سؤال ۱:

( ۲

$$V_k^n(s) \leftarrow (1-\alpha)V^n(s) + \alpha\left[R(s, \Pi(s), s') + \gamma V^n(s')\right]$$

$$V(A) \leftarrow 0/9(0) + 0/1\left[-۳+1(0)\right] = -0/۳$$

$$V(B) \leftarrow 0/9(0) + 0/1\left[۴+1(-0/۳)\right] = 0/۳V$$

$$V(A) \leftarrow 0/9(-0/۳) + 0/1\left[\underset{-0/۶۳}{-۳}+1(-0/۳)\right] = -0/۵V$$

$$V(A) \leftarrow 0/9(-0/۵V) + 0/1\left[-۳+1(0/۳V)\right] = -0/۸۹۳$$

$$V(B) \leftarrow 0/9(0/۳V) + 0/1\left[1+1(0)\right] = 0/۴۳۳$$

( ب

$$Q_1(A,1) = 0/9\left(Q_0(A,1)\right) + 0/1\left(-۳ + \max Q(B)\right) =$$
$$0/9(0) + 0/1(-3+0) = -0/۳$$

$$Q(B,1) = 0/9\left(Q_0(B,1)\right) + 0/1(4+0) = 0/4$$

$$Q(A,2) = 0/9(0) + 0/1(-4+0) = -0/4$$

$$Q_2(A,1) = 0/9(-0/۳) + 0/1(-3+0/4) = -(0/۲V + 0/۲۴) = -0/۵۳$$

$$Q(B,2) = 0/9(0) + 0/1(1+0) = 0/1$$

ج) با توجه به مقدار Qvalue ها و بیشینه آن برای هر state در بارد Action ۲ و در A اکشن ۱ را انتخاب می کنیم.

د) استفاده از Qvalue نسبت به مرکت تصادفی می تواند دقیق تر باشد و پی از اعمال این شرط احتمالاً به حالت های دقیق تری برسیم.

سؤال ۲: $\qquad Q_1(A^r, \rightarrow) = \frac{1}{2} Q_0(A, \rightarrow) + \frac{1}{2}(2 + \gamma \max_{ci} Q(B, ci))$

$\gamma = 0/7$

$Q \text{ values} = 1$

الف)

$a = 0/5$

$= \frac{1}{2}(1) + \frac{1}{2}(2 + 0/7(-2)) = 0/3 + 0/5 = 0/8$

$Q_1(C, \leftarrow) = \frac{1}{2} Q_0(C, \leftarrow) + \frac{1}{2}(2 + 0/7(-2)) = ((2 - 0/7)) = 0/8$

$Q_1(B, \rightarrow) = \frac{1}{2} Q_0(B, \rightarrow) + \frac{1}{2}(-2 + 0/7(2)) = 0/5 - 0/3 = 0/2$

$Q_2(A, \rightarrow) = \frac{1}{2} Q_1(A, \rightarrow) + \frac{1}{2}(4 + 0/7(0)) = 2,4$

$Q(A, \rightarrow) = 2,4$

$Q(B, \rightarrow) = 0/2$

ج) در این مثال با توجه به مقادیر $Q$ های بدست آمده از A به راست از B به راست و از C به چپ حرکت خواهیم کرد.

در سیاست حریصانه با توجه به مقدار $Q$ های بدست می آوریم که انتخاب در حریصانه انجام می دهیم و حرکت بعدی را دنبال می کنیم. اشکال اینجا است که ممکن است تخمین های ما مناسب نبوده و ما با پیشروی بر این آن حریصانه زودتر آن بی ترک reward و امتیازها را بیابیم در حالی که اگر کمی explore کنیم دیر تر به پاداش های بزرگتر بیابیم بشد.

می توان از روش $\varepsilon$ greedy استفاده کرد که در بار با احتمال کمی explore و با بقیه آن راه حریصی را انتخاب می کند.

