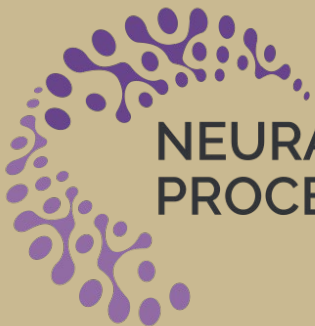


Learning General Parameterized Policies for Infinite Horizon Average Reward Constrained MDPs via Primal-Dual Policy Gradient Algorithm

Qinbo Bai¹, Washim Uddin Mondal², Vaneet Aggarwal¹

1 Purdue University, 2 Indian Institute of Technology Kanpur



NEURAL INFORMATION
PROCESSING SYSTEMS

Formulation: Gain Function

- CMDP(S, A, P, r, c, ρ), Gain function (long-term average reward)

$$J_{g,\rho}^{\pi} \triangleq \lim_{T \rightarrow \infty} \frac{1}{T} \mathbf{E} \left[\sum_{t=0}^{T-1} g(s_t, a_t) \middle| s_0 \sim \rho, \pi \right]$$

- We consider a parametrized class of policies π_{θ} where $\theta \in \Theta \subset R^d$
- The original problem is defined as

$$\max_{\theta \in \Theta} J_r^{\pi_{\theta}} \quad \text{s.t.} \quad J_c^{\pi_{\theta}} \geq 0$$

- Regret and violation are defined as

$$\text{Reg}_T(\mathbb{A}, \mathcal{M}) \triangleq \sum_{t=0}^{T-1} \left(J_r^{\pi^*} - r(s_t, a_t) \right) \quad \text{Vio}_T(\mathbb{A}, \mathcal{M}) \triangleq - \sum_{t=0}^{T-1} c(s_t, a_t)$$

Formulation: Ergodicity

- Assumption 1: The MDP is ergodic.
- Unique stationary distribution

$$d^{\pi_\theta}(s) = \lim_{T \rightarrow \infty} \frac{1}{T} \left[\sum_{t=0}^{T-1} \Pr(s_t = s | s_0 \sim \rho, \pi_\theta) \right]$$

- Define the mixing time

Definition 1. *The mixing time of an MDP \mathcal{M} with respect to a policy parameter θ is defined as,*

$$t_{\text{mix}}^\theta \triangleq \min \left\{ t \geq 1 \mid \left\| (P^{\pi_\theta})^t(s, \cdot) - d^{\pi_\theta} \right\| \leq \frac{1}{4}, \forall s \in \mathcal{S} \right\}$$

- One key property of mixing time under ergodic assumption is

$$\left\| (P^{\pi_\theta})^t(s, \cdot) - d^{\pi_\theta} \right\| \leq 2 \cdot 2^{-\frac{t}{t_{\text{mix}}}}$$

Method

- Define the Lagrange function $J_L(\theta, \lambda) = J_r(\theta) + \lambda J_c(\theta)$
- The problem becomes a saddle-point problem

$$\max_{\theta \in \Theta} \min_{\lambda \geq 0} J_L(\theta, \lambda)$$

- Use primal-dual method with policy gradient

$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} J_L(\theta_k, \lambda_k), \quad \lambda_{k+1} = \mathcal{P}_{[0, \frac{2}{\delta}]}[\lambda_k - \beta J_c(\theta_k)]$$

- It is well known that the gradient can be written as

$$\nabla_{\theta} J_L(\theta, \lambda) = \mathbf{E}_{s \sim d^{\pi_{\theta}}, a \sim \pi_{\theta}(s)} [A_{L, \lambda}^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

