

# Average-Constrained Policy Optimization

**Akhil Agnihotri**

**Rahul Jain**

**Haipeng Luo**

*University of Southern California*

AKHIL.AGNIHOTRI@USC.EDU

RAHUL.JAIN@USC.EDU

HAIPENGL@USC.EDU

## Abstract

Reinforcement Learning (RL) with constraints is becoming an increasingly important problem for various applications. Often, the average criterion is more suitable than the discounted criterion. Yet, RL for average criterion-constrained MDPs remains a challenging problem. Algorithms designed for discounted constrained RL problems often do not perform well for the average CMDP setting. In this paper, we introduce a new policy optimization with function approximation algorithm for constrained MDPs with the average criterion. We develop basic sensitivity theory for average MDPs, and then use the corresponding bounds in the design of the algorithm. We provide theoretical guarantees on its performance, and through extensive experimental work in various challenging MuJoCo environments, show the superior performance of the algorithm when compared to other state-of-the-art algorithms adapted for the average CMDP setting.

## 1. Introduction

In recent years, reinforcement learning (RL) techniques have achieved remarkable successes on a variety of complex tasks including the game of Go [30], robotic control [22], and the real-time game StarCraft [36]. In these tasks, the agents are allowed to explore the state space and experiment with every action during training. However, in many realistic applications such as navigation robots [4, 13] and robotic assistants, considerations of safety and cost require the agent to operate with some constraints.

Even if constraints are incorporated using the Constrained Markov Decision Process (CMDP) framework [12], these deep RL (DRL) constrained policy algorithms are only suitable for the discounted setting [3, 32]. Given the need for average-reward algorithms, one might posit that having the discount factor  $\gamma \rightarrow 1$  in algorithms for the discounted setting would obviate the need for policy optimization algorithms specifically for the constrained average-reward setting. However, many algorithms break down when the discount factor gets too close to 1 [5, 17, 20]. Although policy optimization algorithms for the average reward case have been proposed [23, 38, 43], they cannot incorporate constraints.

In this work, we propose such a policy optimization algorithm for an average-CMDP with function approximation called Average-Constrained Policy Optimization (ACPO). Our approach is motivated by theoretical guarantees that bound the difference between the average long-run rewards or costs of different policies. For experimental validation, we use several MuJoCo environments [34], and demonstrate the effectiveness and superior performance of the ACPO algorithm as compared to others.

**Main Contributions.** A brief summary of our work’s contributions:

- We propose ACPO, the first on-policy Deep RL algorithm with function approximation specifically designed for average-CMDPs, inspired by policy optimization algorithms.

- Our empirical results show that ACPO out-performs all state-of-the-art Deep RL algorithms such as CPO [3], PCPO [41], PPO [29], BVF-PPO [26] and ATRPO [43], that while not designed for average-CMDPs could be applied with some tweaks.

## Related Work

Learning constraint-satisfaction policies has been explored in the Deep RL literature [10]. This can either be (1) through expert annotations and demonstrations [9, 25] or, (2) by exploration with constraint satisfaction [3, 32]. While the former approach is not scalable since it requires human interventions, current state-of-the-art algorithms for the latter are not applicable to the average reward setting. Please refer to A.1.1 for a more detailed overview of past literature. We now begin with some preliminaries.

## 2. Preliminaries

A Markov decision process (MDP) is a tuple,  $(S, A, r, P, \mu)$ , where  $S$  is the set of states,  $A$  is the set of actions,  $r : S \times A \times S \rightarrow \mathbb{R}$  is the reward function,  $P : S \times A \times S \rightarrow [0, 1]$  is the transition probability function such that  $P(s'|s, a)$  is the probability of transitioning to state  $s'$  from state  $s$  by taking action  $a$ , and  $\mu : S \rightarrow [0, 1]$  is the initial state distribution. A stationary policy  $\pi : S \rightarrow \Delta(A)$  is a mapping from states to probability distributions over the actions, with  $\pi(a|s)$  denoting the probability of selecting action  $a$  in state  $s$ , and  $\Delta(A)$  is the probability simplex over the action space  $A$ . We denote the set of all stationary policies by  $\Pi$ . For the average setting, we will make the standard assumption that the MDP is ergodic.

In reinforcement learning, we aim to select a policy  $\pi$  which maximizes a performance measure,  $J(\pi)$ , which, for continuous control tasks is either the discounted reward criterion or the average reward approach. Please refer to Appendix A.2 for more details.

### 2.1. Average criterion

The average-reward objective is given by:

$$J(\pi) := \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{N-1} r(s_t, a_t, s_{t+1}) \right] = \mathbb{E}_{s \sim d_\pi, a \sim \pi(\cdot|s), s' \sim P(\cdot|s, a)} [r(s, a, s')], \quad (1)$$

where  $d_\pi(s) := \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{t=0}^{N-1} P_{\tau \sim \pi}(s_t = s)$  is the *stationary state distribution* under policy  $\pi$ . The limits in  $J(\pi)$  and  $d_\pi(s)$  are guaranteed to exist under our ergodic assumption. Since the MDP is aperiodic, it can also be shown that  $d_\pi(s) = \lim_{t \rightarrow \infty} P_{\tau \sim \pi}(s_t = s)$ .

In the average setting, we seek to keep the estimate of the state value function unbiased and hence, introduce the *average-reward bias* and *average-reward action-bias function* respectively as

$$\bar{V}^\pi(s) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} (r(s_t, a_t, s_{t+1}) - J(\pi)) \mid s_0 = s \right]; \quad \bar{Q}^\pi(s, a) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} (r(s_t, a_t, s_{t+1}) - J(\pi)) \mid s_0 = s, a_0 = a \right].$$

Finally, define the *average-reward advantage function* as  $\bar{A}^\pi(s, a) := \bar{Q}^\pi(s, a) - \bar{V}^\pi(s)$ .

### 2.2. Constrained MDPs

A constrained Markov decision process (CMDP) is an MDP augmented with constraints that restrict the set of allowable policies for that MDP. Specifically, we augment the MDP with a set  $C$  of auxiliary cost functions,  $C_1, \dots, C_m$  (with each function  $C_i : S \times A \times S \rightarrow \mathbb{R}$  mapping transition tuples to costs, just like

the reward function), and bounds  $l_1, \dots, l_m$ . Similar to the value functions being defined for the average reward criterion, we define the average cost objective with respect to the cost function  $C_i$  as

$$J_{C_i}(\pi) := \lim_{N \rightarrow \infty} \frac{1}{N} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{N-1} C_i(s_t, a_t, s_{t+1}) \right] = \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi \\ s' \sim P(\cdot | s, a)}} [C_i(s, a, s')]. \quad (2)$$

where  $J_{C_i}$  will be referred to as the *average cost* for constraint  $C_i$ . The set of feasible stationary policies for a CMDP then is given by  $\Pi_C := \{\pi \in \Pi : J_{C_i}(\pi) \leq l_i, \forall i \in \{1, \dots, M\}\}$ . The goal is to find a policy  $\pi^*$  such that  $\pi^* \in \arg \max_{\pi \in \Pi_C} J(\pi)$ .

Lastly, analogous to  $\bar{V}^\pi$ ,  $\bar{Q}^\pi$ , and  $\bar{A}^\pi$ , we define similar quantities for the cost functions  $C_i(\cdot)$ , and denote them by  $\bar{V}_{C_i}^\pi$ ,  $\bar{Q}_{C_i}^\pi$ , and  $\bar{A}_{C_i}^\pi$ .

### 3. ACPO: The Average-Constrained Policy Optimization Algorithm

In this section, we present the main results of our work. For conciseness, we denote by  $d_\pi \in \mathbb{R}^{|S|}$  the column vector whose components are  $d_\pi(s)$  and  $P_\pi \in \mathbb{R}^{|S| \times |S|}$  to be the state transition probability matrix under policy  $\pi$ . For the theoretical buildup to our main result, please see Appendix A.3. Proofs of theorems and lemmas, if not already given, are available in Appendix A.

#### 3.1. Trust Region Based Approach

For large or continuous state and action CMDPs, solving for the exact optimal policy is impractical. However, *trust region*-based policy optimization algorithms have proven to be effective for solving such problems [2, 27–29]. Hence, we now introduce the ACPO algorithm, which is inspired by the trust region formulations above as the following optimization problem:

$$\begin{aligned} & \underset{\pi \in \Pi_\Theta}{\text{maximize}} && \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi}} [\bar{A}^{\pi_{\theta_k}}(s, a)] \\ & \text{s.t.} && J_{C_i}(\pi_{\theta_k}) + \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi}} [\bar{A}_{C_i}^{\pi_{\theta_k}}(s, a)] \leq l_i, \quad \forall i \\ & && \bar{D}_{\text{KL}}(\pi \parallel \pi_{\theta_k}) \leq \delta \end{aligned} \quad (3)$$

where  $\bar{D}_{\text{KL}}(\pi \parallel \pi_{\theta_k}) := \mathbb{E}_{s \sim d_{\pi_{\theta_k}}} [D_{\text{KL}}(\pi \parallel \pi_{\theta_k})[s]]$ ,  $\bar{A}^{\pi_{\theta_k}}(s, a)$  is the average advantage function defined earlier, and  $\delta > 0$  is a step size. We use this form of updates as it is an approximation to the lower bound given in Proposition A.4 and the upper bound given in Corollary A.5.

In most cases, the trust region threshold for formulations like Eq. (3) are heuristically motivated. We now show that it is quantitatively motivated and comes with a worst case performance degradation and constraint violation. Proof is in Appendix A.5.

**Theorem 3.1** *Let  $\pi_{\theta_{k+1}}$  be the optimal solution to Eq. (3) for some  $\pi_{\theta_k} \in \Pi_\Theta$ . Then, we have*

$$J(\pi_{\theta_{k+1}}) - J(\pi_{\theta_k}) \geq -\sqrt{2(\delta + V_{\max})} \nu^{\pi_{\theta_{k+1}}} \quad (4)$$

$$\text{and } J_{C_i}(\pi_{\theta_{k+1}}) \leq l_i + \sqrt{2(\delta + V_{\max})} \nu_{C_i}^{\pi_{\theta_{k+1}}} \quad \forall i, \quad (5)$$

where  $\nu^{\pi_{\theta_{k+1}}} = \sigma^{\pi_{\theta_{k+1}}} \max_s \left| \mathbb{E}_{a \sim \pi_{\theta_{k+1}}} [\bar{A}^{\pi_{\theta_k}}(s, a)] \right|$ ,  $\nu_{C_i}^{\pi_{\theta_{k+1}}} = \sigma^{\pi_{\theta_{k+1}}} \max_{i,s} \left| \mathbb{E}_{a \sim \pi_{\theta_{k+1}}} [\bar{A}_{C_i}^{\pi_{\theta_k}}(s, a)] \right|$ ,  $V_{\max} = \max_i \beta_i^2$ , and  $\beta_i = [J_{C_i}(\pi_{\theta_k}) - l_i]_+$ .

