

MAXIMUM NEXT-STATE ENTROPY FOR EFFICIENT REINFORCEMENT LEARNING

Anonymous authors

Paper under double-blind review

ABSTRACT

Maximum entropy algorithms have demonstrated significant progress in Reinforcement Learning (RL), which offers additional guidance in the form of entropy, which is particularly beneficial in tasks with sparse rewards. Nevertheless, current approaches grounded in policy entropy encourage the agent to explore diverse actions, yet they do not directly help the agent explore diverse states. In this study, we theoretically reveal the challenge of optimizing the next-state entropy of the agent. To address this limitation, we introduce **Maximum Next-State Entropy (MNSE)**, a novel method that maximizes next-state entropy through an action mapping layer following the inner policy. We provide a theoretical analysis demonstrating that MNSE can maximize next-state entropy by optimizing the action entropy of the inner policy. We conduct extensive experiments on various continuous control tasks and show that MNSE can significantly improve the exploration capability of RL algorithms.

1 INTRODUCTION

Maximum entropy RL algorithms have demonstrated remarkable performance across various domains, including games (Gao et al., 2018), robotic control (Haarnoja et al., 2018a;b), and autonomous navigation (Sun et al., 2022). The maximum entropy framework enhances policy exploration and robustness by optimizing both reward and policy entropy simultaneously (Ziebart et al., 2008; 2010). Recent advancements adapt the temperature dynamically (Hu et al., 2021) to improve the trade-off between reward and entropy. Meanwhile, Deep Soft Policy Gradient (DSPG) Shi et al. (2019) integrates soft policy gradients with the soft Bellman equation to address stability issues in off-policy learning.

However, maximizing policy entropy may not directly help the agent explore diverse states due to redundancy in the action space. [For example, in real-world tasks, actuators often exhibit saturation and dead zone issues due to design redundancies or equipment aging \(Galuppini et al., 2018; Bai, 2002\). This leads to a redundant action space, where multiple actions can result in the same state.](#)

If we cannot deal with this issue, the nonlinear actuators can reduce control accuracy, thereby significantly degrading the control performance.

To help the agent directly explore diverse states, current researchers propose novelty-based exploration methods, such as random network distillation (Burda et al., 2018) and pseudo-count techniques (Lobel et al., 2023; Machado et al., 2020). These methods aim to drive agents toward discovering more diverse states by directly evaluating the novelty of a state and incorporating it as an exploration bonus into the extrinsic reward. Current novelty-based and pseudo-count exploration methods have achieved considerable success across various domains. However, it is unclear whether state diversity can provably benefit maximum entropy RL and how to bridge the gap under theoretical guarantees, which naturally leads to the following question:

How can we maximize the next-state entropy of the RL agent in a principled way?

To answer this question, we start with the analysis of next-state entropy. Intuitively, if the policy is non-redundant, we can optimize the next-state entropy by optimizing policy entropy. However, the condition of non-redundant policy is hard to meet in the stochastic transition case. To solve this issue, we first introduce the inner policy and reversible action mapping layer. Then, we rigorously derive the gap between the next-state entropy and the policy entropy of the inner policy. We

demonstrate that by optimizing both the action mapping layer and inverse dynamics model, we can maximize the next-state entropy by optimizing the inner policy entropy. In practice, our method also achieves superior performance. We conduct extensive experiments on various environments, including MuJoCo (Todorov et al., 2012) and Meta-World (Yu et al., 2020). Results show that our algorithm performs better than baselines.

2 RELATED WORK

Maximum Entropy Reinforcement Learning has been widely adopted for improving policy exploration and robustness in RL. Early work introduced Soft Actor-Critic (SAC) (Haarnoja et al., 2018c;d), an off-policy actor-critic algorithm that formalizes MaxEnt RL by balancing the goals of maximizing expected return and policy entropy. Deep Soft Policy Gradient (DSPG) by (Shi et al., 2019) integrates soft policy gradients with the soft Bellman equation to address stability issues in off-policy learning. Count-Based Soft Q-Learning (CBSQL) by (Hu et al., 2021) adapts the temperature dynamically to improve the trade-off between reward and entropy. Additionally, Han & Sung (2021) propose a max-min entropy framework to improve exploration in model-free learning by promoting low-entropy state visitation.

Action Representation. Extensive efforts have been made to effectively represent actions within large action spaces. Zahavy et al. (2018b) propose a method that directly identifies redundant or irrelevant actions using external elimination signals provided by the environment, removing them from the sampling process in text-based games. Tennenholtz & Mannor (2019) adopt a negative sampling procedure, leveraging expert demonstrations to better understand the action space. However, valuable prior information is often scarce and expensive, limiting the scalability of these approaches. Chandak et al. (2019) demonstrates how to learn and utilize action representations without relying on prior knowledge by embedding them within the policy structure to train agents effectively. Similarly, Metz et al. (2017) introduces a novel approach to discretize high-dimensional continuous action spaces by sequentially combining one-dimensional discrete actions.

Exploration is a cornerstone of reinforcement learning, with various strategies enhancing agents' ability to learn from complex environments. For example, the Go-Explore strategy (Ecoffet et al., 2019) advocates for a phased approach to overcoming challenging exploration dilemmas. Count-based methods (Bellemare et al., 2016) capitalize on environmental novelty by employing pseudo-counts. Disagreement-based exploration (Pathak et al., 2019) harnesses the variance in model predictions to propel the agent toward exploration. Curiosity-driven exploration mechanisms, such as ICM (Pathak et al., 2017), utilize prediction errors to incentivize exploration. RND (Burda et al., 2018) employs a novel neural network to generate intrinsic rewards based on the prediction error of environmental dynamics, driving the agent towards unexplored territories. NGU (Badia et al., 2020) integrates intrinsic motivation with an episodic memory mechanism to encourage the revisitation of novel states, promoting long-term exploration.

State Entropy Maximization aims to learn a reward-free policy in which state visitations are uniformly distributed across the state space, thus promoting robust policy initialization and efficient adaptation. Additionally, when task rewards are available, incorporating state entropy as an intrinsic reward has proven to be an effective approach for enhancing exploration. Lee et al. (2019) propose optimizing the state marginal distribution to align with a target distribution, effectively enhancing exploration. Building on this idea, Islam et al. (2019) introduce entropy regularization based on the marginal state distribution, achieving superior state space coverage in complex domains. Further advancements include the work of Guo et al. (2021), who incorporate geometry-aware Shannon entropy of state visitations in both discrete and continuous domains, framing exploration as a computationally tractable problem. Additionally, Hazan et al. (2019) provide a provably efficient algorithm for state entropy maximization, leveraging a black-box planning oracle. Expanding on these methods, Liu & Abbeel (2021) maximize a particle-based entropy in an abstract representation space, demonstrating human-level performance in navigating complex environments.

3 BACKGROUND

Reinforcement Learning. We consider the Markov Decision Processes (MDPs) as the model process, defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \gamma)$, where \mathcal{S} is a state space, \mathcal{A} is an action space, $\gamma \in [0, 1)$

is the discount factor and $\mathcal{P} : \mathcal{S} \times \mathcal{A} \rightarrow \text{Dist}(\mathcal{S})$, $r : \mathcal{S} \rightarrow [r_{\min}, r_{\max}]$ are the transition function and reward function, respectively. We assume a fixed distribution μ_0 as the initial state distribution. The goal of an RL agent is to learn a policy $\pi(a | s)$ under dataset \mathcal{D} , which maximizes the expectation of a discounted cumulative reward: $J(\pi) = \mathbb{E}_{\mu_0, \pi} [\sum_{t=0}^{\infty} \gamma^t r_t]$. For any policy π , the corresponding state-action value function is $Q^\pi(s, a) = \mathbb{E}[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | S_t = s, A_t = a, \pi]$.

Maximum Policy Entropy. Different with standard reinforcement learning, maximum policy entropy reinforcement learning aims to augment the objective with the expected entropy of the policy:

$$J(\pi) = \mathbb{E}_{\mu_0, \pi} \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right], \quad (1)$$

where α is a temperature parameter, determining the relative importance of the entropy term against the reward, and thus controls the stochasticity of the optimal policy.

