\ &= 0.06 + 0.4 = 0.46 \end{aligned}$$

as  $E_1 > E_2$

optimal action must be set  $S_1$ .

↑ Policy 1.

$$\pi_1(a) = \begin{cases} 0.2 & a=1 \\ 0.8 & a=2 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_2(a) = \begin{cases} 0.3 & a=4 \\ 0.7 & a=5 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_3(a) = \begin{cases} 0.5 & a=9 \\ 0.5 & a=2 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_4(a) = \begin{cases} 0.2 & a=4 \\ 0.2 & a=5 \\ 0.6 & a=10 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_5(a) = \begin{cases} 0.8 & a=9 \\ 0.1 & a=4 \\ 0.1 & a=1 \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_6(a) = \begin{cases} 0.1 & a=7 \\ 0.9 & a=5 \\ 0 & \text{otherwise} \end{cases}$$

③

### Initial estimates

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

$$Q(1) \approx 1$$

$$Q(2) \approx 2$$

$$Q(3) \approx 3$$

$$\epsilon = 0.5$$

$$\{0, 1\}$$

t=1

Explore

$$a=2$$

$$R(a=2) = 0$$

$$Q(1) \approx 1$$

$$Q(2) \approx 1$$

$$Q(3) = 1$$

$$Q(3) = 3$$

$$t=2$$

Exploit

$$a=3$$

~~$$R(3) = 0$$~~

$$Q(1) = 1$$

$$Q(3) = 1.5$$

$$Q(2) = 1$$

$$Q(3) = 1.5$$

$$t=3$$

Explore

$$a=1$$

$$R(1) = 1$$

$$Q(1) \approx 1$$

$$Q(1) = 1$$

$$Q(2) \approx 1$$

$$Q(3) \approx 1.5$$

$$t=4$$

Exploit

$$a=3$$

$$R(3) = 0$$

$$Q(1) \approx 1$$

$$Q(3) \approx 1$$

$$Q(2) \approx 1$$

$$t=5$$

Explore

$$a=2$$

$$R(2) \approx 1$$

$$Q(2) \approx 1$$

$$Q(1) \approx 1$$

$$Q(2) \approx 1$$

$$Q(3) \approx 1$$

# t26 Exploit

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

$$a = 1$$

$$R(a=1) = 0$$

$$Q(1) = 0.667$$

$$Q(1) = 0.667$$

$$Q(2) = 1$$

$$Q(3) = 1$$

④

Current state:  $s$

Next state:  $s'$

action taken:  $a$

reward:  $r$ .

$s$	$a$	$s'$	$r$	$p(s', r   s, a)$
h	s	w	$r_s$	$\alpha$
h	s	l	$r_s$	$1-\alpha$
h	w	h	$r_w$	1
h	w	l	-	0
l	s	h	$r_p$	$1-\beta$
l	s	l	$r_s$	$\beta$
l	w	h	-	0
l	w	l	$r_w$	1
l	re	h	0	1
l	re	l	-	0

States: {high, low}

~~Actions~~ A(high) = {search, wait}

Actions = {search, wait, recharge}

ACRONYMS:

high  $\rightarrow$  h

low  $\rightarrow$  l

search  $\rightarrow$  s

wait  $\rightarrow$  w

recharge  $\rightarrow$  re

reward of search $\rightarrow r_s$ reward of wait $\rightarrow r_w$ reward penalty $\rightarrow r_p$
--

⑥

3.15

$$c_{t+1} = \sum_{k=0}^{\infty} \gamma^k R_{k+t+1}$$

~~$c_{t+1} = \sum_{k=0}^{\infty} \gamma^k R_{k+t+1}$~~

Adding a constant  $c$ .