**Remark 3.7.** Note that if the constraints are ignored (by setting  $V_{max} = 0$ ), then this bound is tighter than given in [43] for the unconstrained average-reward setting.

However, the update rule of Eq. (3) is difficult to implement in practice as it takes steps that are too small, which degrades convergence. In addition, it requires the exact knowledge of  $\bar{A}^{\pi_{\theta_k}}(s, a)$  which is computationally infeasible for large-scale problems. In the next section, we will introduce a specific sampling-based practical algorithm to alleviate these concerns.

## 4. Practical Implementation of ACPO

In this section, we introduce a practical version of the ACPO algorithm with a principle recovery method. With a small step size  $\delta$ , we can approximate the reward function and constraints with a first order expansion, and approximate the KL divergence constraint with a second order expansion. This gives us a new optimization problem which can be solved exactly using Lagrangian duality.

### 4.1. An Implementation of ACPO

Since we are working with a parameterized class, we shall now overload notation to use  $\theta_k$  as the policy at iteration  $k$ , i.e.,  $\theta_k \equiv \pi_{\theta_k}$ . In addition, we use  $g$  to denote the gradient of the advantage function objective,  $a_i$  to denote the gradient of the advantage function of the cost  $C_i$ ,  $H$  as the Hessian of the KL-divergence. Formally,

$$g := \nabla_{\theta} \mathbb{E}_{\substack{s \sim d_{\theta_k} \\ a \sim \theta}} [\bar{A}^{\theta_k}(s, a)] \Big|_{\theta=\theta_k}, \quad a_i := \nabla_{\theta} \mathbb{E}_{\substack{s \sim d_{\theta_k} \\ a \sim \theta}} [\bar{A}_{C_i}^{\theta_k}(s, a)] \Big|_{\theta=\theta_k}, \quad H := \nabla_{\theta}^2 \mathbb{E}_{s \sim d_{\theta_k}} [D_{\text{KL}}(\theta \parallel \theta_k)] \Big|_{\theta=\theta_k}.$$

In addition, let  $c_i := J_{C_i}(\theta_k) - l_i$ . The approximation to the problem in Eq. (3) is:

$$\max_{\theta} g^T(\theta - \theta_k) \quad \text{s.t.} \quad c_i + a_i^T(\theta - \theta_k) \leq 0, \quad \forall i \quad \text{and} \quad \frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k) \leq \delta. \quad (6)$$

This is a convex optimization problem in which strong duality holds, and hence it can be solved using a Lagrangian method. The update rule for the dual problem then takes the form

$$\theta_{k+1} = \theta_k + \frac{1}{\lambda^*} H^{-1}(g - A\mu^*). \quad (7)$$

where  $\lambda^*$  and  $\mu^*$  are the Lagrange multipliers satisfying the dual

$$\max_{\substack{\lambda \geq 0 \\ \mu \geq 0}} \frac{-1}{2\lambda} (g^T H^{-1} g - 2r^T \mu + \mu^T S \mu) + \mu^T c - \frac{\lambda \delta}{2}, \quad (8)$$

with  $r := g^T H^{-1} A$ ,  $S := A^T H^{-1} A$ ,  $A := [a_1, \dots, a_m]$ , and  $c := [c_1, \dots, c_m]^T$ .

### 4.2. Feasibility and Recovery

The approximation regime described in Eq. (6) requires  $H$  to be invertible. For large parametric policies,  $H$  is computed using the conjugate gradient method as in [27]. However, in practice, using this approximation along with the associated statistical sampling errors, there might be potential violations of the approximate constraints leading to infeasible policies.

To rectify this, for the case where we only have one constraint, one can recover a feasible policy by applying a recovery step inspired by the TRPO update on the cost surrogate as:

$$\theta_{k+1/2} = \theta_k - \sqrt{2\delta} \left[ t \cdot \frac{H^{-1}a}{\sqrt{a^T H^{-1}a}} + (1-t) \cdot \frac{H^{-1}g}{\sqrt{g^T H^{-1}g}} \right] \quad (9)$$

**Algorithm 1** Average-Constrained Policy Optimization (ACPO)

- 
- 1: **Input:** Initial random policy  $\pi_0 \in \Pi_\theta$
  - 2: **for**  $k = 0, 1, 2, \dots, K$  **do**
  - 3:   Sample a set of trajectories  $\Omega$  using  $\pi_k = \pi_{\theta_k}$
  - 4:   Find estimates of  $g, a, H, c$  using  $\Omega$
  - 5:   **if** a feasible solution to Equation (6) exists **then**
  - 6:     Solve dual problem in Equation (8) for  $\lambda_k^*, \mu_k^*$
  - 7:     Find policy update  $\pi_{k+1}$  with Equation (7)
  - 8:   **else**
  - 9:     Find recovery policy  $\pi_{k+1/2}$  with Equation (9)
  - 10:    Obtain  $\pi_{k+1}$  by linesearch till approximate constraint satisfaction of Equation (6)
  - 11:   **end if**
  - 12: **end for**
- 

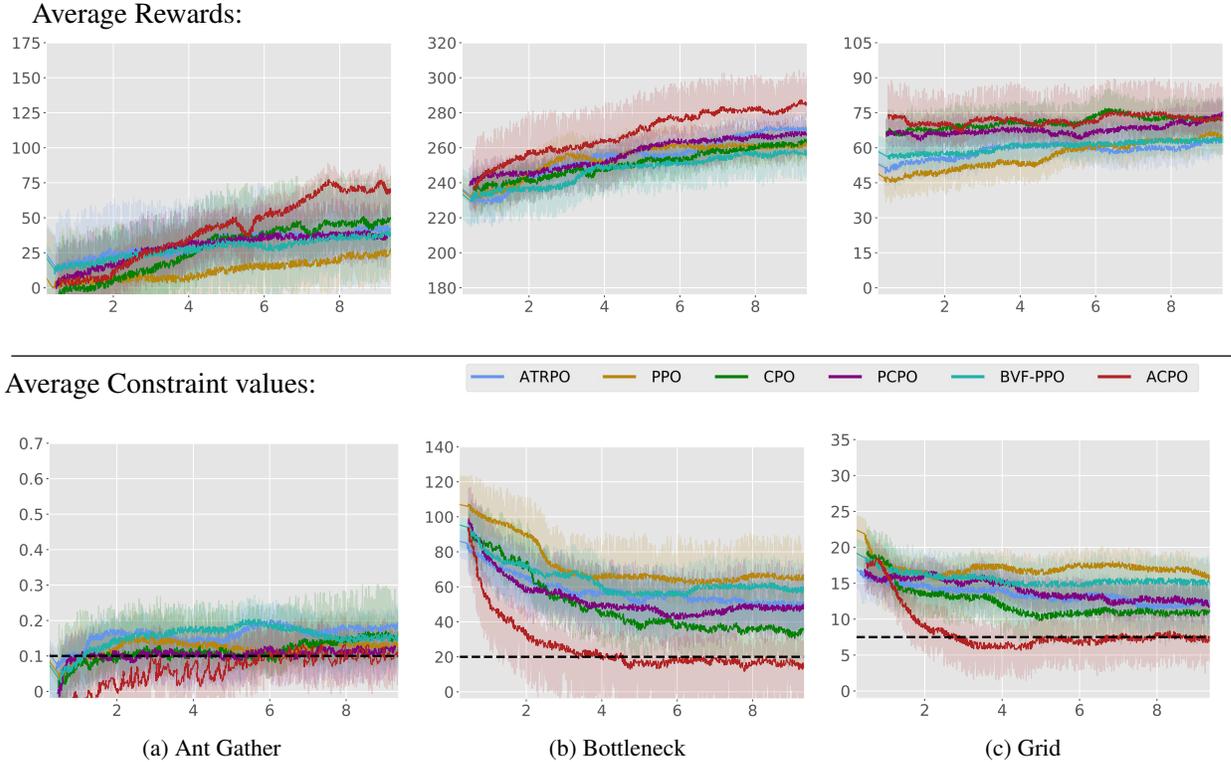


Figure 1: The average reward and constraint cost function values vs iterations (in  $10^4$ ) learning curves for some algorithm-task pairs. Solid lines in each figure are the empirical means, while the shaded area represents 1 standard deviation, all over 5 runs. The dashed line in constraint plots is the constraint threshold  $l$ . ATRPO and PPO are tested with constraints, which are included in their Lagrangian formulation. Additional results are available in Appendix A.9.

where  $t \in [0, 1]$ . Based on this, Algorithm 1 provides a basic outline of ACPO. For more details of the algorithm, see Appendix A.6.

## 5. Empirical Results

We conducted a series of experiments to evaluate the relative performance of the ACPO algorithm and answer the following questions: (i) Does ACPO learn a sequence of constraint satisfying policies while maximizing the average reward in the long run? (ii) How does ACPO compare with the already existing constraint policy optimization algorithms which are applied with a large discount factor? (iii) What are the factors that affect the performance of ACPO?

We work with the MuJoCo environments to train the various learning agent on the following tasks - *Gather*, *Circle*, *Grid*, and *Bottleneck* tasks (see Figure 2 in Appendix A.8.1 for more details on the environments). For our experiments we only work with a single constraint with policy recovery using Eq. (9) (this is only a computational limitation; ACPO in principle can handle multiple constraints). We compare ACPO with the following baseline algorithms: CPO [3], ATRPO [43], PCPO [41] (a close variant of CPO), BVF-PPO [26] and PPO [29].

Although ATRPO and PPO originally do not incorporate constraints, for fair comparison, we introduce constraints using a Lagrangian. Also, CPO, PCPO and PPO are compared with  $\gamma = 0.999$ . Learning curves for the algorithms are compiled in Figure 1 (for Point-Circle, Point-Gather, and Ant-Circle see Appendix A.9). Since there are two objectives (rewards in the objective and costs in the constraints), we show the plots which maximize the reward objective while satisfying the cost constraint. See Appendix A.8 for more details.

