
Revisit Policy Optimization in Matrix Form

Sitao Luan^{1,2}, Xiao-Wen Chang¹, Doina Precup^{1,2,3}
{sitao.luan@mail, chang@cs, dprecup@cs}.mcgill.ca
¹McGill University; ²Mila; ³DeepMind

Abstract

In tabular case, when the reward and environment dynamics are known, policy evaluation can be written as $V_\pi = (I - \gamma P_\pi)^{-1} r_\pi$, where P_π is the state transition matrix given policy π and r_π is the reward signal given π . What annoys us is that P_π and r_π are both mixed with π , which means every time when we update π , they will change together. In this paper, we leverage the notation from [4] to disentangle π and environment dynamics which makes optimization over policy more straightforward. We show that policy gradient theorem [3] and TRPO [2] can be written in a more general notation framework and such framework has good potential to be extended to model-based reinforcement learning.

1 Introduction

1.1 Bellman Equation in Matrix Form

Markov decision process (MDP) is a framework to model the learning process that the agent learns from the interaction with the environment [3]. The interaction happens in discrete time steps, $t = 0, 1, 2, 3, \dots$. At step t , given a state $S_t = s_t \in \mathcal{S}$, the agent picks an action $a_t \in \mathcal{A}(s_t)$ according to a policy $\pi(\cdot|s_t)$, which is a rule of choosing actions given a state. Then, at time $t + 1$, the environmental dynamics $p : \mathcal{S} \times \mathcal{R} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ takes the agent to a new state $S_{t+1} = s_{t+1} \in \mathcal{S}$ and provide a numerical reward $R_{t+1} = r_{t+1}(s_t, a_t, s_{t+1}) \in \mathbb{R}$. Such a sequence of interactions gives us a trajectory $\tau = \{S_0, A_0, R_1, S_1, A_1, R_2, S_2, A_2, R_3, \dots\}$. Our objective is to find an optimal policy to maximize the expected long-term discounted cumulative rewards $V_\pi(s) = E_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$ for each state s or $Q_\pi(s, a) = E_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, a_t = a]$ for each state-action pair (s, a) , where γ is the discount factor. The Bellman equation for $V_\pi(s)$ can be written as follows:

$$\begin{aligned} V_\pi(s) &= \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_\pi(s')] \\ &= \sum_a \pi(a|s) \sum_{s', r} r \cdot p(s', r|s, a) + \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) \gamma V_\pi(s') \\ &= \sum_a \pi(a|s) r(s, a) + \gamma \sum_{s'} \{ \sum_a \pi(a|s) \cdot p(s'|s, a) \} V_\pi(s') \\ &= r_\pi(s) + \gamma \sum_{s'} P_\pi(s'|s) V_\pi(s') \quad \text{for all } s \in \mathcal{S} \end{aligned} \tag{1}$$

Where $r_\pi(s)$ is the expected immediate reward at state s under π and $P_\pi(s'|s)$ is the transition probability of moving from s to s' under π . Notice that the transition probability is a combination of

policy and the environment dynamics. We can write V_π, r_π, P_π in matrix form

$$\mathbf{V}_\pi = \begin{bmatrix} V_\pi(s_1) \\ V_\pi(s_2) \\ \vdots \\ V_\pi(s_n) \end{bmatrix}, \mathbf{r}_\pi = \begin{bmatrix} r_\pi(s_1) \\ r_\pi(s_2) \\ \vdots \\ r_\pi(s_n) \end{bmatrix} = \begin{bmatrix} \sum_a \pi(a|s_1)r(s_1, a) \\ \sum_a \pi(a|s_2)r(s_2, a) \\ \vdots \\ \sum_a \pi(a|s_n)r(s_n, a) \end{bmatrix}, P_\pi = \begin{bmatrix} P_\pi(s_1|s_1) & P_\pi(s_2|s_1) & \cdots & P_\pi(s_n|s_1) \\ P_\pi(s_1|s_2) & P_\pi(s_2|s_2) & \cdots & P_\pi(s_n|s_2) \\ \vdots & \vdots & \ddots & \vdots \\ P_\pi(s_1|s_n) & P_\pi(s_2|s_n) & \cdots & P_\pi(s_n|s_n) \end{bmatrix} \quad (2)$$

Then the Bellman equation can be rewritten in matrix form

$$\mathbf{V}_\pi = \mathbf{r}_\pi + \gamma P_\pi \mathbf{V}_\pi$$

Assume that P_π and \mathbf{r}_π are given and $I - \gamma P_\pi$ is nonsingular, then

$$\mathbf{V}_\pi = (I - \gamma P_\pi)^{-1} \mathbf{r}_\pi \quad (3)$$

For later uses, we define the vector $\mathbf{Q}_\pi = [\mathbf{Q}(s_1, \cdot)^T, \mathbf{Q}(s_2, \cdot)^T, \dots, \mathbf{Q}(s_n, \cdot)^T]^T$.

