

Homework #15

Junwu Zhang

CME 241: Reinforcement Learning for Finance

March 04, 2020

Problem 1.

Write Proof (with precise notation) of the Policy Gradient Theorem

Solution.

For Policy Gradient (PG), we have the following crucial components to the problem:

- *States:* $s_t \in \mathcal{S}$
- *Actions:* $a_t \in \mathcal{A}$
- *Rewards:* $r_t \in \mathbb{R}, \forall t \in \{0, 1, 2, \dots\}$
- State Transition Probabilities $\mathcal{P}_{s,s'}^a = \Pr(s_{t+1} = s' | s_t = s, a_t = a)$
- Initial State Probability Distribution $p_0 : \mathcal{S} \rightarrow [0, 1]$
- Expected Rewards $\mathcal{R}_s^a = E[r_t | s_t = s, a_t = a]$
- Policy Function Approximation $\pi(s, a; \theta) = \Pr(a_t = a | s_t = s, \theta), \theta \in \mathbb{R}^k$
- *Discount Factor:* γ

The Policy Gradient Theorem (PGT) can be written as:

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \nabla_{\theta} \pi(s, a; \theta) \cdot Q^{\pi}(s, a) \cdot da \cdot ds \quad (1)$$

where $\rho^{\pi}(s) = \int_{\mathcal{S}} \sum_{t=0}^{\infty} \gamma^t \cdot p_0(s_0) \cdot p(s_0 \rightarrow s, t, \pi) \cdot ds_0$. We can it *Discounted-Aggregate State-Visitation Measure*, which is a key function for PG. $\rho^{\pi}(s)$ depends on θ . We also note that $\nabla_{\theta} \log \pi(s, a; \theta)$ is the Score Function.

To prove the PGT, we can first write:

$$J(\theta) = \int_{\mathcal{S}} p_0(s_0) \cdot V^{\pi}(s_0) \cdot ds_0 \quad (2)$$

$$= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \pi(s_0, a_0; \theta) \cdot Q^{\pi}(s_0, a_0) \cdot da_0 \cdot ds_0 \quad (3)$$

Then, by calculating $\pi(s_0, a_0; \theta)$ and $Q^{\pi}(s_0, a_0)$, we can write $\nabla_{\theta} J(\theta)$ as:

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_{\theta} \pi(s_0, a_0; \theta) \cdot Q^{\pi}(s_0, a_0) \cdot da_0 \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \pi(s_0, a_0; \theta) \cdot \nabla_{\theta} Q^{\pi}(s_0, a_0) \cdot da_0 \cdot ds_0 \end{aligned} \quad (4)$$

Using the Bellman operator, we can expand the action value function $Q^\pi(s_0, a_0)$ as:

$$Q^\pi(s_0, a_0) = \mathcal{R}_{s_0}^{a_0} + \int_{\mathcal{S}} \gamma \cdot \mathcal{P}_{s_0, s_1}^{a_0} \cdot V^\pi(s_1) \cdot ds_1 \quad (5)$$

Plugging Equation (5) into the original equation, we have:

$$\nabla_\theta J(\theta) = \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a) \cdot da_0 \cdot ds_0 \quad (6)$$

$$+ \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \pi(s_0, a_0; \theta) \cdot \nabla_\theta \left(\int_{\mathcal{S}} \gamma \cdot \mathcal{P}_{s_0, s_1}^{a_0} \cdot V^\pi(s_1) \cdot ds_1 \right) \cdot da_0 \cdot ds_0 \quad (7)$$

since $\nabla_\theta \mathcal{R}_{s_0}^{a_0} = 0$.

Moving ∇_θ into $\int_{\mathcal{S}}$, and further moving the outside $\int_{\mathcal{A}}$, we have:

$$\begin{aligned} \nabla_\theta J(\theta) &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a) \cdot da_0 \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \pi(s_0, a_0; \theta) \int_{\mathcal{S}} \gamma \cdot \mathcal{P}_{s_0, s_1}^{a_0} \cdot \nabla_\theta V^\pi(s_1) \cdot ds_1 \cdot da_0 \cdot ds_0 \\ &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a_0) \cdot da_0 \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} \left(\int_{\mathcal{S}} \gamma \cdot p_0(s_0) \int_{\mathcal{A}} \pi(s_0, a_0; \theta) \cdot \mathcal{P}_{s_0, s_1}^{a_0} \cdot da_0 \cdot ds_0 \right) \cdot \nabla_\theta V^\pi(s_1) \cdot ds_1 \end{aligned} \quad (8)$$

