

## Lecture 7: Approximately optimal approximate RL, TRPO

By Shipra Agrawal

Based on [Kakade and Langford, 2002] and Schulman et al. [2015].

## 1 Examples demonstrating problems with the policy gradient algorithm

Below are two examples of situations when estimating the policy gradient direction is difficult. In particular, the lack of exploration in the policy gradient method translates into large number of samples in order to accurately estimate the gradient direction.

### 1.1 Example 1

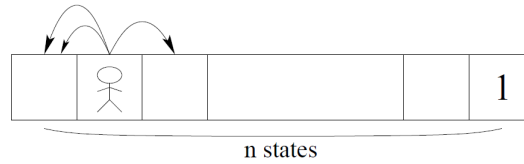


Figure 1: MDP for Example 1

This example from Kakade and Langford [2002] illustrates a scenario where non-zero estimates of policy gradient require observing sample trajectories of exponential length.

We are given an MDP as in Figure 1.1. Each state  $s \in \{n, n-1, \dots, 1\}$  has 3 possible actions. For two of those actions, the next state is deterministically  $s+1$ , and for one action it is deterministically  $s-1$ . The reward for all actions is 0 in all states except state 1 where the reward is 1. Therefore, the goal for discounted reward MDP is to reach the state 1 as soon as possible. And, the optimal policy is to take the third action in all states. Suppose we start in state  $n-1$  with a uniform policy that takes all actions with equal probability. The policy gradient algorithm will try to improve this policy by taking a step in the gradient direction, where gradient

$$\nabla \rho(\pi) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla \pi(s, a)$$

Now  $\frac{1}{1-\gamma} d^\pi(s)$  is the total discounted probability to be in state  $s$  at time  $t = 1, 2, \dots$  according to the uniform policy  $\pi$ , and  $Q^\pi(s, a)$  is the  $Q$ -value (estimate) of this policy. Note that  $Q$ -value estimates will be 0 if we do not observe state 1 in sample trajectories. The given policy takes two steps back and one step forward. By standard random walk analysis, the expected time to reach that goal state 1 from state  $n$  is exponential in  $n$  for such policy. Therefore, intuitively, to obtain non-zero estimates of gradient we need to observe on-policy sample trajectories of exponential length.

More formally, by standard random walk analysis, the probability to reach any state  $s \leq n/2$  is  $(p/q)^{n/2} = \left(\frac{1/3}{2/3}\right)^{n/2} = (1/2)^{n/2}$ . Therefore,  $d^\pi(1) \leq (1/2)^{n/2}$ , i.e., exponentially small in  $n$ . Now, for any state  $s \geq n/2$ ,  $Q^\pi(s, a)$  is at most the probability to reach the state 1, i.e.,  $Q^\pi(s, a) \leq (1/2)^{n/2}$ . Therefore, each term in above summation is at most  $(1/2)^{n/2}$ , giving that the total magnitude of each component of the gradient is at most  $n(1/2)^{n/2}$ . Therefore, all components of the gradient must be estimated within an exponentially small error to get meaningful gradient estimates so that policy improvement can occur, requiring exponentially many samples.

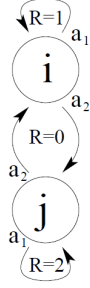


Figure 2: MDP for Example 2

## 1.2 Example 2

The second example illustrates a scenario with a small number of states (2 states), where a policy gradient based approach can still take exponential time to converge.

Consider Figure 1.2, there are 2 states, and 2 actions in each state. Action  $a_1$  in each state loops back to the same state, with reward 1 in state  $i$  and reward 2 in state  $j$ . And, action  $a_2$  transitions to the other state with reward 0. Goal is to maximize average reward. Optimal policy is to play action  $a_2$  in state  $i$  and action  $a_1$  in state  $j$  to obtain an average reward of 2.

Let the stationary distribution of initial policy  $\pi$  is

$$p(i) = 0.8, p(j) = 0.2.$$

An example of such a policy is  $\pi(i, a_1) = 0.8, \pi(i, a_2) = 0.2, \pi(j, a_1) = 0.2, \pi(j, a_2) = 0.8$ . Therefore, initially, the probability of being in state  $i$  is much higher than that of being in state  $j$ . On the other hand the optimal policy has stationary distribution with probability 1 on being in state  $j$ .