1.2 Problem

One problem with (3) is that P_π and \mathbf{r}_π are both dependent on π , which means every time when π changes, we need to reconstruct them. In addition, although \mathbf{V}_π is a function of π , we cannot write the function in terms of π explicitly. What we desire is to have

$$\mathbf{V}_\pi = f(\pi|P, r) \quad \text{or} \quad \mathbf{V}_{\pi_\theta} = f(\pi_\theta|P, r) \quad (4)$$

where π_θ is a parameterized policy and P, r only depend on environment dynamics, instead of $\mathbf{V}_\pi = f(\pi|P_\pi, \mathbf{r}_\pi)$ or $\mathbf{V}_{\pi_\theta} = f(\pi_\theta|P_{\pi_\theta}, \mathbf{r}_{\pi_\theta})$. Then, we can directly write our objective as a function of π , which can make optimization more straightforward and efficient, just as follows,

$$\operatorname{argmax}_\pi \mathbf{V}_\pi = f(\pi|P, r) \quad \text{or} \quad \operatorname{argmax}_\theta \mathbf{V}_{\pi_\theta} = f(\pi_\theta|P, r) \quad (5)$$

and solve it via gradient descent or other techniques. The notations from [4] is helpful to construct (4)(5). We introduce the formulation in the following sections and try to rewrite TRPO [2] with it.

2 Notation and Preliminary

2.1 Notation and Properties

For simplicity, we assume both the number of states in \mathcal{S} and the number of actions in \mathcal{A} are finite. We define the notations as follows [4]:

- $|\mathcal{S}|$ and $|\mathcal{A}|$ denote the number of states in \mathcal{S} and the number of actions in \mathcal{A} , respectively.
- $P \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{S}|}$ is a transition matrix whose entries are $P_{(sa, s')} = p(s'|s, a)$, where $p(s'|s, a) \geq 0$, $\sum_{s'} p(s'|s, a) = 1$, for all s and a , *i.e.*,

$$P \mathbf{1}_{|\mathcal{S}|} = \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|}, \quad \mathbf{1}_{|\mathcal{S}|} = [1, 1, \dots, 1]^T \in \mathbb{R}^{|\mathcal{S}| \times 1} \quad (6)$$

- $\pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times 1}$ is a stationary policy, whose entries are $\pi_{(sa)} = \pi(a|s)$, where $\sum_a \pi_{(sa)} = 1$. It is convenient to rewrite the policy as a matrix $\Pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{A}|}$, where $\Pi_{(s, s', a)} = \pi_{(sa)}$ if $s' = s$, otherwise 0, *i.e.*

$$\Pi = \operatorname{diag}(\pi(\cdot|s_1)^T, \dots, \pi(\cdot|s_n)^T), \quad \pi(\cdot|s_i)^T \in \mathbb{R}^{1 \times |\mathcal{A}|} \quad (7)$$

$\Xi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}| \times |\mathcal{A}|}$ is an auxiliary (marginalization) matrix defined by

$$\Xi = \operatorname{diag}(\mathbf{1}_{|\mathcal{A}|}^T, \dots, \mathbf{1}_{|\mathcal{A}|}^T) \quad (8)$$

It is easy to verify that

$$\Pi \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|} = \mathbf{1}_{|\mathcal{S}|}, \quad \Pi \Xi^T = I_{|\mathcal{S}|} \quad (9)$$

We can easily reconstruct the state-to-state transition matrix by $P_\pi = \Pi P \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$, and state-action-to-state-action transition matrix by $P'_\pi = P \Pi \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}| \times |\mathcal{S}| \times |\mathcal{A}|}$. From (6) and (9),

$$\begin{aligned} P_\pi \mathbf{1}_{|\mathcal{S}|} &= \Pi P \mathbf{1}_{|\mathcal{S}|} = \Pi \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|} = \mathbf{1}_{|\mathcal{S}|} \\ P'_\pi \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|} &= P \Pi \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|} = P \mathbf{1}_{|\mathcal{S}|} = \mathbf{1}_{|\mathcal{S}| \times |\mathcal{A}|} \end{aligned}$$