### 5.1. Performance Analysis

From Figure 1, we can see that ACPO is able to improve the reward objective while having approximate constraint satisfaction on all tasks. In particular, ACPO is the only algorithm that best learns almost-constraint-satisfying maximum average-reward policies across all tasks: in a simple Gather environment, ACPO is able to almost exactly track the cost constraint values to within the given threshold  $l$ ; however, for the high dimensional Grid and Bottleneck environments we have more constraint violations due to complexity of the policy behavior. Regardless, in these environments, ACPO still outperforms all other baselines with less average constraint violation, while maintaining similar if not better average-reward performance. See Appendix A.8.3 for a more detailed discussion.

In Equation (9) we introduced a hyperparameter  $t$ , which provides for an intuitive trade-off as follows: either we purely decrease the constraint violations ( $t = 1$ ), or we decrease the average-reward ( $t = 0$ ), which consequently decreases the constraint violation. The latter formulation is principled in that if we decrease rewards, we are bound to decrease constraints violation due to the nature of the environments. See Appendix A.9.2 for more discussion.

## 6. Conclusion

In this paper, we studied the problem of learning policies that maximize average-rewards for a given CMDP with average-cost constraints. We showed that the current algorithms with state-of-the-art performance and constraint violation bounds for the discounted setting do not generalize to the average setting. We then proposed a new algorithm, the Average-Constrained Policy Optimization (ACPO) that is of course directly inspired by the PPO class of algorithms but based on theoretical sensitivity-type bounds for average-CMDPs we derive, and use in designing the algorithm. Our experimental results on a range of MuJoCo environments (including some high dimensional ones) show the effectiveness of ACPO on average CMDP RL problems, as well as its superior performance vis-a-vis some current alternatives. A possible direction for future work is implementation of ACPO to fully exploit the parallelization potential.

## References

- [1] Jinane Abounadi, D Bertsekas, and Vivek S Borkar. Learning algorithms for markov decision processes with average cost. *SIAM Journal on Control and Optimization*, 40(3):681–698, 2001.
- [2] Joshua Achiam. UC Berkeley CS 285 (Fall 2017), Advanced Policy Gradients, 2017. URL: [http://rail.eecs.berkeley.edu/deeprlcourse-fa17/f17docs/lecture\\_13\\_advanced\\_pg.pdf](http://rail.eecs.berkeley.edu/deeprlcourse-fa17/f17docs/lecture_13_advanced_pg.pdf). Last visited on 2020/05/24.
- [3] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 22–31. JMLR. org, 2017.
- [4] Akhil Agnihotri, Prathamesh Saraf, and Kriti Rajesh Bapnad. A convolutional neural network approach towards self-driving cars. In *2019 IEEE 16th India Council International Conference (INDICON)*, pages 1–4, 2019. doi: 10.1109/INDICON47234.2019.9030307.
- [5] Marcin Andrychowicz, Anton Raichuk, Piotr Stańczyk, Manu Orsini, Sertan Girgin, Raphael Marinier, Léonard Hussenot, Matthieu Geist, Olivier Pietquin, Marcin Michalski, et al. What matters in on-policy reinforcement learning? a large-scale empirical study. *arXiv preprint arXiv:2006.05990*, 2020.
- [6] Liyu Chen, Rahul Jain, and Haipeng Luo. Learning infinite-horizon average-reward Markov decision process with constraints. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, pages 3246–3270. PMLR, 17–23 Jul 2022. URL <https://proceedings.mlr.press/v162/chen22i.html>.
- [7] Grace E Cho and Carl D Meyer. Comparison of perturbation bounds for the stationary distribution of a markov chain. *Linear Algebra and its Applications*, 335(1-3):137–150, 2001.
- [8] Peter G Doyle. The kemeny constant of a markov chain. *arXiv preprint arXiv:0909.2636*, 2009.
- [9] Yang Gao, Huazhe Xu, Ji Lin, Fisher Yu, Sergey Levine, and Trevor Darrell. Reinforcement learning from imperfect demonstrations. *arXiv preprint arXiv:1802.05313*, 2018.
- [10] Javier Garcia and Fernando Fernandez. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.
- [11] Charles Miller Grinstead and James Laurie Snell. *Introduction to probability*. American Mathematical Soc., 2012.
- [12] Moshe Haviv. On constrained markov decision processes. *Operations research letters*, 19(1):25–28, 1996.
- [13] Hanjiang Hu, Zuxin Liu, Sharad Chitlangia, Akhil Agnihotri, and Ding Zhao. Investigating the impact of multi-lidar placement on object detection for autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2550–2559, June 2022.
- [14] Jeffrey J Hunter. Stationary distributions and mean first passage times of perturbed markov chains. *Linear Algebra and its Applications*, 410:217–243, 2005.

- [15] Jeffrey J Hunter. *Mathematical techniques of applied probability: Discrete time models: Basic theory*, volume 1. Academic Press, 2014.
- [16] Jeffrey J Hunter. The computation of the mean first passage times for markov chains. *Linear Algebra and its Applications*, 549:100–122, 2018.
- [17] Nan Jiang, Satinder P Singh, and Ambuj Tewari. On structural properties of mdps that bound loss due to shallow planning. In *IJCAI*, pages 1640–1647, 2016.
- [18] Sham Kakade and John Langford. Approximately optimal approximate reinforcement learning. In *International Conference on Machine Learning*, volume 2, pages 267–274, 2002.
- [19] J.G. Kemeny and I.J. Snell. *Finite Markov Chains*. Van Nostrand, New Jersey, 1960.
- [20] Lucas Lehnert, Romain Laroche, and Harm van Seijen. On value function representation of long horizon problems. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [21] Mark Levene and George Loizou. Kemeny’s constant and the random surfer. *The American mathematical monthly*, 109(8):741–745, 2002.
- [22] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [23] Peng Liao, Zhengling Qi, Runzhe Wan, Predrag Klasnja, and Susan A Murphy. Batch policy learning in average reward markov decision processes. *The Annals of Statistics*, 50(6):3364–3387, 2022.
- [24] Xiaoteng Ma, Xiaohang Tang, Li Xia, Jun Yang, and Qianchuan Zhao. Average-reward reinforcement learning with trust region methods. *arXiv preprint arXiv:2106.03442*, 2021.
- [25] Aravind Rajeswaran, Vikash Kumar, Abhishek Gupta, Giulia Vezzani, John Schulman, Emanuel Todorov, and Sergey Levine. Learning complex dexterous manipulation with deep reinforcement learning and demonstrations. *arXiv preprint arXiv:1709.10087*, 2017.
- [26] Harsh Satija, Philip Amortila, and Joelle Pineau. Constrained Markov decision processes via backward value functions. In Hal Daumé III and Aarti Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 8502–8511. PMLR, 13–18 Jul 2020. URL <https://proceedings.mlr.press/v119/satija20a.html>.
- [27] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region policy optimization. In *International Conference on Machine Learning*, pages 1889–1897, 2015.
- [28] John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *International Conference on Learning Representations (ICLR)*, 2016.
- [29] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- [30] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of go without human knowledge. *nature*, 550(7676):354–359, 2017.

- [31] H Francis Song, Abbas Abdolmaleki, Jost Tobias Springenberg, Aidan Clark, Hubert Soyer, Jack W Rae, Seb Noury, Arun Ahuja, Siqu Liu, Dhruva Tirumala, et al. V-mpo: on-policy maximum a posteriori policy optimization for discrete and continuous control. *International Conference on Learning Representations*, 2020.
- [32] Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. *International Conference on Learning Representation (ICLR)*, 2019.
- [33] Ryan J Tibshirani. Dykstra’s algorithm, admm, and coordinate descent: Connections, insights, and extensions. *Advances in Neural Information Processing Systems*, 30, 2017.
- [34] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- [35] Eugene Vinitsky, Aboudy Kreidieh, Luc Le Flem, Nishant Kheterpal, Kathy Jang, Fangyu Wu, Richard Liaw, Eric Liang, and Alexandre M. Bayen. Benchmarks for reinforcement learning in mixed-autonomy traffic. In *Proceedings of Conference on Robot Learning*, pages 399–409, 2018.
- [36] Oriol Vinyals, Igor Babuschkin, Wojciech M Czarnecki, Michaël Mathieu, Andrew Dudzik, Junyoung Chung, David H Choi, Richard Powell, Timo Ewalds, Petko Georgiev, et al. Grandmaster level in starcraft ii using multi-agent reinforcement learning. *Nature*, 575(7782):350–354, 2019.
- [37] Quan Vuong, Yiming Zhang, and Keith W Ross. Supervised policy update for deep reinforcement learning. In *International Conference on Learning Representation (ICLR)*, 2019.
- [38] Yi Wan, Abhishek Naik, and Richard S Sutton. Learning and planning in average-reward markov decision processes. In *International Conference on Machine Learning*, pages 10653–10662. PMLR, 2021.
- [39] Chen-Yu Wei, Mehdi Jafarnia Jahromi, Haipeng Luo, and Rahul Jain. Learning infinite-horizon average-reward mdps with linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 3007–3015. PMLR, 2021.
- [40] Yuhuai Wu, Elman Mansimov, Roger B Grosse, Shun Liao, and Jimmy Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In *Advances in neural information processing systems (NIPS)*, pages 5285–5294, 2017.
- [41] Tsung-Yen Yang, Justinian Rosca, Karthik Narasimhan, and Peter J Ramadge. Projection-based constrained policy optimization. In *International Conference on Learning Representation (ICLR)*, 2020.
- [42] Shangdong Zhang, Yi Wan, Richard S Sutton, and Shimon Whiteson. Average-reward off-policy policy evaluation with function approximation. *arXiv preprint arXiv:2101.02808*, 2021.
- [43] Yiming Zhang and Keith Ross. Average reward reinforcement learning with monotonic policy improvement. *Preprint*, 2020.

## Appendix A. Appendix

### A.1. Preliminaries

#### A.1.1. RELATED WORK REVISITED

Previous work on RL with the average reward criterion has mostly attempted to extend stochastic approximation schemes for the tabular setting, such as Q-learning [1, 38], to the non-tabular setting with function approximation [39, 42]. [6] deals with online learning in a constrained MDP setting, but their aim is regret minimization or exploration, both in tabular settings. [42] provide bounds on the performance of a trust region algorithm for the the average reward setting but do not incorporate constraints.

