# Trust Region-Guided Proximal Policy Optimization

Yuhui Wang, Hao He , Xiaoyang Tan, Yaozhong Gan

College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, China

## Problem

- **Reinforcement Learning**
  Find a policy $\pi$ which could maximize the accumulative reward
  $$\eta(\pi) = \mathbb{E}\left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t)|\pi\right]$$
  .

- **Policy Gradient**
  $$L_{\pi_{old}}(\pi) = \mathbb{E}_{s \sim \rho^{\pi_{old}}, a \sim \pi_{old}}\left[\frac{\pi(a|s)}{\pi_{old}(a|s)}A_{s,a}^{\pi_{old}}\right] + \eta(\pi_{old}),$$
  where $A_{s,a}^{\pi} = \mathbb{E}[R_t^{\gamma}|s_t = s, a_t = a; \pi] - \mathbb{E}[R_t^{\gamma}|s_t = s; \pi]$ is the advantage function value and $R_t^{\gamma} = \Sigma_{k=0}^{\infty} \gamma^k c(s_{t+k}, a_{t+k})$;
  $\rho^{\pi}(s) = (1-\gamma)\Sigma_{t=1}^{\infty} \gamma^{t-1}\rho_t^{\pi}(s)$, $\rho_t^{\pi}$ is the visitation probability of $\pi$ at $t$.

- **Trust Region Policy Optimization (TRPO)**
  $$\max_{\pi} L_{\pi_{old}}(\pi) \text{ subject to } \max_{s \in \mathcal{S}} D_{KL}^s(\pi_{old}, \pi) \leq \delta$$
  Maximizing $L_{\pi_{old}}(\pi)$ within the trust region guarantee non-decreasing of the performance.
  **Theorem 1.** Let $M_{\pi_{old}}(\pi) = L_{\pi_{old}}(\pi) - C \max_{s \in \mathcal{S}} D_{KL}^s(\pi_{old}, \pi)$,
  $$\eta(\pi) \geq M_{\pi_{old}}(\pi), \eta(\pi_{old}) = M_{\pi_{old}}(\pi_{old}).$$

- **Proximal Policy Optimization (PPO)**
  $$L_{\pi_{old}}^{CLIP}(\pi) = \mathbb{E}\left[\min\left(\frac{\pi(a|s)}{\pi_{old}(a|s)}A_{s,a}^{\pi_{old}}, clip\left(\frac{\pi(a|s)}{\pi_{old}(a|s)}, l_{s,a}, u_{s,a}\right)A_{s,a}^{\pi_{old}}\right)\right]$$
  PPO attempts to enforce restriction on the policy
  $$l_{s,a} \leq \frac{\pi(a|s)}{\pi_{old}(a|s)} \leq u_{s,a}$$
  $$\Downarrow$$
  $$-\pi_{old}(a|s)(1 - l_{s,a}) \leq \pi(a|s) - \pi_{old}(a|s) \leq \pi_{old}(a|s)(u_{s,a} - 1)$$
  **clipping range**: $(l_{s,a}, u_{s,a})$;
  **feasible range**: $(-\pi_{old}(a|s)(1 - l_{s,a}), \pi_{old}(a|s)(u_{s,a} - 1))$.

## Motivation

**Original PPO:** **clipping range**: $(1 - \epsilon, 1 + \epsilon)$;
**feasible range**: $(-\pi_{old}(a|s)\epsilon, \pi_{old}(a|s)\epsilon)$.

Consider an optimal action $a_{opt}$ and a sub-optimal one $a_{subopt}$. If $\pi_{old}(a_{opt}|s) < \pi_{old}(a_{subopt}|s)$, then $|(-\pi_{old}(a_{opt}|s)\epsilon, \pi_{old}(a_{opt}|s)\epsilon)| < |(-\pi_{old}(a_{subopt}|s)\epsilon, \pi_{old}(a_{subopt}|s)\epsilon)|$. This means that the feasible range of $\pi(a_{opt}|s)$ is limited than that of $\pi(a_{subopt}|s)$.

Note that $\pi(a_{opt}|s)$ and $\pi(a_{subopt}|s)$ are in a zero-sum competition. Such unequal preference may continuously weaken the likelihood of the optimal action and make the policy trapped in local optima.

## Contributions

- We propose an enhanced PPO, TRGPPO, to improve the exploration ability while not harming the learning stability.
- Theoretically prove that PPO is prone to suffer from the risk of lack of exploration and TRGPPO has better exploration ability.
- Extensive experiments demonstrate the effectiveness of TRGPPO on benchmark tasks.