Consider policy improvement method using parametric policy of form  $\pi(s, a) \propto e^{\theta^\top \phi_{s,a}}$ , with  $\phi_{s,a} = \mathbf{1}_{s,a}$ . The optimal policy can be encoded by setting  $\theta_{i,a_2} \gg \theta_{i,a_1}$  and  $\theta_{j,a_1} \gg \theta_{j,a_2}$ . Policy gradient is given by

$$\begin{aligned} \nabla_{\theta} \rho(\pi_{\theta}) &= \mathbb{E}_{s \sim d^{\pi}, a \sim \pi(s)} [Q^{\pi}(s, a) \nabla_{\theta} \log(\pi_{\theta}(s, a))] \\ &= \mathbb{E}_{s \sim d^{\pi}, a \sim \pi(s)} [Q^{\pi}(s, a) (\phi_{s,a} - \sum_{a'} \pi_{\theta}(s, a') \phi_{s,a'})] \\ &= [d^{\pi}(s) Q^{\pi}(s, a) \pi_{\theta}(s, a) (1 - \pi_{\theta}(s, a))]_{s=i,j, a=a_1, a_2} \end{aligned}$$

Update of  $\theta_{s,a}$  by gradient ascent step:

$$\theta_{s,a} \leftarrow \theta_{s,a} + \alpha d^{\pi}(s) Q^{\pi}(s, a) \pi_{\theta}(s, a) (1 - \pi_{\theta}(s, a))$$

In the given MDP, since there are only two actions,  $\pi(s, a_1) = 1 - \pi(s, a_2)$ , and therefore for component  $(s, a_1)$  and  $(s, a_2)$  of gradient are the same. Further for the given policy

$$\pi(i, a_1)(1 - \pi(i, a_1)) = \pi(i, a_2)(1 - \pi(i, a_2)) = \pi(j, a_2)(1 - \pi(j, a_2)) = \pi(j, a_1)(1 - \pi(j, a_1)) = 0.8(0.2)$$

This means that the gradient only depends on  $d^{\pi}(s), s = i, j$ , which is the stationary distribution of the current policy  $\pi$ , and  $Q^{\pi}(s, a)$ . For state  $i$ ,  $Q^{\pi}(i, a_1)$  is higher than  $Q^{\pi}(i, a_2)$ . Therefore,  $\theta_{i,a_1}$  increases compared to  $\theta_{i,a_2}$ , which means the updated policy will favor looping on  $i$  even more rather than transitioning to  $j$ . For  $j$ , again  $\theta_{j,a_1}$  increases more, but since the stationary probability  $d^{\pi}(i) = p(i)$  is higher for state  $i$  than  $p(j)$ , the state  $i$  gets more magnitude of update. The stationary probability is even worse in the next policy because of increased probability of taking action  $a_1$  in state  $i$ . Due to these reasons, initially the stationary probability of  $j$  decreases, and becomes exponentially small (see Figure 3.3 in Kakade and Langford [2002]), before it comes back to the correct policy in exponential time steps.

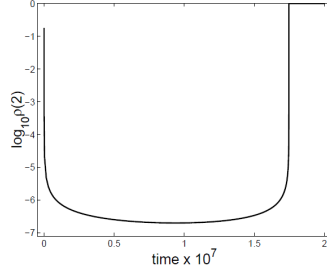


Figure 3: Stationary probability of state  $j$ : Figure 3.3 from Kakade and Langford [2002]

## 2 Conservative greedy algorithm for guaranteed policy improvement

Above examples demonstrate that the improvement in policy in every step of the policy gradient can either be very small (example 1), or may even be negative, i.e., the policy may become worse (example 2, where initially, the new policy has increased probability of looping on state  $i$ ). One underlying reason is that the gradient may not be an accurate measure of policy improvement (i.e., change in the gain of policy) when the policy is changed substantially. Kakade and Langford [2002] design a “conservative greedy” policy improvement method with guaranteed improvement in every step. Therefore, one can precisely quantify the number of steps required for this algorithm to converge,

This algorithm updates the policy in a lazy manner to ensure improvement with every update. Let  $\pi'$  be a new policy. (For example, this could be the policy obtained after gradient ascent from the old policy  $\pi$ ). The conservative greedy algorithm makes the following lazy update to the policy:

$$\pi^{new}(s, a) := (1 - \alpha)\pi(s, a) + \alpha\pi'(s, a) \quad (1)$$

Kakade and Langford [2002] provide a lower bound on increase in gain for such an update for any  $\pi', \pi$ . Following, they suggest that the policy  $\pi'$  and step size  $\alpha$  should be picked in such a way that this lower bound is maximized. The conservative greedy algorithm (with right choice of  $\pi'$  and step size  $\alpha$ ) is guaranteed to converge to a policy where no further local improvement can be made to the policy. Below we describe these results in detail.

### 2.1 Quantifying policy improvement

Following lemma lower bounds the gain from policy improvement.

