



Deep Reinforcement Learning

Professor Mohammad Hossein Rohban

Homework 8:

Policy-Based Theory

By:

Ayeen Poostforoushan

401105742



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:

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

where $P r^\pi(s_t = s | s_0)$ 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:

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

As we know,

$$V^{\pi_\theta}(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 \right].$$

So

$$\nabla_\theta V^{\pi_\theta}(s_0) = \nabla_\theta \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 \right].$$

we do the log trick:

$$\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \nabla_{\theta} \log p_{\theta}(\tau) \mid s_0 \right].$$

also we know that by factorizing the probability of the trajectory, conditioned on s_0 , we have:

$$p_{\theta}(\tau \mid s_0) = \prod_{t=0}^{\infty} \pi_{\theta}(a_t \mid s_t) P(s_{t+1} \mid s_t, a_t).$$

so we have:

$$\nabla_{\theta} \log p_{\theta}(\tau \mid s_0) = \sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_t \mid s_t).$$

therefore:

$$\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \sum_{u=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(a_u \mid s_u) \mid s_0 \right].$$

As we studied before, we know that in the double sum, each $\nabla_{\theta} \log \pi(a_u | s_u)$ only affects rewards from time u onward, so

$$\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \times \sum_{u=0}^{\infty} \nabla_{\theta} \log \pi(a_u | s_u) = \sum_{u=0}^{\infty} \nabla_{\theta} \log \pi(a_u | s_u) \times \sum_{k=0}^{\infty} \gamma^{u+k} r(s_{u+k}, a_{u+k}).$$

showing the equivalence and then to $Q^{\pi_{\theta}}$:

$$\sum_{k=0}^{\infty} \gamma^k r(s_{u+k}, a_{u+k}) = Q^{\pi_{\theta}}(s_u, a_u).$$

now we open the expectation with summation, collect each (s, a) and factor out $Q^{\pi_{\theta}}(s, a)$, then note that each (s, a) appears with weight $(1 - \gamma)\gamma^u$ under the expectation, giving the discounted visitation. Here's the proof:

$$\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \mathbb{E} \left[\sum_{u=0}^{\infty} \gamma^u \nabla_{\theta} \log \pi_{\theta}(a_u \mid s_u) Q^{\pi_{\theta}}(s_u, a_u) \mid s_0 \right].$$

open the expectation with the sum:

$$\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \sum_{u=0}^{\infty} \gamma^u \mathbb{E} \left[\nabla_{\theta} \log \pi_{\theta}(a_u \mid s_u) Q^{\pi_{\theta}}(s_u, a_u) \mid s_0 \right].$$

now write that inner expectation as a sum over all (s, a) :

$$\mathbb{E} \left[\nabla_{\theta} \log \pi_{\theta}(a_u \mid s_u) Q^{\pi_{\theta}}(s_u, a_u) \mid s_0 \right] = \sum_s \sum_a \Pr(s_u = s, a_u = a \mid s_0) \nabla_{\theta} \log \pi_{\theta}(a \mid s) Q^{\pi_{\theta}}(s, a).$$

factor out the $Q^{\pi_{\theta}}(s, a)$ since it does not depend on u :

$$\nabla_{\theta} V^{\pi_{\theta}}(s_0) = \sum_s \sum_a \nabla_{\theta} \log \pi_{\theta}(a \mid s) Q^{\pi_{\theta}}(s, a) \sum_{u=0}^{\infty} \gamma^u \Pr(s_u = s, a_u = a \mid s_0).$$

but $\Pr(s_u = s, a_u = a \mid s_0) = \Pr(s_u = s \mid s_0) \pi_{\theta}(a \mid s)$. So

$$\sum_{u=0}^{\infty} \gamma^u \Pr(s_u = s, a_u = a \mid s_0) = \pi_{\theta}(a \mid s) \sum_{u=0}^{\infty} \gamma^u \Pr(s_u = s \mid s_0).$$

define the discounted visitation, $d_{s_0}^{\pi_\theta}(s) = (1 - \gamma) \sum_{u=0}^{\infty} \gamma^u \Pr(s_u = s \mid s_0)$, so $\sum_{u=0}^{\infty} \gamma^u \Pr(s_u = s \mid s_0) = \frac{d_{s_0}^{\pi_\theta}(s)}{1 - \gamma}$. plug back in:

$$\nabla_\theta V^{\pi_\theta}(s_0) = \sum_s \sum_a \nabla_\theta \log \pi_\theta(a \mid s) Q^{\pi_\theta}(s, a) \pi_\theta(a \mid s) \frac{d_{s_0}^{\pi_\theta}(s)}{1 - \gamma}.$$

recognize the double sum as an expectation:

$$\nabla_\theta V^{\pi_\theta}(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d_{s_0}^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot \mid s)} [\nabla_\theta \log \pi_\theta(a \mid s) Q^{\pi_\theta}(s, a)].$$

And now we have reached the required form of the policy gradient theorem.

