

The optimization problem can be formulated as

$$\min_{\substack{\lambda_0 \\ \lambda_k \geq 0, 1 \leq k \leq K}} \left[ \max_{\pi} \mathbf{E} \sum_{t=0}^{T-1} \gamma^t \hat{\mathcal{R}}_{\lambda_0, \dots, \lambda_K}(s_t, \tilde{a}_t) + \sum_{k=1}^K \lambda_k \gamma^T \delta_k \right], \quad (0.1)$$

where the penalized reward function

$$\hat{\mathcal{R}}_{\lambda_0, \dots, \lambda_K}(s_t, \tilde{a}_t) = \mathcal{R}(s_t, a_t) - \lambda_0 \sum_{j=1}^J \mathcal{H}_j(s_t) - \sum_{k=1}^K \lambda_k \mathbf{1}(\mathcal{C}_{k,t}(\tilde{a}_t) = 1). \quad (0.2)$$

The parameters we need to choose before running:

- (i) *hard-constraints* (step-wise) :  $h_j, 1 \leq j \leq J$
- (ii) *soft-constraint* (episode-wise) :  $c_k, 1 \leq k \leq K$
- (iii) *soft-constraint-tolerance* (optional):  $t_k, 1 \leq k \leq K$
- (iv) *soft-violation-risk-threshold* :  $r_k, 1 \leq k \leq K$
- (v) *hard-lambda-learning-rate* :  $\alpha$
- (vi) *soft-lambda-learning-rate* :  $\beta_k, 1 \leq k \leq K$
- (vii) *number of policy exploitation* (outer loop) :  $n_p$
- (viii) *number of updates on lambda* (inner loop)  $n_\lambda$
- (ix) *exploitation-steps* :  $s_e$
- (x) *lambda-update-steps* :  $s_\lambda$
- (xi) *last-training-steps* :  $s_l$
- (xii) *lambda-sample-size* :  $m_\lambda$

---

**Algorithm 1:** PPO algorithm on updating policy with fixed  $\lambda = (\lambda_0, \dots, \lambda_K)$

---

Input:  $h_j, 1 \leq j \leq J; c_k, t_k, r_k, 1 \leq k \leq K$ ; number of steps  $s$ ; policy parameter  $\theta$ ;  
value function parameter  $\phi$ .

**for**  $k = 0, 1, \dots, s - 1$  **do**

    Generate a random environment (task) parameter  $\sigma$ .

    Collect set of episode(s)  $D$  by sampling using policy  $\pi_\theta$  in environment  $\mathcal{E}_\sigma$ .

    Compute penalized rewards-to-go  $\tilde{\mathcal{R}}_t$  w.r.t.  $\lambda_k, 0 \leq k \leq K$ .

    Compute advantage estimates  $A_t$  based on the current value function  $V_\phi$ .

    Update the policy parameter by maximizing the PPO-Clip objective:

$$\theta \leftarrow \arg \max_{\theta'} \frac{1}{|D|T} \sum_{\tau \in D} \sum_{t=0}^{T-1} \min \left( \frac{\pi_{\theta'}(a_t|s_t)}{\pi_\theta(a_t|s_t)} A_t, g(\varepsilon, A_t) \right), \quad (0.3)$$

        using MLP.

    Update the value function parameter  $\phi$  by minimizing the mean square error:

$$\phi \leftarrow \arg \min_{\phi'} \frac{1}{|D|T} \sum_{\tau \in D} \sum_{t=0}^{T-1} \min \left( V_{\phi'}(s_t) - \tilde{\mathcal{R}}_t \right)^2, \quad (0.4)$$

        using MLP.

**end for**

---

---

**Algorithm 2:** Algorithm on updating on  $\lambda = (\lambda_0, \dots, \lambda_K)$

---

Input: Numbers of steps, sample size, and constraints.

Initialize  $\lambda = 0$ , policy parameter  $\theta$ , and value function parameter  $\phi$ .

**for**  $i = 0, 1, \dots, n_p - 1$  **do**

    Apply Algorithm 1 with  $\theta, \phi, \lambda$ , and  $s = s_e$ .

**for**  $n = 0, 1, \dots, n_\lambda - 1$  **do**

        Apply Algorithm 1 with  $\theta, \phi, \lambda$ , and  $s = s_\lambda$ .

**for**  $m = 0, 1, \dots, m_\lambda - 1$  **do**

            Generate a random environment (task) parameter  $\sigma$ .

            Collect an episode by sampling using policy  $\pi_\theta$  in environment  $\mathcal{E}_\sigma$ .

            Compute the number of hard violations  $H^m$  and the indicator  $C_k^m$  of soft violations,  $1 \leq k \leq K$ .

**end for**

        Compute the gradient w.r.t.  $\lambda_k$ ,  $0 \leq k \leq K$ :

$$\nabla_\lambda \leftarrow -\frac{1}{m_\lambda} \sum_{m=0}^{m_\lambda-1} (H^m, C_1^m - \delta_1, \dots, C_K^m - \delta_K). \quad (0.5)$$

$\lambda \leftarrow \lambda - \text{diag}(\alpha, \beta_1, \dots, \beta_K) \nabla_\lambda$ .

**end for**

**end for**

Apply Algorithm 1 with  $\theta, \phi, \lambda$ , and  $s = s_l$ .

---