4 ANALYSIS OF NEXT-STATE ENTROPY

How do we maximize the state entropy of the agent? Current approaches have been to encourage exploration by adding bonus rewards related to the new state (Burda et al., 2018; Badia et al., 2020; Zhang et al., 2021). Adding an exploration bonus has achieved considerable success, while theoretical analysis of the state entropy has not been explored well. To bridge this gap, in this section, we theoretically analyze the optimization objective of the next-state entropy. Firstly, we define the next-state entropy as $\mathcal{H}(S_{t+1} | S_t = s, \pi)$, which represents the entropy of the next state after executing policy π in state s .

Definition 4.1 (Next-State Entropy). *We define the next-state entropy under policy π following state s by*

$$\mathcal{H}(S_{t+1} | S_t = s, \pi) = -\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim P(\cdot | s, a)} \log [P^\pi(s' | s)] \quad (2)$$

where $P^\pi(s' | s) = \int_{a \in \mathcal{A}} \pi(a | s) P(s' | s, a)$. The **next-state entropy**, as defined above, measures the diversity of the subsequent states under the policy π . In classical entropy-regularized reinforcement learning, the **policy entropy** is often used to encourage diversity in the actions taken. However, these two concepts are not generally equivalent.

In this section, we will illustrate the discrepancy between policy entropy and next-state entropy during policy updates in an illustrative example (Section 4.1) and then reveal their relationship under both deterministic (Section 4.2) and stochastic transitions (Section 4.3) with theoretical analysis.

4.1 TOY EXAMPLE

The maximum entropy RL framework is often credited with improving exploration efficiency and promoting more diverse state visitation, particularly in sparse reward settings. However, through the following illustrative example, we demonstrate that optimizing policy entropy alone can be inefficient in certain cases, especially when there is redundancy in the action space.

Consider a one-step MDP with a deterministic transition, where $s_{t+1} = \max(a_t, 0)$, which means actions less than zero are redundant. We assume an initial policy π_{policy} ¹ is a Gaussian policy with a mean less than zero: $\pi_{policy}(a) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(a-\mu)^2}{2\sigma^2}\right)$, $\mu < 0$. The corresponding policy entropy is given by:

$$\mathcal{H}(\pi_{policy}(\cdot | s_t)) = \log(\sqrt{2\pi e}\sigma), \quad (3)$$

which is independent of the mean μ . As shown in the Figure 1 (Left), updating the policy by maximizing policy entropy alone leads to increased variance σ . However, since the mean of the Gaussian distribution remains unchanged, there is still more than a 50% probability (the gray area) that actions sampled from the updated policy will be less than zero, leading to the same next state.

¹In Section 4.1, to avoid confusion between the policy π and the mathematical constant π , we denote the policy by π_{policy} .

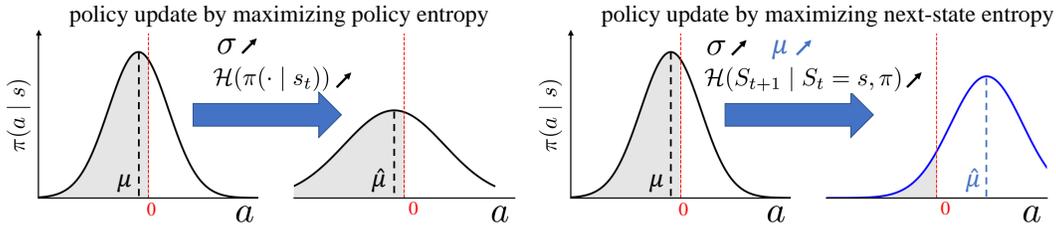


Figure 1: Probability distribution of policy π . For a one-step MDP with a deterministic transition $s_{t+1} = \max(a_t, 0)$, the gray area shows the probability of actions leading to the same next state. **Left:** Policy update under maximum policy entropy by increasing variance σ . **Right:** Policy update under maximum next-state entropy by increasing both variance σ and mean μ .

In contrast, the next-state entropy can be mathematically expressed as:

$$\mathcal{H}(S_{t+1} | S_t = s, \pi_{policy}) = \underbrace{\log(\sqrt{2\pi e}\sigma)}_{\text{policy entropy}} + \left\{ \Phi\left(\frac{\mu}{\sigma}\right) \log \sigma - \int_{-\infty}^{-\frac{\mu}{\sigma}} -\varphi(z) \log \varphi(z) dz \right\}, \quad (4)$$

where $\varphi(z) = \frac{e^{-z^2/2}}{\sqrt{2\pi}}$ is the probability density function of the standard normal distribution, and $\Phi(x) = \int_{-\infty}^x \varphi(z) dz$ is its cumulative distribution function. The Equation 4 highlights that next-state entropy is influenced not only by the policy entropy but also by an additional term associated with μ . Notably, when we assume $\sigma > 1$ in the toy example, the variable μ positively correlates with next-state entropy. A detailed proof and numerical analysis are provided in Appendix A. As illustrated in the Figure 1 (Right), policy updates driven by next-state entropy increase both the variance σ and the mean μ , significantly reducing the probability of sampling actions less than zero, which decreases redundancy in the next states and enhance state diversity.

4.2 DETERMINISTIC CASE

Section 4.1 demonstrates the superiority of next-state entropy over policy entropy in terms of exploration. To better understand the relationship between next-state entropy and policy entropy, we start with an MDP with deterministic transitions. As discussed by Baram et al. (2021), the following corollary provides insights into the equivalent relationship between next-state entropy and policy entropy under deterministic case:

Corollary 4.2. *Let M be a deterministic MDP with a transition function $T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$. If $T(s, a) \neq T(s, a'), \forall s \in \mathcal{S}, \forall a, a' \in \mathcal{A}$, then*

$$\mathcal{H}(S_{t+1} | S_t = s, \pi) = \mathcal{H}(\pi(\cdot | s_t)) \quad , \quad \forall s \in \mathcal{S}$$

This corollary indicates that, in the deterministic case, if there are no redundant actions in the action space \mathcal{A} , next-state entropy and policy entropy are equivalent. Therefore, recent studies (Zahavy et al., 2018a; Tennenholtz & Mannor, 2019; Zhong et al., 2024) have focused on eliminating redundant actions in the action space \mathcal{A} and significantly enhance the exploration and performance in various domains.

4.3 STOCHASTIC CASE

Similarly, in the stochastic case, we can define a non-redundant policy and reveal the relationship between next-state entropy and policy entropy. Firstly, we define the non-redundant policy as follows:

Definition 4.3 (Non-Redundant Policy). *Given a stochastic MDP with a transition dynamics P , π is a non-redundant policy if*

$$\forall s, s' \in \mathcal{S}, \quad \forall a_i, a_j \in \{a \in \mathcal{A} | \pi(a | s) > 0\}, \quad P(s' | s, a_i) * P(s' | s, a_j) = 0.$$

In the stochastic case, we define π as a non-redundant policy only if there is an absolutely non-intersection between probability distribution $P(s' | s, a)$ of possible actions sampling from π . In other words, for any reachable next state s' , only one action can lead to it under a non-redundant policy. Under the above strong assumption on policy π , we can reveal the relationship between next-state entropy and policy entropy in the stochastic case:

Theorem 4.4. *If π is a non-redundant policy, then*

$$\mathcal{H}(S_{t+1} | S_t = s, \pi) = \mathcal{H}(\pi(\cdot | s)) + \mathcal{H}_{\text{model}}, \quad (5)$$

where $\mathcal{H}_{\text{model}} = \mathbb{E}_{a \sim \pi(\cdot | s)} \mathcal{H}(S_{t+1} | S_t = s, a)$ is the entropy of the dynamics model.

Proof. Please refer to Appendix B for the detailed proof. \square

The entropy of the dynamics model $\mathcal{H}_{\text{model}}$ represents an inherent property of the system. Suppose we treat this term as a constant c . In that case, the above equation reveals that in the case of stochastic transitions, as long as the policy is non-redundant, next-state entropy is equivalent to policy entropy.