د) الگوریتم های TD

اگر بر در مدل سازی مورد استفاده است و بطور پیوسته states ها update می شوند. در هر اپیزود TD(0) تنها یک قدم جلوتر نگاه می کند برای مثال episodic خوب است.

سؤال ۳: با استفاده از @Learning

الف)

$Q_1(A, Down) = \frac{1}{2}(0) + \frac{1}{2}(2 + \frac{1}{2}(0)) = 1$

$Q_1(B, Down) = \frac{1}{2}(0) + \frac{1}{2}(-4 + \frac{1}{2}(0)) = -2$

$Q_1(B, up) = \frac{1}{2}(0) + \frac{1}{2}(0 + \frac{1}{2}(0)) = 0$

$Q_2(B, up) = \frac{1}{2}(0) + \frac{1}{2}(3 + \frac{1}{2}(1)) = \frac{3}{2} + \frac{1}{4} = \frac{6+1}{4} = \frac{7}{4}$

$Q_1(A, Down) = 1$

$Q_1(B, up) = \frac{V}{F}$

ب)

$\hat{T}(A, up, A) = \frac{1}{1}$ $\qquad$ $\hat{T}(A, up, B) = 0$ $\qquad$ $\hat{T}(B, up, A) = \frac{1}{2}$

$\hat{T}(B, up, A) = \frac{1}{2}$ $\qquad$ $\hat{R}(A, up, A) = -1$ $\qquad$ $\hat{R}(A, up, B) = n/a$

$\hat{R}(B, up, A) = 3$ $\qquad$ $\hat{R}(B, up, B) = 0$

با توجه به تعداد تکرار هر کدام از شرایط $\hat{T}$ ها به دست آمده جایزبین از اول جدول به دست آمده.

$$V^{\Pi}(15) = E_{\Pi} \left\{ r_{t+1} + V^{\Pi}(s_{t+1}) \mid s_t = s \right\}$$

$$= -1 + \frac{1}{4} V^{\Pi}(12) + \frac{1}{4} V^{\Pi}(13) + \frac{1}{4} V^{\Pi}(14) + \frac{1}{4} V^{\Pi}(15)$$

$$= -1 - \frac{1}{4} \times 22 - \frac{1}{4} \times 20 - \frac{1}{4} \times 14 - \frac{1}{4} V^{\Pi}(15)$$

$$\frac{3}{4} V^{\Pi}(15) = -15 \implies V^{\Pi}(15) = -20$$

سؤال ۵:

الف) exploration یک راه سودمند و پیدا نیست است که به agent اجازه می دهد که دانش خود را در مورد action
هایی که در طولانی مدت به سودی انجامد بیشتر کند.

در مقابل exploitation از ارزش تخمینی فعلی Agent استفاده کرده و راه حل سریعتر را انتخاب می کند تا
بیشترین reward دست پیدا کند. آن جایی که Agent تخمین های مربوطه ی زیاد دارد و راضی نیست عملکر
بیشترین reward و بهترین راه باشد.

ب) با کاهش دادن مقدار E در E-greedy اکنون و میل دادن آن به صفر به نوعی حریان را احتمال
explore کردن را کم می کنیم و به آنویم greedy نزدیک تری می شویم که خیلی امید ی مناسب است.

V value چ) برای اندازه گیری میزان خوب بودن یک state خاص از نظر پاداش کلی مورد انتظار برای agent است و طبق یک policy عمل است و می‌شود  * $V^*$ حداکثر مقدار قابل دسترس برای هر حالت در فضای حالت دارویشن، می‌دهد

درحالیکه Q value ـ دانش هردو کریک یک خاص با توجه وضعیت برای عامل که از یک policy پیروی کند مقدار خوب است و  $Q^*$ حداکثر پاداش قابل دستی از یک جفت حالت عمل معین توسط خواهش بدهان می‌دهد

Q-value  function  transition  برای دسترسی سیستم مدل نری اینار چای یک 'action space کری دارم اما Q ما ضی راحت نیست برای استفاده

آ transition function مترس دارم گاهی اوقات V مناسب است