

---

# Distributional Pareto-Optimal Multi-Objective Reinforcement Learning

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Multi-objective reinforcement learning (MORL) has been proposed to learn control  
2 policies over multiple competing objectives with each possible preference over  
3 returns. However, current MORL algorithms fail to account for distributional  
4 preferences over the multi-variate returns, which are particularly important in real-  
5 world scenarios such as autonomous driving. To address this issue, we extend the  
6 concept of Pareto-optimality in MORL into distributional Pareto-optimality, which  
7 captures the optimality of return distributions, rather than the expectations. Our  
8 proposed method, called Distributional Pareto-Optimal Multi-Objective Reinforce-  
9 ment Learning (DPMORL), is capable of learning distributional Pareto-optimal  
10 policies that balance multiple objectives while considering the return uncertainty.  
11 We evaluated our method on several benchmark problems and demonstrated its  
12 effectiveness in discovering distributional Pareto-optimal policies and satisfying  
13 diverse distributional preferences compared to existing MORL methods.

## 14 1 Introduction

15 Multi-Objective Reinforcement Learning (MORL) has recently received extensive attention in the  
16 artificial intelligence realm due to its adeptness at managing intricate decision-making issues with  
17 multiple conflicting objectives [1, 2, 3, 4]. In many multi-objective tasks, the relative preferences of  
18 users over different objectives are typically indeterminate a priori. Consequently, MORL’s primary  
19 aim is to learn a variety of optimal policies under different preferences to approximate the Pareto  
20 frontier of optimal solutions. It has been demonstrated that MORL can significantly reduce the  
21 reliance on scalar reward design for objective combinations and dynamically adapt to the varying  
22 preferences of different users.

23 However, numerous real-world situations involve not only unknown relative preferences across multi-  
24 ple objectives, but also uncertain preferences in return distributions. These may include preference in  
25 risk [5, 6], safety conditions [7], and non-linear user’s utility [8]. Consider, for instance, autonomous  
26 driving scenarios where agents need to strike a balance between safety and efficiency objectives.  
27 Different users might possess varying levels of risk tolerance. Some may demand high safety perfor-  
28 mance, tolerating lower expected efficiency, while others may seek a more balanced performance  
29 between safety and efficiency. Current MORL methods, by focusing exclusively on expected values  
30 with linear preferences, may not adequately capture multi-objective risk-sensitive preferences, hence  
31 being unable to deliver diverse policies catering to users with varied risk preferences.

32 In this work, we broaden the concept of Pareto-optimality in MORL to encompass distributional  
33 Pareto-optimality, which prioritizes the optimality of return distributions over mere expectations.  
34 From a theoretical perspective, we define Distributional Pareto-Optimal (DPO) policies, which capture  
35 the optimality of multivariate distributions through stochastic dominance [9, 10]. Distributional

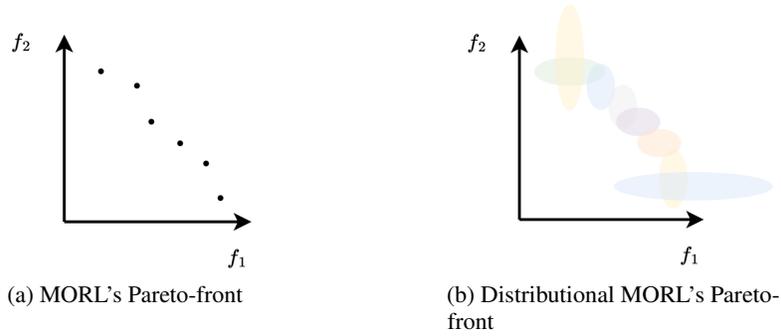


Figure 1: The comparison of learning targets between the traditional MORL tasks and distributional MORL tasks, in which  $f_1$  and  $f_2$  are two (conflicting) objectives.

36 Pareto-Optimality delineates the set of policies with optimal return distributions, which we formally  
 37 establish as an extension of Pareto-optimality in MORL.