Since our work is related to CPO [3] and ATRPO [43], we highlight our work’s novelty as follows. ACPO is designed to find the optimal policy of the *undiscounted* average-CMDP problem directly, with an average reward objective and constraints on average cost functions. While CPO is designed for discounted CMDPs (setting  $\gamma = 1$  in CPO does not yield the ACPO algorithm), and though we could use it for the average case by setting  $\gamma \rightarrow 1$ , we shall see that empirically it has inferior performance to ACPO. Moreover, theory for CPO and ATRPO cannot be extended trivially to the ACPO algorithm, but since they are all trust region methods, they share some analysis. Several of our techniques are still unique, e.g. in Lemma A.3, we use eigenvalues of the transition matrix to relate total variation (TV) of stationary distributions with the TV of the policies, and in Lemma A.11, we use the sublevel sets of constraints and projection inequality of Bregman divergence.

### A.2. Warmup on MDPs

#### A.2.1. DISCOUNTED CRITERION

For a given discount factor  $\gamma \in (0, 1)$ , the discounted reward objective is defined as

$$J_\gamma(\pi) := \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \right] = \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d_{\pi, \gamma} \\ a \sim \pi \\ s' \sim P(\cdot|s, a)}} [r(s, a, s')]$$

where  $\tau$  refers to a sample trajectory of  $(s_0, a_0, s_1, \dots)$  generated when following a policy  $\pi$ , that is,  $a_t \sim \pi(\cdot|s_t)$  and  $s_{t+1} \sim P(\cdot|s_t, a_t)$ ;  $d_{\pi, \gamma}$  is the *discounted occupation measure* that is defined by  $d_{\pi, \gamma}(s) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P_{\tau \sim \pi}(s_t = s)$ , which essentially refers to the discounted fraction of time spent in state  $s$  while following policy  $\pi$ .

#### A.2.2. POLICY IMPROVEMENT FOR DISCOUNTED CMDPS

In many on-policy constrained RL problems, we improve policies iteratively by maximizing a predefined function within a local region of the current best policy [3, 31, 32, 41]. [3] derived a policy improvement bound for the discounted CMDP setting as:

$$J_\gamma(\pi_{k+1}) - J_\gamma(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{\substack{s \sim d^{\pi_k} \\ a \sim \pi_{k+1}}} \left[ A_\gamma^{\pi_k}(s, a) - \frac{2\gamma\epsilon^{\pi_{k+1}}}{1-\gamma} D_{TV}(\pi_{k+1}||\pi_k)[s] \right], \quad (10)$$

where  $A_\gamma^{\pi_k}$  is the discounted version of the advantage function,  $\epsilon^{\pi_{k+1}} := \max_s |\mathbb{E}_{a \sim \pi_{k+1}} [A_\gamma^{\pi_k}(s, a)]|$ , and  $D_{TV}(\pi_{k+1}||\pi_k)[s] = (1/2) \sum_a |\pi_{k+1}(a|s) - \pi_k(a|s)|$  is the total variational divergence between  $\pi_{k+1}$  and  $\pi_k$  at  $s$ . These results laid the foundations for on-policy constrained RL algorithms [37, 40].

However, Equation (10) does not generalize to the average setting ( $\gamma \rightarrow 1$ ) (see Appendix A.6). In the next section, we will derive a policy improvement bound for the average case and present an algorithm based on trust region methods, which will generate almost-monotonically improving iterative policies.

### A.3. Build-up to Main Result

#### A.3.1. POLICY IMPROVEMENT FOR THE AVERAGE-CMDP

Let  $\pi'$  be the policy obtained via some update rule from the current policy  $\pi$ . Analogous to the discounted setting of a CMDP, we would like to characterize the performance difference  $J(\pi') - J(\pi)$  by an expression which depends on  $\pi$  and some divergence metric between the two policies. [43] give such a lemma.

**Lemma A.1** [43] *Under the unichain assumption of the underlying Markov chain, for any stochastic policies  $\pi$  and  $\pi'$ :*

$$J(\pi') - J(\pi) = \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)]. \quad (11)$$

Note that this difference depends on the stationary state distribution obtained from the *new* policy,  $d_{\pi'}$ . This is computationally impractical as we do not have access to this  $d_{\pi'}$ . Fortunately, by use of the following lemma we can show that if  $d_\pi$  and  $d_{\pi'}$  are “close” with respect to some metric, we can approximate Eq. (11) using samples from  $d_\pi$ .

**Lemma A.2** *Under the unichain assumption, for any stochastic policies  $\pi$  and  $\pi'$  we have:*

$$\left| J(\pi') - J(\pi) - \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \right| \leq 2\epsilon D_{TV}(d_{\pi'} \parallel d_\pi) \quad \text{where} \quad \epsilon = \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] \right|. \quad (12)$$

See Appendix A.4 for proof. Lemma A.2 implies  $J(\pi') \approx J(\pi) + \mathbb{E}[\bar{A}^\pi(s, a)]$  when  $d_\pi$  and  $d_{\pi'}$  are “close”. Now that we have established this approximation, we need to study the relation of how the actual change in policies affects their corresponding stationary state distributions. For this, we turn to standard analysis of the underlying Markov chain of the CMDP.

Under the ergodic assumption, we have that  $P_\pi$  is irreducible and hence its eigenvalues  $\{\lambda_{\pi,i}\}_{i=1}^{|S|}$  are such that  $\lambda_{\pi,1} = 1$  and  $\lambda_{\pi,i \neq 1} < 1$ . For our analysis, we define  $\sigma^\pi = \max_{i \neq 1} (1 - \lambda_{\pi,i})^{-1/2}$ , and from [21] and [8], we connect  $\{\lambda_{\pi,i}\}_{i=1}^{|S|}$  to the sensitivity of the stationary distributions to changes in the policy using the result below.

**Lemma A.3** *Under the ergodic assumption, the divergence between the stationary distributions  $d_\pi$  and  $d_{\pi'}$  is upper bounded as:*

$$D_{TV}(d_{\pi'} \parallel d_\pi) \leq \sigma^\star \mathbb{E}_{s \sim d_\pi} [D_{TV}(\pi' \parallel \pi)[s]] \quad \text{where} \quad \sigma^\star = \max_\pi \sigma^\pi. \quad (13)$$

See Appendix A.4 for proof. This bound is tighter and easier to compute than the one given by [43], which replaces  $\sigma^\star$  by  $\kappa^\star = \max_\pi \kappa^\pi$ , where  $\kappa^\pi$  is known as *Kemeny’s constant* [11, 16, 19]. It is interpreted as the expected number of steps to get to any goal state, where the expectation is taken with respect to the stationary-distribution of those states.

Combining the bounds in Lemma A.2 and Lemma A.3 gives us the following result:

**Proposition A.4** *Under the ergodic assumption, the following bounds hold for any stochastic policies  $\pi$  and  $\pi'$ :*

$$L_{\pi}^{-}(\pi') \leq J(\pi') - J(\pi) \leq L_{\pi}^{+}(\pi') \quad (14)$$

where

$$L_{\pi}^{\pm}(\pi') = \mathbb{E}_{\substack{s \sim d_{\pi} \\ a \sim \pi'}} [\bar{A}^{\pi}(s, a)] \pm 2\nu \mathbb{E}_{s \sim d_{\pi}} [D_{TV}(\pi' \parallel \pi)[s]] \quad \text{and} \quad \nu = \sigma^* \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^{\pi}(s, a)] \right|.$$

It is interesting to compare the inequalities of Equation (14) to Equation (11). The term  $\mathbb{E}[\bar{A}^{\pi}(s, a)]$  in Prop. A.4 is somewhat of a *surrogate* approximation to  $J(\pi') - J(\pi)$ , in the sense that it uses  $d_{\pi}$  instead of  $d_{\pi'}$ . As discussed before, we do not have access to  $d_{\pi'}$  since the trajectories of the new policy are not available unless the policy itself is updated. This surrogate is a first order approximation to  $J(\pi') - J(\pi)$  in the parameters of  $\pi'$  in a neighborhood around  $\pi$  [18]. Hence, Eq. (14) can be viewed as bounding the worst-case approximation error.

Extending this discussion to the cost function of our CMDP, similar expressions follow immediately.

**Corollary A.5** *For any policies  $\pi', \pi$ , and any cost function  $C_i$ , the following bound holds:*

$$M_{\pi}^{-}(\pi') \leq J_{C_i}(\pi') - J_{C_i}(\pi) \leq M_{\pi}^{+}(\pi') \quad (15)$$

where

$$M_{\pi}^{\pm}(\pi') = \mathbb{E}_{\substack{s \sim d_{\pi} \\ a \sim \pi'}} [\bar{A}_{C_i}^{\pi}(s, a)] \pm 2\nu_{C_i} \mathbb{E}_{s \sim d_{\pi}} [D_{TV}(\pi' \parallel \pi)[s]] \quad \text{and} \quad \nu_{C_i} = \sigma^* \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}_{C_i}^{\pi}(s, a)] \right|.$$

Until now, we have been dealing with bounds given with regards to the TV divergence of the policies. However, in practice, bounds with respect to the KL divergence of policies is more commonly used [24, 27, 28]. From Pinsker's and Jensen's inequalities, we have that

$$\mathbb{E}_{s \sim d_{\pi}} [D_{TV}(\pi' \parallel \pi)[s]] \leq \sqrt{\mathbb{E}_{s \sim d_{\pi}} [D_{KL}(\pi' \parallel \pi)[s]]/2}. \quad (16)$$

We can thus use Eq. (16) in the bounds of Proposition A.4 and Corollary A.5 to make policy improvement guarantees, i.e., if we find updates such that  $\pi_{k+1} \in \arg \max_{\pi} L_{\pi_k}^{-}(\pi)$ , then we will have monotonically increasing policies as, at iteration  $k$ ,  $\mathbb{E}_{\substack{s \sim d_{\pi_k}, a \sim \pi}} [\bar{A}^{\pi_k}(s, a)] = 0$ ,  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{KL}(\pi \parallel \pi_k)[s]] = 0$  for  $\pi = \pi_k$ , implying that  $J(\pi_{k+1}) - J(\pi_k) \geq 0$ . However, this sequence does not guarantee constraint satisfaction at each iteration, so we now turn to trust region methods in Section 3.1 to incorporate constraints, do policy improvement and provide safety guarantees.

