

# HW1 - Part 2

## Problem 1

$$\begin{array}{lll} r(s,a) = 100 & \theta_a = 0 & \pi_{\theta}(a|s) = 0.1 \\ r(s,b) = 98 & \theta_b = \ln 5 & \pi_{\theta}(b|s) = 0.5 \\ r(s,c) = 95 & \theta_c = \ln 4 & \pi_{\theta}(c|s) = 0.4 \end{array}$$

a) What are the  $E[\hat{V}]$  and the covariance matrix of  $\hat{V}$

$$E[\hat{V}] = E \begin{bmatrix} r(s_t, a_t) \frac{\partial \log \pi_{\theta}(a_t|s_t)}{\partial \theta(s, a)} \\ r(s_t, a_t) \frac{\partial \log \pi_{\theta}(a_t|s_t)}{\partial \theta(s, b)} \\ r(s_t, a_t) \frac{\partial \log \pi_{\theta}(a_t|s_t)}{\partial \theta(s, c)} \end{bmatrix}$$

$$= \begin{bmatrix} \pi_{\theta}(a|s) \cdot r(s,a) \cdot (1 - \pi_{\theta}(s,a)) - \pi_{\theta}(b|s) \cdot r(s,b) \cdot \pi_{\theta}(s,a) - \pi_{\theta}(c|s) \cdot r(s,c) \cdot \pi_{\theta}(s,a) \\ -\pi_{\theta}(a|s) \cdot r(s,a) \cdot \pi_{\theta}(s,b) + \pi_{\theta}(b|s) \cdot r(s,b) \cdot (1 - \pi_{\theta}(s,b)) - \pi_{\theta}(c|s) \cdot r(s,c) \cdot \pi_{\theta}(s,b) \\ -\pi_{\theta}(a|s) \cdot r(s,a) \cdot \pi_{\theta}(s,c) - \pi_{\theta}(b|s) \cdot r(s,b) \cdot \pi_{\theta}(s,c) + \pi_{\theta}(c|s) \cdot r(s,c) \cdot (1 - \pi_{\theta}(s,c)) \end{bmatrix}$$

$$= \begin{bmatrix} 0.1 \cdot 100 \cdot 0.9 - 0.5 \cdot 98 \cdot 0.1 - 0.4 \cdot 95 \cdot 0.1 \\ -0.1 \cdot 100 \cdot 0.5 + 0.5 \cdot 98 \cdot 0.5 - 0.4 \cdot 95 \cdot 0.5 \\ -0.1 \cdot 100 \cdot 0.4 - 0.5 \cdot 98 \cdot 0.4 + 0.4 \cdot 95 \cdot 0.6 \end{bmatrix} = \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}$$

covariance matrix of  $\hat{V}$

$$E[(\hat{V} - E[\hat{V}])(\hat{V} - E[\hat{V}])^T]$$

$$= E[\hat{V}\hat{V}^T] - \nabla_{\theta} V^{\pi_{\theta}} \nabla_{\theta} V^{\pi_{\theta}^T}$$

$$= E[\hat{V}\hat{V}^T] - \begin{bmatrix} 0.09 & 0.15 & -0.24 \\ 0.15 & 0.25 & -0.4 \\ -0.24 & -0.4 & 0.64 \end{bmatrix}$$

$$\hat{V} = r(S_0, a_0) \begin{bmatrix} \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, a)} \\ \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, b)} \\ \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, c)} \end{bmatrix}$$

$$\hat{V}\hat{V}^T = r(S_0, a_0)^2 \begin{bmatrix} \left( \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, a)} \right)^2 & \frac{\partial \dots}{\partial \theta(S, a)} \frac{\partial \dots}{\partial \theta(S, b)} & \frac{\partial \dots}{\partial \theta(S, a)} \frac{\partial \dots}{\partial \theta(S, c)} \\ \frac{\partial \dots}{\partial \theta(S, a)} \frac{\partial \dots}{\partial \theta(S, b)} & \left( \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, b)} \right)^2 & \frac{\partial \dots}{\partial \theta(S, b)} \frac{\partial \dots}{\partial \theta(S, c)} \\ \frac{\partial \dots}{\partial \theta(S, a)} \frac{\partial \dots}{\partial \theta(S, c)} & \frac{\partial \dots}{\partial \theta(S, b)} \frac{\partial \dots}{\partial \theta(S, c)} & \left( \frac{\partial \log \pi_{\theta}(a_0 | S_0)}{\partial \theta(S, c)} \right)^2 \end{bmatrix}$$

$$E[\hat{V}\hat{V}^T] = \sum \pi(a_0 | S_0) \cdot r(a_0, S_0)^2 \cdot \begin{bmatrix} \dots & \dots & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

$$= \begin{bmatrix} 894.12 & -509.6 & -384.52 \\ -509.6 & 2353 & -1843.4 \\ -384.52 & -1843.4 & 2227.92 \end{bmatrix}$$

$$\text{所求} = E[\hat{V}\hat{V}^T] - \nabla_{\theta} V^{\pi_{\theta}} \nabla_{\theta} V^{\pi_{\theta}^T} = \begin{bmatrix} 894.03 & -509.75 & -384.28 \\ -509.75 & 2352.75 & -1843 \\ -384.28 & -1843 & 2227.28 \end{bmatrix} \quad \#$$