38 On the practical side, we propose a novel method named Distributional Pareto-Optimal Multi-  
 39 Objective Reinforcement Learning (DPMORL). It aims to learn a set of DPO policies. We demonstrate  
 40 that a policy achieving the highest expected utility for a given utility function is a DPO policy. We  
 41 suggest an iterative process to learn diverse non-linear utility functions and optimize policies under  
 42 them [11]. Experimental outcomes on several benchmark problems attest to DPMORL’s effectiveness  
 43 in optimizing policies that meet preferences on multivariate return distributions. We posit that our  
 44 proposed method significantly addresses the challenge of managing multiple conflicting objectives  
 45 with unknown preferences on multivariate return distributions in complex decision-making situations.

46 Our main contributions are listed as follows:

- 47 1. We introduce the concept of Distributional Pareto-Optimal (DPO) policies. This expands  
 48 the notion of Pareto optimality in MORL to include preferences over the entire distribution  
 49 of returns, not just their expected values. This enables agents to express more nuanced  
 50 preferences over policies and align their decisions more closely with their actual objectives.
- 51 2. We propose DPMORL, a new algorithm for learning DPO policies under non-linear utility  
 52 functions. This algorithm accommodates a broad range of distributional preferences, thus  
 53 offering a more flexible and expressive approach to MORL.
- 54 3. We execute extensive experiments on various MORL tasks, demonstrating the effectiveness  
 55 of our approach in learning DPO policies. Our findings show that our algorithm consistently  
 56 surpasses existing MORL methods in terms of optimizing policies for multivariate expected  
 57 and distributional preferences, underscoring its practical benefits.

## 58 2 Related Work

59 **Multi-Objective Reinforcement Learning.** MORL has emerged as a vibrant research area in the AI  
 60 community owing to its adeptness in managing intricate decision-making scenarios involving multiple  
 61 conflicting objectives [1, 2]. A plethora of MORL algorithms have been put forth in the literature.  
 62 For instance, [12] proposed the utilization of Generalized Policy Improvement (GPI) to infer a set of  
 63 policies, employing Optimistic Linear Support (OLS) for dynamic reward weight exploration [13].  
 64 This GPI-based learning process was further enhanced by [3] through the incorporation of a world  
 65 model to augment sample efficiency. [4] proposed an evolutionary approach to learn policies with  
 66 diverse weights. However, our proposed approach diverges from these by leveraging utility functions  
 67 to guide policy learning [14, 15]. The utility-based paradigm, inspired by axiomatic approaches,  
 68 accentuates the user’s utility in decision-making problems, capitalizing on known information about  
 69 the user’s utility function and permissible policy types. Such methods scalarize the multi-dimensional  
 70 objectives into a single reward function, enabling traditional RL algorithms to infer desirable policies,  
 71 while also respecting axiomatic principles where necessary [2]. Several studies have utilized non-  
 72 linear utility functions to guide policy learning, catering to more intricate preferences compared to  
 73 linear ones [8, 16]. However, conventional MORL methodologies concentrate on optimizing the  
 74 expected returns, neglecting the distributional characteristics of the returns. In contrast, our work

75 extends the utility-based MORL framework by learning a set of utility functions with distributional  
76 attributes to guide policy learning.

77 **Distributional Reinforcement Learning.** Distributional RL extends traditional RL by modeling  
78 the entire distribution of returns, rather than just their expected values [17]. This approach has been  
79 shown to improve both learning efficiency and policy quality in a variety of single-objective RL  
80 tasks [18, 19]. Key algorithms in this area include Categorical DQN (C51) [17], Quantile Regression  
81 DQN (QR-DQN) [18], and Distributional MPO [19]. Despite the success of these distributional RL  
82 algorithms in single-objective settings, their direct extension to MORL has been limited due to the  
83 added complexity of handling multiple conflicting objectives and expressing preferences over the  
84 distribution of returns. On the other hand, the insights and techniques developed in distributional  
85 RL can provide a valuable foundation for incorporating distributional preferences into MORL, as  
86 demonstrated by our proposed approach.

