الف)

$$\text{utility\_up} = -50 + \gamma + \gamma^2 + \gamma^3 + \ldots + \gamma^{100}$$

$$\text{utility\_down} = 50 - \gamma - \gamma^2 - \ldots - \gamma^{100}$$

ب)

$$\text{utility\_up} > \text{utility\_down} \longrightarrow \gamma + \gamma^2 + \gamma^3 + \ldots + \gamma^{100} > 50$$
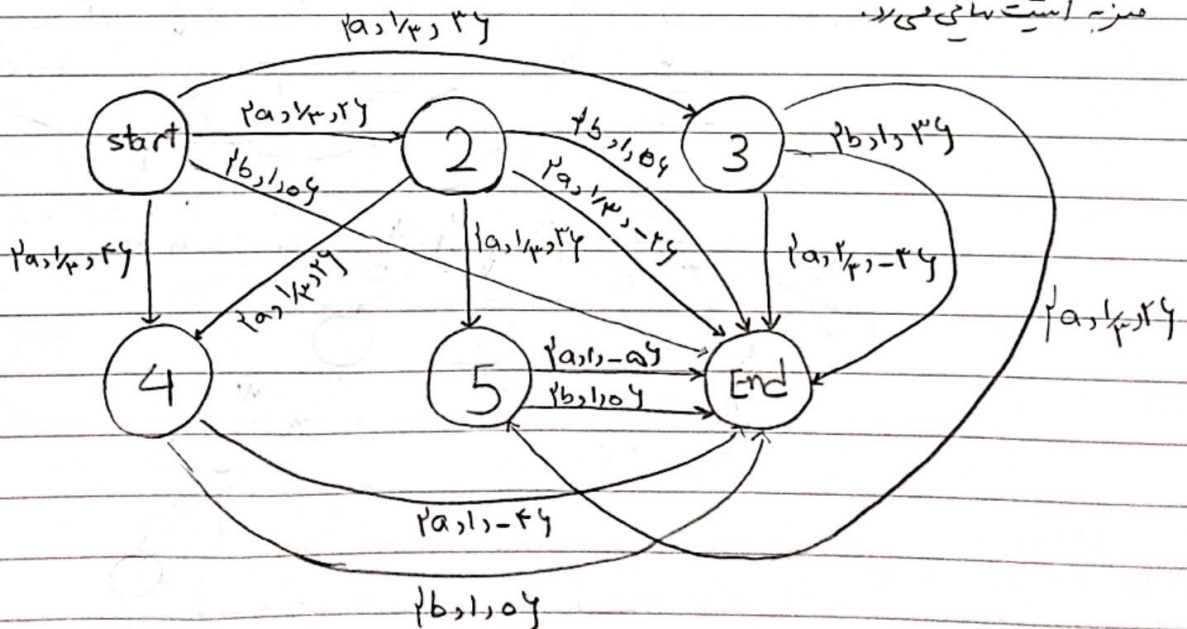
$$\frac{\gamma(1-\gamma^{100})}{1-\gamma} > 50 \longrightarrow \gamma - \gamma^{101} > 50 - 50\gamma$$

$$\longrightarrow 51\gamma - \gamma^{101} - 50 > 0 \longrightarrow \boxed{0.99 < \gamma < 1}$$

به ازای این مقادیر γ که در این معادله صدق کند اینت‌ها بالا می‌رود.

اگر در ( اکشن های up و down و left و right ) باشد که a و b پرداکشن ثابت و طا اکشن را نشی است (اگر مدل ۳ چیز آمده که اولی اکشن) در احتمال رفتن یک است و دسترسی رفتن از یک است به استیت بعدی است برای اینکه اگر از x بیشتر شد صورت می‌با reward منفی تنذار است و استیت بعدی می‌برد در غیر این صورت اگر انتقال دهد با reward صفر استیت بنای می‌برد.

ج) در صورتی که سیاست بهینه optimal policy را پیدا کنیم و policy بهتری پیدا نشود Converge کرده ایم.