**Lemma 1** (Lemma 4.1 of Kakade and Langford [2002]).

$$\rho(\pi^{new}) - \rho(\pi) \geq \alpha A_\pi(\pi') - \frac{2\alpha^2\gamma\epsilon}{(1 - \gamma(1 - \alpha))}$$

where  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ ,

$$A_\pi(\pi') := \sum_s d^\pi(s) \sum_a \pi'(s, a) A^\pi(s, a),$$

$$\epsilon := \frac{1}{(1 - \gamma)} \left( \max_s \sum_a \pi'(s, a) A^\pi(s, a) \right).$$

*Intuitively,  $A_\pi(\pi')$ , measures to what extent advantage can increase if a different action (according to  $\pi'$ ) was chosen in every visited state under  $\pi$ . And,  $\epsilon$  is maximum of such advantage over different states.*

*Proof.* We can use the policy gradient theorem to get an intuition for this. For a fixed  $\pi'$ , there is one policy  $\pi_{new}$  for every value of  $\alpha \in [0, 1]$ . Therefore, we can consider  $\pi^{new}$  as a parametric policy  $\pi_\alpha^{new}$  parametrized by  $\alpha$ ; with

$\pi_\alpha^{new} = \pi$  for  $\alpha = 0$ , and  $\pi_\alpha^{new} = \pi'$  for  $\alpha = 1$ . Then, using policy gradient theorem, we get that gradient of  $\rho(\pi_\alpha^{new})$  at  $\alpha = 0$  is given by

$$\begin{aligned}
\nabla_\alpha \rho(\pi_\alpha^{new}) \Big|_{\alpha=0} &= \sum_s d^\pi(s) \sum_a (\nabla_\alpha \pi_\alpha^{new}(s, a)) A^{\pi_\alpha^{new}}(s, a) \Big|_{\alpha=0} \\
&= \sum_s d^\pi(s) \sum_a (\pi'(s, a) - \pi(s, a)) A^\pi(s, a) \Big|_{\alpha=0} \\
&= \sum_s d^\pi(s) \sum_a (\pi'(s, a) - \pi(s, a)) A^\pi(s, a) \\
&= \sum_s d^\pi(s) \sum_a \pi'(s, a) A^\pi(s, a) - \sum_s d^\pi(s) \pi(s, a) A^\pi(s, a) \\
&= \sum_s d^\pi(s) \sum_a \pi'(s, a) A^\pi(s, a) \\
&= A_\pi(\pi')
\end{aligned}$$

The second last step follows because  $\sum_s d^\pi(s) \pi(s, a) A^\pi(s, a) = \sum_s d^\pi(s) \pi(s, a) (Q^\pi(s, a) - V^\pi(s)) = \sum_s d^\pi(s) (V^\pi(s) - V^\pi(s)) = 0$ . In fact, by the same insight,  $A_\pi(\pi^{new}) = \alpha A_\pi(\pi')$ . Therefore, using Taylor expression, a lower bound on the improvement is given by

$$\rho(\pi_\alpha^{new}) - \rho(\pi) \geq \alpha A_\pi(\pi') - O(\alpha^2) = A_\pi(\pi^{new}) - O(\alpha^2)$$

To get a precise lower bound expression, we use a stronger lemma below (Lemma 2). This lemma shows that improvement in gain is given precisely by the following expression:

$$\rho(\pi^{new}) - \rho(\pi) = \sum_s d^{\pi^{new}}(s) \sum_a \pi^{new}(s, a) A^\pi(s, a)$$

Compare this to

$$\alpha A_\pi(\pi') = A_\pi(\pi^{new}) := \sum_s d^\pi(s) \sum_a \pi^{new}(s, a) A^\pi(s, a)$$

The first expression above uses state distribution under  $\pi^{new}$  instead of  $\pi$ . To get a precise lower bound we need to bound the error due to this measure mismatch.

To compare the two, a coupling argument is used. In any given state  $s$ ,  $\pi^{new}$  picks actions according to  $\pi'$  with probability  $\alpha$  and according to  $\pi$  with probability  $1 - \alpha$ . Now, for any fixed time  $t$ , let  $\eta_t$  be the number of steps before time  $t$  where  $\pi^{new}$  did not take action according to  $\pi$ . Then, conditional on event  $\eta_t = 0$ , the distribution of states before time  $t$  is same for trajectories generated from  $\pi^{new}$  and  $\pi$ . More precisely,

$$\Pr_{\tau \sim \pi^{new}}(s_t = s | \eta_t = 0) = \Pr_{\tau \sim \pi}(s_t = s)$$

