Qs ① Limiting ODE of the TD(0) algorithm with linear function approximation:

$$\dot{w}(t) = b - A w(t)$$

where,

$$b = \phi^T D R$$
$$A = \phi^T D [I - \gamma P] \phi$$

Global Lyapunov function for the limiting ODE

$$V(x) = \| w - w_* \|_P$$

Here,

$$w_* = A^{-1} b$$

Verification:

① To show: $V(x) \geq 0 \quad \forall x \quad$ with equality iff
$$x = x_*$$

Here, $V(w) = \| w - w_* \|_P$

$$= \left( \sum_{i=1}^{d} \lceil w_i \rceil^P_{-w_*} \right)^{1/p}$$

w.k.t $\quad w_i \geq w_{*i}$

$\therefore \quad V(w) \geq 0$

with $V(w) = 0 \quad$ iff $\quad w_i = w_{*i} \quad \forall i$
$$\Rightarrow w = w_*$$

② To show: $\lim\limits_{\|w\| \to \infty} V(w) = +\infty$

$$\lim\limits_{\|w\| \to \infty} \|w - w_* \|_p = \lim\limits_{\|w\| \to \infty} \left( \sum_{i=1}^{d} |w_i - w_{*i}|^p \right)^{1/p}$$

$$\leq \lim\limits_{\|w\| \to \infty} \left( \sum_{i=1}^{d} \|w_i\|^p \|w_*\|^p \right)^{1/p} \quad \text{[using Cauchy Swanchz Inequality]}$$

$$\leq \lim\limits_{\|w\| \to \infty} \left( \|w\|_p - \|w_*\|_p \right)$$

$$\leq +\infty$$

$$\Rightarrow \lim\limits_{\|w\| \to \infty} \|w - w_*\|_p \leq +\infty$$

$$\Rightarrow \lim\limits_{\|w\| \to \infty} \|w - w_*\|_p = +\infty$$

Qs ② Given: $\phi(1) = 1$      $\phi(2) = 2$

$$\Rightarrow \phi = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

Target policy $\mu$ and Behavior policy $\eta$

* $\eta(\text{right} | 1) = 0.5$

$\therefore \eta(\text{left} | 1) = 1 - \eta(\text{right} | 1) = 1 - 0.5 = 0.5$
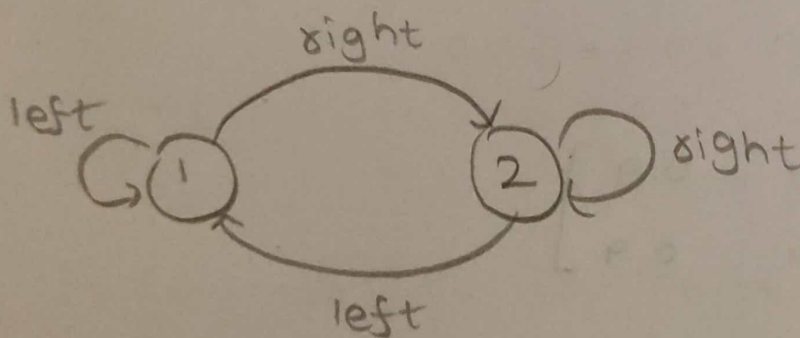
* $\eta(\text{right} | 2) = 0.5$

$\Rightarrow \eta(\text{left} | 2) = 1 - \eta(\text{right} | 2) = 1 - 0.5 = 0.5$

* $\mu(\text{right} | 1) = 0.9$

$\Rightarrow \mu(\text{left} | 1) = 1 - \mu(\text{right} | 1) = 1 - 0.9 = 0.1$

* $\mu(\text{right} | 2) = 0.9$

$\Rightarrow \mu(\underset{\text{left}}{\text{right}} | 2) = 1 - \mu(\text{right} | 2) = 1 - 0.9 = 0.1$

Qs ② ⓐ To find: $P^M$

$P^M(S_{t+1}=1|S_t=1) = \mu(\text{right}|S_t=1) \, P^M(S_{t+1}=1|S_t=1, a_t=\text{right})$
$\qquad\qquad + \mu(\text{left}|S_t=1) \, P^M(S_{t+1}=1|S_t=1, a_t=\text{left})$