87 **Risk-Sensitive and Safe Reinforcement Learning.** Risk-sensitive and safe RL are special cases  
88 of MORL that accentuate specific facets of reward distributions [2]. Risk-sensitive RL primarily  
89 considers reward variance, striving to optimize policies that negotiate the trade-off between expected  
90 return and risk [20, 21, 22]. Conversely, safe RL prioritizes constraints on the agent’s behavior,  
91 ensuring adherence to specified safety criteria throughout the learning process [23, 24, 25]. Some  
92 studies have proposed the integration of constraints into the state space, constructing new Markov  
93 Decision Processes [7]. Others have explored the distributional aspects of rewards, investigating  
94 the implications of distributional RL on risk-sensitive and safety-oriented decision-making [18, 17].  
95 However, our proposed setting is more general in terms of objectives, as it considers a broader range of  
96 user preferences and captures the entire distribution of rewards, rather than focusing solely on specific  
97 aspects such as risk, variance, or constraints. By extending MORL to incorporate distributional  
98 properties, our approach enables the learning of distributional Pareto-optimal policies that cater to  
99 diverse user preferences and offer better decision-making in a wide range of real-world applications.

### 100 3 Preliminaries: Multi-Objective Reinforcement Learning

101 In Multi-Objective Reinforcement Learning (MORL), the agent interacts with an environment  
102 modeled as a Multi-Objective Markov Decision Process (MOMDP) with multi-dimensional reward  
103 functions. A Multi-Objective MDP is defined by a tuple  $(S, A, P, P_0, \mathbf{R}, \gamma, T)$ , where  $S$  is the state  
104 space,  $A$  is the action space,  $P$  is the state transition function,  $P_0(s)$  is the initial state distribution,  
105  $\mathbf{R} : S \times A \times S \rightarrow \mathbb{R}^K$  is a vectored reward function and  $K$  is the number of objectives,  $\gamma$  is the  
106 discount factor,  $T$  is the total timesteps<sup>1</sup>. The goal of the agent is to learn a set of Pareto-optimal  
107 policies, which represent the best possible trade-offs among the different objectives.

108 One popular methodology for MORL problems is the utility-based method, which combines the  
109 multi-dimensional reward functions into a single scalar reward function using a weighted sum or  
110 another aggregation method [26]. The intuition is to map the agent’s preferences over different  
111 objectives to a scalar value for the agent training. Given a weight vector  $\mathbf{w} = (w_1, w_2, \dots, w_K)$ ,  
112 with  $w_i$  representing the importance of the  $i$ -th objective, the scalarized reward function is defined  
113 as  $r_{\text{scalar}}(s, a) = \sum_{i=1}^K w_i r_i(s, a)$ . The agent then solves the scalarized MDP by optimizing its  
114 policy to maximize the expected scalar reward, using standard reinforcement learning algorithms  
115 like Q-learning or policy gradient methods. This approach can be straightforward to implement and  
116 has been shown to be effective in various MORL settings under expected preferences [27]. However,  
117 such a linear combination of each dimension of the reward cannot deal with the preferences with  
118 distributional considerations.

### 119 4 Distributionally Pareto-Optimal Multi-Objective Reinforcement Learning

120 In this section, we introduce our proposed theoretical framework and algorithms for extending MORL  
121 to handle distributional preferences.

#### 122 4.1 Distributionally Pareto-Optimal Policies

123 In this work, we consider the reward as a  $K$ -dimensional vector, where each element represents the  
124 reward for a specific objective. Given the MOMDP and a policy  $\pi$ , we define the random variable of

---

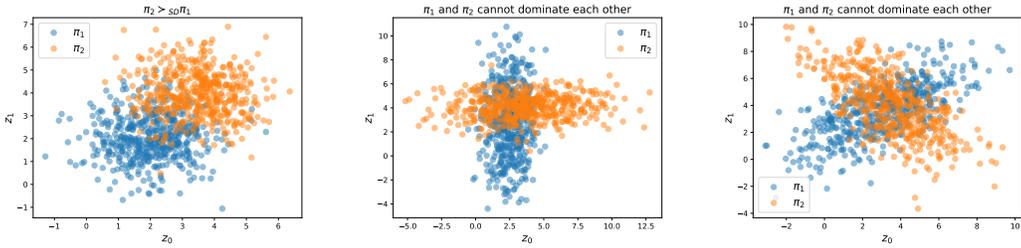
<sup>1</sup>We assume that current time  $t$  is a part of state to accommodate for the finite horizon setting.