$$c_{t+1} = \sum_{k=0}^{\infty} \gamma^k (R_{k+t+1} + c)$$

$$c_t = \sum_{k=0}^{\infty} \gamma^k R_{k+t+1} + \sum_{k=0}^{\infty} \gamma^k c \quad ; \quad \gamma < 1$$

$$\hat{c}_t = c_t + \frac{c}{1-\gamma}$$

$$V_{\pi}(s) = \mathbb{E}_{\pi}[c_t | s_t = s]$$

$$\begin{aligned} V_{\pi}(s) &= \mathbb{E}_{\pi}[\hat{c}_t | s_t = s] \\ &= \mathbb{E}_{\pi}[c_t + \frac{c}{1-\gamma} | s_t = s] \\ &\geq V_{\pi}(s) + \frac{c}{1-\gamma} \end{aligned}$$

Value function will change but new policy will stay same  
 when a constant gets added to all the rewards.

3.16

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

for An episodic will end at some finite time T.

~~$$C_{t+T} = \sum_{k=0}^T \gamma^k R_{k+t+1}$$~~

Let  $\gamma = 1$

$$C_{t+T} = \sum_{k=0}^T R_{k+t+1}$$

assuming  $T = 1$

$$C_{t+1} = \sum_{k=0}^T R_{k+t+1}$$

$$\begin{aligned} C'_{t+1} &= \sum_{k=0}^T (R_{k+t+1} + c) \\ &= \sum_{k=0}^T R_{k+t+1} + \sum_{k=0}^T c \end{aligned}$$

$$C'_{t+1} = C_t + cT$$

$$V_\pi(s) = \mathbb{E}_\pi[C_t | s_t = s]$$

$$\begin{aligned} V'_\pi(s) &= \mathbb{E}_\pi[C'_{t+1} | s_t = s] \\ &= \mathbb{E}_\pi[C_t + cT | s_t = s] \end{aligned}$$

$$= \mathbb{E}_\pi[C_t | s_t = s] + c \mathbb{E}[T | s_t = s]$$

Addition of  $c$  a constant in Rewards of an episodic task would result in drastic changes in value function.

As  $T$  is a random variable, its expectation would change from problem to problem & episode to episode.

This might even result in change policy.

$$(8) V_*(s) = \max_a q_*(s, a)$$

(11)  $R_{t+2}$  is conditional on  $s_t, A_t$  if  $s_{t+1}$  is not taken into consideration.

$$\begin{aligned} & P(R_{t+2} = r'' | s_t = s, A_t = a) \\ &= \sum_{s'} P(R_{t+2} = r'', s_{t+1} = s' | s_t = s, A_t = a) \\ &= \sum_{s'} P(R_{t+2} = r'' | s_{t+1} = s', s_t = s, A_t = a) \\ &\quad P(s_{t+1} = s' | s_t = s, A_t = a) \\ &= \sum_{s'} P(R_{t+2} = r'') | s_{t+1} = s' \\ &\quad P(s_{t+1} = s' | s_t = s, A_t = a) \end{aligned}$$

(MDP Assumption)

$$P(R_{t+2} = r'' | s_t = s, A_t = a)$$

$$= \sum_{s''} P(R_{t+2} = r'', s_{t+2} = s'' | s_t = s, A_t = a)$$

$$\begin{aligned} & P(R_{t+2} = r'', s_{t+2} = s'' | s_t = s, A_t = a) = \\ & \sum_{s'} \sum_{a'} P(R_{t+2} = r'', s_{t+2} = s'' | s_{t+1} = s', A_{t+1} = a', s_t = s, A_t = a) \end{aligned}$$

$$\begin{aligned} &= \sum_{s'} \sum_{a'} P(R_{t+2} = r'', s_{t+2} = s'' | s_{t+1} = s', A_{t+1} = a', s_t = s, A_t = a) \\ &\quad \pi(A_{t+1} = a' | s_{t+1} = s') P(s_{t+1} = s' | s_t = s, A_t = a) \end{aligned}$$

$$= \sum_{s'} \sum_{a'} P(r'', s'' | s', a') \pi(a' | s') P(s' | s, a)$$

$$\frac{P(R_{t+2}=1)}{P(R_{t+2}=r'')} | s_t = s, a_t = a = \sum_{s''} \sum_{s'} \sum_{a'} p(r'', s'') | s', a') \pi(a' | s') p(s' | s, a)$$

