

$$\left\{ \begin{array}{l} \text{پایین} \\ \text{متوسط} \\ \text{بالا} \end{array} \right\} \left\{ \begin{array}{l} 3 + \text{بالا} \\ 3 + \text{متوسط} \\ 0 + \text{پایین} \end{array} \right\} \left\{ \begin{array}{l} -1 \\ -1 \\ -1 \end{array} \right\}$$

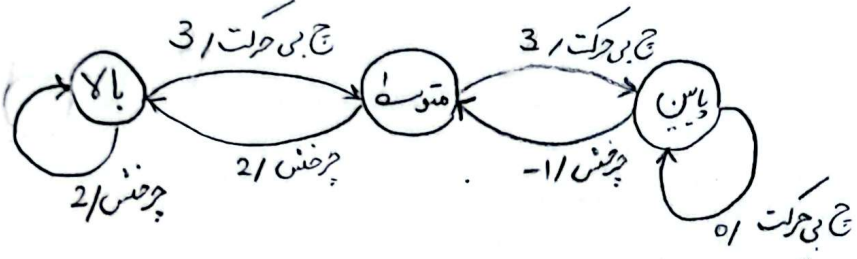
تبدیل هر حرکت

نودها ← عمل قرارگیری بابت  
یال ها ← تفسیر محل قرارگیری

نام یال ←  $act_{max}$  /  $act_{min}$  برآیند یا  $act_{diff}$  و  $act_{diff}$

گران ماشین حالت

صیغ بی حرکت / جرخش



امتیاز اولیه =

جدول گرانها را بر حسب شرایط اولیه و حل می کنیم.  
خانه ها را به عبار حرکت با  $max$  امتیاز همراه با مقدار بعدی است.

محاسبه قرارگیری / محاسبه حرکت بعدی

1 2 3

پایین	جرخش	متوسط , $-1 - 1 = -2$	$LS2 = 1.4$	بالا	$LS3 = 2.52$
	بی جرخش	پایین , $0 + 0 = 0$	$LD2 = 0.4$	پایین	$LD3 = 1.12$
متوسط	جرخش	بالا , $0 + 2 = 2$	$MS2 = 4.4$	بالا	$MS3 = 6.32$
	بی جرخش	پایین , $0 + 3 = 3$	$MD2 = 3.4$	پایین	$MD3 = 4.12$
بالا	جرخش	بالا , $0 + 2 = 2$	$HS2 = 4.4$	بالا	$HS3 = 6.32$
	بی جرخش	متوسط , $0 + 3 = 3$	$HD2 = 5.4$	متوسط	$HD3 = 6.52$

$LS2: \begin{matrix} 3 \rightarrow 2 \\ 0.8 \times 3 - 1 = 1.4 \end{matrix}$   $LD2: \begin{matrix} 2 \rightarrow 2 \\ 0.8 \times 0 + 0 = 0 \end{matrix}$

$MS2: 2 \rightarrow 4 : 0.8 \times 3 + 2 = 4.4$   $MD2: 2 \rightarrow 3 : 0.8 \times 0 + 3 = 3$

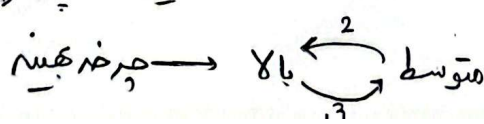
$HS2: 2 \rightarrow 2 : 0.8 \times 3 + 2 = 4.4$   $HD2: 2 \rightarrow 3 : 0.8 \times 3 + 3 = 5.4$

$LS3: 2 \rightarrow 2 : 0.8 \times 4.4 - 1 = 2.52$   $LD3: 2 \rightarrow 2 : 0.8 \times 1.4 + 0 = 1.12$

$MS3: 2 \rightarrow 2 : 0.8 \times 5.4 + 2 = 6.32$   $MD3: 2 \rightarrow 2 : 0.8 \times 1.4 + 3 = 4.12$

$HS3: 2 \rightarrow 2 : 0.8 \times 5.4 + 2 = 6.32$   $HD3: 2 \rightarrow 2 : 0.8 \times 1.4 + 3 = 4.12$

سیاست بهینه این است که اگر خانه پایین هستیم ابتدا به خانه متوسط برویم. سپس جرخه زیر را تکرار کنیم



الف) نادرست. در الگوریتم Q-L اصل سیاست برای این است که exploration و exploitation همواره با هم باشد و با احتمال بزرگتر از صفر هر دو exploration و exploitation با هم یافت شود.

ب) درست. چون مقدار تنبیه برای عامل بسیار زیاد است؛ عامل ناچاراً مجبور است بیشتر میزبانی را بیاموزد و در صورت وجود جواب.

ج) درست. این ضریب اگر بین 0 تا 1 باشد، به عنوان exploration و exploitation همواره در برابری.

چون برای است عملاً exploitation و exploration با هم یافت شود. چون همواره آشنایی انتخاب می شود تماماً به جواب خواص می رسد.