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.,*

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

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

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

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

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

b) Prove theorem 1.1.

We begin with the exact policy gradient:

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

Since ϕ minimizes $\epsilon = \mathbb{E}_{s,a} [(Q^{\pi_\theta}(s, a) - Q_\phi(s, a))^2]$, its gradient on ϕ becomes zero:

$$0 = \nabla_\phi \epsilon = \mathbb{E}_{s,a} [2 (Q_\phi(s, a) - Q^{\pi_\theta}(s, a)) \nabla_\phi Q_\phi(s, a)].$$

By the equation we assumed, $\nabla_\phi Q_\phi(s, a) = \nabla_\theta \log \pi_\theta(a \mid s)$, this becomes

$$\mathbb{E}_{s,a} [(Q_\phi(s, a) - Q^{\pi_\theta}(s, a)) \nabla_\theta \log \pi_\theta(a \mid s)] = 0.$$

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

Rearrange:

$$\mathbb{E}_{s,a} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q^{\pi_{\theta}}(s, a)] = \mathbb{E}_{s,a} [\nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\phi}(s, a)].$$

So we have:

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

and the policy gradient using the learned critic is exact.

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 :

$$Q_\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}, a_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right]$$

$$V_\pi(s_t) = \mathbb{E}_{a_t, s_{t+1}, \dots} \left[\sum_{l=0}^{\infty} \gamma^l r(s_{t+l}) \right]$$

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

a) Prove the following identity:

$$\eta(\pi') = \eta(\pi) + \mathbb{E}_{s_0, a_0, \dots \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right] \quad (3)$$

We start from

$$\eta(\pi') = \mathbb{E}_{s_0 \sim \rho_0} [V^{\pi'}(s_0)] = \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) \right],$$

Since ρ_0 is the same initial-state distr. under any π' , we replaced the outer s_0 -expectation by the full trajectory expectation under.

We can replace expectation over s_0 distribution with the whole trajectory distribution over any policy π' because the s_0 distribution is still the same and the other parts of the distribution can simply be marginalized out.

Likewise,

$$\eta(\pi) = \mathbb{E}_{s_0 \sim \rho_0} [V^\pi(s_0)] = \mathbb{E}_{\tau \sim \pi'} [V^\pi(s_0)].$$

But

$$V^\pi(s_0) = \sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) - \sum_{t=1}^{\infty} \gamma^t V^\pi(s_t),$$

so

$$\eta(\pi) = \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) - \sum_{t=1}^{\infty} \gamma^t V^\pi(s_t) \right].$$

Subtracting from $\eta(\pi')$ gives

$$\begin{aligned} \eta(\pi') - \eta(\pi) &= \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t) - \left(\sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) - \sum_{t=1}^{\infty} \gamma^t V^\pi(s_t) \right) \right] \\ &= \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t) + \gamma V^\pi(s_{t+1}) - V^\pi(s_t)) \right] \\ &= \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t) \right]. \end{aligned}$$

Rearranging gives $\eta(\pi') = \eta(\pi) + \mathbb{E}_{\tau \sim \pi'} [\sum_{t=0}^{\infty} \gamma^t A_\pi(s_t, a_t)]$, as claimed.

Equation 3 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:

$$\eta(\pi') = \eta(\pi) + \sum_s \rho_{\pi'}(s) \sum_a \pi'(a|s) A_\pi(s, a) \quad (4)$$

Starting from

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

we open the trajectory expectation into an expectation over states and then actions:

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

Then we swap the order of summation over t and the sum over states:

$$\sum_{t=0}^{\infty} \gamma^t \sum_s \Pr(s_t = s | \pi') \mathbb{E}_{a \sim \pi'(\cdot|s)} [A_\pi(s, a)] = \sum_s \left(\sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi') \right) \left(\sum_a \pi'(a | s) A_\pi(s, a) \right),$$

Notice that $\sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s \mid \pi') = \rho_{\pi'}(s)$ by definition. So we finally have:

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

which is exactly (4).