(12)

$$E[R_{t+2} | s_t = s, a_t = a] =$$

$$\sum_{s''} \sum_{s'} \sum_{a'} p(r'', s'') | s', a') \pi(a' | s') p(s' | s, a) [R_{t+2}]$$

(13)

$$v_\pi(s) \doteq E[a_t | s_t = s]$$

$$= E[R_{t+1} + r_{t+1} | s_t = s]$$

$$= \sum_{s'} \sum_{a_g} [R_{t+1} + r_{t+1}] P(s' | s) P(a_{t+1} = g | s_t = s)$$

$$= \sum_s \sum_g [s + r_g] P(s, g | s)$$

$$= \sum_{s', s, a, g} [s + r_g] P(s, g, s' | s, a)$$

$$= \sum_{s', s, a, g} [s + r_g] P(s, g, s' | s, a) \pi(a | s)$$

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

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

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

(14)

$$R_1 = 2, R_2 = -1, R_3 = 10, R_4 = -3$$

$$\gamma = 0.5$$

$$h_4 = 0$$

$$h_3 = -3 + 0.5 \times 0 = -3$$

$$h_2 = 10 + 0.5(-3) = 8.5$$

$$h_1 = -1 + 0.5 \times 8.5 = 3.25$$

$$h_0 = 2 + 0.5 \times 3.25 = 3.625$$

$$R_{t+k+1} = c \quad (\text{constant reward}).$$

$$\begin{aligned} a_{\infty} &= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1}) \\ &= \sum_{k=0}^{\infty} \gamma^k (c) = c \sum_{k=0}^{\infty} \gamma^k \\ &= \frac{c}{1-\gamma} \end{aligned}$$

(15)

 $v_k(s) \leftarrow \text{given } \forall s \in S$ 

$$\pi_k(a|s) = \begin{cases} 1 & \text{if } a = \arg \max_a \sum_{s', r} p(s', r|s, a) (r + \gamma v_k(s')) \\ 0 & \text{otherwise} \end{cases}$$

$$\pi_k^*(a|s) = \begin{cases} 1, & \text{if } a = \arg \max_a q^*(s, a) \\ 0, & \text{otherwise} \end{cases}$$

⑯ Current state:  $s$

Next State:  $s'$

Reward:  $r$

Action:  $a$

States = {Fresh, Stale}

Action = {query, silent}

~~Actions~~

~~Sense~~

DATE \_\_\_\_\_  
PAGE \_\_\_\_\_

$s$	$a$	$s'$	$r$	$P(s', r   s, a)$
Fresh	query	Fresh	-4	0.9
Fresh	query	Stale	-4	0.1
Fresh	silent	Fresh	4	0.5
Fresh	silent	Stale	4	0.5
Stale	query	Fresh	-8	0.8
Stale	query	Stale	-8	0.2
Stale	silent	Fresh	4	0
Stale	silent	Stale	4	1

$$V_F(\text{Fresh}) \equiv V_F$$

$$V_F(\text{Stale}) \equiv V_S$$

$V_F(\text{Stale}) =$

$$V_F = \frac{1}{2} [0.9(-4 + 0.5V_F) + 0.1(-8 + 0.5V_S)]$$

$$+ \frac{1}{2} [0.5(4 + 0.5V_F) + 0.5(4 + V_S)]$$

$$= \frac{1}{2} [-4 + 0.45V_F + 0.05V_S + 4 + 0.025V_F + 0.25V_F + 0.25V_S]$$

$$= \frac{1}{2} [0.7V_F + 0.3V_S] \Rightarrow 2V_F = 0.7V_F + 0.3V_S$$

$$\Rightarrow 1.3V_F = 0.3V_S$$

$$V_S = 0.5 [0.8(-8 + 0.5V_F) + 0.2(-8 + 0.5V_S) + 1(4 + 0.5V_S)]$$

$$2V_S = -8 + 0.4V_F + 0.1V_S + 0.5V_S$$

$$2V_S = -8 + 0.4V_F + 0.6V_S$$

$$1.4V_S = -8 + 0.4V_F = -8 + 0.092V_S$$

$$V_S = (-8)/(1.367) = -6.1176, V_F = -1.4117$$

(17)

Policy improvement theorem.