125 multi-objective return as  $\mathbf{Z}(\pi) = \sum_{t=0}^T \gamma^t \mathbf{r}_t$ , where  $\mathbf{r}_t$  is the multi-dimensional reward at step  $t$ , and  
 126 the states  $s_0, s_1, \dots, s_T$  and actions  $a_0, a_1, \dots, a_T$  are sampled from the MOMDP and the policy  $\pi$   
 127 respectively. The return distribution of policy  $\pi$ , denoted as  $\boldsymbol{\mu}(\pi)$ , represents the joint distribution of  
 128 the returns  $\mathbf{Z}(\pi)$  following policy  $\pi$ . The utility function  $f$  is a non-decreasing function that maps a  
 129  $K$ -dimensional return into a scalar utility value, capturing the user’s distributional preferences over  
 130 the different objectives. The expected utility of policy  $\pi$  under the utility function  $f$ , represented as  
 131  $\mathbb{E}_{\mathbf{z} \sim \boldsymbol{\mu}(\pi)} f(\mathbf{z})$ , is the expected value of applying  $f$  to the return distribution  $\boldsymbol{\mu}(\pi)$ .

132 Our goal is to learn a set of policies that are distributionally Pareto-optimal, which means that their  
 133 return distributions of each policy cannot dominate that of another. To measure such a relationship,  
 134 here we first introduce the concept of stochastic dominance:

135 **Definition 1** (Stochastic Dominance for Multivariate Distribution). *A multivariate distribution*  
 136  $\boldsymbol{\mu}_1$  *dominates another distribution*  $\boldsymbol{\mu}_2$ , *denoted as*  $\boldsymbol{\mu}_1 \succ_{SD} \boldsymbol{\mu}_2$ , *if and only if*  $\boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$  *and for*  
 137 *any non-decreasing utility function*  $f : \mathbb{R}^K \rightarrow \mathbb{R}$ ,  $\boldsymbol{\mu}_1$  *has greater expected utility than*  $\boldsymbol{\mu}_2$ , *i.e.*  
 138  $\mathbb{E}_{\mathbf{z} \sim \boldsymbol{\mu}_1} f(\mathbf{z}) \geq \mathbb{E}_{\mathbf{z} \sim \boldsymbol{\mu}_2} f(\mathbf{z})$ .

139 **Definition 2** (Stochastic Dominance for Policies). *A policy*  $\pi_1$  *stochastically dominates another*  
 140 *policy*  $\pi_2$ , *denoted as*  $\pi_1 \succ_{SD} \pi_2$ , *if and only if*  $\boldsymbol{\mu}(\pi_1) \succ_{SD} \boldsymbol{\mu}(\pi_2)$ , *indicating*  $\pi_1$  *has greater*  
 141 *expected utility than*  $\pi_2$  *under any non-decreasing utility function*  $f$ .



(a) Stochastic Dominance Example 1 (b) Stochastic Dominance Example 2 (c) Stochastic Dominance Example 3

Figure 2: Three 2D examples of stochastic dominance given synthetic returns of two policies. (a)  $\pi_2 \succ_{SD} \pi_1$ , since any non-decreasing utility function satisfies  $\mathbb{E}f(\mathbf{Z}(\pi_2)) \geq \mathbb{E}f(\mathbf{Z}(\pi_1))$ . (b)  $\pi_1$  and  $\pi_2$  cannot dominate each other. For example, there exists  $f(\mathbf{z}) = \text{ReLU}(z_0 - 5.0)$ , such that  $\mathbb{E}f(\mathbf{Z}(\pi_2)) > \mathbb{E}f(\mathbf{Z}(\pi_1))$ , thus  $\pi_1$  cannot dominate  $\pi_2$ . Similarly,  $\pi_2$  cannot dominate  $\pi_1$ . Thus there is no dominant relationship between the two policies. (c)  $\pi_1$  and  $\pi_2$  cannot dominate each other, which is similar to (b).