5 METHOD

Based on the analysis in Section 4, we can directly optimize state entropy by maximizing the entropy of the non-redundant policy. However, constructing a non-redundant policy in the stochastic transition setting is extremely challenging since we must strictly satisfy the Definition 4.3. To solve this issue, we decompose the overall policy into an inner policy π_i and a parameterized, reversible action mapping layer f . Based on this framework, we rigorously derive the gap between the next-state entropy of the overall policy and the policy entropy of the inner policy. Meanwhile, we demonstrate that by optimizing both the action mapping layer and an inverse dynamics model to minimize the gap term, the inner policy’s entropy will be equivalent to the next-state entropy of the overall policy.

5.1 ACTION MAPPING

We consider a new action space \mathcal{E} , where action $a \in \mathcal{A}$ is a function of the inner action $e \in \mathcal{E}$:

$$a = f(e; \theta), \quad \forall e \in \mathcal{E}, \quad (6)$$

where f is the action mapping parameterized by θ and it has an invertible function:

$$e = f^{-1}(a; \theta), \quad \forall a \in \mathcal{A}. \quad (7)$$

Subsequently, we define policy based on the action mapping as follows: for a given state s , the inner policy $\pi_i(\cdot | s)$ outputs the action e in new action space \mathcal{E} and then f transforms e back to the original action space \mathcal{A} . For the policy π , we have the following conclusion:

$$\pi(a | s) = \pi_i(e | s) * \left| \frac{\partial f^{-1}(a; \theta)}{\partial a} \right| = \pi_i(e | s) * \left| \frac{\partial f(e; \theta)}{\partial e} \right|^{-1}, \quad (8)$$

where $\left| \frac{\partial f(e; \theta)}{\partial e} \right|$ is the Jacobian determinant of the function f . We denote the entropy of the internal policy as $\mathcal{H}(\pi_i(\cdot | s))$ and denote the entropy of the state as $\mathcal{H}(S_{t+1} | S_t = s, \pi)$. Following the above policy framework and notations, we can derive the entropy of the state as follows:

Theorem 5.1. *For any inner policy $e \sim \pi_i(\cdot | s)$ and invertible action mapping layer $a = f(e; \theta)$, we have*

$$\mathcal{H}(S_{t+1} | S_t = s, \pi) = \mathcal{H}(\pi_i(\cdot | s)) + \underbrace{\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, a)} [\log [p_{\text{inv}}(e | s, s')]]}_{\text{Gap Term}} + \mathcal{H}_{\text{model}}$$

where $\mathcal{H}_{\text{model}} = \mathbb{E}_{e \sim \pi_i(\cdot | s)} [\mathcal{H}(S_{t+1} | S_t = s, e)]$ is the entropy of the dynamics model, $p_{\text{inv}}(e | s, s')$ is the inverse dynamic of inner policy, *which is a function that predicts the inner action e required to transition from a current state s to a next state s' .*

Algorithm 1 Maximum Next-State Entropy RL

```

1: Inputs: Initialize inner policy  $\pi_i^\psi$ , inverse dynamics model  $p_{\text{inv}}^\phi$ , action mapping layer  $f_\theta$ .
2: for iteration  $t = 0, 1, 2, \dots$  do
3:   Observe state  $s$  and select internal action  $e \sim \pi_i^\psi(\cdot | s)$ 
4:   Execute action  $a = f_\theta(e)$  in the environment
5:   Obtain next state  $s'$ , reward  $r$ , and done signal  $d$ 
6:   Store  $(s, e, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
7:   if it's time to update then
8:     Update inverse dynamics model parameters  $\phi$  based on Equation 11
9:     Update the action mapping layer parameters  $\theta$  based on Theorem 5.2
10:    Update inner policy  $\pi_i^\psi$  based on the standard maximum policy entropy RL algorithms
11:  end if
12: end for
13: Return: Policy parameters  $\psi$  and action mapping layer parameters  $\theta$ 

```

Proof. Please refer to Appendix C for the detailed proof. □

Theorem 5.1 suggests that we do not require π_i is the non-redundant policy. In addition, it is noteworthy that the entropy of the dynamics model $\mathcal{H}_{\text{model}}$ is an inherent characteristic of the system. We can regard this term as a constant c . Therefore, if we can minimize the gap term, we can maximize the next-state entropy by optimizing the policy entropy of the inner policy π_i .

5.2 MAXIMIZE NEXT-STATE ENTROPY

Based on the analysis in Theorem 5.1, we can maximize next-state entropy as follows:

$$\begin{aligned}
J_{\text{MNSE}}(\pi) &= \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \alpha \mathcal{H}(S_{t+1} | S_t = s_t, \pi)) \right] \\
&= \mathbb{E}_{(s_t, a_t, s_{t+1}) \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t \left(r_t + \alpha \underbrace{\mathcal{H}(\pi_i^\psi(\cdot | s_t))}_{\text{Entropy of Inner Policy}} \right. \right. \\
&\quad \left. \left. + \alpha \underbrace{\mathbb{E}_{a \sim \pi(\cdot | s_t)} \mathbb{E}_{s_{t+1} \sim P(\cdot | s, a)} \left[\log \left[p_{\text{inv}}^\phi(e = f^{-1}(a_t; \theta) | s_t, s_{t+1}) \right] \right]}_{\text{Gap Term}} + c \right) \right], \tag{9}
\end{aligned}$$

where the inverse dynamics p_{inv}^ϕ is parameterized with ϕ , and the inner policy π_i^ψ is parameterized with ψ . It is noteworthy that $\pi(\cdot | s)$ in the above equation implicitly includes the inner policy π_i^ψ and the action mapping layer f_θ defined in the Equation 6. Therefore, we need to optimize these three parameters ϕ, ψ, θ simultaneously.

Based on the Equation 9, we can conclude that maximizing $J_{\text{MNSE}}(\pi)$ is equivalent to maximizing the entropy of the inner policy π_i^ψ if and only if the *Gap Term* is zero. Further, with an appropriate model to estimate the inverse dynamics, since $\log p_{\text{inv}}^\phi \leq 0$, the *Gap Term* is always ≤ 0 . Meanwhile, in the optimization of the gap term, we set $\gamma = 1$ to facilitate sampling and training, so we need to optimize ϕ, θ by maximizing the gap term as follows:

$$\begin{aligned}
\phi^*, \theta^* &= \arg \max J_{\text{Gap Term}}(\phi, \theta) \\
&= \arg \max_{\phi, \theta} \mathbb{E}_{\substack{(s_t, a_t, s_{t+1}) \sim \pi \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} \left[p_{\text{inv}}^\phi(e = f^{-1}(a_t; \theta) | s_t, s_{t+1}) \right] \tag{10}
\end{aligned}$$

Specifically, we use the iterative optimization mechanism to optimize ϕ and θ for the given inner policy π_i . Let ϕ^k and θ^k denote the learned parameters after iteration k , then:

Step 1. Given $\theta = \theta^k$, we optimize the objective function in Equation 10 by optimizing ϕ :

$$\begin{aligned}\phi^{k+1} &= \arg \max_{\phi} J_{\text{Gap Term}}(\phi, \theta^k) \\ &= \arg \min_{\phi} \mathbb{E}_{(s, s', e) \in \mathcal{D}} - \log \left[p_{\text{inv}}^{\phi}(e | s, s') \right]\end{aligned}\quad (11)$$

where \mathcal{D} denotes the dataset collected by policy π . Intuitively, minimizing the objective in Equation 11 amounts to maximum likelihood estimation of actions.