$\qquad = 0.9 * 0 + 0.1 * 1 = 0.1$

$P^M(S_{t+1}=2|S_t=1) = \mu(\text{right}|1) \, P^M(2|1,\text{right}) +$
$\qquad\qquad\qquad \mu(\text{left}|1) \, P^M(2|1,\text{left})$

$\qquad = 0.9 * 1 + 0.1 * 0 = 0.9$

$P^M(S_{t+1}=1|S_t=2) = \mu(\text{right}|2) \, P^M(1|2,\text{right}) +$
$\qquad\qquad\qquad \mu(\text{left}|2) \, P^M(1|2,\text{left})$

$\qquad = 0.9 * 0 + 0.1 * 1 = 0.1$

$P^M(S_{t+1}=2|S_t=2) = \mu(\text{right}|2) \, P^M(2|2,\text{right}) +$
$\qquad\qquad\qquad \mu(\text{left}|2) \, P^M(2|2,\text{left})$

$\qquad = 0.9 * 1 + 0.1 * 0 = 0.9$

$\therefore \ P^M = \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix}$

To find: $d^M$

we know that $\qquad d^{M^T} P^M = d^{M^T}$

$\Rightarrow \begin{bmatrix} d^M(1) & d^M(2) \end{bmatrix} \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix} = \begin{bmatrix} d^M(1) & d^M(2) \end{bmatrix}$

$\Rightarrow 0.1 \, d^M(1) + 0.1 \, d^M(2) = d^M(1)$

$\Rightarrow 0.1 \, d^M(2) = 0.9 \, d^M(1)$

$\Rightarrow d^M(2) = 9 \, d^M(1) \quad \longrightarrow \text{①}$

we also
know that $\qquad d^M(1) + d^M(2) = 1$

Subs Eq ①

$\Rightarrow d^M(1) + 9 \, d^M(1) = 1$

$\Rightarrow 10 \, d^M(1) = 1$

$\Rightarrow d^M(1) = 0.1$

From ① $\qquad d^M(2) = 9 * 0.1 = 0.9$

$\therefore d^M = \begin{bmatrix} 0.1 \\ 0.9 \end{bmatrix}$

To find : A

w.k.t $\qquad A = \phi^T D (I - \gamma P^M) \phi$

$= \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 0.1 & 0 \\ 0 & 0.9 \end{bmatrix} \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \gamma \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} 0.1 & 1.8 \end{bmatrix} \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - \begin{bmatrix} 0.09 & 0.81 \\ 0.09 & 0.81 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} 0.1 & 1.8 \end{bmatrix} \begin{bmatrix} 0.91 & -0.81 \\ -0.09 & 0.19 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} -0.071 & 0.261 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} -0.071 + 0.522 \end{bmatrix}$

$= \begin{bmatrix} 0.451 \end{bmatrix}$

As the entry in A. is positive, here, A is a positive definite matrix.

∴. The limiting ODE associated with the on-policy TD(0) has a globally asymptotically stable equilibrium.

Qs ② ⓑ To find: $P\pi$?

$$P\pi = \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix}$$

To find: $d\pi$

w.k.t $\quad d\pi^T P\pi = d\pi^T$

$\Rightarrow \begin{bmatrix} d\pi(1) & d\pi(2) \end{bmatrix} \begin{bmatrix} 0.5 & 0.5 \\ 0.5 & 0.5 \end{bmatrix} = \begin{bmatrix} d\pi(1) & d\pi(2) \end{bmatrix}$

$\Rightarrow \quad 0.5 d\pi(1) + 0.5 d\pi(2) = d\pi(1)$

$\Rightarrow \quad d\pi(1) = d\pi(2) \quad \longrightarrow ②$

Also, $\quad d\pi(1) + d\pi(2) = 1$

Subs Eq ②

$\Rightarrow \quad d\pi(1) + d\pi(1) = 1$

$\Rightarrow \quad d\pi(1) = d\pi(2) = 0.5$

∴ $\quad d\pi = \begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$

To find : A'

w.k.t $A' = \phi^T D \eta (I - \gamma P^\mu) \phi$

$= \begin{bmatrix} 1 & 2 \end{bmatrix} \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix} \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} 0.5 & 1 \end{bmatrix} \begin{bmatrix} 0.91 & -0.81 \\ -0.09 & 0.19 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} 0.365 & -0.215 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$

$= \begin{bmatrix} 0.365 & -0.43 \end{bmatrix}$

$= \begin{bmatrix} -0.065 \end{bmatrix}$

The entry in A' is negative.