#### A.4. Proofs

**Lemma A.6 (Trivialization of Discounted Criterion Bounds)** *Consider the policy performance bound of [3], which says that for any two stationary policies  $\pi$  and  $\pi'$ :*

$$J_{\gamma}(\pi') - J_{\gamma}(\pi) \geq \frac{1}{1 - \gamma} \left[ \mathbb{E}_{\substack{s \sim d_{\pi, \gamma} \\ a \sim \pi'}} [A_{\gamma}^{\pi}(s, a)] - \frac{2\gamma\epsilon^{\gamma}}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi, \gamma}} [D_{TV}(\pi' \parallel \pi)[s]] \right] \quad (17)$$

where  $\epsilon^{\gamma} = \max_s \left| \mathbb{E}_{a \sim \pi'} [A_{\gamma}^{\pi}(s, a)] \right|$ . Then, the right hand side times  $1 - \gamma$  goes to negative infinity as  $\gamma \rightarrow 1$ .

**Proof** Since  $d_{\pi,\gamma}$  approaches the stationary distribution  $d_\pi$  as  $\gamma \rightarrow 1$ , we can multiply the right hand side of (17) by  $(1 - \gamma)$  and take the limit which gives us:

$$\begin{aligned} & \lim_{\gamma \rightarrow 1} \left( \mathbb{E}_{\substack{s \sim d_{\pi,\gamma} \\ a \sim \pi'}} [A_\gamma^\pi(s, a)] \pm \frac{2\gamma\epsilon^\gamma}{1 - \gamma} \mathbb{E}_{s \sim d_{\pi,\gamma}} D_{\text{TV}}(\pi' \parallel \pi)[s] \right) \\ &= \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] - 2\epsilon \mathbb{E}_{s \sim d_\pi} [D_{\text{TV}}(\pi' \parallel \pi)[s]] \cdot \lim_{\gamma \rightarrow 1} \frac{\gamma}{1 - \gamma} \\ &= -\infty \end{aligned}$$

Here  $\epsilon = \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] \right|$ . The first equality is a standard result of  $\lim_{\gamma \rightarrow 1} A_\gamma^\pi(s, a) = \bar{A}^\pi(s, a)$ . ■

**Lemma A.7** [43] *Under the unichain assumption of the underlying Markov chain, for any stochastic policies  $\pi$  and  $\pi'$ :*

$$J(\pi') - J(\pi) = \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)]. \quad (11)$$

**Proof** We directly expand the right-hand side using the definition of the advantage function and Bellman equation, which gives us:

$$\begin{aligned} \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] &= \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{Q}^\pi(s, a) - \bar{V}^\pi(s)] \\ &= \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi' \\ s' \sim P(\cdot|s,a)}} [r(s, a, s') - J(\pi) + \bar{V}^\pi(s') - \bar{V}^\pi(s)] \\ &= J(\pi') - J(\pi) + \underbrace{\mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi' \\ s' \sim P(\cdot|s,a)}} [\bar{V}^\pi(s')] - \mathbb{E}_{s \sim d_{\pi'}} [\bar{V}^\pi(s)]}_A \end{aligned}$$

Analyzing  $A$ , since  $d_{\pi'}(s)$  is a stationary distribution:

$$\begin{aligned} \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi' \\ s' \sim P(\cdot|s,a)}} [\bar{V}^\pi(s')] &= \sum_s d_{\pi'}(s) \sum_a \pi'(a|s) \sum_{s'} P(s'|s, a) \bar{V}^\pi(s') \\ &= \sum_s d_{\pi'}(s) \sum_{s'} P_{\pi'}(s'|s) \bar{V}^\pi(s') = \sum_{s'} d_{\pi'}(s') \bar{V}^\pi(s') \end{aligned}$$

Therefore,  $A = \sum_{s'} d_{\pi'}(s') \bar{V}^\pi(s') - \mathbb{E}_{s \sim d_{\pi'}} [\bar{V}^\pi(s)] = 0$ . Hence, proved. ■

**Lemma A.8** *Under the unichain assumption, for any stochastic policies  $\pi$  and  $\pi'$  we have:*

$$\left| J(\pi') - J(\pi) - \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \right| \leq 2\epsilon D_{\text{TV}}(d_{\pi'} \parallel d_\pi) \quad \text{where} \quad \epsilon = \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] \right|. \quad (12)$$

**Proof**

$$\begin{aligned}
 \left| J(\pi') - J(\pi) - \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \right| &= \left| \mathbb{E}_{\substack{s \sim d_{\pi'} \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] - \mathbb{E}_{\substack{s \sim d_\pi \\ a \sim \pi'}} [\bar{A}^\pi(s, a)] \right| && \text{(from Lemma A.1)} \\
 &= \left| \sum_s \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] (d_{\pi'}(s) - d_\pi(s)) \right| \\
 &\leq \sum_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] (d_{\pi'}(s) - d_\pi(s)) \right| \\
 &\leq \max_s \left| \mathbb{E}_{a \sim \pi'} [\bar{A}^\pi(s, a)] \right| \|d_{\pi'} - d_\pi\|_1 && \text{(Holder's inequality)} \\
 &= 2\epsilon D_{\text{TV}}(d_{\pi'} \parallel d_\pi)
 \end{aligned}$$

■

**Lemma A.9** *Under the ergodic assumption, the divergence between the stationary distributions  $d_\pi$  and  $d_{\pi'}$  is upper bounded as:*

$$D_{\text{TV}}(d_{\pi'} \parallel d_\pi) \leq \sigma^* \mathbb{E}_{s \sim d_\pi} [D_{\text{TV}}(\pi' \parallel \pi)[s]] \quad \text{where} \quad \sigma^* = \max_\pi \sigma^\pi. \quad (13)$$

**Proof** This proof takes ideas from Markov chain perturbation theory [7, 14, 43]. Firstly we state a standard result with  $P_\pi^* = \mathbf{1}d_\pi^T$

$$(d_{\pi'} - d_\pi)^T (I - P_{\pi'} + P_{\pi'}^*) = d_{\pi'}^T - d_\pi^T - d_{\pi'}^T + d_\pi^T P_{\pi'} = d_\pi^T P_{\pi'} - d_{\pi'}^T = d_\pi^T (P_{\pi'} - P_\pi).$$

Denoting the fundamental matrix of the Markov chain  $Z^{\pi'} = (I - P_{\pi'} + P_{\pi'}^*)^{-1}$  and the mean first passage time matrix  $M^{\pi'} = (I - Z^{\pi'} + E Z_{\text{dg}}^{\pi'}) D^{\pi'}$ , and right multiplying the above by  $(Z^{\pi'})^{-1}$  we have,

$$d_{\pi'}^T - d_\pi^T = d_\pi^T (P_{\pi'} - P_\pi) (I - M^{\pi'} (D^{\pi'})^{-1}) \Rightarrow d_{\pi'} - d_\pi = (I - M^{\pi'} (D^{\pi'})^{-1})^T (P_{\pi'}^T - P_\pi^T) d_\pi \quad (18)$$

$$\text{i.e.} \quad \|d_{\pi'} - d_\pi\|_1 \leq \left\| (I - M^{\pi'} (D^{\pi'})^{-1})^T (P_{\pi'}^T - P_\pi^T) d_\pi \right\|_1$$

(submultiplicative property)

$$\|d_{\pi'} - d_\pi\|_1 \leq \underbrace{\left\| (I - M^{\pi'} (D^{\pi'})^{-1}) \right\|_\infty}_{T_1} \underbrace{\left\| (P_{\pi'}^T - P_\pi^T) d_\pi \right\|_1}_{T_2}$$

(Holder's inequality)

We know that  $\kappa^\pi = \text{Tr}(Z^\pi)$  and from [15], we can write  $T_1$  using the eigenvalues  $\{\lambda_{\pi,i}\}_{i=1}^{|S|}$  of the underlying  $P_\pi$  as

$$T_1 \leq \frac{1}{|S|} \sum_{i=2}^{|S|} \frac{1}{(1 - \lambda_{\pi,i})^{1/2}} \leq \max_i (1 - \lambda_{\pi,i})^{-1/2} = \sigma^\pi \leq \max_\pi \sigma^\pi = \sigma^*.$$

For  $T_2$ , we refer to the result by [43], and provide the proof for completeness below.

$$\begin{aligned}
T_2 &= \sum_{s'} \left| \sum_s \left( \sum_a P(s'|s, a) \pi'(a|s) - P(s'|s, a) \pi(a|s) \right) d_\pi(s) \right| \\
&\leq \sum_{s', s} \left| \sum_a P(s'|s, a) (\pi'(a|s) - \pi(a|s)) \right| d_\pi(s) \\
&\leq \sum_{s, s', a} P(s'|s, a) |\pi'(a|s) - \pi(a|s)| d_\pi(s) \\
&\leq \sum_{s, a} |\pi'(a|s) - \pi(a|s)| d_\pi(s) = 2 \mathbb{E}_{s \sim d_\pi} [D_{\text{TV}}(\pi' \parallel \pi)[s]]
\end{aligned}$$

Combining these inequalities of  $T_1$  and  $T_2$ , we get the desired result.  $\blacksquare$

### A.5. Performance and Constraint Bounds of Trust Region Approach

Consider the trust region formulation in Equation (3). To prove the policy performance bound when the current policy is infeasible (i.e., constraint-violating), we prove the KL divergence between  $\pi_k$  and  $\pi_{k+1}$  for the KL divergence projection, along with other lemmas. We then prove our main theorem for the worst-case performance degradation.

