

---

**Algorithm 1: PAC-Bayesian Actor-Critic (PBAC)**

---

**Input:** Polyak coefficient  $\tau \in (0, 1)$ , batch size  $n$ , bootstrap rate  $\kappa$ , posterior sampling rate  $\text{PSR}$ , prior variance  $\sigma_0^2$ , ensemble size  $K$   
**Initialization:** Replay buffer  $\mathcal{D} \leftarrow \emptyset$ ; critic parameters  $\{\theta_k\}_{k=1}^K$  and targets  $\bar{\theta}_k \leftarrow \theta_k$ ; actor trunk  $g$  and heads  $h_1, \dots, h_K$

1 Initialize state  $s \leftarrow \text{env.reset}()$ , interaction counter  $e \leftarrow 0$ ;  
2 **while** *training* **do**  
3     **if**  $e \bmod \text{PSR} = 0$  **then**  
4         Sample index  $j \sim \mathcal{U}(\{1, \dots, K\})$ ;  
5         Set active policy:  $\pi \leftarrow \pi_{g \circ h_j}$ ;  
6     **end**  
7     Sample action  $a \sim \pi(s)$ ;  
8     Execute  $a$  in environment, observe reward  $r$ , next state  $s'$ ;  
9     Store transition:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s, a, r, s')\}$ ;  
10     $e \leftarrow e + 1$ ;  
11    Sample batch  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^n \sim \mathcal{D}$ ;  
12    Sample bootstrap mask:  $b_{ik} \sim \text{Bernoulli}(1 - \kappa)$  for all  $i \in [n], k \in [K]$ ;  
13    **foreach**  $(s_i, a_i, r_i, s'_i)$  *in batch* **do**  
14         
$$\bar{\mu}_\pi(s'_i) \leftarrow \frac{1}{K} \sum_{k=1}^K b_{ik} \cdot \bar{X}_k(s'_i, \pi(s'_i))$$
$$\mu_\pi(s_i) \leftarrow \frac{1}{K} \sum_{k=1}^K b_{ik} \cdot X_k(s_i, \pi(s_i))$$
$$\sigma_\pi^2(s_i) \leftarrow \frac{1}{K-1} \sum_{k=1}^K b_{ik} \cdot (X_k(s_i, \pi(s_i)) - \mu_\pi(s_i))^2$$
  
15     **end**  
16     Update critic parameters  $\{\theta_k\}$  by minimizing:  
$$\frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K b_{ik} (r_i + \gamma \bar{X}_k(s'_i, \pi(s'_i)) - X_k(s_i, \pi(s_i)))^2$$
$$+ \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \frac{b_{ik} (r_i + \gamma \bar{\mu}_\pi(s'_i) - X_k(s_i, \pi(s_i)))^2}{2\gamma^2 \sigma_0^2} - \frac{\gamma^2 + 1/2}{n} \sum_{i=1}^n \log \sigma_\pi^2(s_i)$$
  
17     Update actor:  
$$(g, h_1, \dots, h_K) \leftarrow \arg \max_{g, h_1, \dots, h_K} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K X_k(s_i, \pi_{g \circ h_k}(s_i))$$
  
18     Update target critics:  
$$\bar{\theta}_k \leftarrow \tau \theta_k + (1 - \tau) \bar{\theta}_k, \quad \text{for all } k \in [K]$$
  
19     **if** *episode ended* **then**  
20          $s \leftarrow \text{env.reset}()$ ;  
21     **else**  
22          $s \leftarrow s'$ ;  
23     **end**  
24 **end**

---