## Method

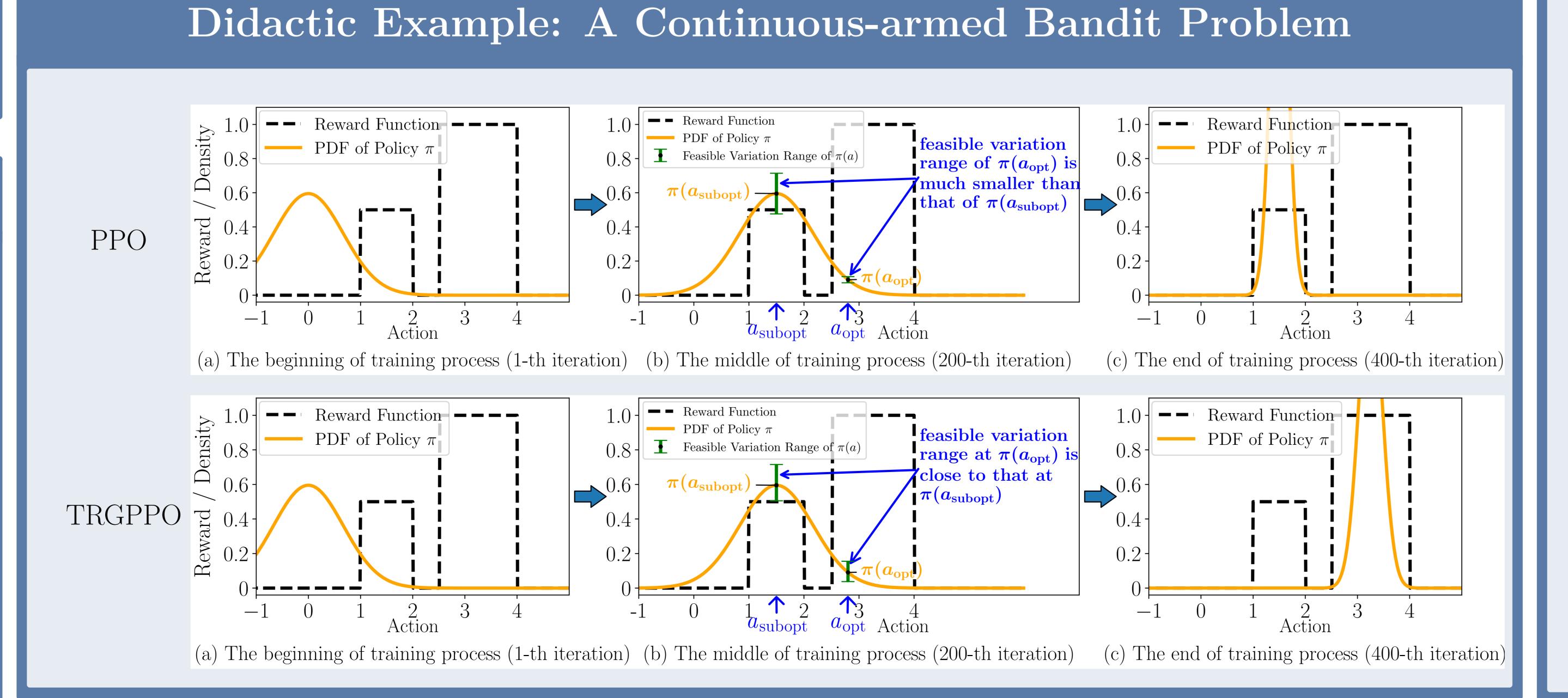- **Improve exploration:** relax the restrictions on actions which are not preferred by $\pi_{old}$.
- **Harmless to the learning stability:** make relaxation guided by trust region-based criterion.

**The original PPO method:**
$$l_{s,a} = 1 - \epsilon, \text{ for any } (s, a)$$
$$u_{s,a} = 1 + \epsilon, \text{ for any } (s, a)$$

**Our TRGPPO method:**
$$l_{s,a}^{\delta} = \min_{\pi}\left\{\frac{\pi(a|s)}{\pi_{old}(a|s)} : D_{KL}^s(\pi_{old}, \pi) \leq \delta\right\}$$
$$u_{s,a}^{\delta} = \max_{\pi}\left\{\frac{\pi(a|s)}{\pi_{old}(a|s)} : D_{KL}^s(\pi_{old}, \pi) \leq \delta\right\}$$



(a) clipping range (discrete action space)    (b) feasible range (discrete action space)    (c) clipping range (continuous action space)

## Didactic Example: A Continuous-armed Bandit Problem



(a) The beginning of training process (1-th iteration)    (b) The middle of training process (200-th iteration)    (c) The end of training process (400-th iteration)

