

ENGG6600 Reinforcement Learning Assignment-3 Fall 2022

Yaowen Mei
1177855

Problem 1- Variance of return G_t

Prove that $Var(G_{t+1}) \geq Var(G_t)$ is true if we assume R_{t+1} is, on average, correlated with the previous rewards.

Given that:

- $G_t = \sum_{i=0}^t R_i$
- $\frac{1}{t+1} \sum_{i=0}^t Cov(R_i, R_{t+1}) > 0$

From statistic, we also have the following properties/definitions:

$$\begin{aligned} Cov(X, Y) &= E[XY] - E[X]E[Y] \\ Var(X) &= E[X^2] - E[X]^2 \geq 0 \end{aligned} \tag{1}$$

$$\sum_{i=0}^t Cov(R_i, R_{t+1}) = Cov\left(\sum_{i=0}^t R_i, R_{t+1}\right) = Cov(G_t, R_{t+1}) > 0 \tag{2}$$

SOLUTION-1:

$$\begin{aligned} Var(G_{t+1}) &= Var(G_t + R_{t+1}) \\ &= E[(G_t + R_{t+1})^2] - E[G_t + R_{t+1}]^2 \\ &= E[G_t^2 + R_{t+1}^2 + 2G_t R_{t+1}] - E[G_t]^2 - E[R_{t+1}]^2 - 2E[G_t]E[R_{t+1}] \\ &= E[G_t^2] - E[G_t]^2 + E[R_{t+1}^2] - E[R_{t+1}]^2 + 2E[G_t R_{t+1}] - 2E[G_t]E[R_{t+1}] \\ &= Var(G_t) + Var(R_{t+1}) + 2Cov(G_t, R_{t+1}) \end{aligned}$$

Therefore, by using Eq. 1 and Eq. 2, the equation above can be rewrite as:

$$Var(G_{t+1}) - Var(G_t) = \underbrace{Var(R_{t+1})}_{\geq 0} + \underbrace{2Cov(G_t, R_{t+1})}_{> 0} \geq 0 \quad \square$$

Problem 2—Variance Reduce in Policy Gradient Method

Potentially, at the cost of increased bias, the variance in policy gradient methods could be reduced.

Let us consider an infinite horizon MDP $\langle S, A, Pr, \gamma \rangle$, let us define:

The advantage function:

$$A_{\pi_\theta}(S_t, A_t) = q_{\pi_\theta}(S_t, A_t) - v_{\pi_\theta}(S_t)$$

The policy gradient function:

$$\nabla_\theta J(\theta) = E_{\pi_\theta} [A_{\pi_\theta}(S_t, A_t) \nabla_\theta \log \pi(A_t | S_t; \theta)]$$

In practice, we do not have access to the true advantage function $A_{\pi_\theta}(S_t, A_t)$, so we would like to consider using the general form of an estimator $\hat{A}_{\pi_\theta}(S_t, A_t)$ that can be a function of the entire trajectory.

Questions List:

1. ***Policy Gradient with Baseline:***

Given that $\hat{q}_{\pi_\theta}(S_{t:\infty}, A_{t:\infty})$ is an unbiased estimator of the true $q_{\pi_\theta}(S_t, A_t)$, and b_t is an arbitrary function of the actions and states sampled before A_t , **prove** that by using this estimate of A_t , we obtain an unbiased estimate of the policy gradient.

2. ***TD error as unbiased estimator of the advantage function:***

Let's look at another variants of \hat{A}_{π_θ} .

Recall that TD error:

$$\delta_{\pi_\theta} = R_{t+1} + \gamma \hat{v}(S_{t+1}) - \hat{v}(S_t) \quad (3)$$

Prove that δ_{π_θ} is an unbiased estimator of π_θ when $\hat{v} = v_{\pi_\theta}$