142 The definition of stochastic dominance extends univariate first-order stochastic dominance [9] into  
 143 multivariate cases. Figure 2 provides some examples of stochastic dominance given 2D returns of two  
 144 policies. The definition of stochastic dominance in policies allows for comparing the optimality of  
 145 distributions between policies, which allows us to define the Distributionally Pareto-optimal policy:

146 **Definition 3** (Distributionally Pareto-Optimal Policy). *Formally,  $\pi_1$  is Distributionally Pareto-*  
 147 *Optimal Policy if there does not exist a policy  $\pi_2$  such that  $\boldsymbol{\mu}(\pi_2) \succ_{SD} \boldsymbol{\mu}(\pi_1)$ .*

148 In other words,  $\pi_1$  is considered Distributionally Pareto-Optimal (DPO) if it is not stochastically  
 149 dominated by any other policy. This implies that  $\pi_1$  offers the highest expected utility under some  
 150 non-decreasing utility function  $f$  and cannot be outperformed by any other policy in the problem.  
 151 DPO policies are essential in our framework as they represent the most desirable policies for users  
 152 with different distributional preferences. To find such a policy, we have the following theorem:

153 **Theorem 1.** *If a policy  $\pi$  has optimal expected utility under some non-decreasing utility func-*  
 154 *tion  $f$ , then it is a Distributionally Pareto-Optimal policy. Also, any Pareto-Optimal policy is a*  
 155 *Distributionally Pareto-Optimal policy.*

156 This theorem guarantees that a policy achieving the highest expected utility for a given utility function  
 157 will be a DPO policy, making it a suitable candidate for deployment in MORL problems. In the  
 158 next section, we base on the result of this theorem to find optimal policies for a diverse set of utility  
 159 functions, in order to learn the set of DPO policies. This result also shows that our definition of  
 160 Distributional Pareto-Optimal policies is an extension of Pareto-optimal policies, allowing for a more  
 161 diverse set of optimal policies to be captured.

162 We also formally prove that optimal policies for risk-sensitive and safe constraint objectives belong  
 163 to the set of Pareto-Optimal policies. This proves that DPO policies can successfully cover policies  
 164 with diverse distributional preferences. The detailed proof is provided by Theorem 2 in Appendix A.

## 165 4.2 Distributionally Pareto-Optimal Multi-Objective Reinforcement Learning

166 We now present our proposed algorithm, Distributionally Pareto-Optimal Multi-Objective Reinforce-  
 167 ment Learning (DPMORL). The main idea of DPMORL is to learn a set of non-linear utility functions  
 168 that can guide the agent to discover distributional Pareto-optimal policies. The algorithm proceeds in  
 169 two stages: (1) generating the utility functions and (2) training the policies.

### 170 4.2.1 Utility Function Generation with Diversity-based Objective

171 The first component of our algorithm focuses on generating a diverse set of plausible utility functions,  
 172 upon which we find the optimal policies to find a diverse set of optimal policies. This is essential  
 173 to ensure our method can accommodate various distributional preferences and adapt to different  
 174 problem settings. To achieve this, we propose (1) Non-decreasing Neural Network for parameterizing  
 175 a diverse set of non-linear utility functions (2) an objective function of minimum distance which  
 176 encourages generating a diverse set of utility functions.

177 **Non-decreasing Neural Network.** We employ non-decreasing neural network to parameterize  
 178 the utility function. The use of a neural network allows us to represent complex, non-linear, and  
 179 arbitrary continuous utility functions, while the non-decreasing constraint ensures that the utility  
 180 function satisfies the desired properties for multi-objective problems. We ensure the non-decreasing  
 181 property in neural networks by constraining the weight matrix to be non-negative and the activation  
 182 function to be non-decreasing following existing work in convex neural networks [28] and QMIX  
 183 [29], which can approximate any multivariate non-decreasing function with arbitrary small errors  
 184 [30]. The implementation details of the Non-decreasing Neural Network are provided in Appendix B.

185 **Diversity-based Objective Function.** We propose an objective function based on diversity for  
 186 learning a diverse set of plausible utility functions. Specifically, we define  $f_{\theta_1}, f_{\theta_2}, \dots, f_{\theta_M}$  to be  
 187 the set of candidate utility functions, parameterized by  $\theta_1, \dots, \theta_M$ . For any given utility function  
 188  $f_{\theta_i}$ , the learning objective is defined as

$$J^{\text{value}}(\theta_i) = \min_{j \neq i} \mathbb{E}_{\mathbf{z} \sim \mathcal{U}([0,1]^K)} [f_{\theta_i}(\mathbf{z}) - f_{\theta_j}(\mathbf{z})]^2 \quad (1)$$

