**Algorithm 1:** Randomized Ensembled Double Q-Learning (REDQ)

---

**Input:** Initial policy parameters $\theta$, ensemble of $M$ Q-functions $\{\phi_i\}_{i=1}^{M}$, target Q-function parameters $\{\bar{\phi}_i \leftarrow \phi_i\}_{i=1}^{M}$, polyak coefficient $\tau$, number of critics sampled $N < M$, batch size $B$, number of critic updates per iteration $G$, replay buffer $\mathcal{D}$

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

**2**      Reset environment and observe initial state $s_0$;

**3**      **while** *not terminal* **do**

**4**          Select action $a_t \sim \pi_\theta(\cdot|s_t)$;

**5**          Execute $a_t$ in environment;

**6**          Observe reward $r_t$, next state $s_{t+1}$, done signal $d_t$;

**7**          Store $(s_t, a_t, r_t, s_{t+1}, d_t)$ in $\mathcal{D}$;

**8**          **if** $d_t$ *is True* **then**

**9**              Reset environment and observe new initial state $s_{t+1}$;

**10**          **end**

**11**          $s_t \leftarrow s_{t+1}$;

**12**      **end**

**13**      **for** $g = 1$ **to** $G$ **do**

**14**          Sample a batch of $B$ transitions $(s, a, r, s', d)$ from $\mathcal{D}$;

**15**          Sample a random subset $\mathcal{I} \subset \{1, \ldots, M\}$ of size $N$;

**16**          Compute target Q-values:

$$y = r + \gamma(1-d) \left( \min_{i \in \mathcal{I}} Q_{\bar{\phi}_i}(s', a') - \alpha \log \pi_{\bar{\theta}}(a'|s') \right), \quad a' \sim \pi_{\bar{\theta}}(\cdot|s')$$

**17**          **for** $i = 1$ **to** $M$ **do**

**18**              Update critic $i$ parameters by gradient descent:

$$\phi_i \leftarrow \phi_i - \lambda_Q \nabla_{\phi_i} \frac{1}{B} \sum \left( Q_{\phi_i}(s, a) - y \right)^2$$

**19**              Update target critic network $i$:

$$\bar{\phi}_i \leftarrow \tau \phi_i + (1 - \tau) \bar{\phi}_i$$

**20**          **end**

**21**      **end**

**22**      Update actor parameters by gradient ascent:

$$\theta \leftarrow \theta + \lambda_\pi \nabla_\theta \frac{1}{B} \sum \left( \frac{1}{M} \sum_{i=1}^{M} Q_{\phi_i}(s, \pi_\theta(s)) - \alpha \log \pi_\theta(\pi_\theta(s)|s) \right)$$

**23** **end**

---