## Theoretical Analysis

1: Initialize a policy $\pi_0$, $t \leftarrow 0$.
2: **repeat**
3:    Sample an action $\hat{a}_t \sim \pi_t$.
4:    Get the new policy $\pi_{t+1}$ by optimizing the surrogate objective function of PPO based on $\hat{a}_t$:
$$\hat{\pi}_{t+1}(a) = \begin{cases} \pi_t(a)u_a & a = \hat{a}_t \text{ and } c(a) > 0 \\ \pi_t(a)l_a & a = \hat{a}_t \text{ and } c(a) < 0 \\ \pi_t(a) - \frac{\pi_t(\hat{a}_t)u_{\hat{a}_t} - \pi_t(\hat{a}_t)}{|\mathcal{A}| - 1} & a \neq \hat{a}_t \text{ and } c(\hat{a}_t) > 0 \\ \pi_t(a) + \frac{\pi_t(\hat{a}_t)(1 - l_{\hat{a}_t})}{|\mathcal{A}| - 1} & a \neq \hat{a}_t \text{ and } c(\hat{a}_t) < 0 \\ \pi_t(a) & c(\hat{a}_t) = 0 \end{cases}$$
5:    $\pi_{t+1} = Normalize(\hat{\pi}_{t+1})$. $t \leftarrow t + 1$.
6: **until** $\pi_t$ converge

**Algorithm 1:** Simplified PPO for bandit problem

**Notation**:
$a_{opt} = argmax_a c(a)$: the optimal action;
$a_{subopt} \in \{a \in \mathcal{A}|c(a) > 0, a \neq a_{opt}\}$: a sub-optimal action;
$\pi^*$: the optimal policy, where $\pi^*(a_{opt}) = 1$, $\pi^*(a) = 0$ for $a \neq a_{opt}$;
$\Delta_{\pi_0,t} \triangleq \mathbb{E}_{\pi_t}[\|\pi_t - \pi^*\|_{\infty}|\pi_0]$: the expected distance between $\pi_t$ and $\pi^*$;

- **Convergence of PPO**
  If the policy is initialized sufficiently far from the optimal one, PPO is expected to diverge from the optimal policy.
  **Theorem 2.**
  If $\pi_0^2(a_{opt}) \cdot |\mathcal{A}| < \Sigma_{a_{subopt} \in \mathcal{A}_{subopt}} \pi_0^2(a_{subopt}) - \Sigma_{a^- \in \mathcal{A}} \pi_0^2(a^-)$, then $\Delta_{\pi_0,0} < \Delta_{\pi_0,1}^{PPO} < \cdots < \Delta_{\pi_0,t}^{PPO}$.