where random variable  $\tau = (s_1, s_2, \dots, s_t, \dots)$  denotes a trajectory. Further, the probability that  $\pi^{new}$  differed from  $\pi$  in some time step before  $t$  is given by  $p_t := \Pr(\eta_t > 0) = 1 - (1 - \alpha)^{t-1}$ . Therefore,

$$\begin{aligned}
\rho(\pi^{new}) - \rho(\pi) &= \sum_s d^{\pi^{new}}(s) \sum_a \pi^{new}(s, a) A^\pi(s, a) \\
&= \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \pi^{new}} \left[ \sum_t \gamma^{t-1} \sum_a \pi^{new}(s_t, a) A^\pi(s_t, a) \right] \\
&= \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \pi^{new}} \left[ \sum_t \gamma^{t-1} \sum_a \alpha \pi'(s_t, a) A^\pi(s_t, a) \right] \\
&= \alpha \sum_t (1 - p_t) \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots, s_t) \sim \pi^{new}} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) | \eta_t = 0 \right] \\
&\quad + \alpha \sum_t p_t \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \pi^{new}} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\
&= \alpha \sum_t (1 - p_t) \gamma^{t-1} \mathbb{E}_{(s_1, s_2, \dots, s_t) \sim \pi} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) \right] \\
&\quad + \alpha \sum_t p_t \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \pi^{new}} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\
&= \alpha A_\pi(\pi') - \alpha \sum_t p_t \gamma^{t-1} \mathbb{E}_{(s_1, s_2, \dots, s_t) \sim \pi} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) \right] \\
&\quad + \alpha \sum_t \gamma^{t-1} p_t \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \pi^{new}} \left[ \sum_a \pi'(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\
&\geq \alpha A_\pi(\pi') - 2\alpha \sum_t (1 - (1 - \alpha)^{t-1}) \gamma^{t-1} \left( \max_s \sum_a \pi'(s, a) A^\pi(s, a) \right) \\
&= \alpha A_\pi(\pi') - 2\alpha \epsilon (1 - \gamma) \left( \frac{1}{1 - \gamma} - \frac{1}{1 - (1 - \alpha)\gamma} \right) \\
&= \alpha A_\pi(\pi') - 2\alpha \epsilon \frac{\alpha \gamma}{(1 - (1 - \alpha)\gamma)}
\end{aligned}$$

where  $\epsilon = \frac{1}{1-\gamma} (\max_s \sum_a \pi'(s, a) A^\pi(s, a))$ .

□

**Lemma 2** (Lemma 6.1 of Kakade and Langford [2002]). *For any two policies  $\tilde{\pi}, \pi$ ,*

$$\rho(\tilde{\pi}) - \rho(\pi) = \sum_s d^{\tilde{\pi}}(s) \sum_a \tilde{\pi}(s, a) A^\pi(s, a)$$

Here,  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ .

(Compare this to policy advantage expression:  $A_{\tilde{\pi}}(\pi) := \sum_s d^\pi(s) \sum_a (\tilde{\pi}(s, a)) A^\pi(s, a)$ . Note that the policy advantage uses the state distribution under  $\pi$  instead of  $\tilde{\pi}$ )

*Proof.*

$$\begin{aligned}
\rho(\tilde{\pi}) = V^{\tilde{\pi}}(s_1) &= \mathbb{E}_{s_1, a_1, s_2, a_2, \dots \sim \tilde{\pi}} \left[ \sum_{t=1}^{\infty} \gamma^{t-1} R(s_t, a_t) | s_1 \right] \\
&= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{s_t, a_t \sim \tilde{\pi}} [R(s_t, a_t) + V^{\pi}(s_t) - V^{\pi}(s_t) | s_1] \\
&= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{s_t, a_t \sim \tilde{\pi}, s_{t+1} \sim P_{s_t, a_t}} [R(s_t, a_t) + \gamma V^{\pi}(s_{t+1}) - V^{\pi}(s_t) | s_1] + V^{\pi}(s_1) \\
&= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{s_t, a_t, s_{t+1} \sim \tilde{\pi}} [Q^{\pi}(s_t, a_t) - V^{\pi}(s_t) | s_1] + V^{\pi}(s_1) \\
&= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{s_t, a_t, s_{t+1} \sim \tilde{\pi}} [A^{\pi}(s_t, a_t) | s_1] + V^{\pi}(s_1) \\
&= \sum_s d^{\tilde{\pi}}(s) \sum_a \tilde{\pi}(s, a) A^{\pi}(s, a) + V^{\pi}(s_1) \\
&= \sum_s d^{\tilde{\pi}}(s) \sum_a \tilde{\pi}(s, a) A^{\pi}(s, a) + \rho(\pi)
\end{aligned}$$

□

## 2.2 Algorithm design: selecting $\pi'$ and $\alpha$

Lemma 1 provides a lower bound on the improvement for a given  $\alpha$  and  $\pi'$

**Choice of  $\alpha$ .** For  $\alpha = 1$  (greedy update), the improvement is lower bounded as:

$$\rho(\pi^{new}) - \rho(\pi) \geq A_{\pi}(\pi') - 2\gamma\epsilon$$

However, the second quantity can be larger than the first in above, as  $\epsilon$  is an upper bound on  $A_{\pi}(\pi')$ , and therefore, the above improvement may be negative.