Step 2. Given $\phi = \phi^{k+1}$, we optimize the objective function in Equation 10 by optimizing θ :

$$\begin{aligned}\theta^{k+1} &= \arg \max_{\theta} J_{\text{Gap Term}}(\phi^{k+1}, \theta) \\ \text{where } a &= f(e; \theta), \quad \pi(a | s) = \pi_i(e | s) * \left| \frac{\partial f(e; \theta)}{\partial e} \right|^{-1}.\end{aligned}\quad (12)$$

It is noteworthy that the action mapping layer f_{θ} is implicitly included in π , preventing direct optimization of θ . To solve this issue, we adopt the gradient descent method as follows:

Theorem 5.2. *Given the inverse dynamic $p_{\text{inv}}^{\phi^{k+1}}(e | s, s')$, the gradient of $J_{\text{Gap Term}}(\theta)$ can be derived as:*

$$\nabla_{\theta} J_{\text{Gap Term}}(\theta) = \mathbb{E}_{\substack{s_0 \in \mathcal{S}, a_t \sim \pi(\cdot | s_t) \\ s_{t+1} \sim P(\cdot | s_t, a_t)}} \left[\nabla_{\theta} \log \left| \frac{\partial f(e; \theta)}{\partial e} \right|^{-1} \Big|_{e=e_t} \log p_{\text{inv}}^{\phi^{k+1}}(e_t | s_t, s_{t+1}) \right], \quad (13)$$

where $\left| \frac{\partial f(e; \theta)}{\partial e} \right|$ is the Jacobian determinant of the function f and $e_t = f^{-1}(a_t; \theta)$.

Proof. Please refer to Appendix D for the detailed proof. \square

Based on Theorem 5.2, we can perform gradient updates on θ to train the action mapping layer f_{θ} .

5.3 PRACTICAL IMPLEMENTATION

The overall framework of our algorithm is illustrated in Algorithm 1. After interacting with the environment, we iteratively train the inverse dynamics network $p_{\text{inv}}^{\phi}(e | s, s')$ and the action mapping layer $f_{\theta}(e)$ using the collected data based on the Equation 11 and Theorem 5.2. For the inner policy π_i , we use the standard maximum policy entropy RL methods, such as SAC (Haarnoja et al., 2018c).

Specifically, we construct the invertible action mapping function f using a piecewise linear function with N parameters ($\vec{\theta} \in \mathbb{R}^N$), defined as follows:

$$f(x) = \sum_{j=1}^{i-1} k_j \cdot \frac{1}{N} + k_i \left(x - \frac{i-1}{N} \right) \quad \text{for } x \in \left[\frac{i-1}{N}, \frac{i}{N} \right]$$

where $k_i = N \cdot \frac{\exp(\theta_i)}{\sum_{j=1}^N \exp(\theta_j)}$ for $i = 1, 2, \dots, N$. For each $i \in [1, N]$, k_i represents the slope of the linear function in the interval $\left[\frac{i-1}{N}, \frac{i}{N} \right]$. For environments with multidimensional action spaces, we construct $|\mathcal{A}|$ action mapping functions, each applying an independent transformation to its respective dimension.

For the inverse dynamics model, rather than using Gaussian distributions to predict the distribution of continuous actions, we discretize the actions in the dataset and employ discrete multinomial distributions. These multinomial distributions output an M -dimensional vector, where each dimension corresponds to the probability that the predicted action lies within the interval $\left[\frac{j-1}{M}, \frac{j}{M} \right]$, for $j \in [1, M]$. This approach enables the inverse dynamics model to capture complex and multimodal behaviors effectively.

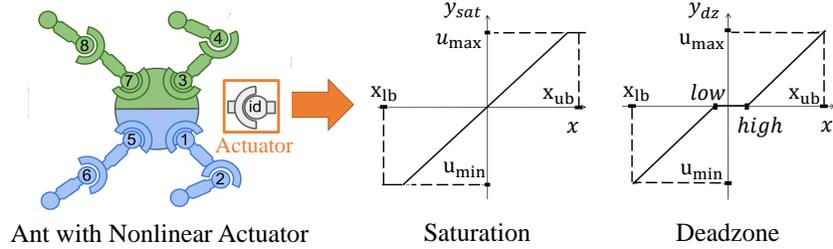


Figure 2: The control of systems with actuators demonstrating input nonlinearities (e.g., saturation, deadzone).

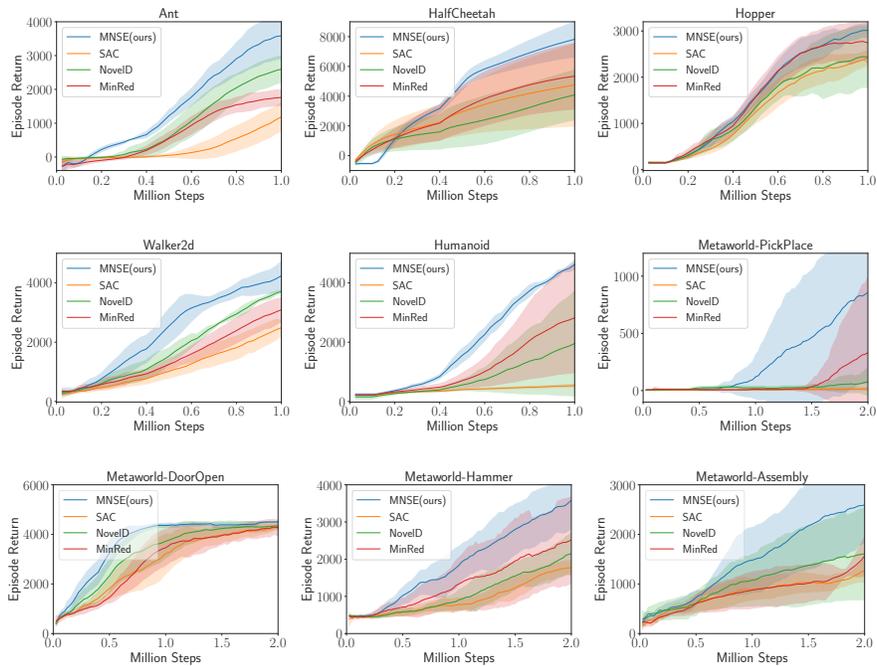


Figure 3: The experimental results in MuJoCo and MetaWorld with nonlinear actuators under five random seeds.

6 EXPERIMENTS

In this section, we aim to address the following questions: (1) How do traditional RL methods perform when the policy entropy does not accurately reflect the diversity of the state? (2) How does MNSE compare with other state-of-the-art approaches for maximizing entropy? (3) How sensitive is MNSE to the hyperparameters of the algorithm?

6.1 EXPERIMENTAL SETTING

In this section, we examine actuators with input nonlinearities. In industrial applications, nonlinear actuators are a common challenge due to wear-and-tear or inaccuracies in mechanical components. In this work, we consider two types of input nonlinearities: saturation and deadzone, as shown in Figure 2. Both saturation and deadzone can reduce control accuracy, severely influencing the control

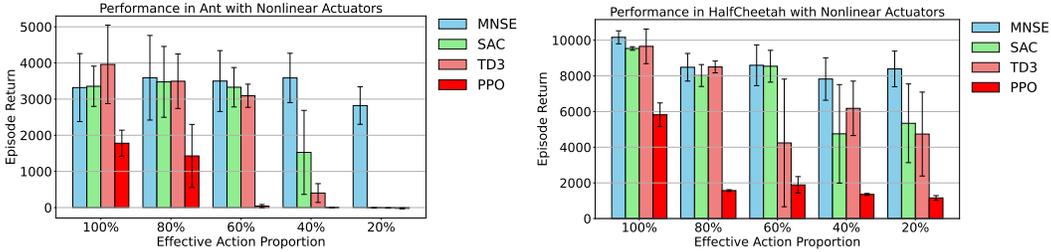


Figure 4: **(Left)**: Performance of MNSE, SAC, TD3, PPO in MuJoCo Ant with nonlinear actuators under different effective action proportions. **(Right)**: Performance of MNSE, SAC, TD3, PPO in MuJoCo HalfCheetah with nonlinear actuators under different effective action proportions.