Algorithm

Algorithm 1 Primal-Dual Parameterized Policy Gradient

- 1: **Input:** Episode length H , learning rates α, β , initial parameters θ_1, λ_1 , initial state $s_0 \sim \rho(\cdot)$,
- 2: $K = T/H$
- 3: **for** $k \in \{1, \dots, K\}$ **do**
- 4: $\mathcal{T}_k \leftarrow \phi$
- 5: **for** $t \in \{(k-1)H, \dots, kH-1\}$ **do**
- 6: Execute $a_t \sim \pi_{\theta_k}(\cdot|s_t)$
- 7: Observe $r(s_t, a_t), c(s_t, a_t)$ and s_{t+1}
- 8: $\mathcal{T}_k \leftarrow \mathcal{T}_k \cup \{(s_t, a_t)\}$
- 9: **end for**
- 10: **for** $t \in \{(k-1)H, \dots, kH-1\}$ **do**
- 11: Obtain $\hat{A}_{L, \lambda_k}^{\pi_{\theta_k}}(s_t, a_t)$ via Algorithm 2 and \mathcal{T}_k
- 12: **end for**
- 13: Compute ω_k using (15) $\omega_k \triangleq \hat{\nabla}_{\theta} J_L(\theta_k, \lambda_k) = \frac{1}{H} \sum_{t=t_k}^{t_{k+1}-1} \hat{A}_{L, \lambda_k}^{\pi_{\theta_k}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta_k}(a_t|s_t)$
- 14: Update the parameters:

$$\theta_{k+1} = \theta_k + \alpha \omega_k,$$

$$\lambda_{k+1} = \mathcal{P}_{[0, \frac{2}{\delta}]}[\lambda_k - \beta \hat{J}_c(\theta_k)]$$

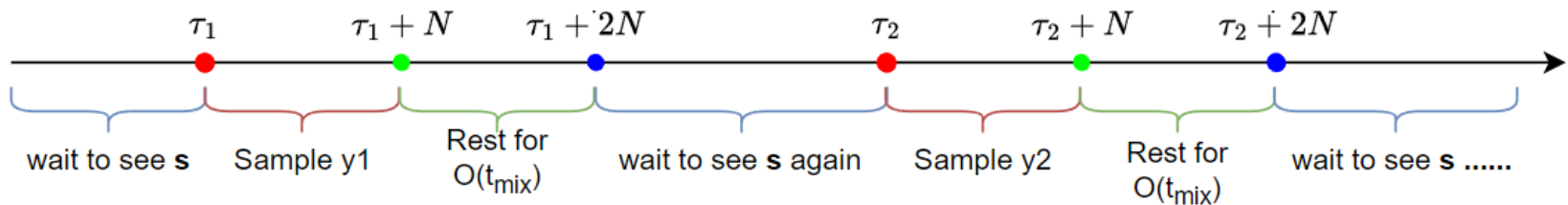
$$\text{where } \hat{J}_c(\theta_k) = \frac{1}{H - N} \sum_{t=(k-1)H+N}^{kH-1} c(s_t, a_t)$$

- 15: **end for**
-

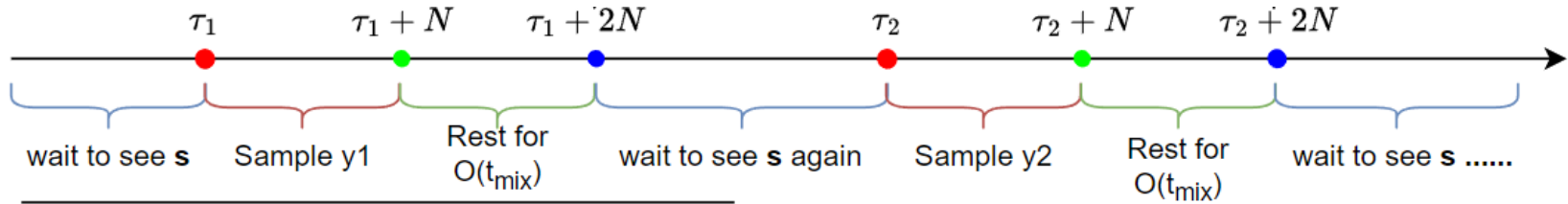
Challenge

$$\nabla_{\theta} J_L(\theta, \lambda) = \mathbf{E}_{s \sim d^{\pi_{\theta}}, a \sim \pi_{\theta}(s)} [A_{L, \lambda}^{\pi_{\theta}}(s, a) \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