However, we can select a step size to ensure positive improvement as long as  $A_{\pi}(\pi') > 0$ . Let  $R$  be an upper bound on rewards, so that  $A_{\pi}(\pi') \leq \frac{R}{1-\gamma}$ ,  $\epsilon \leq \frac{R}{(1-\gamma)}$ . Then, setting

$$\alpha = \frac{A_{\pi}(\pi')(1-\gamma)^2}{4R}, \quad (2)$$

and substituting in Lemma 1, we get

$$\rho(\pi^{new}) - \rho(\pi) \geq \alpha A_{\pi}(\pi') - \frac{2\alpha^2\epsilon}{(1-\gamma)} \geq \frac{A_{\pi}(\pi')^2(1-\gamma)^2}{8R} \quad (3)$$

**Remark:** Note that since the ‘distribution’  $d^{\pi}(s)$  in the definition of  $A_{\pi}(\pi')$  is not normalized to 1, and instead sums  $1/(1-\gamma)$ ,  $A_{\pi}(\pi')$  in the above expression is upper bounded by  $R/(1-\gamma)$ . So, the above expression is of order  $R$ .

**Choice of  $\pi'$ .** From above quantification, to ensure maximum improvement,  $\pi'$  should maximize  $A_{\pi}(\pi')$ . Intuitively,  $A_{\pi}(\pi')$  measures to what extent advantage can increase if a different action (according to  $\pi'$ ) was chosen in every visited state under  $\pi$ . Clearly,

$$\max_{\pi'} A_{\pi}(\pi') = \sum_s d^{\pi}(s) \max_a A^{\pi}(s, a)$$

Therefore, the policy that maximizes policy advantage is given by

$$\pi'(s) = \arg \max_a A^\pi(s, a).$$

However, to use this policy the advantage  $A^\pi(s, a)$  needs to be estimated. The following algorithm allows approximate estimation of  $A^\pi(s, a)$ :

**Algorithm.** Initialize  $\pi$ .

1. Set  $\hat{A} = \sum_s d^\pi(s) \max_a \hat{A}^\pi(s, a)$ , where  $\hat{A}^\pi(s, a)$  are estimates of  $A^\pi(s, a)$  for every  $s, a$  and such that

$$(1 - \gamma)\hat{A} \geq (1 - \gamma) \max_{\pi'} A_\pi(\pi') - \frac{\delta}{3}$$

Here  $\max_{\pi'} A_\pi(\pi') = \sum_s d^\pi(s) \max_a A^\pi(s, a)$ . This can be done for example by estimating  $A^\pi(s, a)$  as function approximation  $\hat{A}^\pi = f_\omega(s, a)$  where parameter  $\omega$  is set through sample estimation with loss function

$$(1 - \gamma) \sum_s d^\pi(s) \max_a |A^\pi(s, a) - f_\omega(s, a)|$$

Since this is expected error over state distribution under current policy  $\pi$ , this loss can be approximated using trajectory samples from the current policy. Roughly  $\frac{R^2}{\delta^2} \log \frac{R^2}{\delta^2}$  samples are required to ensure a  $\delta$  with probability  $1 - \delta$ .

2. If  $(1 - \gamma)\hat{A} < \frac{2\delta}{3}$ , STOP.

3. Otherwise, update policy:

$$\pi \leftarrow (1 - \alpha)\pi + \alpha\pi'$$

where

$$\begin{aligned} \pi'(s) &:= \arg \max_a f_\omega(s, a). \\ \alpha &= \left( \hat{A} - \frac{\delta}{3(1 - \gamma)} \right) \frac{(1 - \gamma)^2}{4R} \end{aligned}$$

4. Go back to Step 1.

In above procedure, in every iteration  $\hat{A}(1 - \gamma) - \frac{\delta}{3} \geq \frac{2\delta}{3} - \frac{\delta}{3} = \frac{\delta}{3}$ , therefore, from (3), the increase in gain is at least

$$\left( \hat{A} - \frac{\delta}{3(1 - \gamma)} \right)^2 \frac{(1 - \gamma)^2}{8R} \geq \frac{\delta^2}{72R}$$

Since the total improvement to be made is at most maximum value of gain, i.e.,  $R/(1 - \gamma)$ , the procedure terminates in at most  $\frac{R}{1 - \gamma} \frac{72R}{\delta^2} = \frac{72R^2}{\delta^2(1 - \gamma)}$  steps. That is, we have the following result.