We can further expand the above equations to:

$$\begin{aligned} \nabla_\theta J(\theta) &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a_0) \cdot da_0 \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} \left(\int_{\mathcal{S}} \gamma \cdot p_0(s_0) \cdot p(s_0 \rightarrow s_1, 1, \pi) \cdot ds_0 \right) \cdot \nabla_\theta V^\pi(s_1) \cdot ds_1 \end{aligned} \quad (9)$$

based on the fact that:

$$\int_{\mathcal{A}} \pi(s_0, a_0; \theta) \cdot \mathcal{P}_{s_0, s_1}^{a_0} \cdot da_0 = p(s_0 \rightarrow s_1, 1, \pi) \quad (10)$$

Next, expand the above equations to:

$$\begin{aligned} \nabla_\theta J(\theta) &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a_0) \cdot da \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} \left(\int_{\mathcal{S}} \gamma \cdot p_0(s_0) p(s_0 \rightarrow s_1, 1, \pi) \cdot ds_0 \right) \cdot \nabla_\theta \left(\int_{\mathcal{A}} \pi(s_1, a_1; \theta) Q^\pi(s_1, a_1) \cdot da_1 \right) \cdot ds_1 \end{aligned} \quad (11)$$

based on the fact that:

$$V^\pi(s_1) = \int_{\mathcal{A}} \pi(s_1, a_1; \theta) \cdot Q^\pi(s_1, a_1) \cdot da_1 \quad (12)$$

We can then split πQ^π and calculate the gradient of Q^π using Bellman expansion, we have:

$$\begin{aligned}\nabla_\theta J(\theta) &= \int_{\mathcal{S}} p_0(s_0) \int_{\mathcal{A}} \nabla_\theta \pi(s_0, a_0; \theta) \cdot Q^\pi(s_0, a_0) \cdot da_0 \cdot ds_0 \\ &\quad + \int_{\mathcal{S}} \int_{\mathcal{S}} \gamma p_0(s_0) p(s_0 \rightarrow s_1, 1, \pi) ds_0 \left(\int_{\mathcal{A}} \nabla_\theta \pi(s_1, a_1; \theta) Q^\pi(s_1, a_1) da_1 + \dots \right) ds_1 \\ &= \sum_{t=0}^{\infty} \int_{\mathcal{S}} \int_{\mathcal{S}} \gamma^t \cdot p_0(s_0) \cdot p(s_0 \rightarrow s_t, t, \pi) \cdot ds_0 \int_{\mathcal{A}} \nabla_\theta \pi(s_t, a_t; \theta) \cdot Q^\pi(s_t, a_t) \cdot da_t \cdot ds_t\end{aligned}\tag{13}$$

Since we know $\int_{\mathcal{A}} \nabla_\theta \pi(s_t, a_t; \theta) \cdot Q^\pi(s_t, a_t) \cdot da_t$ is independent of t , we can move $\sum_{t=0}^{\infty}$ inside the integrals and write:

$$\nabla_\theta J(\theta) = \int_{\mathcal{S}} \int_{\mathcal{S}} \sum_{t=0}^{\infty} \gamma^t \cdot p_0(s_0) \cdot p(s_0 \rightarrow s, t, \pi) \cdot ds_0 \int_{\mathcal{A}} \nabla_\theta \pi(s, a; \theta) \cdot Q^\pi(s, a) \cdot da \cdot ds\tag{14}$$

Since we also have:

$$\int_{\mathcal{S}} \sum_{t=0}^{\infty} \gamma^t \cdot p_0(s_0) \cdot p(s_0 \rightarrow s, t, \pi) \cdot ds_0 \stackrel{\text{def}}{=} \rho^\pi(s)\tag{15}$$

We can finally write the overall equation as:

$$\nabla_\theta J(\theta) = \int_{\mathcal{S}} \rho^\pi(s) \int_{\mathcal{A}} \nabla_\theta \pi(s, a; \theta) \cdot Q^\pi(s, a) \cdot da \cdot ds\tag{16}$$

which is the same as Equation (1). ■

Problem 2.

