



Deep Reinforcement Learning

Professor Mohammad Hossein Rohban

Homework 8:

Policy-Based Theory

By:

Ashkan Majidi
400109984



Spring 2025

Contents

1	Policy Gradient Theorem	1
1.1	Notations	1
1.2	Proving the Policy Gradient Theorem	1
1.3	Compatible Function Approximation Theorem.....	3
2	Trust Region Policy Optimization	5
2.1	Notations and Preliminaries	5
2.2	Monotonic Improvement Guarantee for General Stochastic Policies	7

Grading

The grading will be based on the following criteria, with a total of 100 points:

Task	Points
Policy Gradient - Part (a)	20
Policy Gradient - Part (b)	10
Trust Region Policy Optimization - Part (a)	10
Trust Region Policy Optimization - Part (b)	5
Trust Region Policy Optimization - Part (c)	10
Trust Region Policy Optimization - Part (d)	20
Trust Region Policy Optimization - Part (e)	20
Trust Region Policy Optimization - Part (f)	5
Bonus: Writing your report in Latex	5

1 Policy Gradient Theorem

In this question, we will prove the policy gradient theorem and provide a set of sufficient conditions that allow us to use function approximations as a critic for the Q -value function so that the policy gradient using our function approximation remains exact.

1.1 Notations

Consider a normal finite MDP with bounded rewards. $P(s'|s, a)$ represents the transition model, which corresponds to the probability of transitioning from state s to s' due to action a . Also, the reward model is represented by $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ where $r(s, a)$ is the immediate reward associated with taking action a in state s . Parameter $\gamma \in [0, 1]$ corresponds to the discount factor, and s_0 indicates the starting state of our MDP.

A parametrized policy π_θ induces a distribution over trajectories $\tau = (s_t, a_t, r_t)_{t=0}^\infty$ where s_0 is the starting state, and for all subsequent timesteps t , $a_t \sim \pi(\cdot|s_t)$, $s_{t+1} \sim P(\cdot|s_t, a_t)$. The state value function and the state-action value (Q -value) functions are defined as follows by the Bellman operator:

$$\begin{aligned} V^{\pi_\theta}(s) &= \mathbb{E}_{a \sim \pi_\theta(\cdot|s)}[Q^{\pi_\theta}(s, a)] \\ Q^{\pi_\theta}(s, a) &= r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)}[V^{\pi_\theta}(s')] \end{aligned}$$

We also define the discounted state visitation distribution $d_{s_0}^\pi$ of a policy π as:

$$d_{s_0}^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t Pr^\pi(s_t = s | s_0), \quad (1)$$

where $Pr^\pi(s_t = s | s_0)$ is the state visitation probability that $s_t = s$, after we execute π starting at state s_0 .

1.2 Proving the Policy Gradient Theorem

The objective function of our RL problem is defined as $J(\theta) = V^{\pi_\theta}(s_0)$. The policy gradient method uses the gradient ascent algorithm to optimize θ . This can be done by the direct differentiation of the objective function.

a) Prove the following identity, which is known as the Policy Gradient Theorem:

$$\nabla_\theta J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a)] \quad (2)$$

Solution:

First we prove the following equation:

$$\nabla_\theta J(\theta) = \mathbb{E}_{s_0 \sim \rho^{\pi_\theta}} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t|s_t) Q^\pi(s_t, a_t) \right],$$

where $J(\theta)$ is the objective function, π_θ is the policy parameterized by θ , $Q^\pi(s_t, a_t)$ is the state-action value function, and ρ^{π_θ} is the discounted state visitation distribution.

Also we know that:

$$\begin{aligned} J(\theta) &= V^{\pi_\theta}(s_0) \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \end{aligned}$$

which represents the expected return starting from the initial state s_0 when following the policy π_θ . Now we focus on the differentiation of $J(\theta)$:

$$\begin{aligned} \nabla_\theta J(\theta) &= \nabla_\theta \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right] \\ &= \int \nabla_\theta p_\theta(\tau) r(\tau) d\tau \\ &= \int p_\theta(\tau) \nabla_\theta \log p_\theta(\tau) r(\tau) d\tau \\ &= \mathbb{E}_{\tau \sim p_\theta(\tau)} [\nabla_\theta \log p_\theta(\tau) r(\tau)] \\ &= \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\nabla_\theta \left[\log p(s_1) + \sum_{t=1}^T \log \pi_\theta(a_t | s_t) + \log p(s_{t+1} | s_t, a_t) \right] r(\tau) \right] \\ &= \mathbb{E}_{\tau \sim p_\theta(\tau)} \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) \left(\sum_{t=1}^T \gamma^t r(s_t, a_t) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t \nabla_\theta \log \pi_\theta(a_t | s_t) Q_\theta^\pi(s_t, a_t) \right]. \end{aligned}$$