$$J^{\text{grad}}(\theta_i) = \min_{j \neq i} \mathbb{E}_{\mathbf{z}_1, \mathbf{z}_2 \sim \mathcal{U}([0,1]^K)} \left[ \frac{f_{\theta_i}(\mathbf{z}_2) - f_{\theta_i}(\mathbf{z}_1)}{\|\mathbf{z}_2 - \mathbf{z}_1\|} - \frac{f_{\theta_j}(\mathbf{z}_2) - f_{\theta_j}(\mathbf{z}_1)}{\|\mathbf{z}_2 - \mathbf{z}_1\|} \right]^2 \quad (2)$$

$$J(\theta_i) = \alpha J^{\text{value}}(\theta_i) + (1 - \alpha) J^{\text{grad}}(\theta_i) \quad (3)$$

189 Optimizing this objective function can encourage diversity among the generated utility functions  
 190 within range  $[0, 1]^K$  in both value and derivative, leading to more comprehensive coverage of potential  
 191 user preferences. We maximize the objective function by gradient descent on non-decreasing neural  
 192 networks to generate a set of  $N$  utility functions.  
 193

### 194 4.2.2 Optimizing Policy with Utility-based Reinforcement Learning

195 Once we have generated a diverse set of utility functions, the second component of our algorithm  
 196 focuses on optimizing policies to maximize the expected utility. This process, known as utility-based  
 197 RL, leverages the generated utility functions to guide the optimization of policies. By focusing on the  
 198 expected utility, our method can efficiently balance the trade-offs between multiple objectives and  
 199 distributions, ultimately yielding policies that are more likely to align with user preferences.

200 We show that the following utility-based reinforcement learning algorithm can effectively optimize  
 201 the policy with respect to a given utility function.

---

**Algorithm 1** Utility-based Reinforcement Learning

---

**Input:** policy  $\pi$ , an environment  $\mathcal{M} = (S, A, P, P_0, \mathbf{R}, \gamma, T)$ , utility function  $f$

**Output:** new policy  $\pi'$

- 1: Augment state space with  $\tilde{\mathcal{S}} = \mathcal{S} \times \mathcal{Z}$ , where  $\mathcal{Z}$  is the space of cumulative multivariate returns.
  - 2: Let transition  $\tilde{P}_0(\cdot)$  and  $\tilde{P}(\cdot|(s_t, \mathbf{z}_t), a_t)$  with  $s_0 \sim P(s_0)$ ,  $\mathbf{z}_0 = 0$ , and  $s_{t+1} \sim P(\cdot|s_t, a_t)$ ,  $\mathbf{z}_{t+1} = \mathbf{z}_t + \gamma^t \mathbf{r}_t$ .
  - 3: Let scalar reward function  $R$  as  $R((s_t, \mathbf{z}_t), a_t, (s_{t+1}, \mathbf{z}_{t+1})) = \gamma^{-t}[f(\mathbf{z}_{t+1}) - f(\mathbf{z}_t)]$ .
  - 4: Optimize policy  $\pi$  under environment  $\tilde{\mathcal{M}} = (\tilde{\mathcal{S}}, A, \tilde{P}, \tilde{P}_0, R, \gamma, T)$  under off-the-shelf RL algorithm (such as PPO or SAC) to  $\pi'$ .
- 

202 Briefly, Algorithm 1 augments the state space with the cumulative multi-objective returns, and  
203 transforms the multi-dimensional rewards into a scalar reward by the difference in the utility function  
204  $f$ . The following result shows that the new scalar-reward environment  $\tilde{\mathcal{M}}$  generated by Algorithm 1  
205 has the same optimal policy as the optimal policy under utility function  $f$ :

206 **Theorem 2.** *The optimal policy  $\pi^*$  under environment  $\tilde{\mathcal{M}} = (\tilde{\mathcal{S}}, A, \tilde{P}, \tilde{P}_0, R, \gamma, T)$ , with scalar  
207 reward function*

$$R((s_t, \mathbf{z}_t), a_t, (s_{t+1}, \mathbf{z}_{t+1})) = \gamma^{-t}[f(\mathbf{z}_{t+1}) - f(\mathbf{z}_t)]$$

208 *is the optimal policy in the utility function  $f$ , i.e.*

$$\mathbb{E}_{\mathbf{z} \sim \mu(\pi^*)} [f(\mathbf{z})] = \max_{\pi} \mathbb{E}_{\mathbf{z} \sim \mu(\pi)} [f(\mathbf{z})].$$

209 An advantage of Algorithm 1 is that we can directly utilize off-the-shelf RL algorithms, such as  
210 PPO [31] and SAC [32] without any modification, which makes the algorithm easy to implement  
211 using widespread existing implementations of online RL algorithms.