Derive the score function for softmax policy (for finite set of actions)

Solution. Since we know that θ is n -vector and features vector is:

$$\phi(s, a) = (\phi_1(s, a), \dots, \phi_n(s, a)) \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}\tag{17}$$

we can weight the actions by doing a linear combination of all the features, namely $\theta^T \cdot \phi(s, a)$. Then, we can see that action probabilities are proportional to weights exponentiated, which can be written as:

$$\pi(s, a; \theta) = \frac{e^{\theta^T \cdot \phi(s, a)}}{\sum_b e^{\theta^T \cdot \phi(s, b)}} \text{ for all } s \in \mathcal{S}, a \in \mathcal{A}\tag{18}$$

By the definition of score function, we can take the logarithm of the above equation, and have:

$$\nabla_\theta \log \pi(s, a; \theta) = \phi(s, a) - \sum_b \pi(s, b; \theta) \cdot \phi(s, b) = \phi(s, a) - \mathbb{E}_\pi[\phi(s, \cdot)]\tag{19}$$

■

Problem 3.

Derive the score function for gaussian policy (for continuous actions)

Solution. Slightly differs from the last question, since the action space is now continuous, we have the feature vector as:

$$\phi(s) = (\phi_1(s), \dots, \phi_n(s)) \text{ for all } s \in \mathcal{S} \quad (20)$$

Similar to the last question, we have Gaussian Mean as $\theta^T \cdot \phi(s)$ and with the Gaussian policy $a \sim \mathcal{N}(\theta^T \cdot \phi(s), \sigma^2)$ for all $s \in \mathcal{S}$, we have the score function for continuous action space as:

$$\nabla_{\theta} \log \pi(s, a; \theta) = \frac{(a - \theta^T \cdot \phi(s)) \cdot \phi(s)}{\sigma^2} \quad (21)$$

■

Problem 4.

Write code for the REINFORCE Algorithm (Monte-Carlo Policy Gradient Algorithm, i.e., no Critic)

Solution. Code is attached separately.

■

Problem 5.

Write Proof (with proper notation) of the Compatible Function Approximation Theorem

Solution. The *Compatible Function Approximation Theorem* states: If the following two conditions are satisfied:

- Critic gradient is compatible with the Actor score function

$$\nabla_w Q(s, a; w) = \nabla_{\theta} \log \pi(s, a; \theta) \quad (22)$$

- Critic parameters w minimize the following mean-squared error:

$$\epsilon = \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \pi(s, a; \theta) (Q^{\pi}(s, a) - Q(s, a; w))^2 \cdot da \cdot ds \quad (23)$$

Then the Policy Gradient using critic $Q(s, a; w)$ is exactly:

$$\nabla_{\theta} J(\theta) = \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \nabla_{\theta} \pi(s, a; \theta) \cdot Q(s, a; w) \cdot da \cdot ds \quad (24)$$

To prove this, we can first use the second point in the theorem and know that:

$$\begin{aligned} \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot (Q^{\pi}(s, a) - Q(s, a; w)) \cdot \nabla_w Q(s, a; w) \cdot da \cdot ds &= 0 \\ \forall w \text{ that minimizes } \epsilon &= \int_{\mathcal{S}} \rho^{\pi}(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot (Q^{\pi}(s, a) - Q(s, a; w))^2 \cdot da \cdot ds \end{aligned} \quad (25)$$

Using the first bullet point of the theorem which is:

$$\nabla_w Q(s, a; w) = \nabla_\theta \log \pi(s, a; \theta), \quad (26)$$

we can replace the subscript and have:

$$\int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot (Q^\pi(s, a) - Q(s, a; w)) \cdot \nabla_\theta \log \pi(s, a; \theta) \cdot da \cdot ds = 0 \quad (27)$$

Therefore, we can see that Equation (25) and Equation (27) are equal:

$$\begin{aligned} & \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot Q^\pi(s, a) \cdot \nabla_\theta \log \pi(s, a; \theta) \cdot da \cdot ds \\ &= \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot Q(s, a; w) \cdot \nabla_\theta \log \pi(s, a; \theta) \cdot da \cdot ds \end{aligned} \quad (28)$$

Since we also know:

$$\nabla_\theta J(\theta) = \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot Q^\pi(s, a) \cdot \nabla_\theta \log \pi(s, a; \theta) \cdot da \cdot ds \quad (29)$$

Combining Equation (28) and Equation (29), we have:

$$\begin{aligned} \nabla_\theta J(\theta) &= \int_S \rho^\pi(s) \int_{\mathcal{A}} \pi(s, a; \theta) \cdot Q(s, a; w) \cdot \nabla_\theta \log \pi(s, a; \theta) \cdot da \cdot ds \\ &= \int_S \rho^\pi(s) \int_{\mathcal{A}} \nabla_\theta \pi(s, a; \theta) \cdot Q(s, a; w) \cdot da \cdot ds \end{aligned} \quad (30)$$

and we can see that this is the same as Equation (24), therefore the theorem is proved. ■