Now if we prove following statement we will get our desired identity.

$$\sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{(s_t, a_t)} [f(s_t, a_t)] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [f(s, a)]$$

Proof:

$$\begin{aligned}
\mathbb{E}_{\tau \sim \pi_\theta} [f(\tau)] &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{(s_t, a_t)} [f(s_t, a_t)] \\
&= \sum_{t=0}^{\infty} \sum_{s_t, a_t} \gamma^t P(s_t, a_t) f(s_t, a_t) \\
&= \sum_{s, a} f(s, a) \sum_{t=0}^{\infty} \gamma^t P(s_t = s, a_t = a) \\
&= \sum_{s, a} f(s, a) \sum_{t=0}^{\infty} \gamma^t P(s_t = s) \pi_\theta(a|s) \\
&= \sum_{s, a} f(s, a) \frac{d_{s_0}^\pi(s)}{1 - \gamma} \pi_\theta(a|s) \\
&= \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [f(s, a)] \quad \blacksquare
\end{aligned}$$

Putting all together we get our desired identity.

1.3 Compatible Function Approximation Theorem

Now, consider the case in which Q^{π_θ} is approximated by a learned function approximator. If the approximation is sufficiently good, we might hope to use it in place of Q^{π_θ} in equation 2. If we use the function approximator $Q_\phi(s, a)$, the convergence of our method is not necessarily maintained due to the fact that our gradient will not be exact anymore. The following theorem provides sufficient conditions for our function approximator so that our gradient using the approximator remains exact.

Theorem 1.1 (Compatible Function Approximation). *If the following two conditions are satisfied for any function approximator with parameter ϕ :*

1. Critic gradient is compatible with the Actor score function, i.e.,

$$\nabla_\phi Q_\phi(s, a) = \nabla_\theta \log \pi_\theta(a|s)$$

2. Critic parameters ϕ minimize the following mean-squared error¹:

$$\epsilon = \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [(Q^{\pi_\theta}(s, a) - Q_\phi(s, a))^2]$$

Then, the policy gradient using critic $Q_\phi(s, a)$ is exact, i.e.,

$$\nabla_\theta J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^\pi} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} [\nabla_\theta \log \pi_\theta(a|s) Q_\phi(s, a)]$$

b) Prove theorem 1.1.

Solution:

The true policy gradient is given by:

$$\nabla_\theta J(\theta) = \frac{1}{1 - \gamma} \mathbb{E}_{\substack{s \sim d_{s_0}^\pi \\ a \sim \pi_\theta(\cdot|s)}} [\nabla_\theta \log \pi_\theta(a|s) Q^{\pi_\theta}(s, a)]. \quad (3)$$

¹ Assume that the mean-squared error has only one critical point which corresponds to its minimum.

When using the function approximator $Q_\phi(s, a)$, the approximate gradient becomes:

$$\tilde{\nabla}_\theta J(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d_{s_0}^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [\nabla_\theta \log \pi_\theta(a|s) Q_\phi(s, a)]. \quad (4)$$

The difference between the true and approximate gradients is:

$$\nabla_\theta J(\theta) - \tilde{\nabla}_\theta J(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d_{s_0}^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [\nabla_\theta \log \pi_\theta(a|s) (Q^{\pi_\theta}(s, a) - Q_\phi(s, a))]. \quad (5)$$

We now use the two conditions in the theorem:

Condition 1: $\nabla_\phi Q_\phi(s, a) = \nabla_\theta \log \pi_\theta(a|s)$.

Condition 2: ϕ minimizes the mean-squared error $\epsilon = \mathbb{E}_{s,a} [(Q^{\pi_\theta}(s, a) - Q_\phi(s, a))^2]$.