**Lemma 3.** *Above conservative greedy algorithm terminates in at most  $\frac{72R^2}{\delta^2(1 - \gamma)}$  steps to find a policy  $\pi$  such that*

$$(1 - \gamma) \max_{\pi'} A_\pi(\pi') \leq \delta$$

### 3 How good is the policy found?

As demonstrated in the last section, the conservative greedy algorithm (with right choice of  $\pi'$  and step size  $\alpha$ ) is guaranteed to terminate. It terminates at the policy  $\pi$  such that  $(1 - \gamma) \max_{\pi'} A_\pi(\pi') \leq \delta$ . This can be interpreted as the condition that there is no (or very little) advantage increase on changing the policy under the state distribution of the current policy. But, how does this policy compare to the “optimal policy”, which may have completely different state distribution.

The following theorem shows that the gap can be large if the stationary distribution over states for the chosen policy is very different from the stationary distribution over states for the optimal policy. This can happen if there isn't enough exploration over states. The following theorem also provides a way to ensure exploration. It states that one could start from a different starting state distribution (e.g. uniform) than the target starting state distribution, and then the gap depends only on how the stationary distribution of optimal policy differs from the uniform distribution.

**Theorem 4** (Theorem 6.2 of Kakade and Langford [2002]). *Let  $d^{\pi, \mu}$  denotes the (non-normalized) stationary distribution over states for policy  $\pi$  when starting state distribution is  $\mu$ . Also, let  $\rho(\pi; \mu)$  denote the gain of policy  $\pi$  when starting state distribution is given by  $\mu$ .*

*Suppose we have a policy  $\pi$ , with  $(1 - \gamma) \max_{\pi'} A_{\pi}(\pi') \leq \delta$ , where  $A_{\pi}(\pi') = \sum_s d^{\pi, \mu}(s) \sum_a \pi'(s, a) A^{\pi}(s, a)$ . Then, for any policy  $\pi^*$  and starting state distribution  $\mu^*$ ,*

$$\begin{aligned} \rho(\pi^*; \mu^*) - \rho(\pi; \mu^*) &\leq \frac{\delta}{1 - \gamma} \left\| \frac{d^{\pi^*, \mu^*}}{d^{\pi, \mu}} \right\|_{\infty} \\ &\leq \frac{\delta}{(1 - \gamma)^2} \left\| \frac{(1 - \gamma) d^{\pi^*, \mu^*}}{\mu} \right\|_{\infty} \end{aligned}$$

*Proof.*

$$\begin{aligned} \frac{\delta}{1 - \gamma} \geq \max_{\pi'} A_{\pi}(\pi') &= \sum_s d^{\pi, \mu}(s) \max_a A^{\pi}(s, a) \\ &= \sum_s \frac{d^{\pi, \mu}(s)}{d^{\pi^*, \mu^*}(s)} d^{\pi^*, \mu^*}(s) \max_a A^{\pi}(s, a) \\ &\geq \left( \min_s \frac{d^{\pi, \mu}(s)}{d^{\pi^*, \mu^*}(s)} \right) \sum_s d^{\pi^*, \mu^*}(s) \pi^*(s, a) A^{\pi}(s, a) \\ &= \left\| \frac{d^{\pi, \mu}}{d^{\pi^*, \mu^*}} \right\|_{\infty}^{-1} (\rho(\pi^*; \mu^*) - \rho(\pi; \mu^*)) \end{aligned}$$

where the last step follows from Lemma 2. □

Therefore, if we can choose starting state distribution to be uniform distribution, we can bound the gap from optimal policy by  $\frac{n\delta}{(1-\gamma)^2}$ , where  $n$  is the number of states. (Some applications may not have the flexibility of choosing the starting state distribution).

