**Algorithm 1:** Optimistic Actor-Critic (OAC)

---

**Input:** $w_1, w_2, \theta$      `// Initial parameters of the critics and`
         `target policy` $\pi_{\mathcal{T}}$

**Initialization:** $\tilde{w}_1 \leftarrow w_1, \tilde{w}_2 \leftarrow w_2, \mathcal{D} \leftarrow \emptyset$      `// Initialize target`
         `networks and replay pool`

**1** **for** *each iteration* **do**

**2**    **for** *each environment step* **do**

**3**      Sample action $a_t \sim \pi_E(a_t|s_t)$     `// Exploration policy from`
       `Eq. (9)`;

**4**      Sample transition $s_{t+1} \sim p(s_{t+1}|s_t, a_t)$;

**5**      Store transition: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, R(s_t, a_t), s_{t+1})\}$;

**6**    **end**

**7**    **for** *each training step* **do**

**8**      **for** $i \in \{1, 2\}$ **do**

**9**        Update $w_i$ with:

$$\nabla_{w_i} \left\| \hat{Q}_{\mathrm{LB}}^i(s_t, a_t) - R(s_t, a_t) - \gamma \min(\hat{Q}_{\mathrm{LB}}^1(s_{t+1}, a), \hat{Q}_{\mathrm{LB}}^2(s_{t+1}, a)) \right\|_2^2$$

**10**      **end**

**11**      Update $\theta$ with:

$$\nabla_\theta \hat{J}_{\hat{Q}_{\mathrm{LB}}}^\alpha$$

**12**      Update target critics:

$$\tilde{w}_1 \leftarrow \tau w_1 + (1 - \tau)\tilde{w}_1, \quad \tilde{w}_2 \leftarrow \tau w_2 + (1 - \tau)\tilde{w}_2$$

**13**    **end**

**14** **end**

**Output:** $w_1, w_2, \theta$          `// Optimized parameters`

---