**Lemma A.10** *For a closed convex constraint set, if we have a constraint satisfying policy  $\pi_k$  and the KL divergence  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1/2} \parallel \pi_k)[s]]$  of the ‘Improve’ step is upper bounded by step size  $\delta$ , then after KL divergence projection of the ‘Project’ step we have*

$$\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \parallel \pi_k)[s]] \leq \delta.$$

#### Proof

We make use of the fact that Bregman divergence (hence, KL divergence) projection onto the constraint set ( $\in \mathbb{R}^d$ ,  $d \in \mathbb{N}$ ) exists and is unique. Since  $\pi_k$  is safe, we have  $\pi_k$  in the constraint set, and  $\pi_{k+1}$  is the projection of  $\pi_{k+\frac{1}{2}}$  onto the constraint set. Using the projection inequality, we have

$$\begin{aligned}
\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_k \parallel \pi_{k+1})[s]] + \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \parallel \pi_{k+\frac{1}{2}})[s]] &\leq \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_k \parallel \pi_{k+\frac{1}{2}})[s]] \\
\Rightarrow \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_k \parallel \pi_{k+1})[s]] &\leq \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_k \parallel \pi_{k+\frac{1}{2}})[s]] \leq \delta.
\end{aligned}$$

( $D_{\text{KL}}(\cdot \parallel \cdot) \geq 0$ )

Since KL divergence is asymptotically symmetric when updating the policy within a local neighbourhood ( $\delta \ll 1$ ), we have

$$\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \parallel \pi_k)[s]] \leq \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+\frac{1}{2}} \parallel \pi_k)[s]] \leq \delta.$$

$\blacksquare$

**Lemma A.11**

For a closed convex constraint set, if we have a constraint violating policy  $\pi_k$  and the KL divergence  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1/2} \| \pi_k)[s]]$  of the first step is upper bounded by step size  $\delta$ , then after KL divergence projection of the second step we have

$$\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_k)[s]] \leq \delta + V_{\max},$$

where  $V_{\max} = \max_i \alpha_i \beta_i^2$ ,  $\beta_i = [J_{C_i}(\pi_k) - l_i]_+$ ,  $\alpha_i = \frac{1}{2a_i^T H^{-1} a_i}$ , with  $a_i$  as the gradient of the cost advantage function corresponding to constraint  $C_i$ , and  $H$  as the Hessian of the KL divergence constraint.<sup>1</sup>

**Proof**

Let the sublevel set of cost constraint function for the current infeasible policy  $\pi_k$  be given as:

$$L_{\pi_k} = \{\pi \mid J_{C_i}(\pi) + \mathbb{E}_{\substack{s \sim d_{\pi_k} \\ a \sim \pi}} [\bar{A}_{C_i}^{\pi_k}(s, a)] \leq J_{C_i}(\pi_k) \quad \forall i\}.$$

This implies that the current policy  $\pi_k$  lies in  $L_{\pi_k}$ . The constraint set onto which  $\pi_{k+\frac{1}{2}}$  is projected onto is given by:  $\{\pi \mid J_{C_i}(\pi) + \mathbb{E}_{\substack{s \sim d_{\pi_k} \\ a \sim \pi}} [\bar{A}_{C_i}^{\pi_k}(s, a)] \leq l_i \quad \forall i\}$ . Let  $\pi_{k+1}^L$  be the projection of  $\pi_{k+\frac{1}{2}}$  onto  $L_{\pi_k}$ .

Note that the Bregman inequality of Lemma A.10 holds for any convex set in  $\mathbb{R}^d$ ,  $d \in \mathbb{N}$ . This implies  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1}^L \| \pi_k)[s]] \leq \delta$  since  $\pi_k$  and  $\pi_{k+1}^L$  are both in  $L_{\pi_k}$ , which is also convex since the constraint functions are convex. Using the Three-point Lemma<sup>2</sup>, for policies  $\pi_k$ ,  $\pi_{k+1}$ , and  $\pi_{k+1}^L$ , with  $\varphi(\mathbf{x}) := \sum_i x_i \log x_i$ , we have

$$\begin{aligned} \delta &\geq \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1}^L \| \pi_k)[s]] = \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_k)[s]] \\ &\quad - \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_{k+1}^L)[s]] \\ &\quad + \mathbb{E}_{s \sim d_{\pi_k}} [(\nabla \varphi(\pi_k) - \nabla \varphi(\pi_{k+1}^L))^T (\pi_{k+1} - \pi_{k+1}^L)[s]] \\ \Rightarrow \mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_k)[s]] &\leq \delta + \underbrace{\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_{k+1}^L)[s]]}_{T_1} \\ &\quad - \underbrace{\mathbb{E}_{s \sim d_{\pi_k}} [(\nabla \varphi(\pi_k) - \nabla \varphi(\pi_{k+1}^L))^T (\pi_{k+1} - \pi_{k+1}^L)[s]]}_{T_2}. \end{aligned} \quad (19)$$

If the constraint violations of the current policy  $\pi_k$  are small, i.e.,  $J_{C_i}(\pi_k) - l_i = b_i$  is small for all  $i$ , then  $T_1$  can be approximated by a second order expansion. We analyze  $T_1$  for any constraint  $C_i$  and then bound it over all the constraints. As before we overload the notation with  $\pi_k = \pi_{\theta_k} = \theta_k$  to write. For any constraint  $C_i$ , we can write  $T_1^i$  as the expected KL divergence if projection was onto the constraint set of  $C_i$  i.e.

1. For any  $x \in \mathbb{R}$ ,  $[x]_+ := \max(0, x)$

2. For any  $\phi$ , the Bregman divergence identity:  $D_\phi(x, y) + D_\phi(y, z) = D_\phi(x, z) + \langle \nabla \phi(z) - \nabla \phi(y), x - y \rangle$

$$\begin{aligned}
T_1^i &\approx \frac{1}{2}(\pi_{k+1} - \pi_{k+1}^L)^T H (\pi_{k+1} - \pi_{k+1}^L) = \frac{1}{2} \left( \frac{\beta_i}{a_i^T H^{-1} a_i} H^{-1} a_i \right)^T H \left( \frac{\beta_i}{a_i^T H^{-1} a_i} H^{-1} a_i \right) \\
&= \frac{\beta_i^2}{2a_i^T H^{-1} a_i} = \alpha_i \beta_i^2,
\end{aligned}$$

where the second equality is a result of the trust region guarantee (see [27] for more details). Finally we invoke the projection result from [2] which uses Dykstra's Alternating Projection algorithm [33] to bound this projection, i.e.,  $T_1 \leq \max_i T_1^i \approx \max_i \alpha_i \beta_i^2$ .

And since  $\delta$  is small, we have  $\nabla \varphi(\pi_k) - \nabla \varphi(\pi_{k+1}^L) \approx 0$  given  $s$ . Thus,  $T_2 \approx 0$ . Combining all of the above, we have  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_k)[s]] \leq \delta + V_{\text{max}}$ . ■

**Theorem 3.1** *Let  $\pi_{\theta_{k+1}}$  be the optimal solution to Eq. (3) for some  $\pi_{\theta_k} \in \Pi_{\Theta}$ . Then, we have*

$$J(\pi_{\theta_{k+1}}) - J(\pi_{\theta_k}) \geq -\sqrt{2(\delta + V_{\text{max}})} \nu^{\pi_{\theta_{k+1}}} \quad (4)$$

$$\text{and } J_{C_i}(\pi_{\theta_{k+1}}) \leq l_i + \sqrt{2(\delta + V_{\text{max}})} \nu_{C_i}^{\pi_{\theta_{k+1}}} \quad \forall i, \quad (5)$$

where  $\nu^{\pi_{\theta_{k+1}}} = \sigma^{\pi_{\theta_{k+1}}} \max_s |_{a \sim \pi_{\theta_{k+1}}} \mathbb{E} [\bar{A}^{\pi_{\theta_k}}(s, a)]|$ ,  $\nu_{C_i}^{\pi_{\theta_{k+1}}} = \sigma^{\pi_{\theta_{k+1}}} \max_{i,s} |_{a \sim \pi_{\theta_{k+1}}} \mathbb{E} [\bar{A}_{C_i}^{\pi_{\theta_k}}(s, a)]|$ ,  $V_{\text{max}} = \max_i \beta_i^2$ , and  $\beta_i = [J_{C_i}(\pi_{\theta_k}) - l_i]_+$ .

### Proof

Since  $\bar{D}_{\text{KL}}(\pi_{\theta_k} \| \pi_{\theta_k}) = 0$ ,  $\pi_{\theta_k}$  is feasible. The objective value is 0 for  $\pi_{\theta} = \pi_{\theta_k}$ . The bound follows from Equation (14) and Equation (16) where the average KL i.e.  $\mathbb{E}_{s \sim d_{\pi_k}} [D_{\text{KL}}(\pi_{k+1} \| \pi_k)[s]]$  is bounded by  $\delta + V_{\text{max}}$  from Lemma A.11.

Similar to Corollary A.5, expressions for the auxiliary cost constraints also follow immediately as the second result.

**Remark A.12** *Remark If we look at proof as given by [43] in Section 5 of their paper, with the distinction now that  $\delta$  is replaced by  $\delta + V_{\text{max}}$ , we have the same result. Our worse bound is due to the constrained nature of our setting, which is intuitive in the sense that for the sake of satisfying constraints, we undergo a worse worst-case performance degradation.* ■

## A.6. Approximate ACPO

### A.6.1. POLICY RECOVERY ROUTINE

As described in Section 4.2, we need a recovery routine in case the updated policy  $\pi_{k+1/2}$  is not approximate constraint satisfying. In this case, the optimization problem is inspired from a simple trust region approach by [27]. Since we only deal with one constraint in the practical implementation of ACPO, the recovery rule is obtained by solving the following problem:

$$\begin{aligned}
\min_{\theta} \quad & c + a^T (\theta - \theta_k) \\
\text{s.t.} \quad & \frac{1}{2} (\theta - \theta_k)^T H (\theta - \theta_k) \leq \delta.
\end{aligned}$$

Let  $x = \theta - \theta_k$ , then the dual function  $L(x, \lambda)$  is given by:  $L(x, \lambda) = c + a^T x + \lambda (\frac{1}{2} x^T H x - \delta)$ . Now,

$$\frac{\partial L}{\partial x} = a + \lambda(Hx) = 0 \implies x = -\frac{1}{\lambda} H^{-1} a.$$

$x$  obtained above should satisfy the trust-region constraint:

$$\begin{aligned} \frac{1}{2} \left( -\frac{1}{\lambda} H^{-1} a \right)^T H \left( -\frac{1}{\lambda} H^{-1} a \right) &\leq \delta \\ \implies \frac{1}{2} \cdot \frac{1}{\lambda^2} \cdot a^T H^{-1} a &\leq \delta \\ \implies \sqrt{\frac{a^T H^{-1} a}{2\delta}} &\leq \lambda. \end{aligned}$$

Therefore, the update rule in case of infeasibility takes the form  $\theta = \theta_k - \sqrt{\frac{2\delta}{a^T H^{-1} a}} H^{-1} a$ . We augment this rule with the gradient of the reward advantage function as well, so the final recovery is

$$\theta_{k+1/2} = \theta_k - \sqrt{2\delta} \left[ t \cdot \frac{H^{-1} a}{\sqrt{a^T H^{-1} a}} + (1-t) \cdot \frac{H^{-1} g}{\sqrt{g^T H^{-1} g}} \right] \quad ; \quad t \in [0, 1]$$

#### A.6.2. LINE SEARCH

Because of approximation error, the proposed update may not satisfy the constraints in Eq. (3). Constraint satisfaction is enforced via line search, so the final update is

$$\theta_{k+1} = \theta_k + s^j (\theta_{k+1/2} - \theta_k),$$

where  $s \in (0, 1)$  is the backtracking coefficient and  $j \in \{0, \dots, L\}$  is the smallest integer for which  $\pi_{k+1}$  satisfies the constraints in Equation 3. Here,  $L$  is a finite backtracking budget; if no proposed policy satisfies the constraints after  $L$  backtracking steps, no update occurs.