SOLUTION-2.1:

$$\begin{aligned} LHS &= \nabla_\theta J(\theta) = E_{\pi_\theta} [A_{\pi_\theta}(S_t, A_t) \nabla_\theta \log \pi(A_t | S_t; \theta)] \\ &= E_{\pi_\theta} [(q_{\pi_\theta}(S_t, A_t) - v_{\pi_\theta}(S_t)) \nabla_\theta \log \pi(A_t | S_t; \theta)] \\ &= \underbrace{E_{\pi_\theta} [q_{\pi_\theta}(S_t, A_t) \nabla_\theta \log \pi(A_t | S_t; \theta)]}_{(2.1)} - \underbrace{E_{\pi_\theta} [v_{\pi_\theta}(S_t) \nabla_\theta \log \pi(A_t | S_t; \theta)]}_{(2.2)} \end{aligned} \quad (4)$$

$$\begin{aligned} RHS &= E_{\pi_\theta} [\hat{A}_{\pi_\theta}(S_t, A_t) \nabla_\theta \log \pi(A_t | S_t; \theta)] = E_{\pi_\theta} [(\hat{q}_{\pi_\theta}(S_{t:\infty}, A_{t:\infty}) - b_t(S_{0:t-1}, A_{0:t-1}) \nabla_\theta \log \pi(A_t | S_t; \theta))] \\ &= \underbrace{E_{\pi_\theta} [\hat{q}_{\pi_\theta}(S_{t:\infty}, A_{t:\infty}) \nabla_\theta \log \pi(A_t | S_t; \theta)]}_{(2.3)} - \underbrace{E_{\pi_\theta} [b_t(S_{0:t-1}, A_{0:t-1}) \nabla_\theta \log \pi(A_t | S_t; \theta)]}_{(2.4)} \end{aligned} \quad (5)$$

First, we will proof the (2.2) component in LHS equals the (2.4) componenet in RHS.

Theorem 1. *The (2.2) component in LHS equals the (2.4) componenet in RHS equals 0.*

$$E_{\pi_{\theta}} [v_{\pi_{\theta}} (S_t) \nabla_{\theta} \log \pi (A_t | S_t; \theta)] = E_{\pi_{\theta}} [b_t (S_{0:t-1}, A_{0:t-1}) \nabla_{\theta} \log \pi (A_t | S_t; \theta)] = 0$$

Proof.

$$\begin{aligned} E_{\pi_{\theta}} [v_{\pi_{\theta}} (S_t) \nabla_{\theta} \log \pi (A_t | S_t; \theta)] &= E_{S_t \sim \pi_{\theta}} \left[\sum_a \pi (a | S_t; \theta) v_{\pi_{\theta}} (S_t) \frac{\nabla_{\theta} \pi (a | S_t; \theta)}{\pi (a | S_t; \theta)} \right] \\ &= E_{S_t \sim \pi_{\theta}} \left[v_{\pi_{\theta}} (S_t) \sum_a \pi (a | S_t; \theta) \frac{\nabla_{\theta} \pi (a | S_t; \theta)}{\pi (a | S_t; \theta)} \right] \\ &= E_{S_t \sim \pi_{\theta}} \left[v_{\pi_{\theta}} (S_t) \sum_a \nabla_{\theta} \pi (a | S_t; \theta) \right] \\ &= E_{S_t \sim \pi_{\theta}} \left[v_{\pi_{\theta}} (S_t) \sum_a \nabla_{\theta} 1 \right] = 0 \end{aligned}$$

Also:

$$\begin{aligned} E_{\pi_{\theta}} [b_t (S_{0:t-1}, A_{0:t-1}) \nabla_{\theta} \log \pi (A_t | S_t; \theta)] &= E_{S_{t:\infty} \sim \pi_{\theta}} \left[\sum_a \pi (a | S_t; \theta) b_t (S_{0:t-1}, A_{0:t-1}) \frac{\nabla_{\theta} \pi (a | S_t; \theta)}{\pi (a | S_t; \theta)} \right] \\ &= E_{\pi_{\theta}} \left[E_{S_{t+1:\infty}, A_{t+1:\infty} \sim \pi_{\theta}} \left[\sum_a \pi (a | S_t; \theta) b_t (S_{0:t-1}, A_{0:t-1}) \frac{\nabla_{\theta} \pi (a | S_t; \theta)}{\pi (a | S_t; \theta)} \middle| S_t, A_t \right] \right] \\ &= E_{\pi_{\theta}} \left[E_{S_{t+1:\infty}, A_{t+1:\infty} \sim \pi_{\theta}} \left[\sum_a \pi (a | S_t; \theta) \frac{\nabla_{\theta} \pi (a | S_t; \theta)}{\pi (a | S_t; \theta)} \middle| S_t, A_t \right] b_t (S_{0:t-1}, A_{0:t-1}) \right] \\ &= E_{\pi_{\theta}} \left[E_{S_{t+1:\infty}, A_{t+1:\infty} \sim \pi_{\theta}} \left[\sum_a \nabla_{\theta} \pi (a | S_t; \theta) \middle| S_t, A_t \right] b_t (S_{0:t-1}, A_{0:t-1}) \right] \\ &= E_{\pi_{\theta}} \left[E_{S_{t+1:\infty}, A_{t+1:\infty} \sim \pi_{\theta}} \left[\sum_a \nabla_{\theta} 1 \middle| S_t, A_t \right] b_t (S_{0:t-1}, A_{0:t-1}) \right] \\ &= 0 \end{aligned}$$