The second part of the theorem statement follows from the observation that

$$\mu(s) \leq d^{\pi, \mu}(s) := \sum_{t=1}^{\infty} \gamma^{t-1} \Pr(s_t = s) = (I - \gamma P)^{-1} \mu(s) \leq \frac{1}{1 - \gamma} \mu(s)$$

so that

$$\left\| \frac{d^{\pi^*, \mu^*}}{d^{\pi, \mu}} \right\|_{\infty} = \left\| \frac{(1 - \gamma) d^{\pi^*, \mu^*}}{(1 - \gamma) d^{\pi, \mu}} \right\|_{\infty} \leq \frac{1}{1 - \gamma} \left\| \frac{(1 - \gamma) d^{\pi^*, \mu^*}}{\mu} \right\|_{\infty}$$

## 4 Trust Region Policy Optimization Schulman et al. [2015]

The TRPO algorithm from Schulman et al. [2015] can be explained as a simple (and useful!) extension of the above ideas. The paper explores if we need to use mixed policies of form  $\pi^{\text{new}} = (1 - \alpha)\pi + \alpha\pi'$ . Note that every iteration of the conservative greedy algorithm stated above adds a new policy to this mixture, making the collection of policies to maintain potentially quite large and inconvenient. Instead can we just maintain policy parameter  $\theta$ , and simply update the policy parameter  $\theta$  to obtain a new policy?

Schulman et al. [2015] provide following lower bound on improvement of gain when policy  $\pi$  is updated to arbitrary policy  $\tilde{\pi}$ , which is very similar to the quantification in Lemma 1 (even has essentially the same proof).



But, this way of quantifying the gap allows finding a new policy  $\tilde{\pi}$  in conservative manner through parameter search in a small ‘trust’ region around the parameter for  $\pi$ , instead of restricting to mixed policy as was the case in the above conservative greedy algorithm.

**Lemma 5** (Theorem 1 of Schulman et al. [2015]). *For any two policies  $\tilde{\pi}$  and  $\pi$ ,*

$$\begin{aligned}\rho(\tilde{\pi}) - \rho(\pi) &\geq A_\pi(\tilde{\pi}) - \frac{4\epsilon\gamma}{(1-\gamma)} D_{TV}^{\max}(\pi, \tilde{\pi})^2 \\ &\geq A_\pi(\tilde{\pi}) - \frac{4\epsilon\gamma}{(1-\gamma)} D_{KL}^{\max}(\pi, \tilde{\pi})\end{aligned}$$

with  $\epsilon = \frac{1}{(1-\gamma)} \max_s \max_a |A^\pi(s, a)|$ ; and  $D_{TV}^{\max}(\pi, \tilde{\pi}) = \max_s D_{TV}(\pi(s, \cdot) || \tilde{\pi}(s, \cdot))$  is the maximum total variation distance between the two policies, and  $D_{KL}^{\max}(\pi, \tilde{\pi}) = \max_s D_{KL}(\pi(s, \cdot) || \tilde{\pi}(s, \cdot))$  is the maximum KL divergence, over all states.

*Proof.* The proof is almost exactly the same as the proof of Lemma 1. Consider again the proof of Lemma 1. We observe that all the steps of the proof of Lemma 1 follow for arbitrary policies  $\pi$  and  $\pi^{\text{new}} := \tilde{\pi}$ , by setting  $\alpha$  such that in any state  $s$ ,  $\Pr_{a \sim \pi(s), \tilde{a} \sim \tilde{\pi}(s)}(a \neq \tilde{a} | s) \leq \alpha$ .

Below we provide the proof details. As before, we couple the trajectories of  $\pi$  and  $\tilde{\pi}$ . Let  $p_t$  be the probability that  $\tilde{\pi}$  suggests a different action than  $\pi$  in some time step before  $t$ . We denote this event as  $\eta_t > 0$ . Then  $p_t := \Pr(\eta_t > 0) \leq 1 - (1 - \alpha)^{t-1}$ . Now,

$$\begin{aligned}\rho(\tilde{\pi}) - \rho(\pi) &= \sum_s d^{\tilde{\pi}}(s) \sum_a \tilde{\pi}(s, a) A^\pi(s, a) \\ &= \mathbb{E}_{\tau=(s_1, s_2, \dots) \sim \tilde{\pi}} \left[ \sum_t \gamma^{t-1} \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) \right] \\ &= \sum_t (1 - p_t) \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots, s_t) \sim \tilde{\pi}} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) | \eta_t = 0 \right] \\ &\quad + \sum_t p_t \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots, s_t) \sim \tilde{\pi}} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\ &= \sum_t (1 - p_t) \gamma^{t-1} \mathbb{E}_{(s_1, s_2, \dots, s_t) \sim \pi} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) \right] \\ &\quad + \sum_t p_t \gamma^{t-1} \mathbb{E}_{\tau=(s_1, s_2, \dots, s_t) \sim \tilde{\pi}} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\ &= A_\pi(\tilde{\pi}) - \sum_t p_t \gamma^{t-1} \mathbb{E}_{(s_1, s_2, \dots, s_t) \sim \pi} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) \right] \\ &\quad + \sum_t \gamma^{t-1} p_t \mathbb{E}_{\tau=(s_1, s_2, \dots, s_t) \sim \tilde{\pi}} \left[ \sum_a \tilde{\pi}(s_t, a) A^\pi(s_t, a) | \eta_t > 0 \right] \\ &\geq A_\pi(\tilde{\pi}) - 2 \sum_t p_t \gamma^{t-1} \left( \max_s \sum_a \tilde{\pi}(s, a) A^\pi(s, a) \right) \\ &\geq A_\pi(\tilde{\pi}) - 2\tilde{\epsilon} (1 - \gamma) \sum_t (1 - (1 - \alpha)^{t-1}) \gamma^{t-1} \\ &= A_\pi(\tilde{\pi}) - 2\tilde{\epsilon} \left( \frac{1}{1 - \gamma} - \frac{1}{1 - (1 - \alpha)\gamma} \right) \\ &= A_\pi(\tilde{\pi}) - 2\tilde{\epsilon} \frac{\alpha\gamma}{(1 - (1 - \alpha)\gamma)}\end{aligned}$$