- $\rho_0 \in \mathbb{R}^{|\mathcal{S}| \times 1}$ is the initial state distribution. $\mu_\pi \in \mathbb{R}^{|\mathcal{S}| \times 1}$ is the steady-state distribution of π given environment P which satisfies $\mu_\pi^T \Pi P = \mu_\pi^T$. $\rho_\pi^T = \rho_0^T \sum_{i=0}^{\infty} \gamma^i (\Pi P)^i = \rho_0^T (I - \gamma \Pi P)^{-1}$ is the discounted visitation frequency vector under π and

$$\rho_\pi^T \mathbf{1} = \rho_0^T \sum_{i=0}^{\infty} \gamma^i (\Pi P)^i \mathbf{1} = \rho_0^T (\sum_{i=0}^{\infty} \gamma^i) \mathbf{1} = 1/(1 - \gamma)$$

- $\mathbf{r} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$ is the average reward vector whose entries are $\mathbf{r}_{(sa)} = r(s, a) = E[r|s, a] = \sum_{s', r} r \cdot p(s', r|s, a)$, which specify the average reward obtained when taking action a in state s . From the definition of \mathbf{r}_π in (2) and Π in (7), it is easy to see that

$$\mathbf{r}_\pi = \Pi \mathbf{r} \quad (10)$$

Now we can rewrite the Bellman equation. From (2), we have

$$\mathbf{V}_\pi = (I - \gamma P_\pi)^{-1} \mathbf{r}_\pi = (I - \gamma \Pi P)^{-1} \Pi \mathbf{r} = \sum_{i=0}^{\infty} \gamma^i (\Pi P)^i \Pi \mathbf{r} = \Pi \mathbf{r} + \gamma \Pi P \mathbf{V}_\pi \quad (11)$$

Similarly, we have $\mathbf{Q}_\pi = \sum_{i=0}^{\infty} \gamma^i (P \Pi)^i \mathbf{r} = \mathbf{r} + \gamma P \Pi \mathbf{Q}_\pi$. One can easily check that the relations between \mathbf{V}_π and \mathbf{Q}_π are as follows,

$$\mathbf{V}_\pi = \Pi \mathbf{Q}_\pi, \quad \mathbf{Q}_\pi = \sum_{i=0}^{\infty} \gamma^i (P \Pi)^i \mathbf{r} = \mathbf{r} + \gamma P [\sum_{i=0}^{\infty} \gamma^i (P \Pi)^i \Pi \mathbf{r}] = \mathbf{r} + \gamma P \mathbf{V}_\pi \quad (12)$$

2.2 A Recap of TRPO

We do a simple recap of TRPO [2] in this subsection. We will show how these results can easily be obtained when written in matrix form and how to extend them.

For policy-based reinforcement learning algorithms, if we use policy gradient with a fixed learning rate to update the policy, it always happens that the learning rate is sometimes too large that we will get a worse policy. This oscillation makes the training unstable and the convergence slow. So we want to find a way that our policy is guaranteed to improve after each policy update even the improvement might be quite small at some time. The objective $\eta(\pi)$ is defined as follows,

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right], \text{ where } s_0 \sim \rho_0(s_0), a_t \sim \pi(a_t|s_t), s_{t+1} \sim p(s_{t+1}|s_t, a_t) \quad (13)$$

With the definition of the advantage function $\mathbf{A}_\pi(s, a)$:

$$\mathbf{A}_\pi(s, a) = \mathbf{Q}_\pi(s, a) - \mathbf{V}_\pi(s) \quad (14)$$

$\eta(\tilde{\pi})$ can be split into two parts,

$$\eta(\tilde{\pi}) = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \tilde{\pi}} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{A}_\pi(s_t, a_t) \right] = \eta(\pi) + \sum_s \rho_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) \mathbf{A}_\pi(s, a) \quad (15)$$

where $\rho_{\tilde{\pi}}(s) = \sum_{i=0}^{\infty} \gamma^i p(s_i = s | \tilde{\pi}, P, \rho_0^T)$ is the discounted visitation frequency of s under $\tilde{\pi}$. A local approximation to $\eta(\tilde{\pi})$ is constructed as follows:

$$L_\pi(\tilde{\pi}) = \eta(\pi) + \sum_s \rho_\pi(s) \sum_a \tilde{\pi}(a|s) \mathbf{A}_\pi(s, a) \quad (16)$$

It satisfies two important properties when π_θ is parameterized by θ :

$$L_{\pi_{\theta_0}}(\pi_{\theta_0}) = \eta(\pi_{\theta_0}), \quad \nabla_\theta L_{\pi_{\theta_0}}(\pi_\theta) \Big|_{\theta=\theta_0} = \nabla_\theta \eta(\pi_\theta) \Big|_{\theta=\theta_0} \quad (17)$$