∴ Here, A' is not a positive definite matrix

∴ The limiting ODE associated with the naive off-policy algorithm here does not have a globally asymptotically stable equilibrium.

Qs ② ⓒ Given:

$$f(s) = d_\eta(s) \lim_{t \to \infty} \mathbb{E}_\eta[F_t \mid s_t = s]$$

Here, $F_t = P_{t-1} \Upsilon K_{t-1} + 1$

$\Rightarrow f(s) = d_\eta(s) + \Upsilon \sum_{s' f} d_e(s') P^\mu_{s', s}$      [Derived in Lecture 9]

$\Rightarrow f(s) = d_\eta(s) + \Upsilon \alpha f^T P_{., s}$

$\Rightarrow f^T(I - \Upsilon P^\mu) = d_\eta^T$

$\Rightarrow [f(1) \quad f(2)] \left[ \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix} \right] = [0.5 \quad 0.5]$

$\Rightarrow [f(1) \quad f(2)] \begin{bmatrix} 0.91 & -0.81 \\ -0.09 & 0.19 \end{bmatrix} = [0.5 \quad 0.5]$

$\Rightarrow 0.91 f(1) - 0.09 f(2) = 0.5 \longrightarrow ③$

$-0.81 f(1) + 0.19 f(2) = 0.5 \longrightarrow ④$

Multiplying ③ by 0.19 and ④ by 0.09

$\Rightarrow 0.1729 f(1) - 0.0171 f(2) = 0.095$

$\Rightarrow -0.0729 f(1) + 0.0171 f(2) = 0.045$

Adding the two equations

$\Rightarrow 0.1 \alpha f(1) = 0.14$

$\Rightarrow f(1) = 1.4$

Subs in ③

$$0.91 + 1.4 - 0.09 \, f(2) = 0.5$$

$$\Rightarrow 1.274 - 0.09 \, f(2) = 0.5$$

$$\Rightarrow -0.09 \, f(2) = -0.774$$

$$\Rightarrow f(2) = 8.6$$

$$\therefore f(1) = 1.4 \qquad\qquad f(2) = 8.6$$

To find : $A''$

w. k. t $\quad A'' = \phi^T D^e (I - \gamma P M) \phi$

$$= [1 \ \ 2] \begin{bmatrix} 1.4 & 0 \\ 0 & 8.6 \end{bmatrix} \left( \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} - 0.9 \begin{bmatrix} 0.1 & 0.9 \\ 0.1 & 0.9 \end{bmatrix} \right) \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= [1.4 \ \ 17.2] \begin{bmatrix} 0.91 & -0.81 \\ -0.09 & 0.19 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= [0.274 \ \ \ast 2.134] \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$= [-0.274 + 4.268]$$

$$A'' = [3.994]$$

AS the entry in $A''$ is positive. Here, $A''$ is a positive definite matrix.

∴ The limiting ODE associated with emphatic TD(0) has a globally asymptotically stable equilibrium.

Qs ③ To find : $\oint \sum_s f(s)$

w.k.t $\quad f(s) = d_\eta(s) + \gamma f^T P_{\cdot,s}$

$\therefore \sum_s f(s) = \sum_s d_\eta(s) + \gamma f^T P_{\cdot,s}$

$\qquad = \sum_s d_\eta(s) + \gamma \sum_s f^T P_{\cdot,s}$

$\qquad = 1 + \gamma \sum_s f^T P_{\cdot,s} \qquad [\because \sum_s d_\eta(s) = 1]$

$\qquad = 1 + \gamma \sum_s \sum_{s'} f(s') P_{s',s}$

$\qquad = 1 + \gamma \sum_s f(s) \sum_{s'} P_{s,s'}$

$\qquad = 1 + \gamma \sum_s f(s) \cdot 1 \qquad [\because$ sum of row vector of PM is 1]