- **Characteristic of TRGPPO Compared to PPO**
  **1** **Better exploration:**
    larger clipping range for less explored actions;
    closer to the optimal policy.
  **2** **More sample efficient:** larger clipping range.
  **3** **Do not harm the learning stability:** not enlarging the policy divergence of the new policy with the old policy.
  **4** **Better performance guarantee:** larger empirical lower performance bound.
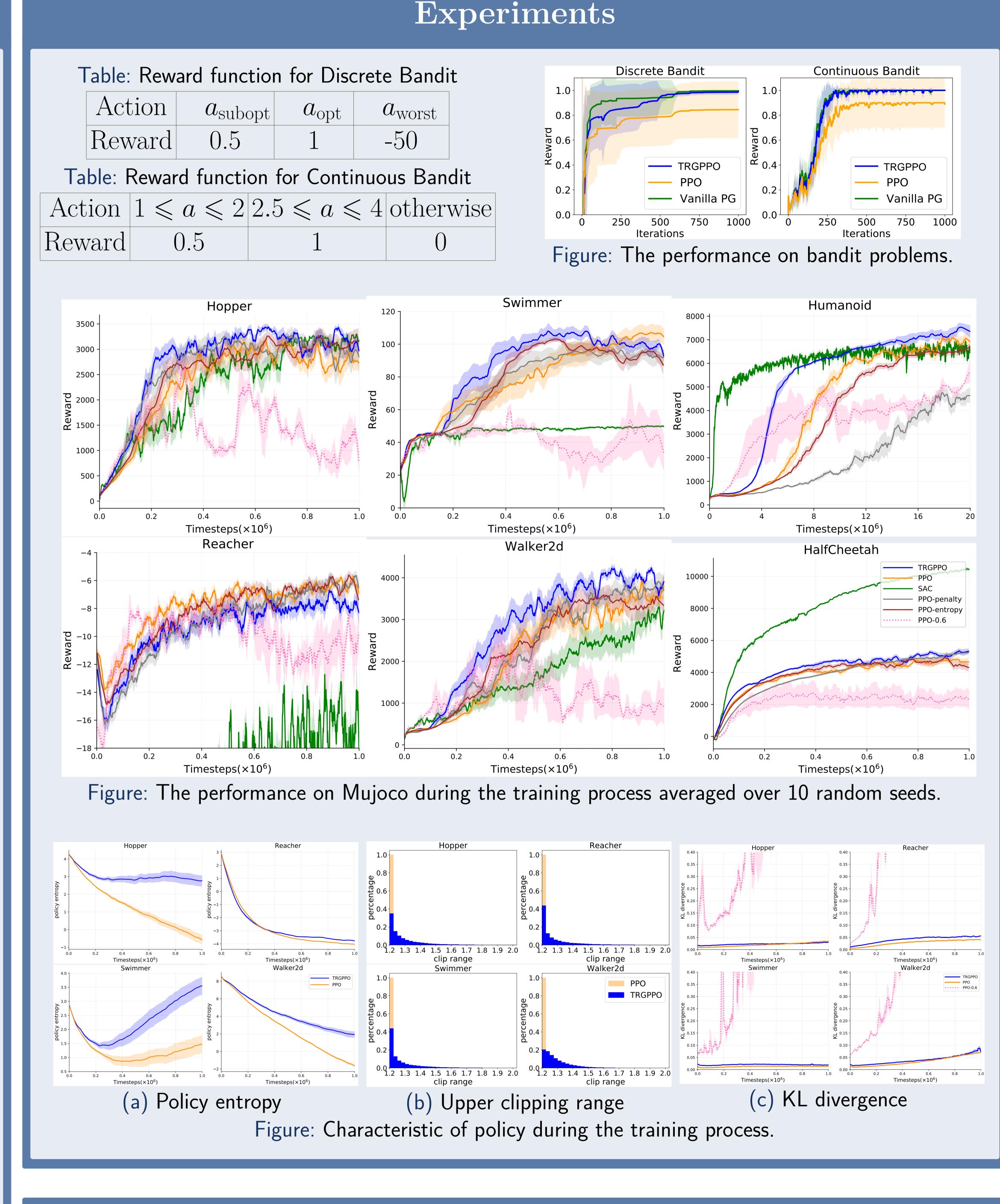  **Lemma 3.** $\frac{du_{s,a}^{\delta}}{d\pi_{old}(a|s)} < 0$, $\frac{dl_{s,a}^{\delta}}{d\pi_{old}(a|s)} > 0$.
  **Theorem 3.** If $\delta \leq g(\max_{a \in \mathcal{A}_{subopt}} \pi_t(a), 1 + \epsilon)$ for all $t$, then $\Delta_{\pi_0,t}^{TRGPPO} \leq \Delta_{\pi_0,t}^{PPO}$ for any $t$.
  **Theorem 4.** If TRGPPO-$\epsilon$ and PPO have the same hyperparameter $\epsilon$, then:
  (i) $u_{s_t,a_t}^{\delta} \geq 1 + \epsilon$ and $l_{s_t,a_t}^{\delta} \leq 1 - \epsilon$ for all $(s_t, a_t)$;
  (ii) $\max_t D_{KL}^{s_t}(\pi_{old}, \pi_{new}^{TRGPPO}) = \max_t D_{KL}^{s_t}(\pi_{old}, \pi_{new}^{PPO})$;
  (iii) $\hat{M}_{\pi_{old}}(\pi_{new}^{TRGPPO}) \geq \hat{M}_{\pi_{old}}(\pi_{new}^{PPO})$.

## Experiments

Table: Reward function for Discrete Bandit

| Action | $a_{subopt}$ | $a_{opt}$ | $a_{worst}$ |
|--------|--------------|-----------|-------------|
| Reward | 0.5 | 1 | -50 |

Table: Reward function for Continuous Bandit

| Action | $1 \leq a \leq 2$ | $2.5 \leq a \leq 4$ | otherwise |
|--------|-------------------|---------------------|-----------|
| Reward | 0.5 | 1 | 0 |



Figure: The performance on bandit problems.



Figure: The performance on Mujoco during the training process averaged over 10 random seeds.



(a) Policy entropy    (b) Upper clipping range    (c) KL divergence

Figure: Characteristic of policy during the training process.

## References

J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In International Conference on Machine Learning, pages 1889-1897,2015.
J. Schulman, F. Wolski, P. Dhariwal, A Radford, and O Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
R. S. Sutton and A. G. Barto. Reinforcement learning: An introduction. MIT press, 2018.

## Contact Information

Yuhui Wang **Email:** y.wang@nuaa.edu.cn **Github:** https://github.com/wangyuhuix
Hao He    **Email:** hugo@nuaa.edu.cn    **Github:** https://github.com/RebornHugo