# HW1 - Part 2

## Problem 1

(b) mean vector of  $\hat{\theta}V$ :

$$V^{\pi_0}(s) = 0.1 \cdot 100 + 0.5 \cdot 98 + 0.4 \cdot 95 = 97$$

$$\nabla_{\theta} V^{\pi_0}(u) \approx \sum_{t=0}^{\infty} \gamma^t (Q_t - V^{\pi_0}(s)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

$$\Rightarrow E[\hat{\nabla} V] = E[(r(s_t, a_t) - V^{\pi_0}(s)) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]$$

$$= \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix}$$

covariance matrix of  $\hat{\nabla} V$

$$\begin{aligned} \hat{\Sigma} &= E[\hat{\nabla} V \hat{\nabla} V^T] = \nabla_{\theta} V^{\pi_0} \nabla_{\theta} V^{\pi_0 T} \\ &= \sum \pi(a_t | s_t) (r(a_t, s_t) - V^{\pi_0}(s))^2 \begin{bmatrix} \left( \frac{\partial \log \pi_{\theta}(a_t | s_t)}{\partial \theta(s, a)} \right)^2 \frac{\partial^{...}}{\partial (s, a)} \frac{\partial^{...}}{\partial (s, b)} \frac{\partial^{...}}{\partial (s, a)} \frac{\partial^{...}}{\partial (s, c)} \\ \frac{\partial^{...}}{\partial \theta(s, b)} \frac{\partial^{...}}{\partial \theta(s, a)} \left( \frac{\partial \log \pi_{\theta}(a_t | s_t)}{\partial \theta(s, b)} \right)^2 \frac{\partial^{...}}{\partial (s, b)} \frac{\partial^{...}}{\partial (s, c)} \\ \frac{\partial^{...}}{\partial \theta(s, a)} \frac{\partial^{...}}{\partial \theta(s, c)} \frac{\partial^{...}}{\partial \theta(s, b)} \frac{\partial^{...}}{\partial \theta(s, c)} \left( \frac{\partial \log \pi_{\theta}(a_t | s_t)}{\partial (s, c)} \right)^2 \end{bmatrix} \end{aligned}$$

$$= \begin{bmatrix} 0.3 \\ 0.5 \\ -0.8 \end{bmatrix} \begin{bmatrix} 0.3 & 0.5 & -0.8 \end{bmatrix}$$

$$= \begin{bmatrix} 0.09 & -0.15 & -0.16 \\ -0.15 & 0.25 & 0 \\ -0.16 & 0 & 0.64 \end{bmatrix} \#$$

Problem 2

(a) Show that  $E_{\tau \sim P_{\pi_0}} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \frac{1}{1-\gamma} E_{s \sim d} E_{a \sim \pi_0(\cdot|s)} [f(s, a)]$

$$\begin{aligned} \text{RHS} &= \frac{1}{1-\gamma} \sum_s d^{\pi_0}(s) \cdot \sum_a \pi_0(a|s) f(s, a) \\ &= \frac{1}{1-\gamma} \sum_s \left( \sum_{s_0} u(s_0) (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \cdot p(s_t=s|s_0, \pi_0) \right) \cdot \sum_a \pi_0(a|s) \cdot f(s, a) \\ &= \sum_{s_0} u(s_0) \cdot \sum_s \sum_a \pi_0(a|s) \sum_{t=0}^{\infty} \gamma^t p(s_t=s|s_0, \pi_0) f(s, a) \\ &= \sum_{s_0} u(s_0) \cdot \sum_s \sum_a \sum_{t=0}^{\infty} \gamma^t p(s_t=s, a_t=a|s_0, \pi_0) \cdot f(s, a) \\ &= \sum_{s_0} u(s_0) \cdot \sum_{\tau} P(\tau|s_0) \cdot \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \\ &= \sum_{\tau} P(\tau) \cdot \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \\ &= E_{\tau \sim P_{\pi_0}} \left[ \sum_{t=0}^{\infty} \gamma^t f(s_t, a_t) \right] = \text{LHS} \quad \# \end{aligned}$$

(b) show that  $\nabla_{\theta} V^{\pi_0}(u) = E_{\tau \sim P_{\pi_0}} \left[ \sum_{t=0}^{\pi(\tau)-1} \gamma^t A^{\pi_0}(s_t, a_t) \nabla_{\theta} \log \pi_0(a_t|s_t) \right]$

for episodic environments.

Pf: The major difference between episodic and continuing environments is the existence of "terminal state".

Suppose  $S_*$  be the terminal state of an episodic environments.

Once the agent reaches  $S_*$ , it will stay at  $S_*$  forever.

Let  $T(\tau)$  be the episodic length of a trajectory  $\tau$ ,

Then,

$$\begin{aligned} \nabla_{\theta} V^{\pi_0}(u) &= E_{\tau \sim P_{\pi_0}} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_0}(s_t, a_t) \nabla_{\theta} \log \pi_0(a_t|s_t) \right] \quad \text{--- (P5)} \\ &= E_{\tau \sim P_{\pi_0}} \left[ \sum_{t=0}^{\pi(\tau)-1} \gamma^t A^{\pi_0}(s_t, a_t) \nabla_{\theta} \log \pi_0(a_t|s_t) \right] \end{aligned}$$

**Problem 3 :**

寫程式用數值方法逼近求得最佳的 **baseline** 約為 **97.137931**

將式子對 **baseline** 微分找微分為 **0** 之地方為 **baseline** 亦為相近之值

附檔為使用數值方法逼近 **baseline** 之 **code**