بنابراین از این حالت شروع می کنیم:

| S | 2 | 3 | 4 | 5 | E |
|---|---|---|---|---|---|
| $\Pi_0 =$ a | a | b | b | b | - |
| p | 0 | 0 | 0 | 0 | 0 |
| ٣ | ١ | 0 | 0 | 0 | 0 |
| ١٠/٣ | ١ | 0 | 0 | 0 | 0 |

این اعداد را از قبل داریم نوشتیم

$$\Pi_{i+1} = \arg\max E[R(S,a,s') + V^\pi(S_i)]$$

$\Pi_1(S) = \arg\max_a$ $\begin{cases} b \to 0 \\ a \to \frac{1}{٣}(٢+١) + \frac{1}{٣}(٣+0) + \frac{1}{٣}(٤+0) = ١٠/٣ \end{cases}$ $\longrightarrow$ a

$\Pi_1(2) = \arg\max_a$ $\begin{cases} b \to 0 \\ a \to \frac{1}{٣}(٢+0) + \frac{1}{٣}(٣+0) + \frac{1}{٣}(-٢+0) = ١ \end{cases}$ $\longrightarrow$ a

$\Pi_1(3) = \arg\max_a$ $\begin{cases} b \to 0 \\ 0 \to \frac{1}{٣}(٢+0) + \frac{1}{٣}(-٣+0) + \frac{1}{٣}(-٣+0) = -٤/٣ \end{cases}$ $\longrightarrow$ b

$\Pi_1(٤) = \arg\max_a$ $\begin{cases} b \to 0 \\ a \to \frac{٣}{٣}(-٤+0) = -٤ \end{cases}$ $\longrightarrow$ b

$\Pi_1(٥) = \arg\max_a$ $\begin{cases} b \to 0 \\ a \to \frac{٣}{٣}(-٥) = -٥ \end{cases}$ $\longrightarrow$ b

| $\Pi_1 =$ | a | a | b | b | b | - |
|---|---|---|---|---|---|---|

چون $\Pi_0, \Pi_1$ مشابه Converge کرده است و $\Pi_1 = \Pi^*$ است.

s.a.m

$$S \quad 1 \quad 2 \quad 3 \quad 4 \quad \omega \quad E$$

$$V_0 = \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \quad 0 \qquad\qquad V = \max_a \sum_{s'} T(s,a,s')\,[R(s,a,s') + \gamma V_k(\cdot)]$$

$$k = 1$$

$$\max_a (S) \begin{cases} b \neq 0 \\ a \Rightarrow \frac{1}{\digamma}(r) + \frac{1}{\digamma}(\digamma) - \frac{1}{\digamma}(\varepsilon) = \digamma \end{cases} \to \digamma$$

$$\max_a (1) = \begin{cases} b = 0 \\ a = \frac{1}{\digamma}(r) + \frac{1}{\digamma}(r) + \frac{1}{\digamma}(-r) = 1 \end{cases} \to 1 \qquad \max_a (\varepsilon) \begin{cases} b = 0 \\ a = -2 \end{cases}$$

$$\max_a (r) = \begin{cases} b = 0 \\ a = \frac{1}{\digamma}(r) + \frac{r}{\digamma}(-r) = -\frac{\varepsilon}{\digamma} \end{cases} \to 0 \qquad \max_\omega (\omega) \begin{cases} b = 0 \\ a = -\partial \end{cases}$$

$$V_1 = \quad \digamma \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$$

$$k = r$$

$$\max_a (S) \begin{cases} b = 0 \\ a = \frac{1}{\digamma}(r+1) + \frac{1}{\digamma}(\digamma + 0) + \frac{1}{\digamma}(\varepsilon + 0) = \frac{1b}{\digamma} \end{cases} \to \frac{1b}{\digamma}$$

$$\max_a (r) = \begin{cases} b = 0 \\ a = \frac{1}{\digamma}(\digamma + 0) + \frac{1}{\digamma}(\digamma + 0) + \frac{1}{\digamma}(-r + 0) \end{cases} \to 1 \qquad \max_a (\varepsilon) \begin{cases} b = 0 \\ a = -\frac{1}{2} \end{cases} \to 0$$

$$\max_a (r) = \begin{cases} b = 0 \\ a = \frac{1}{\digamma}(r+0) + \frac{r}{\digamma}(-r) = -\frac{\varepsilon}{\digamma} = 0 \end{cases} \to 0 \qquad \max_a (\omega) \begin{cases} b = 0 \\ a = -\partial \end{cases} \to 0$$

$$V_\digamma = \quad \frac{1b}{\digamma} \quad 1 \quad 0 \quad 0 \quad 0 \quad 0$$

<span dir="rtl">کندو تکرار در است که مقادیر 1k تا Δ بیشتر منفی خواهد converge شود و k = ۳</span>

سوال ۳: برای ... سیاست پایان ۴ است ... state ← start ... ۱، ۲، ۳، ۴ ... R ...