Given the total variation divergence for two discrete distributions $D_{TV}(p||q) = \frac{1}{2} \sum_i |p_i - q_i|$, we define $D_{TV}^{\max}(\pi, \tilde{\pi}) = \max_s D_{TV}(\pi(\cdot|s) || \tilde{\pi}(\cdot|s))$. It can be proved that

$$\eta(\tilde{\pi}) \geq L_\pi(\tilde{\pi}) - \frac{4\epsilon\gamma}{(1-\gamma)^2} \alpha^2, \text{ where } \alpha = D_{TV}^{\max}(\pi, \tilde{\pi}), \epsilon = \max_{s,a} |\mathbf{A}_\pi(s, a)| \quad (18)$$

With $D_{TV}(p||q)^2 \leq D_{KL}(p||q)$, we define $D_{KL}^{\max}(\pi, \tilde{\pi}) = \max_s D_{KL}(\pi(\cdot|s) || \tilde{\pi}(\cdot|s))$. Then we have

$$\eta(\tilde{\pi}) \geq L_\pi(\tilde{\pi}) - CD_{KL}^{\max}(\pi, \tilde{\pi}), \text{ where } C = \frac{4\epsilon\gamma}{(1-\gamma)^2} \quad (19)$$

Let $M_i(\pi) = L_{\pi_i}(\pi) - CD_{KL}^{\max}(\pi_i, \pi)$, we have

$$\eta(\pi_{i+1}) \geq M_i(\pi_{i+1}), \eta(\pi_i) = M_i(\pi_i) \Rightarrow \eta(\pi_{i+1}) - \eta(\pi_i) \geq M_i(\pi_{i+1}) - M_i(\pi_i)$$

So if we define π_{i+1} as $\pi_{i+1} = \underset{\pi}{\operatorname{argmax}} M_i(\pi)$, then

$$\eta(\pi_0) \leq \eta(\pi_1) \leq \eta(\pi_2) \leq \dots \quad (20)$$

Therefore, $M_i(\pi)$ becomes a surrogate function that we want to maximize.

2.3 Relations Between Several Distance Measures

The total variation distance between two probability measures μ and ν on a sigma-algebra \mathcal{F} of subsets of the sample space Ω is defined as $\delta(\mu, \nu) = \sup_{A \in \mathcal{F}} |\mu(A) - \nu(A)|$. The total variation distance is related to the Kullback–Leibler divergence by Pinsker’s inequality:

$$\delta(\mu, \nu) \leq \sqrt{\frac{1}{2} D_{\text{KL}}(\mu, \nu)}$$

And from [1], we have

$$\delta(\mu, \nu) = \frac{1}{2} \|\mu - \nu\|_1 = \frac{1}{2} \sum_{\omega \in \Omega} |\mu(\omega) - \nu(\omega)|$$

Thus,

$$\|\mu - \nu\|_1^2 \leq 2D_{\text{KL}}(\mu, \nu)$$

3 TRPO in Matrix Form

In this section, we first write TRPO in matrix form and we will see there are more ways to find the local approximation as (16). We derive several interesting properties of these approximations. The norm $\|\cdot\|$ we use in this section is 1-norm. Note that Π (matrix) is just a rewriting of π (vector). They represent the same policy. This relation is the same for parameterized Π_θ and π_θ .

With notations introduced in Section 2.1, η can be written as

$$\eta(\pi) = \rho_0^\top \mathbf{V}_\pi = \rho_0^\top \sum_{i=0}^{\infty} \gamma^i (\Pi P)^i \Pi \mathbf{r}$$

From (14), we see that the advantage function \mathbf{A}_π , the vector form of $\mathbf{A}_\pi(s, a)$, can be written as

$$\mathbf{A}_\pi = \mathbf{Q}_\pi - \Xi^T \mathbf{V}_\pi$$

From (9),(11) and (12), it is easy to verify that $\Pi \mathbf{A}_\pi = 0$. Then

$$(I - \gamma P \Pi)^{-1} \mathbf{A}_\pi = \sum_{i=0}^{\infty} \gamma^i (P \Pi)^i \mathbf{A}_\pi = \mathbf{A}_\pi + \sum_{i=1}^{\infty} \gamma^i (P \Pi)^i \mathbf{A}_\pi = \mathbf{A}_\pi \quad (21)$$