#### A.7. Practical ACPO

As explained in Section 4, we use the below problem formulation, which uses first-order Taylor approximation on the objective and second-order approximation on the KL constraint<sup>3</sup> around  $\theta_k$ , given small  $\delta$ :

$$\begin{aligned} \max_{\theta} \quad & g^T (\theta - \theta_k) \\ \text{s.t.} \quad & c_i + a_i^T (\theta - \theta_k) \leq 0, \quad \forall i \quad ; \quad \frac{1}{2} (\theta - \theta_k)^T H (\theta - \theta_k) \leq \delta. \end{aligned} \quad (20)$$

where

$$\begin{aligned} g &:= \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi_{\theta_k}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \bar{A}^{\pi_{\theta_k}}(s, a) \right] \quad ; \quad c_i := J_{C_i}(\theta_k) - l_i \quad \forall i \\ a_i &:= \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi_{\theta_k}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \bar{A}_{C_i}^{\pi_{\theta_k}}(s, a) \right] \quad ; \quad H := \mathbb{E}_{\substack{s \sim d_{\pi_{\theta_k}} \\ a \sim \pi_{\theta_k}}} \left[ \nabla_{\theta} \log \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \nabla_{\theta} \log \pi_{\theta}(a|s) \Big|_{\theta=\theta_k}^T \right] \end{aligned}$$

---

3. The gradient and first-order Taylor approximation of  $\bar{D}_{\text{KL}}(\pi_{\theta} \parallel \pi_{\theta_k})$  at  $\theta = \theta_k$  is zero.

Similar to the work of [3],  $g$ ,  $a_i$ , and  $H$  can be approximated using samples drawn from the policy  $\pi_{\theta_k}$ . The Hessian  $H$  is identical to the Hessian  $H$  used by [3] and [43]. However, the definitions  $g$  and  $a_i$ 's are different since they include the average reward advantage functions,  $\bar{A}^{\pi_{\theta_k}}(s, a) = \bar{Q}^{\pi_{\theta_k}}(s, a) - \bar{V}^{\pi_{\theta_k}}(s)$ .

Since rewards and cost advantage functions can be approximated independently, we use the framework of [43] to do so. We describe the process of estimation of rewards advantage function, and the same procedure can be used for the cost advantage functions as well. Specifically, first approximate the average-reward bias  $\bar{V}^{\pi_{\theta_k}}(s)$  and then use a one-step TD backup to estimate the action-bias function. Concretely, using the average reward Bellman equation gives

$$\bar{A}^{\pi_{\theta_k}}(s, a) = r(s, a) - J(\pi_{\theta_k}) + \mathbb{E}_{s' \sim P(\cdot|s,a)} [\bar{V}^{\pi_{\theta_k}}(s')] - \bar{V}^{\pi_{\theta_k}}(s) \quad (21)$$

This expression involves the average-reward bias  $\bar{V}^{\pi_{\theta_k}}(s)$ , which we can approximate using the standard critic network  $\bar{V}_{\phi_k}(s)$ . However, in practice we use the average-reward version of the Generalized Advantage Estimator (GAE) from [28], similar to [43]. Hence, we refer the reader to that paper for detailed explanation, but provide an overview below for completeness.

Let  $\hat{J}_\pi = \frac{1}{N} \sum_{t=1}^N r_t$  denote the estimated average reward. The Monte Carlo target for the average reward value function is  $\bar{V}_t^{\text{target}} = \sum_{t'=t}^N (r_{t'} - \hat{J}_\pi)$  and the bootstrapped target is  $\bar{V}_t^{\text{target}} = r_t - \hat{J}_\pi + \bar{V}_\phi^\pi(s_{t+1})$ .

The Monte Carlo and Bootstrap estimators for the average reward advantage function are:

$$\hat{A}_{\text{MC}}^\pi(s_t, a_t) = \sum_{t'=t}^N (r_{t'} - \hat{J}_\pi) - \bar{V}_\phi^\pi(s_t) \quad ; \quad \hat{A}_{\text{BS}}^\pi(s_t, a_t) = r_{i,t} - \hat{J}_\pi + \bar{V}_\phi^\pi(s_{t+1}) - \bar{V}_\phi^\pi(s_t)$$

We can similarly extend the GAE to the average reward setting:

$$\hat{A}_{\text{GAE}}^\pi(s_t, a_t) = \sum_{t'=t}^N \lambda^{t'-t} \delta_{t'} \quad , \quad \delta_{t'} = r_{t'} - \hat{J}_\pi + \bar{V}_\phi^\pi(s_{t'+1}) - \bar{V}_\phi^\pi(s_{t'}). \quad (22)$$

and set the target for the value function to  $\bar{V}_t^{\text{target}} = r_t - \hat{J}_\pi + \bar{V}_\phi^\pi(s_{t+1}) + \sum_{t'=t+1}^N \lambda^{t'-t} \delta_{t'}$ .

## A.8. Experimental Details

### A.8.1. ENVIRONMENTS

All environments tested on are illustrated in Figure 2, along with a detailed description of each.

### A.8.2. EVALUATION DETAILS AND PROTOCOL

For the Gather and Circle tasks we test two distinct agents: a point-mass ( $S \subseteq \mathbb{R}^9$ ,  $A \subseteq \mathbb{R}^2$ ), and an ant robot ( $S \subseteq \mathbb{R}^{32}$ ,  $A \subseteq \mathbb{R}^8$ ). The agent in the Bottleneck task in  $S \subseteq \mathbb{R}^{71}$ ,  $A \subseteq \mathbb{R}^{16}$ , and for the Grid task is  $S \subseteq \mathbb{R}^{96}$ ,  $A \subseteq \mathbb{R}^4$ . We use two hidden layer neural networks to represent Gaussian policies for the tasks. For Gather and Circle, size is (64,32) for both layers, and for Grid and Bottleneck the layer sizes are (16,16) and (50,25). We set the step size  $\delta$  to  $10^{-4}$ , and for each task, we conduct 5 runs to get the mean and standard deviation for reward objective and cost constraint values during training. We train CPO, PCPO, and PPO with the discounted objective, however, evaluation and comparison with BVF-PPO, ATRPO and ACPO<sup>4</sup> is done using the average reward objective (this is a standard evaluation scheme [27, 37, 40]).

For each environment, we train an agent for  $10^5$  steps, and for every  $10^3$  steps, we instantiate 10 evaluation trajectories with the current (deterministic) policy. For each of these trajectories, we calculate

4. Code of the ACPO implementation will be made available on GitHub.

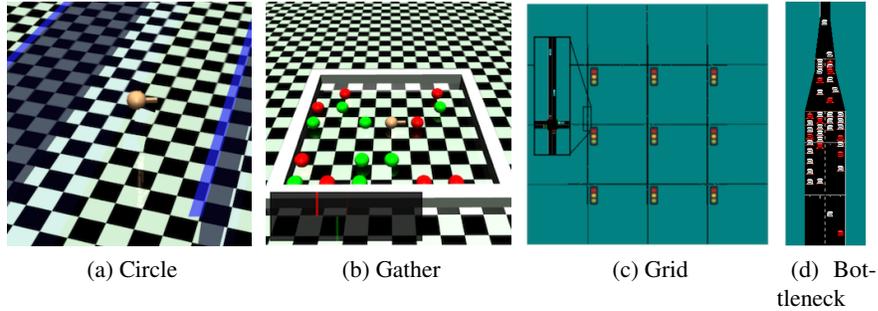


Figure 2: The Circle, Gather, Grid, and Bottleneck tasks. (a) Circle: The agent is rewarded for moving in a specified circle but is penalized if the diameter of the circle is larger than some value [3]. (b) Gather: The agent is rewarded for collecting the green balls while penalized to gather red balls [3]. (c) Grid: The agent controls traffic lights in a 3x3 road network and is rewarded for high traffic throughput but is constrained to let lights be red for at most 5 consecutive seconds [35]. (d) Bottleneck: The agent controls vehicles (red) in a merging traffic situation and is rewarded for maximizing the number of vehicles that pass through but is constrained to ensure that white vehicles (not controlled by agent) have “low” speed for no more than 10 seconds. [35].

the trajectory average reward for the next  $10^3$  steps and finally report the total average-reward as the mean of these 10 trajectories.

For detailed explanation of Point-Circle, Point-Gather, Ant-Circle, and Ant-Gather tasks, please refer to [3]. For detailed explanation of Bottleneck and Grid tasks, please refer to [35]. For the simulations in the Gather and Circle tasks, we use neural network baselines with the same architecture and activation functions as the policy networks. For the simulations in the Grid and Bottleneck tasks, we use linear baselines. For all experiments we use Gaussian neural policies whose outputs are the mean vectors and variances are separate parameters to be learned. Seeds used for generating evaluation trajectories are different from those used for training.

For comparison of different algorithms, we make use of CPO, PCPO, ATRPO, and PPO implementations taken from <https://github.com/rll/rllab> and <https://github.com/openai/safety-starter-agents>. Even the hyperparameters are selected so as to showcase the best performance of other algorithms for fair comparison. The choice of the hyperparameters given below is inspired by the original papers since we wanted to understand the performance of the average reward case.

We use settings which are common in all open-source implementations of the algorithms, such as normalizing the states by the running mean and standard deviation before being fed into the neural network and similarly normalizing the advantage values (for both rewards and constraints) by their batch means and standard deviations before being used for policy updates. Table 1 summarizes the hyperparameters below.