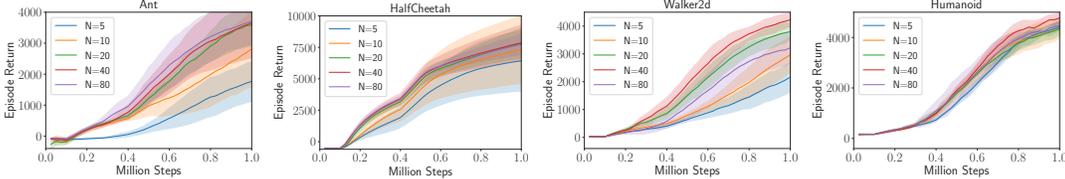


Figure 5: Ablation study of the number of parameters in the piecewise linear function.

system’s performance. Specifically, the saturation and deadzone are formulated as follows:

$$y_{\text{sat}} = \begin{cases} u_{\text{max}} & \text{if } x \geq u_{\text{max}} \\ x & \text{if } u_{\text{min}} < x < u_{\text{max}} \\ u_{\text{min}} & \text{if } x \leq u_{\text{min}} \end{cases} \quad y_{\text{dz}} = \begin{cases} x - \text{high} & \text{if } u \geq \text{high} \\ 0 & \text{if } \text{low} < u < \text{high} \\ x - \text{low} & \text{if } u \leq \text{low} \end{cases}$$

where $x \in [x_{\text{lb}}, x_{\text{ub}}], y \in [u_{\text{min}}, u_{\text{max}}]$. Notably, within a specific environment, u_{min} and u_{max} remain invariant, guaranteeing theoretical optimality across various EAPs.

In the experiments, we employ either saturation or deadzone to each joint of the robot. In addition, we characterize the proportion of effective actions in the action space as Effective Action Proportion (EAP):

$$\text{EAP} = \frac{u_{\text{max}} - u_{\text{min}}}{x_{\text{ub}} - x_{\text{lb}}}. \tag{14}$$

As EAP increases, the actuator nonlinearity increases, making the task more difficult.

Baselines We compare MNSE with baselines using various maximum entropy methods. We first compare our method with the standard maximum policy entropy method, SAC (Haarnoja et al., 2018c;d). In addition, we compare our method with MinRed (Baram et al., 2021), which directly maximizes the next-state entropy to minimize action redundancy. We also compare our method with the strong state-novelty-based exploration method (Zhang et al., 2021). Please refer to Appendix E for the experimental details.

6.2 EXPERIMENTAL RESULTS

Answer of Question 1: To show the impact of the nonlinear actuators on the traditional RL methods, we conduct experiments in MuJoCo environments with various EAP by employing SAC (Haarnoja et al., 2018c;d), TD3 (Silver et al., 2014; Lillicrap, 2015), PPO (Schulman et al., 2017), and our proposed MNSE approach. As shown in Figure 4, as the EAP decreases and actuator nonlinearity increases, the performance of SAC gradually degrades and eventually collapses. In the Ant environment (Figure 4, left), when the EAP decreases to 40%, the performance of the SAC algorithm shows a marked decline, and at 20% EAP or lower, SAC completely fails. Similarly, in the HalfCheetah environment (Figure 4, right), when the EAP drops to 40% or lower, the performance of SAC is only half of that achieved without actuator nonlinearity (EAP = 100%).

486 This decline in performance can be attributed to the fact that as EAP decreases, policy entropy no
 487 longer accurately represents the diversity of the next state, which hinders the effective exploration
 488 guidance within the SAC framework. In contrast, our MNSE approach based on next-state entropy
 489 maintains stable performance across varying EAP environments, significantly outperforming SAC.
 490 Moreover, our experiments indicate that not only SAC, but traditional RL algorithms such as TD3
 491 and PPO, also struggle to handle the reduction in effective action space as actuator nonlinearity
 492 increases. As shown in Figure 4 (right), when the EAP decreases, both TD3 and PPO exhibit
 493 varying degrees of performance degradation, with some cases resulting in complete failure. This
 494 underscores that actuator nonlinearity is a significant challenge across various algorithms in the
 495 field of reinforcement learning.

496 **Answer of Question 2:** To show that MNSE can maximize next-state entropy, we conduct ex-
 497 periments on Mujoco and MetaWorld tasks. As illustrated in Figure 3, our method consistently
 498 outperforms baseline approaches across various experimental environments. NovelD and MinRed
 499 employ bonus-based strategies to promote the exploration of diverse states. Specifically, NovelD
 500 incentivizes agents by evaluating state novelty, while MinRed provides additional rewards based
 501 on transition entropy. In contrast, our method, MNSE, establishes an action mapping layer that
 502 effectively bridges the gap between policy entropy and next-state entropy. By directly promot-
 503 ing exploration through next-state entropy, MNSE demonstrates superior performance compared to
 504 bonus-based methods across various domains, as evidenced by our experimental results.

505 **Answer of Question 3:** To test how the algorithm’s hyperparameters affect the performance of
 506 MNSE, we change the number of parameters N in the piecewise linear function, which significantly
 507 influences the expressive power of the action mapping function f . As shown in Figure 5, we con-
 508 ducted ablation experiments across four MuJoCo environments to evaluate the impact of N . The
 509 results reveal that when N is small, the limited expressive capacity of f leads to suboptimal algo-
 510 rithm performance. As N increases, algorithm performance gradually improves. However, once
 511 $N \geq 20$, the performance stabilizes and shows little variation. To balance expressive power with
 512 computational efficiency, we consistently employed $N = 20$ in all experimental implementations.

514 7 DISCUSSION

515 *Why Maximize Next-State Entropy in Reinforcement Learning?* Entropy regularization is a fun-
 516 damental technique in reinforcement learning. By integrating an entropy maximization term, it
 517 enhances robustness to model and estimation errors (Ziebart et al., 2010), promotes the acquisi-
 518 tion of diverse behaviors (Haarnoja et al., 2017), facilitates broader exploration (Fox et al., 2015;
 519 Haarnoja et al., 2018c;d) and accelerates the learning process by smoothing the optimization land-
 520 scape (Ahmed et al., 2019). However, maximizing policy entropy may not directly promote policy
 521 optimization due to redundancy in the action space. In such cases, next-state entropy extends the
 522 concept of policy entropy more directly. Specifically, next-state entropy measures the entropy of
 523 the next state resulting from the policy, rather than the action itself. This shift allows next-state en-
 524 tropy to capture the diversity of effects induced by actions. By bridging the gap between next-state
 525 and policy entropy, our method retains the benefits of policy entropy while addressing inefficiencies
 526 caused by action redundancy.
 527

528 8 CONCLUSION

529 In this work, we demonstrate a critical problem: the maximum next-state entropy of the RL agent.
 530 We first systematically elucidate the distinctions and interrelationships between next-state entropy
 531 and policy entropy. Then, to bridge the gap between these two concepts, we integrate inverse dynam-
 532 ics with an action mapping layer. We demonstrate that by optimizing both the action mapping layer
 533 and inverse dynamics model, we can maximize the next-state entropy by optimizing the inner policy
 534 entropy. We conduct extensive experiments and demonstrate that our method outperforms baseline
 535 methods across various domains. Future research will focus on extending MNSE to accommodate
 536 more complex action space structures and exploring its potential applications in robotics.
 537
 538
 539