Q3

$$V(C) = \frac{1}{2} (1+1+1+10) + \frac{1}{2} (4+1+10) = 14 \quad \text{الف)}$$

because  $V^\pi(s) = \sum_{s'} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V^\pi(s')]$

ب) آشنایی بهینه \* دارد (همان max است)  $V_{k+1}(s) \leftarrow \max_a \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_k(s')]$

تلاز اول  $V_1: A \rightarrow 1, B \rightarrow 1, C = \frac{1+4}{2} = 2.5, D=1, E=1, F=10, G=0$

تلاز دوم  $V_2: A \rightarrow 1+1=2, \begin{cases} 1+2.5=3.5 \\ 1-2=-1 \end{cases}, \begin{cases} 1+1=2 \\ 1-2=-1 \end{cases}, \begin{cases} 1+10=11 \\ 1-2=-1 \end{cases}, \begin{cases} 10+0=10 \\ 10-2=8 \end{cases}, G=0$

$$\begin{cases} \frac{1}{2}(4+1) + \frac{1}{2}(1+1) = 3.5 \\ 2.5 - 2 = 0.5 \end{cases}$$

$$V_2(B) = \max(3.5, -1) = 3.5 \quad Q_2(B, \text{right}) = 1 + [1 + 2.5] = 3.5 \quad Q_2(B, \text{left}) = 1 + [-2] = -1$$

$$Q_{k+1} = \sum T [R + \gamma \max Q_k]$$

نسخه ج مندرج

	Q(C,L)	Q(C,J)	Q(E,L)	Q(E,R)	Q(F,L)	Q(F,R)
init	0	0	0	0	0	0
trans. 1	0	$\frac{0}{2} + \frac{1}{2}(4+0) = 2$	0	0	0	0
" 2	0	0	$0 \frac{1}{2}(0) + \frac{1}{2}(1) = 0.5$	0	0	0
" 3	0	0	0	$0 \frac{0}{2} + \frac{1}{2}(-2+0.5) = -0.75$	0	0
" 4	0	0	$0 \frac{1}{2}(0) + \frac{1}{2}(1+0) = 0.5$	0	0	0

$$(1-\alpha) Q(s,a) + \alpha [R + \gamma \max_{a'} (Q'_K(s,a'))] = \frac{1}{2} Q_K(s,a) + \frac{1}{2} (R + Q'_K(s,a')) = Q_{K+1}(s,a)$$

Q4  $V_{K+1}^{\pi_i}(s) = \sum_{s'} T(s, \pi_i(s), s') [R(s, \pi_i, s') + \gamma V_K^{\pi_i}(s')]$  ہاں فرمول Q3 سے

$U_0 = 0$  First state: "First Rank"

$$\begin{aligned}
 V_1^{\pi}(FR) &= T(FR, R, FR) [R(FR, R, FR) + \gamma V(FR)] \\
 &+ T(FR, R, BR) [R(FR, R, BR) + \gamma V(BR)] + T(FR, W, FR) [R(FR, W, FR) + \gamma V(FR)] \\
 &+ \dots = \underbrace{1(0+0) + 0(0+0)}_{=0 = Q(FR, R)} + \underbrace{0.9(10+0) + 0.1(-10+0)}_{=8 = Q(FR, W)} + \underbrace{1(1+0) + 0(1+0)}_{=1 = Q(FR, Rest)} \\
 &\quad * V_1^{\pi}(FR) = 8 \text{ FR} \rightarrow W
 \end{aligned}$$

$$\begin{aligned}
 V_1^{\pi}(BR) &= \underbrace{.3(0+0) + .7(0+0)}_{=0 = Q(BR, Run)} + \underbrace{0.05(10+0) + .95(-10+0)}_{=-9 = Q(BR, W)} + \underbrace{0(1+0) + 1(1+0)}_{=1 = Q(BR, Rest)} \\
 &\quad * V_1^{\pi}(BR)
 \end{aligned}$$

$$\begin{aligned}
 V_2^{\pi}(FR) &= \underbrace{1(0 + .5(8)) + 0(0 + .5(1))}_{4.5 = Q(FR, Run)} + \underbrace{.9(10 + .5(8)) + .1(-10 + .5(1))}_{11.65 = Q(FR, W)} + \underbrace{0(1 + .5(8)) + 1(1 + .5(1))}_{1.5 = Q(FR, Rest)} \\
 &\quad * V_2^{\pi}(FR) = 11.65 \text{ FR} \rightarrow W
 \end{aligned}$$

$$\begin{aligned}
 V_2^{\pi}(BR) &= \underbrace{.3(0 + .5(8)) + .7(0 + .5(1))}_{1.55 = Q(BR, Run)} + \underbrace{.05(10 + .5(8)) + .95(-10 + .5(1))}_{-8.325 = Q(BR, W)} + \underbrace{0(1 + .5(8)) + 1(1 + .5(1))}_{1.5 = Q(BR, Rest)} \\
 &\quad * V_2^{\pi}(BR)
 \end{aligned}$$

دراستی



در کنار اول و دوم الویت حرکت « موقعیت First Rank تغییر نمی‌دهند، Water می‌دهند  
اما در کنار دوم، اولویت اول و دوم حرکت برتر موقعیت Bad Rank تغییر می‌دهند اما در کنار اول و دوم هیچ الویت همکار  
شده، والویت‌ها همین ترتیب باقی خواهند ماند.