Since $\eta(\tilde{\pi}) = \rho_0^\top \mathbf{V}_{\tilde{\pi}}$ and $\eta(\pi) = \rho_0^\top \mathbf{V}_\pi$, (15) can easily be shown as follows,

$$\begin{aligned} \mathbb{E}_{s_0, a_0, \dots, \tilde{\pi}} [\sum_{i=0}^{\infty} \gamma^i \mathbf{A}_\pi(s_t, a_t)] &= \rho_0^\top \tilde{\Pi} (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_\pi = \rho_0^\top \tilde{\Pi} (I - \gamma P \tilde{\Pi})^{-1} [\mathbf{Q}_\pi - \Xi^T \mathbf{V}_\pi] \\ &= \rho_0^\top \tilde{\Pi} (I - \gamma P \tilde{\Pi})^{-1} [\mathbf{r} + \gamma P \mathbf{V}_\pi - \Xi^T \mathbf{V}_\pi] \\ &= \rho_0^\top \tilde{\Pi} \sum_{i=0}^{\infty} \gamma^i (P \tilde{\Pi})^i \mathbf{r} + \rho_0^\top \tilde{\Pi} \sum_{i=0}^{\infty} \gamma^i (P \tilde{\Pi})^i \gamma P \mathbf{V}_\pi - \rho_0^\top \tilde{\Pi} \sum_{i=0}^{\infty} \gamma^i (P \Pi)^i \Xi^T \mathbf{V}_\pi \\ &= \rho_0^\top \sum_{i=0}^{\infty} \gamma^i (\tilde{\Pi} P)^i \tilde{\Pi} \mathbf{r} + \rho_0^\top \sum_{i=0}^{\infty} \gamma^{i+1} (\tilde{\Pi} P)^{i+1} \mathbf{V}_\pi - \rho_0^\top \sum_{i=0}^{\infty} \gamma^i (\Pi P)^i \mathbf{V}_\pi \\ &= \rho_0^\top \mathbf{V}_{\tilde{\pi}} - \rho_0^\top \mathbf{V}_\pi = \eta(\tilde{\pi}) - \eta(\pi) \end{aligned}$$

Thus (4) can be written in the following form

$$\eta(\tilde{\pi}) = \eta(\pi) + \rho_0^\top \tilde{\Pi} (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_\pi = \eta(\pi) + \rho_0^\top (I - \gamma \tilde{\Pi} P)^{-1} \tilde{\Pi} \mathbf{A}_\pi \equiv \eta(\pi) + f_\pi(\tilde{\pi})$$

It is easy to see that $f_\pi(\pi) = 0$. Suppose $\tilde{\Pi} = \Pi + d\Pi$ ($d\Pi \rightarrow 0$ and $d\Pi \mathbf{1}_{|S||A| \times 1} = 0$) and if we want $\eta(\tilde{\pi}) \geq \eta(\pi)$, we should have $f_\pi(\tilde{\pi}) = f_\pi(\pi + d\pi) \geq 0$. Note that

$$\begin{aligned} f_\pi(\tilde{\pi}) &= f_\pi(\tilde{\pi}) - f_\pi(\pi) = df_\pi(\tilde{\pi}) \Big|_{\tilde{\pi}=\pi} \\ &= \text{trace} \left(\rho_0^\top (I - \gamma \tilde{\Pi} P)^{-1} \gamma (d\Pi) P (I - \gamma \tilde{\Pi} P)^{-1} \tilde{\Pi} \mathbf{A}_\pi + \rho_0^\top (I - \gamma \tilde{\Pi} P)^{-1} (d\Pi) \mathbf{A}_\pi \right) \Big|_{\tilde{\pi}=\pi} \quad (22) \\ &= \text{trace} \left(\rho_0^\top (I - \gamma \Pi P)^{-1} (d\Pi) \mathbf{A}_\pi \right) = \text{trace} \left(\mathbf{A}_\pi \rho_0^\top (I - \gamma \Pi P)^{-1} d\Pi \right) \end{aligned}$$

Therefore,

$$\left(\nabla_{\tilde{\Pi}} f_\pi(\tilde{\pi}) \Big|_{\tilde{\pi}=\pi} \right)^T = \mathbf{A}_\pi \rho_0^\top (I - \gamma \Pi P)^{-1} \quad (23)$$