540 REFERENCES

- 541 Zafarali Ahmed, Nicolas Le Roux, Mohammad Norouzi, and Dale Schuurmans. Understanding the
542 impact of entropy on policy optimization. In *International conference on machine learning*, pp.
543 151–160. PMLR, 2019.
- 544
- 545 Adrià Puigdomènech Badia, Pablo Sprechmann, Alex Vitvitskyi, Daniel Guo, Bilal Piot, Steven
546 Kapturowski, Olivier Tieleman, Martín Arjovsky, Alexander Pritzel, Andrew Bolt, et al. Never
547 give up: Learning directed exploration strategies. *arXiv preprint arXiv:2002.06038*, 2020.
- 548
- 549 Er-Wei Bai. Identification of linear systems with hard input nonlinearities of known structure. *Auto-*
550 *matica*, 38(5):853–860, 2002.
- 551
- 552 Nir Baram, Guy Tennenholtz, and Shie Mannor. Action redundancy in reinforcement learning. In
553 *Uncertainty in Artificial Intelligence*, pp. 376–385. PMLR, 2021.
- 554
- 555 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
556 Unifying count-based exploration and intrinsic motivation. *Advances in neural information pro-*
557 *cessing systems*, 29, 2016.
- 558
- 559 Yuri Burda, Harrison Edwards, Amos Storkey, and Oleg Klimov. Exploration by random network
560 distillation. In *International Conference on Learning Representations*, 2018.
- 561
- 562 Yash Chandak, Georgios Theodorou, James Kostas, Scott Jordan, and Philip Thomas. Learning
563 action representations for reinforcement learning. In *International conference on machine learn-*
564 *ing*, pp. 941–950. PMLR, 2019.
- 565
- 566 Adrien Ecoffet, Joost Huizinga, Joel Lehman, Kenneth O Stanley, and Jeff Clune. Go-explore: a
567 new approach for hard-exploration problems. *arXiv preprint arXiv:1901.10995*, 2019.
- 568
- 569 Roy Fox, Ari Pakman, and Naftali Tishby. Taming the noise in reinforcement learning via soft
570 updates. *arXiv preprint arXiv:1512.08562*, 2015.
- 571
- 572 Giacomo Galuppini, Lalo Magni, and Davide Martino Raimondo. Model predictive control of sys-
573 tems with deadzone and saturation. *Control Engineering Practice*, 78:56–64, 2018.
- 574
- 575 Yang Gao, Huazhe Xu, Ji Lin, Fisher Yu, Sergey Levine, and Trevor Darrell. Reinforcement learning
576 from imperfect demonstrations. *arXiv preprint arXiv:1802.05313*, 2018.
- 577
- 578 Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Alaa Saade, Shantanu Thakoor, Bilal Piot,
579 Bernardo Avila Pires, Michal Valko, Thomas Mesnard, Tor Lattimore, and Rémi Munos. Geo-
580 metric entropic exploration. *arXiv preprint arXiv:2101.02055*, 2021.
- 581
- 582 Tuomas Haarnoja, Haoran Tang, Pieter Abbeel, and Sergey Levine. Reinforcement learning with
583 deep energy-based policies. In *International conference on machine learning*, pp. 1352–1361.
584 PMLR, 2017.
- 585
- 586 Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning
587 to walk via deep reinforcement learning. *arXiv preprint arXiv:1812.11103*, 2018a.
- 588
- 589 Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine.
590 Composable deep reinforcement learning for robotic manipulation. In *2018 IEEE international*
591 *conference on robotics and automation (ICRA)*, pp. 6244–6251. IEEE, 2018b.
- 592
- 593 Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy
maximum entropy deep reinforcement learning with a stochastic actor. In *International confer-*
ence on machine learning, pp. 1861–1870. PMLR, 2018c.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash
Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and appli-
cations. *arXiv preprint arXiv:1812.05905*, 2018d.

- 594 Seungyul Han and Youngchul Sung. A max-min entropy framework for reinforcement learn-
595 ing. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan (eds.),
596 *Advances in Neural Information Processing Systems*, volume 34, pp. 25732–25745. Curran
597 Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper_files/
598 paper/2021/file/d7b76edf790923bf7177f7ebba5978df-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/d7b76edf790923bf7177f7ebba5978df-Paper.pdf).
- 599 Elad Hazan, Sham Kakade, Karan Singh, and Abby Van Soest. Provably efficient maximum entropy
600 exploration. In *International Conference on Machine Learning*, pp. 2681–2691. PMLR, 2019.
- 601 Dailin Hu, Pieter Abbeel, and Roy Fox. Count-based temperature scheduling for maximum entropy
602 reinforcement learning. In *Deep RL Workshop NeurIPS 2021*, 2021.
- 603 Riashat Islam, Zafarali Ahmed, and Doina Precup. Marginalized state distribution entropy regular-
604 ization in policy optimization. *arXiv preprint arXiv:1912.05128*, 2019.
- 605 Lisa Lee, Benjamin Eysenbach, Emilio Parisotto, Eric Xing, Sergey Levine, and Ruslan Salakhutdi-
606 nov. Efficient exploration via state marginal matching. *arXiv preprint arXiv:1906.05274*, 2019.
- 607 TP Lillicrap. Continuous control with deep reinforcement learning. *arXiv preprint
608 arXiv:1509.02971*, 2015.
- 609 Hao Liu and Pieter Abbeel. Behavior from the void: Unsupervised active pre-training. *Advances in
610 Neural Information Processing Systems*, 34:18459–18473, 2021.
- 611 Sam Lobel, Akhil Bagaria, and George Konidaris. Flipping coins to estimate pseudocounts for
612 exploration in reinforcement learning. In *International Conference on Machine Learning*, pp.
613 22594–22613. PMLR, 2023.
- 614 Marlos C Machado, Marc G Bellemare, and Michael Bowling. Count-based exploration with the
615 successor representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol-
616 ume 34, pp. 5125–5133, 2020.
- 617 Luke Metz, Julian Ibarz, Navdeep Jaitly, and James Davidson. Discrete sequential prediction of
618 continuous actions for deep rl. *arXiv preprint arXiv:1705.05035*, 2017.
- 619 Deepak Pathak, Pulkit Agrawal, Alexei A Efros, and Trevor Darrell. Curiosity-driven exploration
620 by self-supervised prediction. In *International conference on machine learning*, pp. 2778–2787.
621 PMLR, 2017.
- 622 Deepak Pathak, Dhiraj Gandhi, and Abhinav Gupta. Self-supervised exploration via disagreement.
623 In *International conference on machine learning*, pp. 5062–5071. PMLR, 2019.
- 624 Antonin Raffin. RL baselines3 zoo. [https://github.com/DLR-RM/
625 rl-baselines3-zoo](https://github.com/DLR-RM/rl-baselines3-zoo), 2020.
- 626 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
627 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of
628 Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/
629 20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 630 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy
631 optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 632 Wenjie Shi, Shiji Song, and Cheng Wu. Soft policy gradient method for maximum entropy deep
633 reinforcement learning. In *Proceedings of the 28th International Joint Conference on Artificial
634 Intelligence*, pp. 3425–3431, 2019.
- 635 David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller.
636 Deterministic policy gradient algorithms. In *International conference on machine learning*, pp.
637 387–395. Pmlr, 2014.
- 638
- 639
- 640
- 641
- 642
- 643
- 644
- 645
- 646
- 647

- 648 Charles Sun, Jędrzej Orbik, Coline Manon Devin, Brian H. Yang, Abhishek Gupta, Glen Berseth,
649 and Sergey Levine. Fully autonomous real-world reinforcement learning with applications to
650 mobile manipulation. In Aleksandra Faust, David Hsu, and Gerhard Neumann (eds.), *Proceed-*
651 *ings of the 5th Conference on Robot Learning*, volume 164 of *Proceedings of Machine Learn-*
652 *ing Research*, pp. 308–319. PMLR, 08–11 Nov 2022. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v164/sun22a.html)
653 [press/v164/sun22a.html](https://proceedings.mlr.press/v164/sun22a.html).
- 654 Guy Tennenholtz and Shie Mannor. The natural language of actions. *ArXiv*, abs/1902.01119, 2019.
655
- 656 Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control.
657 In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pp. 5026–5033.
658 IEEE, 2012.
- 659 Tianhe Yu, Deirdre Quillen, Zhanpeng He, Ryan Julian, Karol Hausman, Chelsea Finn, and Sergey
660 Levine. Meta-world: A benchmark and evaluation for multi-task and meta reinforcement learning.
661 In *Conference on robot learning*, pp. 1094–1100. PMLR, 2020.
662
- 663 Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J Mankowitz, and Shie Mannor.
664 Learn what not to learn: Action elimination with deep reinforcement learning. In
665 S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett
666 (eds.), *Advances in Neural Information Processing Systems*, volume 31. Curran Asso-
667 ciates, Inc., 2018a. URL [https://proceedings.neurips.cc/paper/2018/file/](https://proceedings.neurips.cc/paper/2018/file/645098b086d2f9e1e0e939c27f9f2d6f-Paper.pdf)
668 [645098b086d2f9e1e0e939c27f9f2d6f-Paper.pdf](https://proceedings.neurips.cc/paper/2018/file/645098b086d2f9e1e0e939c27f9f2d6f-Paper.pdf).
- 669 Tom Zahavy, Matan Haroush, Nadav Merlis, Daniel J Mankowitz, and Shie Mannor. Learn what not
670 to learn: Action elimination with deep reinforcement learning. *Advances in neural information*
671 *processing systems*, 31, 2018b.
- 672 Tianjun Zhang, Huazhe Xu, Xiaolong Wang, Yi Wu, Kurt Keutzer, Joseph E Gonzalez, and Yuan-
673 dong Tian. Noveld: A simple yet effective exploration criterion. *Advances in Neural Information*
674 *Processing Systems*, 34:25217–25230, 2021.
675
- 676 Dianyu Zhong, Yiqin Yang, and Qianchuan Zhao. No prior mask: Eliminate redundant action for
677 deep reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
678 volume 38, pp. 17078–17086, 2024.
- 679 Brian D Ziebart, Andrew L Maas, J Andrew Bagnell, Anind K Dey, et al. Maximum entropy inverse
680 reinforcement learning. In *Aaai*, volume 8, pp. 1433–1438. Chicago, IL, USA, 2008.
681
- 682 Brian D Ziebart, J Andrew Bagnell, and Anind K Dey. Modeling interaction via the principle of
683 maximum causal entropy. 2010.
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701