[Persian handwritten text] (...)



action [Persian handwritten text]

reward [Persian handwritten text]

model [Persian handwritten text]

(ب)

$$V_0 = \quad 1 \quad 2 \quad 3 \quad 4 \quad e$$
$$\quad\quad 0 \quad 0 \quad 0 \quad 0 \quad 0$$

$$\underset{a}{\arg\max}(1) \begin{cases} a \to (-1\times\tfrac{1}{4}) + (\tfrac{1}{4}\times -1) = -1 \\ b \to 1 \end{cases} \longrightarrow 1 \qquad \underset{a}{\arg\max}(2) \begin{cases} a \to (\tfrac{1}{4}\times -1) = -1 \\ b \to 2 \end{cases} \to 2$$

$$V_1 = \quad 1 \quad 2 \quad 3 \quad 4 \quad 0$$

$$\underset{a}{\arg\max}(1) \begin{cases} a \to \tfrac{1}{4}(-1+\tfrac{2}{5}) + \tfrac{1}{4}(-1+\tfrac{2\cdot9}{10}) + \tfrac{1}{4}(-1+\tfrac{3\cdot4}{5}) + \tfrac{1}{4}(-1+\tfrac{4\cdot1}{10}) = \tfrac{5}{4} \\ b \to 1 \end{cases} \longrightarrow \tfrac{5}{4}$$

$$\underset{a}{\arg\max}(2) \begin{cases} a = \tfrac{3}{2} \\ b \to 2 \end{cases} \longrightarrow 2 \qquad \underset{a}{\arg\max}(3) = \begin{cases} a = \tfrac{3}{2} \\ b \to 3 \end{cases} \to 3$$

$$V_2 = \quad \tfrac{5}{4} \quad 2 \quad 3 \quad 4 \quad 0$$

$$\underset{a}{\arg\max}(1) \begin{cases} a \to \tfrac{1}{4}(-1+\tfrac{5\cdot9}{40}) + \tfrac{1}{4}(-1+\tfrac{11}{10}) + \tfrac{1}{4}(-1+\tfrac{3\cdot4}{5}) + \tfrac{1}{4}(-1+\tfrac{4\cdot1}{12}) = 1.409 \\ b \to 1 \end{cases} \longrightarrow 1\tfrac{3}{8}$$

$$V_3 = \quad 1.409 \quad 2 \quad 3 \quad 4 \quad 0 \qquad \underset{s.a.m}{\longrightarrow} \Pi = a \quad b \quad b \quad b \quad 0$$

|  | ١ | ٢ | ٣ | ٤ | E |
|---|---|---|---|---|---|
| $\Pi_0$ = | b | b | a | a | - |

$V_0 = 0 \quad 0 \quad 0 \quad 0 \quad 0$

$V_1 = 1 \quad 2 \quad -1 \quad -1 \quad 0 \qquad +\frac{1}{4}\left(-4 + \frac{1\times 9}{10} + \frac{2\times 9}{10} - \frac{1\times 9}{10} - \frac{1\times 9}{10}\right)$

$V_2 = 1 \quad 2 \quad -0.77 \quad -0.77$

$V_3 = 1 \quad 2 \quad -0.47 \quad -0.47 \qquad \frac{1}{4}\left(-4 + \frac{9}{10}(1 + 2 - \quad)\right)$

$\Pi_1(1) = \underset{a}{\arg\max}\begin{cases} b \to 1 \\ a \to \frac{1}{4}\left(-4 + \frac{9}{10}\left(1 + 2 - \frac{V_{01}}{4\times a}\right)\right) = -0.172 \end{cases} \longrightarrow b$

$\Pi_1(2) = \underset{a}{\arg\max}\begin{cases} b = 2 \\ a \to -0.17 \end{cases} \to b$

$\longrightarrow \Pi_1 = b \quad b \quad b \quad b \quad -$

$V_0^{\Pi_1} \qquad 0 \quad 0 \quad 0 \quad 0 \quad 0$

$V_1^{\Pi_1} \qquad 1 \quad 2 \quad 2 \quad 4 \quad 0$

$\longrightarrow \Pi_2(1) = \begin{cases} b = 1 \\ a = \frac{1}{4}\left(-4 + \frac{9}{10}\times 10\right) = 5/4 \end{cases} \to \frac{5}{4} \to a$

$\Pi_2(2) = b \qquad \Pi_2(3) = b \qquad \Pi_2(4) = b$

$\Pi_2 = a \quad b \quad b \quad b \quad -$

| state | ١ | ٢ | ٣ | ٤ | E |
|---|---|---|---|---|---|
| $\Pi_4$ | a | b | b | b | - |

$\uparrow$

$V_2(\text{تقريب}) = 1.21 \quad 2 \quad 2 \quad 4 \quad 0$

$\Pi_4(1), \Pi_4(2), \Pi_4(4) = b$

$\Pi_3(1) \to \begin{cases} b = 1 \\ a = \frac{1}{4}\left(-4 + \frac{9}{10}(0(10,21))\right) = 1.14 \end{cases} \to a$