Equation 4 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:

$$L_{\pi}(\pi') = \eta(\pi) + \sum_s \rho_{\pi}(s) \sum_a \pi'(a \mid s) A_{\pi}(s, a) \quad (5)$$

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 5 as our objective function, showing that equation 5 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:

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

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:

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

Then we can rewrite equations 4 and 5 as follows:

$$\eta(\pi') = \eta(\pi) + \mathbb{E}_{\tau \sim \pi'} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t) \right] \quad (6)$$

$$L_{\pi}(\pi') = \eta(\pi) + \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \bar{A}(s_t) \right] \quad (7)$$

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 ,

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

We have

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

and also as we know that Value function is the expected of Q function over the policy distribution, the expected value

$$0 = \mathbb{E}_{a \sim \pi(\cdot|s)} [A_{\pi}(s, a)] = \sum_a \pi(a | s) A_{\pi}(s, a).$$

Subtracting gives

$$\bar{A}(s) = \sum_a [\pi'(a | s) - \pi(a | s)] A_{\pi}(s, a).$$

Taking abs of the equation and pulling out the maximum advantage:

$$|\bar{A}(s)| \leq \max_a |A_{\pi}(s, a)| \sum_a |\pi'(a | s) - \pi(a | s)|.$$

By definition of total variation distance,

$$\sum_a |\pi'(a | s) - \pi(a | s)| = 2 D_{\text{TV}}(\pi(\cdot | s), \pi'(\cdot | s)) \leq 2\alpha.$$

So we have the desired result:

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

d) Prove the following lemma:

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

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

We start from

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| = \left| \sum_s (P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)) \bar{A}(s) \right|.$$

Pulling out $\max_{s,a} |\bar{A}(s)|$ gives

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| \leq 2\alpha \max_{s,a} |\bar{A}(s)| \sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)|.$$

Now decompose $P_{\pi'}(s_t = s)$ into the “no-mistake” and “at-least-one-mistake” cases. Let α be the per-step probability of π' choosing a different action than π . Then

$$P_{\pi'}(s_t = s) = (1 - \alpha)^t P_{\text{noMistake}}(s_t = s) + [1 - (1 - \alpha)^t] P_{\text{mistake}}(s_t = s).$$

On the “no-mistake” event, π' and π produce identical state-distributions, so

$$P_{\text{noMistake}}(s_t = s) = P_{\pi}(s_t = s).$$

Therefore

$$P_{\pi'}(s_t = s) - P_{\pi}(s_t = s) = [1 - (1 - \alpha)^t] (P_{\text{mistake}}(s_t = s) - P_{\pi}(s_t = s)).$$

Taking absolute values and summing over s ,

$$\sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)| = [1 - (1 - \alpha)^t] \sum_s |P_{\text{mistake}}(s_t = s) - P_{\pi}(s_t = s)|.$$

Since any two probability distributions have total variation at most 2,

$$\sum_s |P_{\text{mistake}}(s_t = s) - P_{\pi}(s_t = s)| \leq 2.$$

Then we have:

$$\sum_s |P_{\pi'}(s_t = s) - P_{\pi}(s_t = s)| \leq 2[1 - (1 - \alpha)^t].$$

So we finally have:

$$|\mathbb{E}_{s_t \sim \pi'}[\bar{A}(s_t)] - \mathbb{E}_{s_t \sim \pi}[\bar{A}(s_t)]| \leq 2\alpha \max_{s,a} |\bar{A}(s)| \times 2[1 - (1 - \alpha)^t] = 4\alpha(1 - (1 - \alpha)^t) \max_{s,a} |\bar{A}(s)|$$

as required.

e) Prove the following lemma:

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

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

We have

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

Apply the bound from part (d):

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

Thus

$$\begin{aligned}
|\eta(\pi') - L_\pi(\pi')| &\leq \sum_{t=0}^{\infty} \gamma^t 4\alpha(1 - (1 - \alpha)^t) \epsilon \\
&= 4\alpha \epsilon \sum_{t=0}^{\infty} (\gamma^t - \gamma^t(1 - \alpha)^t) \\
&= 4\alpha \epsilon \left(\sum_{t=0}^{\infty} \gamma^t - \sum_{t=0}^{\infty} (\gamma(1 - \alpha))^t \right) \\
&= 4\alpha \epsilon \left(\frac{1}{1 - \gamma} - \frac{1}{1 - \gamma(1 - \alpha)} \right) \\
&= 4\alpha \epsilon \frac{\gamma\alpha}{(1 - \gamma)(1 - \gamma(1 - \alpha))} \\
&= 4\alpha^2 \gamma \epsilon \frac{1}{(1 - \gamma)(1 - \gamma + \gamma\alpha)}.
\end{aligned}$$

Since $0 \leq \alpha \leq 1$ implies $1 - \gamma + \gamma\alpha \geq 1 - \gamma$, we finally have the needed inequality:

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

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 -coupled policy pair (π, π') with appropriate marginals.²

Let $\alpha = D_{TV}^{\max}(\pi, \pi')$. By the hint, there exists an α -coupled policy pair (π, π') . From Lemma (e) we have

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

Since for any two distributions p, q , $D_{TV}(p||q)^2 \leq D_{KL}(p||q)$, it follows that $\alpha^2 \leq D_{KL}^{\max}(\pi, \pi')$. Therefore

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

which immediately gives the claimed bound:

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

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 :

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

²There is no need to prove this hint!

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:

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