From Condition 2, the optimal ϕ satisfies the first-order stationary point:

$$\nabla_\phi \epsilon = -2 \mathbb{E}_{\substack{s \sim d_{s_0}^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [(Q^{\pi_\theta}(s, a) - Q_\phi(s, a)) \nabla_\phi Q_\phi(s, a)] = 0. \quad (6)$$

Substituting Condition 1 into (4):

$$\mathbb{E}_{\substack{s \sim d_{s_0}^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [(Q^{\pi_\theta}(s, a) - Q_\phi(s, a)) \nabla_\theta \log \pi_\theta(a|s)] = 0. \quad (7)$$

This implies:

$$\mathbb{E}_{\substack{s \sim d_{s_0}^{\pi_\theta} \\ a \sim \pi_\theta(\cdot|s)}} [\nabla_\theta \log \pi_\theta(a|s) (Q^{\pi_\theta}(s, a) - Q_\phi(s, a))] = 0. \quad (8)$$

Substituting (6) into (3), we get:

$$\nabla_\theta J(\theta) - \tilde{\nabla}_\theta J(\theta) = 0 \implies \nabla_\theta J(\theta) = \tilde{\nabla}_\theta J(\theta). \quad (9)$$

Thus, the policy gradient computed using $Q_\phi(s, a)$ is exact under the given conditions. ■

2 Trust Region Policy Optimization

In this question, we will dive deep into the mathematical theories behind the TRPO algorithm. As a roadmap, we first prove that minimizing a certain surrogate objective function guarantees policy improvement with non-trivial step sizes. Then, we make a series of approximations to the theoretically justified algorithm, yielding a practical algorithm, which has been called trust region policy optimization (TRPO).

2.1 Notations and Preliminaries

Let π denote a stochastic policy and let $\eta(\pi)$ denote its expected discounted reward:

$$\eta(\pi) = \mathbb{E}_{s_0, a_0, \dots} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]$$

where

$$s_0 \sim \rho_0(s_0), a_t \sim \pi(a_t | s_t), s_{t+1} \sim P(s_{t+1} | s_t, a_t).$$

Also, we will use the following standard definitions of the state-action value function Q_π , the value function V_π , and the advantage function A_π :

$$\begin{aligned} Q_\pi(s_t, a_t) &= \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right] \\ V_\pi(s_t) &= \mathbb{E}_{a_t, s_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right] \\ A_\pi(s, a) &= Q_\pi(s, a) - V_\pi(s) \end{aligned}$$

a) Prove the following identity:

$$\eta(\pi') = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] \quad (10)$$

Solution:

By definition of the advantage function:

$$A_\pi(s, a) = Q_\pi(s, a) - V_\pi(s)$$

Substitute the Bellman equation $Q_\pi(s_t, a_t) = r(s_t) + \gamma \mathbb{E}_{s_{t+1}} [V_\pi(s_{t+1})]$:

$$A_\pi(s_t, a_t) = r(s_t) + \gamma \mathbb{E}_{s_{t+1}} [V_\pi(s_{t+1})] - V_\pi(s_t)$$

For a trajectory generated by π' :

$$\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) = \sum_{t=0}^{\infty} \gamma^t [r(s_t) + \gamma V_{\pi}(s_{t+1}) - V_{\pi}(s_t)]$$

Separate into three sums:

$$= \underbrace{\sum_{t=0}^{\infty} \gamma^t r(s_t)}_{(I)} + \underbrace{\sum_{t=0}^{\infty} \gamma^{t+1} V_{\pi}(s_{t+1})}_{(II)} - \underbrace{\sum_{t=0}^{\infty} \gamma^t V_{\pi}(s_t)}_{(III)}$$

Rewrite (II) with index shift $k = t + 1$:

$$(II) = \sum_{k=1}^{\infty} \gamma^k V_{\pi}(s_k)$$

Subtract (III) from (II):

$$(II) - (III) = \sum_{k=1}^{\infty} \gamma^k V_{\pi}(s_k) - \sum_{t=0}^{\infty} \gamma^t V_{\pi}(s_t) = -V_{\pi}(s_0)$$

Thus:

$$\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) = (I) - V_{\pi}(s_0) = \sum_{t=0}^{\infty} \gamma^t r(s_t) - V_{\pi}(s_0)$$

Take expectations under π'

$$\mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] = \underbrace{\mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right]}_{\eta(\pi')} - \underbrace{\mathbb{E}_{\pi'} [V_{\pi}(s_0)]}_{\eta(\pi)}$$

Rearranging terms gives the desired result:

$$\eta(\pi') - \eta(\pi) = \mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_{\pi}(s_t, a_t) \right] \blacksquare$$

Equation 10 basically shows that the difference between the expected total rewards of any two policies π' and π depends on the advantage function of policy π if the trajectory is sampled by running π' . We will use this equation to derive an optimization scheme further to maximize the expected total reward using the advantage function of policy π to obtain policy π' .