□

Now, the next task is to proof (2.1) component in LHS equals the (2.3) component in RHS:

Theorem 2.

$$E_{\pi_\theta} [q_{\pi_\theta} (S_t, A_t) \nabla_\theta \log \pi (A_t | S_t; \theta)] = E_{\pi_\theta} [\hat{q}_{\pi_\theta} (S_{t:\infty}, A_{t:\infty}) \nabla_\theta \log \pi (A_t | S_t; \theta)]$$

Proof.

$$\begin{aligned} E_{\pi_\theta} [q_{\pi_\theta} (S_t, A_t) \nabla_\theta \log \pi (A_t | S_t; \theta)] &= E_{\pi_\theta} [\hat{q}_{\pi_\theta} (S_{t:\infty}, A_{t:\infty}) \nabla_\theta \log \pi (A_t | S_t; \theta)] \\ &= E_{\pi_\theta} \left[E_{S_{t+1:\infty}, A_{t+1:\infty}, \pi_\theta} [\hat{q}_{\pi_\theta} (S_{t:\infty}, A_{t:\infty}) \nabla_\theta \log \pi (A_t | S_t; \theta) | S_t, A_t] \right] \\ &= E_{\pi_\theta} \left[\underbrace{E_{S_{t+1:\infty}, A_{t+1:\infty}, \pi_\theta} [\hat{q}_{\pi_\theta} (S_{t:\infty}, A_{t:\infty}) | S_t, A_t]}_{q_{\pi_\theta} (S_t, A_t)} \nabla_\theta \log \pi (A_t | S_t; \theta) \right] \\ &= E_{\pi_\theta} [q_{\pi_\theta} (S_t, A_t) \nabla_\theta \log \pi (A_t | S_t; \theta)] = LHS \end{aligned}$$

□

Since we have proofed component (2.1) equals component (2.3), and component (2.2) equals component (2.4), we can claim that Eq. 4 equals Eq. 5. That is:

$$\nabla_\theta J(\theta) = (2.1) + (2.2) = (2.3) + (2.4) = E_{\pi_\theta} [\hat{A}_{\pi_\theta} (S_t, A_t) \nabla_\theta \log \pi (A_t | S_t; \theta)] \quad \square$$

SOLUTION-2.2:

From the TD error equation, Eq. 3, and $\hat{v} = v_{\pi_\theta}$, we have:

$$\delta_{\pi_\theta} = R_{t+1} + \gamma v_{\pi_\theta} (S_{t+1}) - v_{\pi_\theta} (S_t)$$

$$\begin{aligned} E_{\pi_\theta} [\delta_{\pi_\theta} | S_t, A_t] &= E_{\pi_\theta} [\delta_{\pi_\theta} | S_t, A_t] \\ &= E_{\pi_\theta} [R_{t+1} + \gamma v_{\pi_\theta} (S_{t+1}) - v_{\pi_\theta} (S_t) | S_t, A_t] \\ &= E_{\pi_\theta} [R_{t+1} + \gamma v_{\pi_\theta} (S_{t+1}) | S_t, A_t] - v_{\pi_\theta} (S_t) \end{aligned}$$

By definition of q :

$$E_{\pi_\theta} [R_{t+1} + \gamma v_{\pi_\theta} (S_{t+1}) | S_t, A_t] = E_{S_{t+1:\infty}, A_{t+1:\infty}, \pi_\theta} [\hat{q}_{\pi_\theta} (S_{t:\infty}, A_{t:\infty})] = q_{\pi_\theta} (S_t, A_t)$$

Finally, we have proofed that:

$$E_{\pi_\theta} [\delta_{\pi_\theta} | S_t, A_t] = \underbrace{E_{\pi_\theta} [R_{t+1} + \gamma v_{\pi_\theta} (S_{t+1}) | S_t, A_t]}_{q_{\pi_\theta} (S_t, A_t)} - v_{\pi_\theta} (S_t) = A_{\pi_\theta} (S_t, A_t) \quad \square$$