There are 6 ways to set the approximation function as (16):

$$\begin{aligned}
L_{\pi}^1(\tilde{\pi}) &= \eta(\pi) + \rho_0^T \tilde{\Pi} (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi} = \eta(\pi) + \rho_0^T (I - \gamma \tilde{\Pi} P)^{-1} \tilde{\Pi} \mathbf{A}_{\pi} \quad (\text{policy gradient}) \\
L_{\pi}^2(\tilde{\pi}) &= \eta(\pi) + \rho_0^T \Pi (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi} \\
L_{\pi}^3(\tilde{\pi}) &= \eta(\pi) + \rho_0^T \tilde{\Pi} (I - \gamma P \Pi)^{-1} \mathbf{A}_{\pi} \\
L_{\pi}^4(\tilde{\pi}) &= \eta(\pi) + \rho_0^T (I - \gamma \Pi P)^{-1} \tilde{\Pi} \mathbf{A}_{\pi} \quad (\text{TRPO}) \\
L_{\pi}^5(\tilde{\pi}) &= \eta(\pi) + \rho_0^T (I - \gamma \tilde{\Pi} P)^{-1} \Pi \mathbf{A}_{\pi} = \eta(\pi) \quad (\text{trivial}) \\
L_{\pi}^6(\tilde{\pi}) &= \eta(\pi) + \rho_0^T \Pi (I - \gamma P \Pi)^{-1} \mathbf{A}_{\pi} = \eta(\pi) + \rho_0^T (I - \gamma \Pi P)^{-1} \Pi \mathbf{A}_{\pi} = \eta(\pi) \quad (\text{trivial})
\end{aligned}$$

We will discuss $L_{\pi}^2(\tilde{\pi})$, $L_{\pi}^3(\tilde{\pi})$ and $L_{\pi}^4(\tilde{\pi})$ in the following subsections.

3.1 Approximation Function $L_{\pi}^2(\tilde{\pi})$

For $L_{\pi}^2(\tilde{\pi}) = \eta(\pi) + \rho_0^T \Pi (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi}$, it is easy to see that $L_{\pi}^2(\pi) = 0$. For parametrized π_{θ} , to test (17), we have

$$\begin{aligned}
dL_{\pi_{\theta_0}}^2(\pi_{\theta}) \Big|_{\theta=\theta_0} &= d\rho_0^T \Pi_{\theta_0} (I - \gamma P \Pi_{\theta})^{-1} \mathbf{A}_{\pi_{\theta_0}} \Big|_{\theta=\theta_0} \\
&= \text{trace} \left(\rho_0^T \Pi_{\theta_0} (I - \gamma P \Pi_{\theta})^{-1} \gamma P (d\Pi_{\theta}) (I - \gamma P \Pi_{\theta})^{-1} \mathbf{A}_{\pi_{\theta_0}} \right) \Big|_{\theta=\theta_0} \\
&= \text{trace} \left(\mathbf{A}_{\pi_{\theta_0}} \rho_0^T \gamma \Pi_{\theta_0} P (I - \gamma \Pi_{\theta_0} P)^{-1} (d\Pi_{\theta}) \right) \Big|_{\theta=\theta_0}
\end{aligned} \tag{24}$$

It is easy to see

$$\left(\nabla_{\theta} L_{\pi_{\theta_0}}^2(\pi_{\theta}) \Big|_{\theta=\theta_0} \right)^T = \mathbf{A}_{\pi_{\theta_0}} \rho_0^T \gamma \Pi_{\theta_0} P (I - \gamma \Pi_{\theta_0} P)^{-1} \frac{d\Pi_{\theta}}{d\theta} \Big|_{\theta=\theta_0} \neq \left(\nabla_{\theta} \eta(\pi_{\theta}) \Big|_{\theta=\theta_0} \right)^T \tag{25}$$

which means $L_{\pi}^2(\tilde{\pi})$ does not match $\eta(\pi)$ to the first order. To calculate the difference between $L_{\pi}^2(\tilde{\pi})$ and $\eta(\pi)$, we have

$$\eta(\tilde{\pi}) - L_{\pi}^2(\tilde{\pi}) = \rho_0^T (\tilde{\Pi} - \Pi) (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi} = \rho_0^T (d\Pi) (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi}$$

Then we have

$$\begin{aligned}
\|\eta(\tilde{\pi}) - L_{\pi}^2(\tilde{\pi})\| &= \|\rho_0^T (d\Pi) (I - \gamma P \tilde{\Pi})^{-1} \mathbf{A}_{\pi}\| \leq \|\rho_0^T\| \|d\Pi\| \|\sum_{i=0}^{\infty} \gamma^i (P \tilde{\Pi})^i\| \|\mathbf{A}_{\pi}\| \\
&\leq \frac{\sqrt{2D_{\text{KL}}^{\max}(\pi, \tilde{\pi})}}{1 - \gamma} \|\mathbf{A}_{\pi}\|
\end{aligned}$$