- **Challenge: How to get a good estimation of advantage function?**
- Discounted setups typically assume access to simulator to obtain unbiased value estimator, while such estimate is needed from samples.
- Since it is single trajectory, unbiased estimator is not straightforward to obtain.
- **Solution: Divide each episode into following sub-trajectories**



Algorithm



Algorithm 2 Advantage Estimation

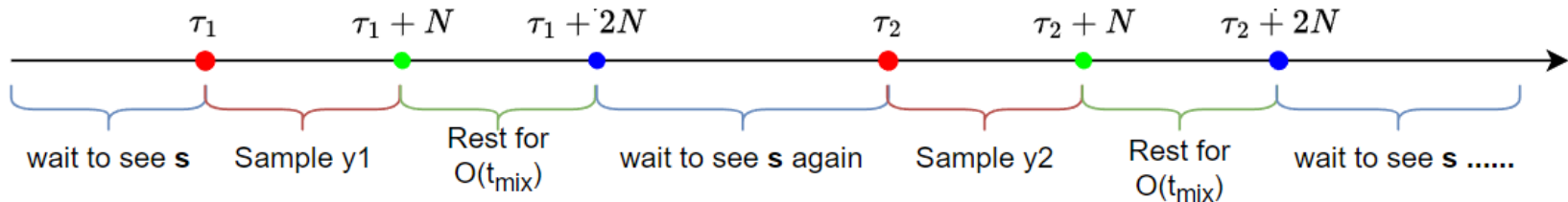
- 1: **Input:** Trajectory $(s_{t_1}, a_{t_1}, \dots, s_{t_2}, a_{t_2})$, state s , action a , Lagrange multiplier λ , and parameter θ
- 2: **Initialize:** $M \leftarrow 0, \tau \leftarrow t_1$
- 3: **Define:** $N = 4t_{\text{mix}} \log_2 T$.
- 4: **while** $\tau \leq t_2 - N$ **do**
- 5: **if** $s_\tau = s$ **then**
- 6: $M \leftarrow M + 1$.
- 7: $\tau_M \leftarrow \tau$
- 8: $g_M \leftarrow \sum_{t=\tau}^{\tau+N-1} g(s_t, a_t), \forall g \in \{r, c\}$
- 9: $\tau \leftarrow \tau + 2N$.
- 10: **else**
- 11: $\tau \leftarrow \tau + 1$.
- 12: **end if**
- 13: **end while**
- 14: **if** $M > 0$ **then**
- 15: $\hat{Q}_g(s, a) = \frac{1}{\pi_\theta(a|s)} \left[\frac{1}{M} \sum_{i=1}^M g_i 1(a_{\tau_i} = a) \right]$,
- 16: $\hat{V}_g(s) = \frac{1}{M} \sum_{i=1}^M g_i, \forall g \in \{r, c\}$
- 17: **else**
- 18: $\hat{V}_g(s) = 0, \hat{Q}_g(s, a) = 0, \forall g \in \{r, c\}$
- 19: **end if**
- 20: **return** $(\hat{Q}_r(s, a) - \hat{V}_r(s)) + \lambda(\hat{Q}_c(s, a) - \hat{V}_c(s))$

$$\hat{V}^{\pi_{\theta_k}}(s) = \frac{1}{i} \sum_{j=1}^i y_j \quad y_j = \sum_{t=\tau_i}^{\tau_i+N} r(s_t, a_t)$$

$$\hat{Q}^{\pi_\theta}(s, a) = \frac{1}{\pi_{\theta_k}(a|s)} \left[\frac{1}{i} \sum_{j=1}^i y_j 1(a_{\tau_j} = a) \right]$$

$$\hat{A}^{\pi_{\theta_k}}(s, a) = \hat{Q}^{\pi_{\theta_k}}(s, a) - \hat{V}^{\pi_{\theta_k}}(s)$$

Algorithm: bound the variance



- Finally, define $H = O(T^{-\xi})$, it can be proved that

$$\mathbf{E} \left[\left(\hat{A}_{L, \lambda_k}^{\pi_{\theta_k}}(s, a) - A_{L, \lambda_k}^{\pi_{\theta_k}}(s, a) \right)^2 \right] \leq \mathcal{O} \left(\frac{t_{\text{hit}} N^3 \log T}{\delta^2 H \pi_{\theta_k}(a|s)} \right) = \mathcal{O} \left(\frac{t_{\text{mix}}^2 (\log T)^2}{\delta^2 T^\xi \pi_{\theta_k}(a|s)} \right)$$