Let ρ_{π} be the unnormalized discounted visitation frequencies:

$$\rho_{\pi}(s) = P(s_0 = s) + \gamma P(s_1 = s) + \gamma^2 P(s_2 = s) + \dots$$

b) Prove the following identity:

$$\eta(\pi') = \eta(\pi) + \sum_s \rho_{\pi'}(s) \sum_a \pi'(a|s) A_\pi(s, a) \quad (11)$$

Solution:

Starting from the result in part (a):

$$\eta(\pi') = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]$$

Expand the expectation over trajectories generated by π' :

$$\mathbb{E}_{\pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] = \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{\substack{s_t \sim \rho_{\pi'}(s_t) \\ a_t \sim \pi'(\cdot|s_t)}} [A_\pi(s_t, a_t)]$$

Using the definition of discounted visitation frequencies $\rho_{\pi'}(s) = \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi')$:

$$= \sum_s \rho_{\pi'}(s) \sum_a \pi'(a|s) A_\pi(s, a)$$

where we:

- Interchanged the sum over t and state s using $\rho_{\pi'}(s)$
- Explicitly wrote the expectation over actions $\mathbb{E}_{a \sim \pi'}[\cdot]$

Substituting back gives:

$$\eta(\pi') = \eta(\pi) + \sum_s \rho_{\pi'}(s) \sum_a \pi'(a|s) A_\pi(s, a) \quad \blacksquare$$

Equation 11 can be used as an optimization objective in reinforcement learning. Note that this equation has been considered difficult to optimize directly due to the complex dependency of $\rho_{\pi'}(s)$ on π' . Instead, the following local approximation of η has been introduced for optimization:

$$L_\pi(\pi') = \eta(\pi) + \sum_s \rho_\pi(s) \sum_a \pi'(a|s) A_\pi(s, a) \quad (12)$$

Note that L_π uses the visitation frequency ρ_π rather than $\rho_{\pi'}$, ignoring changes in state visitation density due to changes in the policy. In the next section, we will derive an algorithm to guarantee a monotonic improvement in our policy using equation 12 as our objective function, showing that equation 12 is good enough in our case.

2.2 Monotonic Improvement Guarantee for General Stochastic Policies

In this section, we build the theoretical foundations to consider the policy optimization problem, assuming that the policy can be evaluated at all states. The ultimate goal of this section is to prove the following theorem:

Theorem 2.1 Let π, π' be two stochastic policies. Then, the following bound holds:

$$\eta(\pi') \geq L_\pi(\pi') - \frac{4\epsilon\gamma}{(1-\gamma)^2} D_{KL}^{\max}(\pi, \pi')$$

where $\epsilon = \max_{s,a} |A_\pi(s, a)|$

During this section, we use the following definitions and inequality for the total variation and KL divergence:

$$\begin{aligned} D_{TV}(p||q) &= \frac{1}{2} \sum_i |p_i - q_i| \\ D_{TV}^{\max}(\pi, \pi') &= \max_s D_{TV}(\pi(.|s)||\pi'(.|s)) \\ D_{KL}^{\max}(\pi, \pi') &= \max_s D_{KL}(\pi(.|s)||\pi'(.|s)) \\ D_{TV}(p||q)^2 &\leq D_{KL}(p||q) \end{aligned}$$

We will prove theorem 2.1 step by step, and you are required to complete the proof as indicated below. To begin the proof, we denote trajectories by τ and define $\bar{A}(s)$ as follows:

$$\bar{A}(s) = \mathbb{E}_{a \sim \pi'(.|s)}[A_\pi(s, a)]$$

Then we can rewrite equations 11 and 12 as follows:

$$\eta(\pi') = \eta(\pi) + \mathbb{E}_{\tau \sim \pi'}[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t)] \quad (13)$$

$$L_\pi(\pi') = \eta(\pi) + \mathbb{E}_{\tau \sim \pi}[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t)] \quad (14)$$