3.2 Approximation Function $L_{\pi}^3(\tilde{\pi})$

For $L_{\pi}^3(\tilde{\pi}) = \eta(\pi) + \rho_0^T \tilde{\Pi} (I - \gamma P \Pi)^{-1} \mathbf{A}_{\pi}$, it is easy to see that $L_{\pi}^3(\pi) = 0$. Moreover, to test (16), we have

$$\begin{aligned}
dL_{\pi_{\theta_0}}^3(\pi_{\theta}) \Big|_{\theta=\theta_0} &= d\rho_0^T \tilde{\Pi}_{\theta_0} (I - \gamma P \Pi_{\theta_0})^{-1} \mathbf{A}_{\pi_{\theta_0}} \Big|_{\theta=\theta_0} \\
&= \text{trace} \left(\rho_0^T d\Pi_{\theta_0} \mathbf{A}_{\pi_{\theta_0}} \right) \Big|_{\theta=\theta_0} = \text{trace} \left(\mathbf{A}_{\pi_{\theta_0}} \rho_0^T d\Pi_{\theta_0} \right) \Big|_{\theta=\theta_0}
\end{aligned} \tag{26}$$

And

$$\left(\nabla_{\theta} L_{\pi_{\theta_0}}^3(\pi_{\theta}) \Big|_{\theta=\theta_0} \right)^T = \mathbf{A}_{\pi_{\theta_0}} \rho_0^T \frac{d\Pi_{\theta}}{d\theta} \Big|_{\theta=\theta_0} \neq \nabla_{\theta} \eta(\pi_{\theta}) \Big|_{\theta=\theta_0} \tag{27}$$

which means $L_{\pi}^3(\tilde{\pi})$ does not match $\eta(\pi)$ to the first order. But it is easy to check,

$$\nabla_{\theta} L_{\pi_{\theta_0}}^2(\pi_{\theta}) \Big|_{\theta=\theta_0} + \nabla_{\theta} L_{\pi_{\theta_0}}^3(\pi_{\theta}) \Big|_{\theta=\theta_0} = \nabla_{\theta} \eta(\pi_{\theta}) \Big|_{\theta=\theta_0} = \nabla_{\theta} L_{\pi_{\theta_0}}^4(\pi_{\theta}) \Big|_{\theta=\theta_0} \tag{28}$$

which means the gradient of $L_{\pi_{\theta_0}}^2(\pi_{\theta})$ and $L_{\pi_{\theta_0}}^3(\pi_{\theta})$ is a gradient decomposition of $\eta(\pi_{\theta})$ at θ_0 . To calculate the difference between $L_{\pi}^3(\tilde{\pi})$ and $\eta(\pi)$, we have

$$\begin{aligned}
\eta(\tilde{\pi}) - L_{\pi}^3(\tilde{\pi}) &= \rho_0^T \tilde{\Pi} \left((I - \gamma P \tilde{\Pi})^{-1} - (I - \gamma P \Pi)^{-1} \right) \mathbf{A}_{\pi} = \rho_0^T \tilde{\Pi} \left((I - \gamma P \Pi)^{-1} (\gamma P d\Pi) (I - \gamma P \Pi)^{-1} \right) \mathbf{A}_{\pi} \\
&= \rho_0^T \tilde{\Pi} \left((I - \gamma P \Pi)^{-1} (\gamma P d\Pi) \right) \mathbf{A}_{\pi}
\end{aligned}$$

Then we have

$$\begin{aligned} \|\eta(\tilde{\pi}) - L_{\pi}^3(\tilde{\pi})\| &= \left\| \rho_0^T \tilde{\Gamma} \left((I - \gamma P \Pi)^{-1} (\gamma P d \Pi) \right) \mathbf{A}_{\pi} \right\| \leq \left\| \rho_0^T \right\| \|d \Pi\| \left\| \sum_{i=0}^{\infty} \gamma^{i+1} (P \Pi)^i P \right\| \|\tilde{\Gamma}\| \|\mathbf{A}_{\pi}\| \\ &= \frac{\gamma \|d \Pi\| \|\mathbf{A}_{\pi}\|}{1 - \gamma} \leq \frac{\gamma \sqrt{2 D_{\text{KL}}^{\max}(\pi, \tilde{\pi})} \|\mathbf{A}_{\pi}\|}{1 - \gamma} \end{aligned}$$

3.3 Approximation Function $L_{\pi}^4(\tilde{\pi})$

If we set

$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + \rho_0^T (I - \gamma P \Pi)^{-1} \tilde{\Gamma} \mathbf{A}_{\pi} = \eta(\pi) + f_{\pi}(\tilde{\pi})$$