A TOY EXAMPLE

Consider a one-step MDP with a deterministic transition, where $s_{t+1} = \max(a_t, 0)$, which means actions less than zero are redundant. We assume an initial policy π is a Gaussian policy with a mean less than zero: $p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$, $\mu < 0$. The corresponding policy entropy is given by:

$$\mathcal{H}(\pi(\cdot | s_t)) = \log(\sqrt{2\pi e}\sigma), \quad (15)$$

which is independent of the mean μ .

Since $P^\pi(s' | s) = 0$ when $s' < 0$, it follows that $P^\pi(s' | s)$ is a rectified Gaussian distribution. Then we have:

$$\begin{aligned} & \mathcal{H}(S_{t+1} | S_t = s, \pi) \\ &= - \int_0^{+\infty} p(x) \log p(x) dx \\ &= - \int_0^{+\infty} \frac{\varphi\left(\frac{x-\mu}{\sigma}\right)}{\sigma} \log \frac{\varphi\left(\frac{x-\mu}{\sigma}\right)}{\sigma} dx \\ &= - \int_{-\frac{\mu}{\sigma}}^{+\infty} \varphi(z) \log \frac{\varphi(z)}{\sigma} dz \\ &= \int_{-\frac{\mu}{\sigma}}^{+\infty} -\varphi(z) \log \varphi(z) dz + \log \sigma \int_{-\frac{\mu}{\sigma}}^{+\infty} \varphi(z) dz \\ &= \int_{-\infty}^{+\infty} -\varphi(z) \log \varphi(z) dz - \int_{-\infty}^{-\frac{\mu}{\sigma}} -\varphi(z) \log \varphi(z) dz + \log \sigma (1 - \Phi(-\frac{\mu}{\sigma})) \\ &= \underbrace{\log(\sqrt{2\pi e}\sigma)}_{\text{policy entropy}} + \left\{ \Phi\left(\frac{\mu}{\sigma}\right) \log \sigma - \int_{-\infty}^{-\frac{\mu}{\sigma}} -\varphi(z) \log \varphi(z) dz \right\}, \end{aligned} \quad (16)$$

where $\varphi(z) = \frac{e^{-z^2/2}}{\sqrt{2\pi}}$ is the probability density function of the standard normal distribution, and $\Phi(x) = \int_{-\infty}^x \varphi(z) dz$ is its cumulative distribution function.

It is challenging to provide a rigorous mathematical analysis of this relationship. Specifically, when we assume $\sigma > 1$ in the toy example, $\log \sigma > 0$, and $\Phi\left(\frac{\mu}{\sigma}\right) \log \sigma$ is positively correlated with μ . Furthermore, since $0 < \varphi(z) \leq \frac{1}{\sqrt{2\pi}} < 1$, the term $-\int_{-\infty}^{-\frac{\mu}{\sigma}} -\varphi(z) \log \varphi(z) dz$ is also positively correlated with μ . Meanwhile, we have conducted numerical experiments to examine the relationship between the μ variable and next-state entropy for different σ values. As shown in the results, when σ values are $\{1.0, 1.5, 2.0\}$, the next-state entropy increases as the μ value increases, demonstrating a positive relationship between μ and entropy.

The above analysis indicates that to increase $\mathcal{H}(S_{t+1} | S_t = s, \pi)$, one effective method is not only to increase σ but also to increase μ .

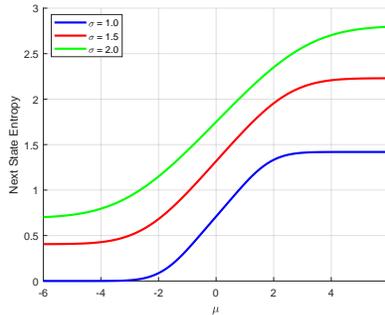


Figure 6: The next-state entropy increases as the mu value increases.

B PROOF OF THEOREM 4.4

Proof. Recall that the entropy of state is:

$$\begin{aligned}\mathcal{H}(S_{t+1} | S_t = s, \pi) &= -\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim P(\cdot | s, a)} \log [P^\pi(s' | s)] \\ &= -\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim P(\cdot | s, a)} \log [\mathbb{E}_{a \sim \pi(\cdot | s)} P(s' | s, a)]\end{aligned}\quad (17)$$

Then, according to the Definition 4.3, we have:

$$\begin{aligned}\mathbb{E}_{s' \sim P^\pi(\cdot | s, a)} \log [\mathbb{E}_{a \sim \pi(\cdot | s)} P(s' | s, a)] &= \mathbb{E}_{s' \sim P^\pi(\cdot | s, a)} \log [\pi(a | s) P(s' | s, a)] \\ &= \mathbb{E}_{s' \sim P^\pi(\cdot | s, a)} [\log \pi(a | s) + \log P(s' | s, a)],\end{aligned}\quad (18)$$

Combining the Equation 17 and Equation 18, we have:

$$\begin{aligned}\mathcal{H}(S_{t+1} | S_t = s, \pi) &= -\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim P^\pi(\cdot | s, a)} [\log \pi(a | s) + \log P(s' | s, a)] \\ &= -\mathbb{E}_{a \sim \pi(\cdot | s)} [\log \pi(a | s)] - \mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim P^\pi(\cdot | s, a)} [\log P(s' | s, a)] \\ &= \mathcal{H}(\pi(\cdot | s)) + \mathbb{E}_{a \sim \pi(\cdot | s)} \mathcal{H}(S_{t+1} | S_t = s, a).\end{aligned}\quad (19)$$

□

C PROOF OF THEOREM 5.1

Proof. Firstly, we can derive the next-state entropy $\mathcal{H}(S_{t+1} | S_t = s, \pi)$ as follows.

$$\begin{aligned}\mathcal{H}(S_{t+1} | S_t = s, \pi) &= -\mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, a)} \log p^\pi(s' | s) \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log p^{\pi_i}(s' | s) \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \{ \log p^{\pi_i}(s' | s) - \log [\pi_i(e | s) p(s' | s, e)] + \log [\pi_i(e | s) p(s' | s, e)] \}\end{aligned}\quad (20)$$

where

$$p^\pi(s' | s) = \mathbb{E}_{\tilde{a} \sim \pi(\cdot | s)} p(s' | s, \tilde{a}), \quad p^{\pi_i}(s' | s) = \mathbb{E}_{\tilde{e} \sim \pi_i(\cdot | s)} p(s' | s, \tilde{e}).$$