- Using Lemma 2 and Concentration inequality for Markov samples, the gradient estimator is close to the true gradient

$$\mathbf{E} \left[\|\omega_k - \nabla_{\theta} J_L(\theta_k, \lambda_k)\|^2 \right] \leq \tilde{\mathcal{O}} \left(\delta^{-2} A G^2 t_{\text{mix}}^2 T^{-\xi} \right)$$

Challenge: Lack of strong duality

- The global convergence of Lagrange function can be bounded as

$$\begin{aligned} \frac{1}{K} \sum_{k=1}^K \mathbf{E} \left(J_L(\pi^*, \lambda_k) - J_L(\theta_k, \lambda_k) \right) &\leq G \left(1 + \frac{1}{\mu_F} \right) \tilde{\mathcal{O}} \left(\sqrt{\beta} + \frac{\sqrt{AG}t_{\text{mix}}}{\delta T^{\xi/2}} + \frac{\sqrt{Lt_{\text{mix}}t_{\text{hit}}}}{\delta T^{(1-\xi)/2}} \right) \\ &+ \frac{B}{L} \tilde{\mathcal{O}} \left(\frac{AG^2t_{\text{mix}}^2}{\delta^2 T^\xi} + \frac{Lt_{\text{mix}}t_{\text{hit}}}{\delta^2 T^{1-\xi}} + \beta \right) + \tilde{\mathcal{O}} \left(\frac{Lt_{\text{mix}}t_{\text{hit}} \mathbf{E}_{s \sim d^{\pi^*}} [KL(\pi^*(\cdot|s) \parallel \pi_{\theta_1}(\cdot|s))] }{T^{1-\xi}\delta} \right) + \sqrt{\epsilon_{\text{bias}}} \end{aligned}$$

- However, the strong duality property does not hold under the general parameterization
- The following Lemma serves as an important tool in disentangling the convergence rates of regret and constraint violation

Lemma . Let Assumption [2](#) hold. For any constant $C \geq 2\lambda^*$, if there exists a $\pi \in \Pi$ and $\zeta > 0$ such that $J_r^{\pi^*} - J_r^\pi + C[-J_c^\pi] \leq \zeta$, then

$$-J_c^\pi \leq 2\zeta/C$$

Global convergence

- The sample complexity of reward and violation are

$$\frac{1}{K} \sum_{k=1}^K \mathbf{E} \left(J_r^{\pi^*} - J_r(\theta_k) \right) \leq \sqrt{\epsilon_{\text{bias}}} + \tilde{\mathcal{O}} \left(T^{-\eta/2} + T^{-\xi/2} + T^{-(1-\xi)/2} \right),$$

$$\mathbf{E} \left[\frac{1}{K} \sum_{k=1}^K -J_c(\theta_k) \right] \leq \delta \sqrt{\epsilon_{\text{bias}}} + \tilde{\mathcal{O}} \left(T^{-(1-\xi-\eta)} + T^{-\eta/2} + T^{-\xi/2} + T^{-(1-\xi)/2} \right)$$

- The best choice for ξ and η can be solved by

$$\max_{(\eta, \xi) \in (0,1)^2} \min \left\{ 1 - \xi - \eta, \frac{\eta}{2}, \frac{\xi}{2}, \frac{1 - \xi}{2} \right\}$$

- Finally, the regret and violation can be achieved as

$$\mathbf{E} [\text{Reg}_T] \leq T \sqrt{\epsilon_{\text{bias}}} + \tilde{\mathcal{O}}(T^{4/5}) + \mathcal{O}(t_{\text{mix}})$$

$$\mathbf{E} [\text{Vio}_T] \leq T \delta \sqrt{\epsilon_{\text{bias}}} + \tilde{\mathcal{O}}(T^{4/5}) + \mathcal{O}(t_{\text{mix}})$$