where  $\tilde{\epsilon} := \frac{1}{1-\gamma} (\max_s \sum_a \tilde{\pi}(s, a) A^\pi(s, a))$ . Now, using  $\Pr_{a \sim \pi, \tilde{a} \sim \tilde{\pi}}(a \neq \tilde{a} | s) \leq \alpha$ , we get that

$$\tilde{\epsilon} = \frac{1}{1-\gamma} \left( \max_s \sum_a \tilde{\pi}(s, a) A^\pi(s, a) - \pi(s, a) A^\pi(s, a) \right) \leq \frac{1}{1-\gamma} \left( \max_s \max_{a, \tilde{a}} \alpha |A^\pi(s, a) - A^\pi(s, \tilde{a})| \right) \leq 2\alpha\epsilon$$

where  $\epsilon = \frac{1}{1-\gamma} \max_{s,a} |A^\pi(s,a)|$ . Substituting,

$$\rho(\tilde{\pi}) - \rho(\pi) \geq A_\pi(\tilde{\pi}) - \frac{4\alpha^2\epsilon\gamma}{(1-\gamma)}$$

Now, setting  $\alpha = D_{TV}^{\max}(\pi, \tilde{\pi})$ , we have  $\Pr(a \neq \tilde{a}|s) \leq \alpha$ . This gives the stated result.  $\square$

Note that the above lemma is very similar in form to Lemma 1 which considered  $\tilde{\pi}$  as mixed policy with  $\alpha$  being the mixing parameter, and provided:

$$\rho(\tilde{\pi}) - \rho(\pi) \geq A_\pi(\tilde{\pi}) - \frac{2\epsilon\gamma}{(1-\gamma)}\alpha^2$$

(For mixed policy  $\tilde{\pi} = (1-\alpha)\pi + \alpha\pi'$ ,  $A_\pi(\tilde{\pi}) = \alpha A_\pi(\pi')$ , and  $\epsilon = \frac{1}{1-\gamma} (\max_s \sum_a \pi'(s,a) A^\pi(s,a))$ ). The main contribution of Lemma 5 is to get rid of mixing through  $\alpha$  and provide bounds in terms of total variation distance for arbitrary updates.

Given the above result, the paper proposes the following way to update the policy parameter. A natural strategy to find a good policy is to search over all policies with small total variation distance or KL-divergence from the policy  $\pi$  and find the one with maximum  $A_\pi(\tilde{\pi})$ . Let  $\pi$  is parameterized by  $\theta$ . And,  $A_\theta(\tilde{\theta})$  denote  $A_{\pi_\theta}(\tilde{\pi}_\theta)$ . Then, TRPO algorithm chooses next policy by solving

$$\begin{aligned} \max_{\tilde{\theta}} \quad & A_\theta(\tilde{\theta}) \\ \text{s.t.} \quad & D_{KL}^{\max}(\pi_\theta || \pi_{\tilde{\theta}}) \leq \eta \end{aligned}$$

The region of parameters with  $D_{KL}^{\max}(\pi_\theta, \pi_{\tilde{\theta}}) \leq \eta$  is referred to as the trust region.

## References

- Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, ICML '02, pages 267–274, San Francisco, CA, USA, 2002. Morgan Kaufmann Publishers Inc. ISBN 1-55860-873-7. URL <http://dl.acm.org/citation.cfm?id=645531.656005>.
- John Schulman, Sergey Levine, Philipp Moritz, Michael Jordan, and Pieter Abbeel. Trust region policy optimization. In *Proceedings of the 32Nd International Conference on International Conference on Machine Learning - Volume 37*, ICML'15, pages 1889–1897. JMLR.org, 2015. URL <http://dl.acm.org/citation.cfm?id=3045118.3045319>.