# PRM with off-policy RL method

Saito Karuha

November 2024

## 1    Pseudo Code

---
**Algorithm 1** DDPG For PRM Traning

---
1: **Initialize:** CriticNet: $Q_\omega(s, a)$ , ActionNet: $\mu_\theta(s)$
2: **Initialize:** $Q_{\omega^-}(s, a) \leftarrow Q_\omega(s, a)$ , $\mu_{\theta^-}(s) \leftarrow \mu_\theta(s)$
3: **Buffer Initialize:** $B \leftarrow \emptyset$
4: **for** $e = 0$ **to** $E$ **do**
5:    Initialize initial state $s_1$ (Randomly pick a question from UCB-Math)
6:    **for** $t = 1$ **to** $T$ **and** $done == True$ **do**
7:       Choose an action $a_t \sim \mu_\theta(s_t)$
8:       $r_t \leftarrow Env(s_t, a_t)$ ; $s_{t+1} \leftarrow [s_t, a_t]$
9:       $B \leftarrow (s_t, a_t, r_t, s_{t+1})$
10:      **if** Buffer is big enough **then**
11:         Randomly pick $N$ touples $\{(s_i, a_i, r_i, s_{i+1})\}_{i=1,...,N}$
12:         Sample $K$ actions: $a_{i+1}^m \sim \mu_{\theta^-}(s_{i+1})$ , $(m = 1, ..., k)$
13:         Calculate for every tuples:

$$y_i = r_i + \gamma * \max_{m \in \{1,...,k\}} Q_{\omega^-}(s_{i+1}, a_{i+1}^m) \qquad (1)$$

$$A_i = r_i + \gamma * Random_j Q_{\omega^-}(s_{i+1}, a_{i+1}^j) - Q_\omega(s_i, a_i) \qquad (2)$$

14:         Compute loss for CriticNet($L$) and ActorNet($J$) respectively:
15:         $L = \frac{1}{N} \sum_{i=1}^N (y_i - Q_\omega(s_i, a_i))^2$
16:         $J = \frac{1}{N} \sum_{i=1}^N min(\frac{\pi_\theta(a_i|s_i)}{\pi_{\theta^-}(a_i|s_i)} A_i, clip(\frac{\pi_\theta(a_i|s_i)}{\pi_{\theta^-}(a_i|s_i)}, 1 - \epsilon, 1 + \epsilon) A_i)$
17:         **Update for Critic and Actor Network**
18:         **Soft update target Network**
19:         $\omega^- \leftarrow \tau\omega + (1 - \tau)\omega$ , $\theta^- \leftarrow \tau\theta + (1 - \tau)\theta$
20:      **end if**
21:   **end for**
22: **end for**

---

# 2 Problems