$\Rightarrow \sum_s f(s) - \gamma \sum_s f(s) = 1$

$\Rightarrow \sum_s f(s) \left[ 1 - \gamma \right] = 1$

$\Rightarrow \sum_s f(s) = \dfrac{1}{1-\gamma}$

Yes, we can identify $\sum_s f(s)$, which is

$$\sum_s f(s) = \dfrac{1}{1-\gamma}$$

Qs $\textcircled{4}$ Optimization Problem : $\min \mathcal{E}(\omega)$

$$\mathcal{E}(\omega) = \frac{1}{2} \| Q^M(s,a) - Q^*(s,a) \|^2_D$$

Intuition : The optimization problem tries to reduce the difference between the true (or target) $Q$ value and the current estimate of the $Q$-value, which is what we want. We want our estimate of $Q$-value to be as close as to value of true $Q$-value.

* But, the problem (or difficulty) in doing this is, we don't know the $Q^*(s,a)$ value

∴ we try to approximate the $Q^*(s,a)$ value by

$$Q^*(s,a) \simeq \phi^T(s,a).\omega$$

where,

     $\phi$ - apriori known feature matrix

     $\omega$ - parameter.

∴ $\min \mathcal{E}(\omega) = \min \frac{1}{2} \| Q^M(s,a) - \phi^T(s,a)\omega \|^2_D$

Convex Optimization:

∴ Differentiate $\mathcal{E}(\omega)$ w.r.t $\omega$

Here, $\mathcal{E}(\omega) = \frac{1}{2} \sum_{a,s} d_\mu(s) \mu(a) (Q^M(s,a) - \phi^T(s,a)\omega)^2$

$$\therefore \nabla \mathcal{E}(\omega) = -\mathbb{E}\left[\left(Q^{\mu}(s_t, a_t) - \phi^T(s_t, a_t)\omega\right)\phi(s_t, a_t)\right]$$

where, $\quad s_t \sim d_\mu$ $\qquad\qquad\qquad\qquad\qquad \hookrightarrow ①$

$\qquad\qquad a_t \sim \mu$

From Bellman Equation, we know that

$$Q^{\mu}(s, a) = \mathbb{E}\left[\gamma(s_t, a_t, s_{t+1}) + \gamma Q^{\mu}(s_{t+1}, a_{t+1}) \mid s_t = s, a_t = a\right]$$

$$= \sum_{a', s'} \mathbb{P}\{a_{t+1} = a', s_{t+1} = s' \mid s_t = s, a_t = a\}$$

$$\left[\gamma(s, a, s') + \gamma Q^{\mu}(s', a')\right]$$

Subs value of $Q^{\mu}(s, a)$ in ①

$$\nabla \mathcal{E}(\omega) = -\sum_{a, s} \mathbb{P}\{s_t = s\} \mu(a_t = a)\underbrace{\left[Q^{\mu}(s, a) - \phi^T(s, a)\omega\right]}_{\phi(s, a)}$$

$$= -\sum_{a, s} \mathbb{P}\{s\} \mu(a) \sum_{a', s'} \mathbb{P}\{s', a' \mid s, a\} (\gamma(s, a, s')$$

$$+ \gamma Q^{\mu}(s', a') - \phi^T(s, a)\omega)\phi(s, a)$$

$$= -\sum_{\substack{a, s, \\ s', a'}} \mathbb{P}\{s_t = s, a_t = a, s_{t+1} = s', a_{t+1} = a'\},$$

$$(\gamma(s, a, s') + \gamma Q^{\mu}(s', a') - \phi^T(s, a)\omega)$$

$$\phi(s, a)$$

$$= -\mathbb{E}\left[\left(\gamma(s, a, s') + \gamma Q^{\mu}(s', a') - \phi^T(s, a)\omega\right)\phi(s, a)\right]$$

where, $\quad s_t \sim d_\mu \qquad\qquad\qquad s_{t+1} \sim \mathbb{P}(\cdot \mid s_t, a_t)$

$\qquad\qquad a_t \sim \mu(\cdot \mid s_t) \qquad\qquad a_{t+1} \sim \mu(\cdot \mid s_{t+1})$

Hence, the SGD algorithm would be:

$$w_{t+1} = w_t + \alpha_t [\gamma(s,a,s') + \gamma Q^M(s',a')$$
$$- \phi^T(s,a)w] \phi(s,a)$$

where, $\alpha_t \to$ step size