Here, we utilize the property of inverse dynamics:

$$p_{\text{inv}}(e | s, s') = \frac{\pi_i(e | s) p(s' | s, e)}{p^{\pi_i}(s' | s)}\quad (21)$$

then Eq. 20 can be derived as follows:

$$\begin{aligned}\mathcal{H}(S_{t+1} | S_t = s, \pi) &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \{ \log [\pi_i(e | s) p(s' | s, e)] + \log p^{\pi_i}(s' | s) - \log [\pi_i(e | s) p(s' | s, e)] \} \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log [\pi_i(e | s) p(s' | s, e)] - \\ &\quad \mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \{ \log p^{\pi_i}(s' | s) - \log [\pi_i(e | s) p(s' | s, e)] \} \quad (\text{Property of Eq. 21}) \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log [\pi_i(e | s) p(s' | s, e)] - \mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \{ -\log p_{\text{inv}}(e | s, s') \} \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log [\pi_i(e | s) p(s' | s, e)] - \mathbb{E}_{e \sim \pi_i(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, e)} \{ -\log p_{\text{inv}}(e | s, s') \} \\ &= -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log [\pi_i(e | s) p(s' | s, e)] + \mathbb{E}_{a \sim \pi(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, a)} [\log [p_{\text{inv}}(e | s, s')]].\end{aligned}$$

The first term is as follows:

$$\begin{aligned}& -\mathbb{E}_{s' \sim p^{\pi_i}(\cdot | s)} \mathbb{E}_{e \sim p_{\text{inv}}(\cdot | s, s')} \log [\pi_i(e | s) p(s' | s, e)] \\ &= -\mathbb{E}_{e \sim \pi_i(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, e)} \log [\pi_i(e | s) p(s' | s, e)] \\ &= -\mathbb{E}_{e \sim \pi_i(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, e)} \log \pi_i(e | s) - \mathbb{E}_{e \sim \pi_i(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, e)} \log p(s' | s, e) \\ &= \mathbb{E}_{e \sim \pi_i(\cdot | s)} [-\log \pi_i(e | s)] + \mathbb{E}_{e \sim \pi_i(\cdot | s)} \mathbb{E}_{s' \sim p(\cdot | s, e)} [-\log p(s' | s, e)] \\ &= \mathcal{H}(\pi_i(\cdot | s)) + \mathcal{H}_{\text{model}}\end{aligned}$$

where $\mathcal{H}_{\text{model}} = \mathbb{E}_{e \sim \pi_i(\cdot|s)} [\mathcal{H}(S_{t+1} | S_t = s, e)]$ is the entropy of the dynamics model.

As a result, we have:

$$\mathcal{H}(S_{t+1} | S_t = s, \pi) = \mathcal{H}(\pi_i(\cdot | s)) + \underbrace{\mathbb{E}_{a \sim \pi(\cdot|s)} \mathbb{E}_{s' \sim p(\cdot|s,a)} [\log [p_{\text{inv}}(e | s, s')]]}_{\text{Gap Term}} + \mathcal{H}_{\text{model}}$$

□

D PROOF OF THEOREM 5.2

Proof. We will derive $\nabla_{\theta} J_{\text{Gap Term}}(\theta)$ in the following.

$$\begin{aligned} \nabla_{\theta} \mathbb{E}_{\substack{s_0 \in \mathcal{S}, a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim P(\cdot|s_t, a_t)}} \left[p_{\text{inv}}^{\phi}(e = f^{-1}(a_t; \theta) | s_t, s_{t+1}) \right] &= \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \frac{\partial \pi(a|s)}{\partial \theta} \mathbb{E}_{s' \sim p(\cdot|s,a)} \log p_{\text{inv}}^{\phi^{k+1}}(e | s, s') \\ &= \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \pi(a|s) \frac{1}{\pi(a|s)} \frac{\partial \pi(a|s)}{\partial \theta} \mathbb{E}_{s' \sim p(\cdot|s,a)} \log p_{\text{inv}}^{\phi^{k+1}}(e | s, s') \\ &= \int_{s \in \mathcal{S}} \int_{a \in \mathcal{A}} \pi(a|s) \frac{\partial \log \pi(a|s)}{\partial \theta} \mathbb{E}_{s' \sim p(\cdot|s,a)} \log p_{\text{inv}}^{\phi^{k+1}}(e | s, s') \\ &= \mathbb{E}_{\substack{s_0 \in \mathcal{S}, a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim P(\cdot|s_t, a_t)}} \frac{\partial \log \pi(a|s)}{\partial \theta} \log p_{\text{inv}}^{\phi^{k+1}}(e | s, s') \end{aligned}$$

when $a = f_{\theta}(e)$, recall that:

$$\pi(a|s) = \pi_i(e | s) * \left| \frac{\partial a}{\partial e} \right|^{-1},$$

According to the inverse function theorem: If $y = f(x)$ and $x = f^{-1}(y)$, we have:

$$\frac{df^{-1}(y)}{dy} = \frac{dx}{dy} = \left(\frac{dy}{dx} \right)^{-1} = \left(\frac{df(x)}{dx} \right)^{-1}$$

then:

$$\frac{\partial \log \pi(a|s)}{\partial \theta} = \frac{\partial \log \left| \frac{\partial f_{\theta}^{-1}(a)}{\partial a} \right|}{\partial \theta} = \nabla_{\theta} \log \left| \frac{\partial f_{\theta}(e)}{\partial e} \right|^{-1}$$

In conclusion,

$$\nabla_{\theta} J_{\text{Gap Term}}(\theta) = \mathbb{E}_{\substack{s_0 \in \mathcal{S}, a_t \sim \pi(\cdot|s_t) \\ s_{t+1} \sim P(\cdot|s_t, a_t)}} \left\{ \nabla_{\theta} \log \left| \frac{\partial f(e; \theta)}{\partial e} \right|^{-1} \Big|_{e=e_t} \log p_{\text{inv}}^{\phi^{k+1}}(e_t | s_t, s_{t+1}) \right\} \quad (22)$$

□

E EXPERIMENTAL DETAILS

Our algorithm, MNSE, is developed based on the SAC algorithm from the RL Baselines3 Zoo (Raffin, 2020; Raffin et al., 2021). The hyperparameters for MNSE are detailed in Table E. In all experimental implementations, we consistently employed $N = 20$ as the number of parameters in the piecewise linear function. The baseline SAC shares the same hyperparameters as those of MNSE.

Hyper-parameter	Value
Shared	
Learning rate	3×10^{-4}
Buffer size	1×10^6
Learning starts	1×10^5
Batch size	64
Soft update coefficient τ	0.005
Discount factor γ	0.99
Activation function	ReLU
Others	
Number of parameters in piecewise linear function N	20

Table 1: MNSE Hyper-parameters sheet

Baseline Hyper-parameter: TD3 and PPO are developed utilizing the RL Baselines3 Zoo (Raffin, 2020), employing the tuned hyperparameters provided by this framework. NovelD (Zhang et al., 2021) and MinRed (Baram et al., 2021) are constructed based on the SAC algorithm within the RL Baselines3 Zoo, with the trade-off coefficient for the additional reward being searched within the range of $5e-3$ to $5e-1$.

Inverse Dynamics: We use discrete multinomial distributions to predict the actions. These multinomial distributions output an M -dimensional vector, where each dimension corresponds to the probability that the predicted action lies within the interval $[\frac{j-1}{M}, \frac{j}{M}]$, for $j \in [1, M]$. We set $M = 20$ in all experimental implementations.

Nonlinear Actuator in Figure 3: We consider two types of input nonlinearities: saturation and deadzone. The specific types of input nonlinearities and the corresponding effective action proportions (EAP) in each environment during our experiments (Figure 3) are presented in Table E.

Environment	Types of Nonlinearities	EAP
MuJoCo Ant	Saturation	40%
MuJoCo HalfCheetah	Saturation	40%
MuJoCo Hopper	Deadzone	20%
MuJoCo Walker2d	Deadzone	20%
MuJoCo Humanoid	Deadzone	20%
Metaworld PickPlace	Saturation	40%
Metaworld Hammer	Saturation	40%
Metaworld DoorOpen	Deadzone	40%
Metaworld Assembly	Deadzone	40%

Table 2: Nonlinear Actuator in Figure 3