For the Lagrangian formulation of ATRPO and PPO, note that the original papers do not provide any blueprint for formulating the Lagrangian, and even CPO and PCPO use *unconstrained* TRPO for benchmarking. However, we feel that this is unfair to these algorithms as they can possibly perform better with a Lagrangian formulation in an average-reward CMDP setting. To this extent, we introduced a Lagrangian parameter  $\ell \in [0, 1]$  that balances the rewards and constraints in the final objective function. More specifically, Equation (6) for a single constraint now becomes

$$\max_{\theta} (1 - \ell)g^T(\theta - \theta_k) - \ell[(c_1 + a_1^T(\theta - \theta_k)) + (\frac{1}{2}(\theta - \theta_k)^T H(\theta - \theta_k) - \delta)]. \quad (23)$$

Table 1: Hyperparameter Setup

Hyperparameter	PPO/ATRPO	CPO/PCPO/ACPO
No. of hidden layers	2	2
Activation	tanh	tanh
Initial log std	-0.5	-1
Batch size	2500	2500
GAE parameter (reward)	0.95	0.95
GAE parameter (cost)	N/A	0.95
Trust region step size $\delta$	$10^{-4}$	$10^{-4}$
Learning rate for policy	$2 \times 10^{-4}$	$2 \times 10^{-4}$
Learning rate for reward critic net	$2 \times 10^{-4}$	$2 \times 10^{-4}$
Learning rate for cost critic net	N/A	$2 \times 10^{-4}$
Backtracking coeff.	0.75	0.75
Max backtracking iterations	10	10
Max conjugate gradient iterations	10	10
Recovery regime parameter $t$	N/A	0.75

**Note.** The authors of the ATRPO and PPO do not suggest any principled approach for finding an optimal  $\ell$ . Hence, the choice of the Lagrangian parameter  $\ell$  is completely empirical and is selected such that these algorithms achieve maximum rewards while satisfying the constraints. Also see in Figure 1, for Ant-Gather, Bottleneck, and Grid environments, where the constraints cannot be satisfied for *any* value of  $\ell$ , we include the results for a specific value of  $\ell$  for illustrative purposes, as detailed in Table 2.

Table 2: Lagrangian parameter  $\ell$  for ATRPO and PPO

Algorithm	Point-Gather	Ant-Circle	Ant-Gather	Bottleneck	Grid
ATRPO	0.50	0.60	0.45	0.50	0.45
PPO	0.55	0.50	0.50	0.50	0.60

### A.8.3. ACPO PERFORMANCE COMMENTARY

**ACPO vs. CPO/PCPO.** For the Point-Gather environment (see Figure 3), we see that initially ACPO and CPO/PCPO give relatively similar performance, but eventually ACPO improves over CPO and PCPO by 52.5% and 36.1% on average-rewards respectively. This superior performance does not come with more constraint violation. The Ant-Gather environment particularly brings out the effectiveness of ACPO where it shows 41.1% and 61.5% improvement over CPO and PCPO respectively, while satisfying the constraint. In the high dimensional Bottleneck and Grid environments, ACPO is particularly quick at optimizing for low constraint violations, while improving over PCPO and CPO in terms of average-reward.

**ACPO vs Lagrangian ATRPO/PPO.** One could suppose to use the state of the art unconstrained policy optimization algorithms with a Lagrangian formulation to solve the average-rewards CMDP problem in consideration, but we see that such an approach, although principled in theory, does not give satisfactory empirical results. This can be particularly seen in the Ant-Circle, Ant-Gather, Bottleneck, and Grid

environments, where Lagrangian ATRPO and PPO give the least rewards, while not even satisfying the constraints. If ATRPO and PPO were used without the Lagrangian (i.e. constraints are ignored), one would see higher rewards but even worse constraint violations, which are not useful for solving the average-reward CMDP problem. Hence, we do not include those results.

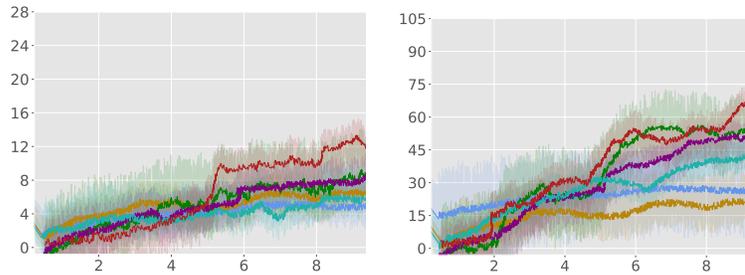
**ACPO vs BVF-PPO.** BVF-PPO is a whole different formulation than the other baselines, as it translates the cumulative cost constraints into state-based constraints, which results in an almost-safe policy improvement method which maximizes returns at every step. However, we see that this approach fails to satisfy the constraints even in the moderately difficult Ant Gather environment, let alone the high dimensional Bottleneck and Grid environments. On the other hand, ACPO performs the best among all baselines in these three environments.

## A.9. Experimental Addendum

### A.9.1. LEARNING CURVES

Due to space restrictions, we present the learning curves for the remaining environments in Figure 3.

Average Rewards:



Average Constraint values:

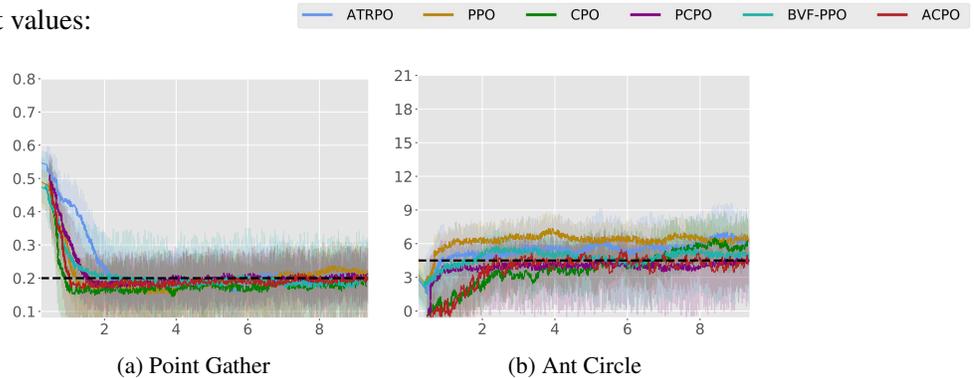


Figure 3: The average reward and constraint cost function values vs iterations (in  $10^4$ ) learning curves for some algorithm-task pairs. Solid lines in each figure are the empirical means, while the shaded area represents 1 standard deviation, all over 5 runs. The dashed line in constraint plots is the constraint threshold  $l$ . ATRPO and PPO are tested with constraints, which are included in their Lagrangian formulation.

## A.9.2. RECOVERY REGIME REVISITED

Figure 4 shows the experiments we conducted with varying  $t \in \{0, 0.25, 0.5, 0.75, 1\}$ . With  $t = 1$ , we obtain the same recovery scheme as that of [3]. Our results show that this scheme does not lead to the best performance, and that  $t = 0.75$  and  $t = 1$  perform the best across all tasks. Figure 5 shows the performance of ACPO with different values of  $t$  in various environments.

Contrasting with the policy recovery update of [3] which only uses the cost advantage function gradient  $a$ , we introduce the reward advantage function gradient  $g$  as well. This choice is to ensure recovery while simultaneously balancing the “regret” of not choosing the best (in terms of the objective value) policy  $\pi_k$ . In other words, we wish to find a policy  $\pi_{k+1/2}$  as close to  $\pi_k$  in terms of their objective function values. We follow up this step with a simple linesearch to find feasible  $\pi_{k+1}$ .

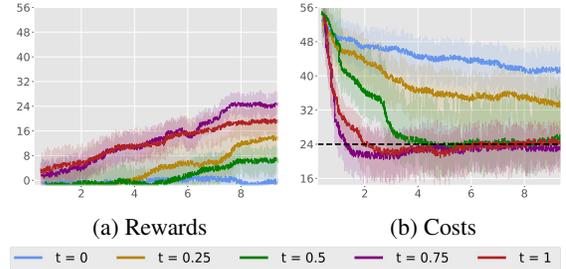


Figure 4: Comparison of performance of ACPO with different values of the hyperparameter  $t$  in the Point-Circle environment. X-axis is iterations in  $10^4$ . See Appendix A.9 for more details.

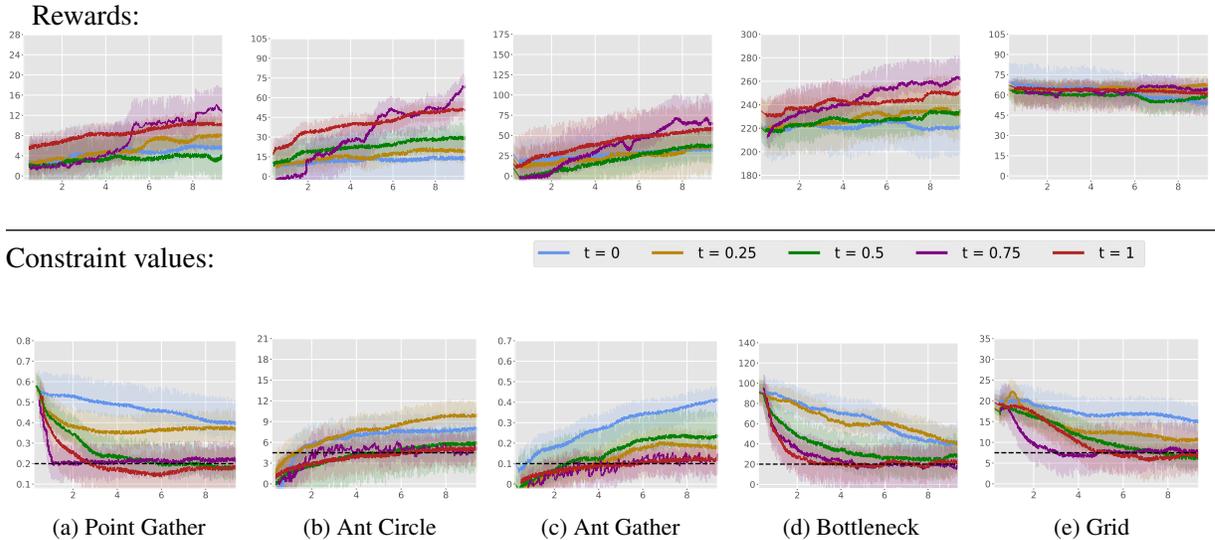


Figure 5: Comparison of performance of ACPO with different values of the hyperparameter  $t$  in various environment. X-axis is iterations in  $10^4$ .