212 It is also important to note that our algorithm simplifies to optimizing the weighted sum of rewards  
213 in MORL when the utility function is linear. This implies that our method is a generalization of  
214 the linear utility function MORL approaches by accommodating a wide range of non-linear utility  
215 functions. This flexibility makes our algorithm particularly suited for problems where the user’s  
216 preferences may not be adequately captured by a linear utility function.

### 217 4.2.3 Iterative Generation of Utility Function and Policies Optimization

218 The range of multi-objective returns may not be known prior to policy optimization, which could  
219 make it challenging to generate relevant utility functions for a specific task. To mitigate this issue, we  
220 introduce an iterative process that alternates between generating currently plausible utility functions  
221 and optimizing policies based on these functions. During each iteration, we gather return samples  
222 throughout the optimization process and update the utility function. This update is particularly aimed  
223 at enhancing the diversity of the utility function with respect to the returns observed during the  
224 optimization process. A detailed outline of the algorithm can be found in Appendix C.

225 In the next experimental section, DPMORL only undergoes a single iteration: we initially generate a  
226 set of  $N$  utility functions as per the methodology detailed in Section 4.2.1. Subsequently, it optimizes  
227 a set of  $N$  policies using these generated utility functions, as outlined in Section 4.2.2. In Appendix  
228 C, we provide a case study demonstrating the application of DPMORL in an iterative training context.

## 229 5 Experiments

230 In this section, we conducted several experiments under the setting of MORL. Through the experi-  
231 ments, we want to investigate the following questions:

- 232 1. Can Utility-based RL effectively learn policies with diverse distributional preferences?
- 233 2. Can DPMORL generate a set of diverse non-linear utility functions?
- 234 3. Can DPMORL obtain promising performance compared with state-of-the-art MORL meth-  
235 ods in view of expected and distributional preferences?

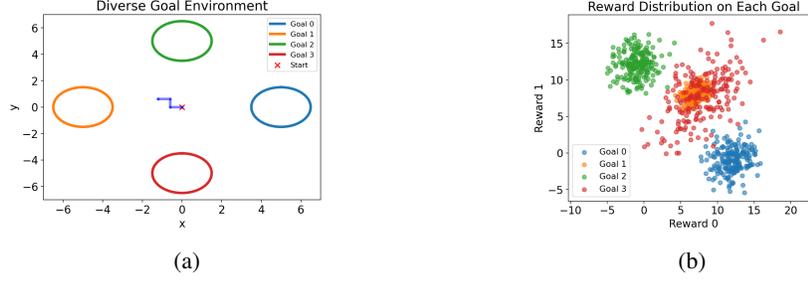


Figure 3: (a) The map of the DiverseGoal environment. (b) Reward distribution within each goal position.

