**Algorithm 1** Priority Experience Replay Constrained Stackelberg Q-Learning with MIP action selection (MIP-PCSQ)

---

1: **Input:** thresholds $d_1, d_2$, discount $\gamma$, soft-update rate $\rho$, exploration $\varepsilon$; PER hyper-parameters $\alpha, \beta_{\text{start}}, \varepsilon_p$; updates per step $K$; buffer capacity $N$

2: **Init:** critics $\phi_i$ (for $Q_i$), cost-critics $\zeta_i$ (for $G_i$), $i \in \{1, 2\}$; targets $\phi_i^{\text{targ}} \leftarrow \phi_i$, $\zeta_i^{\text{targ}} \leftarrow \zeta_i$; replay buffer $D$ with PER; set $\beta \leftarrow \beta_{\text{start}}$

3: **for** $t = 1, 2, \ldots$ **do**

4:     Observe state $s$

5:     Build tables (online): $Q_1^{ij}(s), Q_2^{ij}(s), G_1^{ij}(s), G_2^{ij}(s)$; define safe sets $\mathcal{S}_i(s) = \{j \mid G_1^{ij}(s) \leq d_1, \ G_2^{ij}(s) \leq d_2\}$

6:     **if** rand() $< \varepsilon$ **then**

7:         pick $(a_1, a_2)$ uniformly from $\bigcup_i \{(a_1^i, a_2^j) \mid j \in \mathcal{S}_i(s)\}$

8:     **else**

9:         Solve MIP (1) at $s$ (use (13a)–(13h)) to get $(a_1, a_2)$

10:     **end if**

11:     Execute $(a_1, a_2)$; observe $r_i, c_i, s', d$

12:     Push $(s, a_1, a_2, r_i, c_i, s', d)$ into $D$ with priority $p_{\max}$

13:     **for** $k = 1$ **to** $K$ **do**                     ▷ $K$ critic updates per env-step

14:         $(\text{idx}, B, P) \leftarrow \text{PER.Sample}(D, |B|, \text{stratified} = \text{true})$

15:         $w \leftarrow \left(\frac{1/N}{P}\right)^{\beta}$;   $\tilde{w} \leftarrow w / \max(w)$

16:         Build $Q_1^{ij}(s'), Q_2^{ij}(s'), G_1^{ij}(s'), G_2^{ij}(s')$ (online); define $\mathcal{S}_i(s')$

17:         Solve MIP (1) at $s'$ to get $(a_1', a_2')$

18:         Targets (use target nets): $y_i = r_i + \gamma(1 - d)Q_i^{\text{targ}}(s', a_1', a_2')$,    $g_i = c_i + \gamma(1 - d)G_i^{\text{targ}}(s', a_1', a_2')$

19:         Residuals: $dQ_i \leftarrow Q_i(s, a_1, a_2) - y_i$,   $dG_i \leftarrow G_i(s, a_1, a_2) - g_i$

20:         Losses (weighted mean over $B$): $\mathcal{L}_{Q_i} = \frac{1}{|B|} \sum_{u \in B} \tilde{w}(u) [dQ_i(u)]^2$,   $\mathcal{L}_{G_i} = \frac{1}{|B|} \sum_{u \in B} \tilde{w}(u) [dG_i(u)]^2$

21:         GD steps: $\phi_i \leftarrow \phi_i - \eta_Q \nabla_{\phi_i} \mathcal{L}_{Q_i}$,   $\zeta_i \leftarrow \zeta_i - \eta_G \nabla_{\zeta_i} \mathcal{L}_{G_i}$

22:         Priority update (PER): $\Delta \leftarrow \lambda_{Q1}|dQ_1| + \lambda_{Q2}|dQ_2| + \lambda_{G1}|dG_1| + \lambda_{G2}|dG_2|$; $p_{\text{new}} \leftarrow (\Delta + \varepsilon_p)^{\alpha}$;  PER.Update(idx, $p_{\text{new}}$)

23:         Soft-update targets: $\phi_i^{\text{targ}} \leftarrow \rho \, \phi_i^{\text{targ}} + (1 - \rho)\phi_i$;  $\zeta_i^{\text{targ}} \leftarrow \rho \, \zeta_i^{\text{targ}} + (1 - \rho)\zeta_i$

24:     **end for**

25:     **if** $d = 1$ **then** reset environment

26:     **end if**

27: **end for**

---

**MIP (13a–13h) at state x** ($x \in \{s, s'\}$):

$$\mathcal{S}_i(x) = \{j \mid G_1^{ij}(x) \leq d_1, \ G_2^{ij}(x) \leq d_2\};$$

$$\max_{x,y,v} \quad \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} Q_1^{ij}(x)\, y_{ij} \tag{13a}$$

$$\text{s.t.} \quad \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} y_{ij} = 1 \tag{13b}$$

$$\sum_{j=1}^{m_2} y_{ij} = x_i, \ \forall i \tag{13c}$$

$$x_i = 1 \Rightarrow v_i \geq Q_2^{i\ell}(x), \ \forall i, \ \forall \ell \in \mathcal{S}_i(x) \tag{13d}$$

$$y_{ij} = 1 \Rightarrow v_i \geq Q_2^{ij}(x), \ \forall i, j \tag{13e}$$

$$y_{ij} = 1 \Rightarrow v_i \leq Q_2^{ij}(x), \ \forall i, j \tag{13f}$$

$$y_{ij} = 1 \Rightarrow G_1^{ij}(x) \leq d_1, \ \forall i, j \tag{13g}$$

$$y_{ij} = 1 \Rightarrow G_2^{ij}(x) \leq d_2, \ \forall i, j \tag{13h}$$