The only difference in these two equations is whether the states are sampled using π or π' . To bound the difference between $\eta(\pi')$ and $L_\pi(\pi')$, we first need to introduce a measure of how much π and π' agree. Specifically, we'll couple the policies so that they define a joint distribution over pairs of actions. We use the following definition of α -coupled policy pairs:

Definition 2.2 (π, π') is an α -coupled policy pair if it defines a joint distribution $(a, a')|s$ such that $P(a \neq a'|s) \leq \alpha$ for all s . π and π' will denote the marginal distributions of a and a' , respectively.

c) Prove the following lemma:

Lemma 2.3 Given that π, π' are α -coupled policies, for all s ,

$$|\bar{A}(s)| \leq 2\alpha \max_{s,a} |A_\pi(s, a)|$$

Solution:

For any state s , observe that:

$$\begin{aligned} \bar{A}(s) &= \mathbb{E}_{a' \sim \pi'(.|s)}[A_\pi(s, a')] \\ &= \mathbb{E}_{a' \sim \pi'(.|s)}[A_\pi(s, a')] - \underbrace{\mathbb{E}_{a \sim \pi(.|s)}[A_\pi(s, a)]}_{=0} \end{aligned}$$

The second term vanishes because $\mathbb{E}_{a \sim \pi}[A_\pi(s, a)] = V_\pi(s) - V_\pi(s) = 0$.

Let $\epsilon = \max_{s,a} |A_\pi(s, a)|$. The difference between expectations can be bounded using the Total Variation (TV) distance:

$$|\bar{A}(s)| \leq 2D_{TV}(\pi(\cdot|s), \pi'(\cdot|s)) \cdot \epsilon$$

By the α -coupling assumption and the coupling lemma:

$$D_{TV}(\pi(\cdot|s), \pi'(\cdot|s)) \leq P(a \neq a'|s) \leq \alpha$$

Substitute this into the bound:

$$|\bar{A}(s)| \leq 2\alpha\epsilon \quad \forall s$$

d) Prove the following lemma:

Lemma 2.4 Let (π, π') be an α -coupled policy pair. Then:

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| \leq 4\alpha(1 - (1 - \alpha)^t) \max_{s,a} |A_\pi(s, a)|$$

Solution:

Let $\epsilon = \max_{s,a} |A_\pi(s, a)|$. We analyze the difference in expectations by considering the total variation (TV) distance between state distributions at time t under π and π' .

Let $P_\pi(s_t)$ and $P_{\pi'}(s_t)$ denote the state distributions at time t under policies π and π' , respectively. Since (π, π') are α -coupled, the probability that the trajectories diverge (i.e., take different actions) by time t is bounded by:

$$1 - (1 - \alpha)^t$$

This implies the TV distance between the state distributions satisfies:

$$D_{TV}(P_\pi(s_t), P_{\pi'}(s_t)) \leq 1 - (1 - \alpha)^t$$

For any bounded function $f(s)$ with $\max_s |f(s)| \leq C$, we have:

$$|\mathbb{E}_{s \sim P_{\pi'}}[f(s)] - \mathbb{E}_{s \sim P_\pi}[f(s)]| \leq 2D_{TV}(P_\pi, P_{\pi'}) \cdot C$$

Apply this to $f(s) = \bar{A}(s)$, where $\max_s |\bar{A}(s)| \leq 2\alpha\epsilon$ (from Lemma in part (c)):

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| \leq 2(1 - (1 - \alpha)^t) \cdot 2\alpha\epsilon$$

Simplify the inequality:

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| \leq 4\alpha(1 - (1 - \alpha)^t)\epsilon$$

Substituting $\epsilon = \max_{s,a} |A_\pi(s, a)|$ completes the proof. ■

e) Prove the following lemma:

Lemma 2.5 Let (π, π') be an α -coupled policy pair. Then:

$$|\eta(\pi') - L_\pi(\pi')| \leq \frac{4\alpha^2\gamma\epsilon}{(1 - \gamma)^2}$$

Solution:

Let $\epsilon = \max_{s,a} |A_\pi(s, a)|$. From the definitions:

$$\eta(\pi') - L_\pi(\pi') = \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t) \right] - \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t) \right]$$

This can be rewritten as:

$$\sum_{t=0}^{\infty} \gamma^t (\mathbb{E}_{s_t \sim \pi'} [\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi} [\bar{A}(s_t)])$$