## Q5

الف)  $Q(\text{خوشحال}, \text{ترن هلیکوپتر قدیم}) = w_0 b + w_1 x_1 + w_2 x_2 + w_3 x_3 = 1 + 2 + .5 = 3.5$

$$w^T = [w_0, w_1, w_2, w_3]^T, X = [b, x_1, x_2, x_3]$$

ب) به صورت قطعی از خوشحال به حالت می‌روم ( $T=1$ )، پاداش  $-10.5$  دریافت می‌کنیم؛ داریم: (مطابق وزن هلیکوپتر)

$$-10.5 + \frac{1}{2} \gamma^* = -10.5 + \frac{1}{2} \max \{ 3.5, 4, 2.5, 2, 2 \} = -8.5 \quad Q_{k+1}^{\text{old}}$$

new w:

$$\begin{aligned} w_0 &= 1 + \frac{1}{2} (-12) = -5 & w_2 &= 1 + \frac{1}{2} (-12) = -5 & \leftarrow \text{شاید بزرگترین \gamma^*} \\ w_1 &= 2 + \frac{1}{2} (-12) = -4 & w_3 &= .5 + \frac{1}{2} (-12) = -5.5 \end{aligned} \quad \left. \begin{array}{l} 3.5 = Q_k(\text{old train}) \\ \Delta Q(\text{old train}) = -12 \end{array} \right\}$$

بله ممکن است شما در جدول MDP توصیفیات قبلی یا شقوقی داشته باشید و بزرگترین واریانس نداشت.

ج)

$$\begin{aligned} 1 + 2 + 1 &= 4 & \leftarrow \text{ترن قدیم طولانی‌تر شد} & \leftarrow \text{وزن قدیم} \\ 1 + 1 &= 2 & \leftarrow \text{تقویت کوتاه‌تر شد} & \leftarrow \text{ترن جدید کوتاه‌تر شد} \\ 1 + 1 &= 2 & \leftarrow \text{کوتاه‌تر} & \leftarrow \text{کوتاه‌تر} \end{aligned}$$

با احتمال  $\frac{4\epsilon}{5}$  سایر حرکات و  $1 - \epsilon + \frac{\epsilon}{5} = 1 - \frac{4\epsilon}{5}$  حرکت ترن جدید کوتاه‌تر کند انتخاب می‌شود.

ه) این تابع  $n$  را عنوان تعداد exploitation «نظر می‌دهد» و  $k$  تأثیر exploration فرض می‌دهند تا  $\frac{k}{n}$  بخش exploration

و  $u$  عنوان یادداشت درباره از  $Q$  (همان  $v$ ) باشد.

الف)  $Q((UCLA, UPSET), BUS) = w_0 b + w_1 x_1 + w_2 x_2 + w_3 x_3 = 5$

$\downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow$   
 $1 \quad 1 \quad 2.5 \quad 1 \quad .5 \quad 1 \quad 1$

$$Q((UCLA, UPSET), Taxi) = 1 \times 1 + 2.5 \times 0 + .5 \times 1 + 1 \times 0 = 1.5$$

$\downarrow$   
 nowiki  
 no stops

$$Q((UCLA, UPSET), Monoc) = 1 \times 1 + (w \times \text{zeros}(3, 1)) = 1$$

ب، طبق پرسش. در پرسش انتشار خطای مناسب و بروز رسانی وزن ها جدید، خوشحالی و ناراضی وزن ندارند پس تاثیر ندارند.

ج،  $(1 - \epsilon)$  اتوبوس انتخاب شود و با احتمال  $\epsilon$  رند چون آیام  $\frac{\epsilon}{3}$  ممکن است اتوبوس انتخاب شود پس

$$\boxed{(1 - \frac{2}{3}\epsilon)} \text{ اتوبوس} \quad \boxed{\frac{2}{3}\epsilon} \text{ سایرین} \leftarrow \text{هرکدام } \frac{\epsilon}{3}$$

$$w_i = w_i + \alpha \underbrace{[r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a)]}_{2.5} \frac{\partial}{\partial w_i} Q_w(s, a) \quad (2)$$

$$.5 (-7 + .5 \max_{2.5} \{5, 1.5, 1\}) = -2.25$$

$$[w_0 \ w_1 \ w_2 \ w_3] = [1 \ 2.5 \ .5 \ 1] - 2.25 \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix}^T = [-2 \ 2.5 \ -2.5 \ 1]$$

ه، هنگام تعریف تابع  $Q$  و بروز رسانی  $Q$  ها و با وزن های روش ها مختلف وجود همجنس ممکن است علاوه بر تنوع انتخاب، انتخاب الگو تمرکز همگرا شود و نتایج را ممکن باشد اما مزیت آن این است که «ممکن است بزرگتر باشند و یا امتیاز کمتری داشته باشند» انتخاب آسانتر و آشنایتر باشد.