Then,

$$\begin{aligned} \eta(\tilde{\pi}) - L_{\pi}^4(\tilde{\pi}) &= \rho_0^T \left((I - \gamma P \tilde{\Gamma})^{-1} - (I - \gamma P \Pi)^{-1} \right) \tilde{\Gamma} \mathbf{A}_{\pi} \\ &= \rho_0^T \left((I - \gamma P (\Pi + d \Pi))^{-1} - (I - \gamma P \Pi)^{-1} \right) (\Pi + d \Pi) \mathbf{A}_{\pi} \\ &= \rho_0^T (I - \gamma P \Pi)^{-1} (\gamma P d \Pi) (I - \gamma P \Pi)^{-1} d \Pi \mathbf{A}_{\pi} \end{aligned}$$

Then we have

$$\begin{aligned} \|\eta(\tilde{\pi}) - L_{\pi}^4(\tilde{\pi})\| &= \left\| \rho_0^T (I - \gamma P \Pi)^{-1} (\gamma P d \Pi) (I - \gamma P \Pi)^{-1} d \Pi \mathbf{A}_{\pi} \right\| \\ &\leq \left\| \rho_0^T \right\| \left\| (I - \gamma P \Pi)^{-1} \right\| \left\| \gamma P \right\| \|d \Pi\| \left\| (I - \gamma P \Pi)^{-1} \right\| \|d \Pi\| \|\mathbf{A}_{\pi}\| \\ &= \frac{\gamma \|d \Pi\|^2 \|\mathbf{A}_{\pi}\|}{(1 - \gamma)^2} \leq \frac{2 \gamma D_{\text{KL}}^{\max}(\pi, \tilde{\pi}) \|\mathbf{A}_{\pi}\|}{(1 - \gamma)^2} \end{aligned}$$

which is the same as (19).

Another way to consider this inequality is

$$\begin{aligned} \|\eta(\tilde{\pi}) - L_{\pi}^4(\tilde{\pi})\| &= \left\| \rho_0^T \left((I - \gamma P \tilde{\Gamma})^{-1} - (I - \gamma P \Pi)^{-1} \right) \tilde{\Gamma} \mathbf{A}_{\pi} \right\| \\ &\leq \left\| \rho_0^T \left((I - \gamma P \tilde{\Gamma})^{-1} \right) \tilde{\Gamma} \mathbf{A}_{\pi} \right\| + \left\| \rho_0^T \left((I - \gamma P \Pi)^{-1} \right) \tilde{\Gamma} \mathbf{A}_{\pi} \right\| \\ &\leq \left\| \rho_0^T \right\| \left\| (I - \gamma P \tilde{\Gamma})^{-1} \right\| \|d \Pi\| \|\mathbf{A}_{\pi}\| + \left\| \rho_0^T \right\| \left\| (I - \gamma P \Pi)^{-1} \right\| \|d \Pi\| \|\mathbf{A}_{\pi}\| \\ &= \frac{2 \|d \Pi\| \|\mathbf{A}_{\pi}\|}{1 - \gamma} \leq \frac{2 \sqrt{2 D_{\text{KL}}^{\max}(\pi, \tilde{\pi})} \|\mathbf{A}_{\pi}\|}{1 - \gamma} \end{aligned}$$

Consider these two bounds: if $\sqrt{D_{\text{KL}}^{\max}(\pi, \tilde{\pi})} > \frac{\sqrt{2(1-\gamma)}}{\gamma}$, the second bound is tighter; otherwise, the first bound is tighter.

4 Other Potential Applications

There are some other applications of this set of notations, *e.g.*

- If we know $\nabla_{\pi_{\theta}} \eta(\pi_{\theta})$ and want to get the $\nabla_{\theta} \eta(\pi_{\theta})$, *i.e.* when we know the optimal direction to update policy but we do not know how to control the parameters to make the policy turn to this direction, we can do

$$\min_{\Delta \theta} d(\Pi_{opt}, \Pi_{\theta + \Delta \theta} - \Pi_{\theta}),$$

where $d(\cdot, \cdot)$ is a distance measure, Π_{opt} is the optimal direction to update policy.

- Instead of estimating value function \mathbf{V}_{π} , we can estimate the environment dynamics P , which is independent of π . Each transition information in the trajectory is valuable no matter the reward signal is detected or not. Then, we can use the estimated \hat{P} to update π directly. The estimation (or construction) of \hat{P} is not affected by the update of policy.

References

- [1] D. A. Levin and Y. Peres. *Markov chains and mixing times*, volume 107. American Mathematical Soc., 2017.
- [2] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- [3] R. S. Sutton and A. G. Barto. *Reinforcement learning: an introduction*. 2018.
- [4] T. Wang, M. Bowling, and D. Schuurmans. Dual representations for dynamic programming and reinforcement learning. In *2007 IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning*, pages 44–51. IEEE, 2007.