Using Lemma (d) to bound each term:

$$|\mathbb{E}_{s_t \sim \pi'} [\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi} [\bar{A}(s_t)]| \leq 4\alpha(1 - (1 - \alpha)^t)\epsilon$$

Substitute and sum over all timesteps:

$$|\eta(\pi') - L_\pi(\pi')| \leq 4\alpha\epsilon \sum_{t=0}^{\infty} \gamma^t [1 - (1 - \alpha)^t]$$

Split the summation:

$$= 4\alpha\epsilon \left[\underbrace{\sum_{t=0}^{\infty} \gamma^t}_{(I)} - \underbrace{\sum_{t=0}^{\infty} (\gamma(1 - \alpha))^t}_{(II)} \right]$$

Evaluate geometric series:

$$(I) = \frac{1}{1 - \gamma}, \quad (II) = \frac{1}{1 - \gamma(1 - \alpha)}$$

Simplify the difference:

$$\frac{1}{1 - \gamma} - \frac{1}{1 - \gamma(1 - \alpha)} = \frac{\gamma\alpha}{(1 - \gamma)(1 - \gamma(1 - \alpha))}$$

Since $1 - \gamma(1 - \alpha) \geq 1 - \gamma$, we have:

$$|\eta(\pi') - L_\pi(\pi')| \leq \frac{4\alpha^2\gamma\epsilon}{(1 - \gamma)^2} \blacksquare$$

f) Prove theorem 2.1. Hint: Use the fact that if we have two policies π and π' such that $D_{TV}^{\max}(\pi, \pi') \leq \alpha$, then we can define an α -couples policy pair (π, π') with appropriate marginals.²

Solution:

Given Theorem 2.1, recall from Lemma (e):

$$|\eta(\pi') - L_\pi(\pi')| \leq \frac{4\alpha^2\gamma\epsilon}{(1 - \gamma)^2}$$

From the hint, if $D_{TV}^{\max}(\pi, \pi') \leq \alpha$, then (π, π') can be treated as an α -coupled pair. By the inequality $D_{TV}(p\|q)^2 \leq D_{KL}(p\|q)$, we have:

$$\alpha^2 \leq D_{KL}^{\max}(\pi, \pi')$$

²There is no need to prove this hint!

Substitute $\alpha^2 \leq D_{KL}^{\max}(\pi, \pi')$ into the bound from Lemma (e):

$$\begin{aligned} |\eta(\pi') - L_\pi(\pi')| &\leq \frac{4\gamma\epsilon}{(1-\gamma)^2} \cdot D_{KL}^{\max}(\pi, \pi') \\ \implies \eta(\pi') &\geq L_\pi(\pi') - \frac{4\epsilon\gamma}{(1-\gamma)^2} D_{KL}^{\max}(\pi, \pi') \quad \blacksquare \end{aligned}$$

Note that the inequality in theorem 2.1 becomes an equality in $\pi' = \pi$. Thus, the following optimization problem guarantees a non-decreasing expected return η :

$$\begin{aligned} \pi_{i+1} &= \arg \max_{\pi} L_{\pi_i}(\pi) - CD_{KL}^{\max}(\pi_i, \pi) \\ \text{where } C &= \frac{4\epsilon\gamma}{(1-\gamma)^2} \\ \text{and } L_{\pi_i}(\pi) &= \eta(\pi_i) + \sum_s \rho_{\pi_i}(s) \sum_a \pi(a|s) A_{\pi_i}(s, a) \end{aligned}$$

In practice, if we use the penalty coefficient C as recommended by the theory above, the step sizes would be very small. One way to take larger steps in a robust way is to use a constraint on the KL divergence between the two policies as a trust region:

$$\begin{aligned} \pi_{i+1} &= \arg \max_{\pi} L_{\pi_i}(\pi) \\ \text{subject to } D_{KL}^{\max}(\pi_i, \pi) &\leq \delta \end{aligned}$$

This problem imposes a constraint that the KL divergence is bounded at every point in the state space. While it is motivated by the theory, this problem is impractical to solve due to the large number of constraints. Instead, we can use a heuristic approximation by considering the average KL divergence. The following optimization problem has been proposed as the TRPO algorithm:

$$\begin{aligned} \pi_{i+1} &= \arg \max_{\pi} L_{\pi_i}(\pi) \\ \text{subject to } \mathbb{E}_{s \sim \rho}[D_{KL}(\pi_i(.|s)||\pi(.|s))] &\leq \delta \end{aligned}$$