

➤ A general version of RL policy update forms

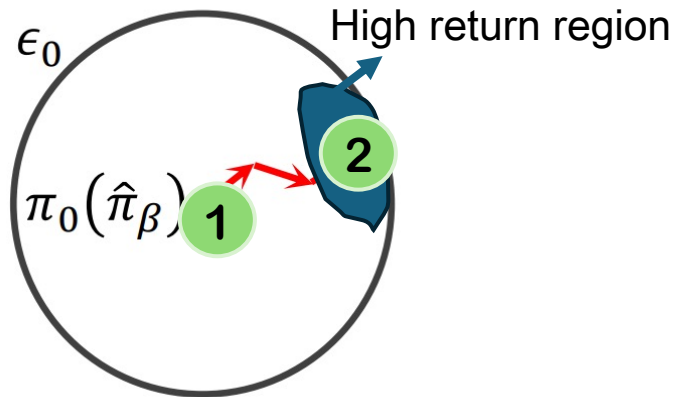
Policy evaluation

$$\hat{Q}^{k+1} = \arg \min_Q \mathbb{E}_{(s,a) \sim w} \left[ \left( r + \gamma \underbrace{\mathbb{E}_{s'|s,a} \mathbb{E}_{a' \sim \pi^k(\cdot|s')} [T^{\pi^k} \hat{Q}]}_{\mathbb{E}_{a' \sim \pi^k(\cdot|s')} [\hat{Q}(s', a')]} - Q(s, a) \right)^2 \right] \dots \textcircled{1}$$

Policy improvement

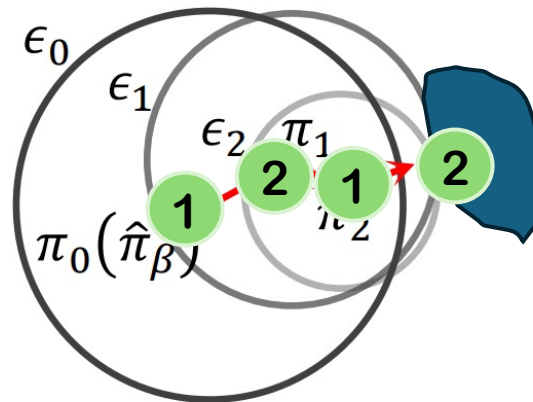
$$\pi^{k+1} = \arg \max_{\pi} \mathbb{E}_{s \sim \rho, a \sim \pi} [\hat{Q}^{k+1}(s, a)] - \beta \mathbb{E}_{s \sim \rho} [D_{\text{KL}}(\pi \parallel \pi_{\text{ref}})] \dots \textcircled{2}$$

➤ One-step RL



- One-step RL/IQL/IDQL...

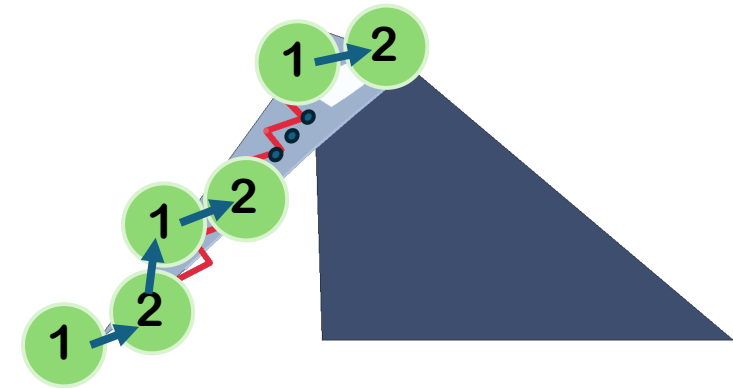
➤ Multi-step RL



- BPPO → Uni-O4/RL-100...

More conservative, but stable

➤ Iterative RL



- PPO/SAC/...

More conservative, but stable