$\pi'(s)$  is the improved policy over  $\pi(s)$ .

$\pi'(s)$  is same as  $\pi(s)$ , except  $\pi'(s) = a \neq \pi(s)$

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s, A_t = \pi'(s)]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | s_t = s]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma \mathbb{E}_{\pi'} [R_{t+2} + \gamma V_{\pi}(s_{t+2}) | s_{t+1}, A_{t+1} = \pi'(s_{t+1})] | s_t = s]$$

$$= \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]$$

$$\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots | s_t = s]$$

$$= V_{\pi'}(s)$$

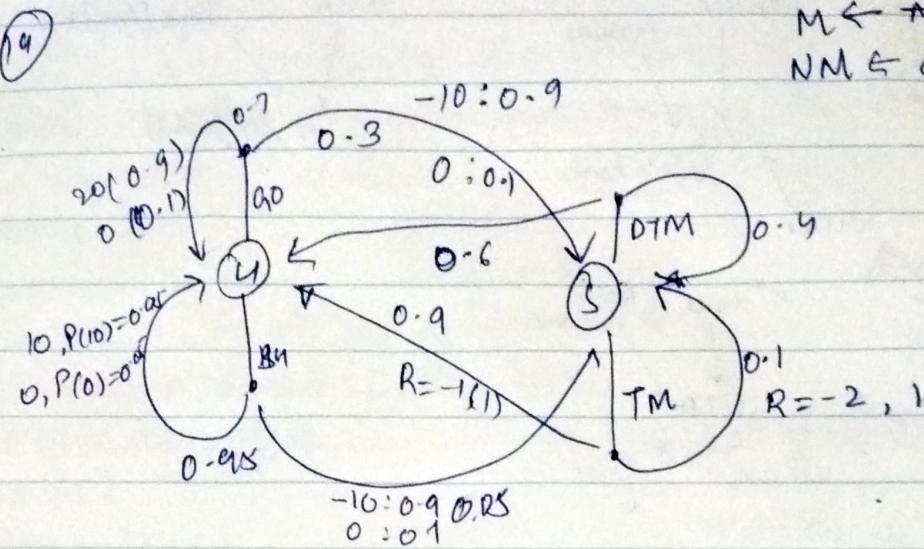
$$\pi'(s) = \operatorname{argmax}_a q_{\pi}(s, a)$$

$$V_{\pi'}(s) = \max_a \mathbb{E}[R_{t+1} + \gamma V_{\pi'}(s_{t+1}) | s_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', s | s, a) (r + \gamma V_{\pi'}(s'))$$

Since if  $\pi'(s)$  is not better than  $\pi(s)$ , then  $\pi'(s)$  must be optimal as it satisfies bellman optimality equation, so  ~~$\pi(s)$~~   $\pi(s)$  will be optimal.

M ← took Medicine  
NM ← didn't take Medicine



S	a	s'	λ	$P(s', s   s, a)$
healthy	Stay home	healthy	10	0.9025
healthy	Stay home	healthy	0	0.0475
healthy	Stay home	sick	-10	0.045
healthy	Stay home	sick	0	0.05
healthy	Go out	healthy	20	0.63
healthy	Go out	healthy	0	0.07
healthy	Go out	sick	-10	0.27
healthy	Go out	sick	0	0.03
sick	to M	healthy	-1	0.9
sick	M	sick	-2	0.1
sick	NM	healthy	0	0.6
sick	NM	<del>sick</del> healthy	0	0.4