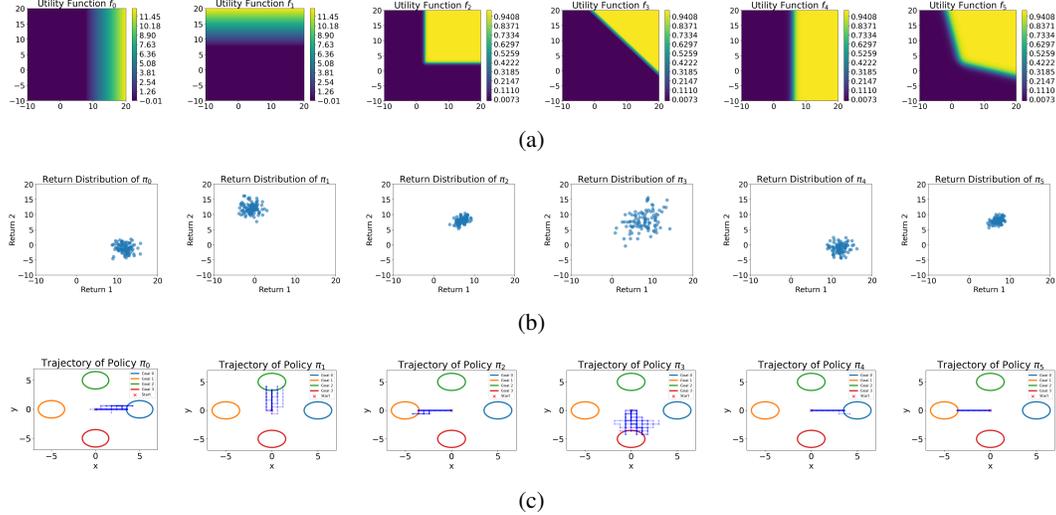


Figure 4: (a) Utility function used for training policies. (b) Return distribution of policy learned by Algorithm 1 to maximize expected utility for each utility function. (c) The trajectory of policies learned by maximizing expected utility for each utility function.

## 236 5.1 Case Study on DiverseGoal Environment

237 To answer the first question, we train policies with a diverse set of utility functions on DiverseGoal,  
 238 an environment that provides a diverse set of reward distributions.

239 DiverseGoal is a MORL environment in Figure 3 with multiple goals that the agent can take multiple  
 240 steps to reach, where each goal has its unique reward distributions. Upon reaching a particular goal,  
 241 the agent secures a 2D reward. This reward is sampled from the specific distribution associated with  
 242 that goal, as illustrated in Figure 3b. Conversely, if the agent reaches the boundary of the map, it  
 243 incurs a negative 2-dimensional reward. The environment requires the agent to navigate the trade-offs  
 244 between various reward distributions, underscoring the complexity and nuance inherent to the task.

245 Here, we focus on showing the effectiveness of the Utility-based RL algorithm (Algorithm 1). We  
 246 select six different non-linear functions  $f_0, f_1, \dots, f_5$ , each in favor of distributions at one goal. We  
 247 train policy  $\pi_i$  with utility function  $f_i$  with Algorithm 1 where  $i = 0, 1, \dots, 5$ , and show the policy's  
 248 trajectory and return samples. Ideally, the learned policy should reach the goal that yields the highest  
 249 expected utility under the utility function. The results are illustrated in Figure 3. Under different type  
 250 of utility function and return distributions, the utility-based RL algorithm is able to find the optimal  
 251 policy with highest expected utility, which shows the effectiveness of the utility-based RL algorithm.

## 252 5.2 Main Experiment

253 In this section, we illustrate the performance of DPMORL on five environments based on MO-  
 254 Gymnasium [33] to answer the latter two questions.

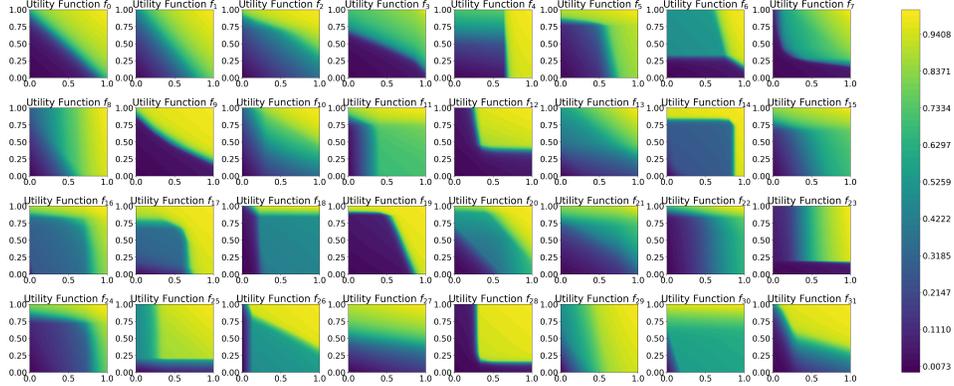


Figure 5: Illustration of 2D utility functions learned by our methods in Section 4.2.1.

255 **Environments.** We conducted experiments across five environments based on MO-Gymnasium [33]  
 256 to evaluate the performance of our proposed method, DPMORL. These environments represent a  
 257 diverse range of tasks, from simple toy problems to more complex continuous control tasks, and  
 258 cover various aspects of multi-objective reinforcement learning:

- 259 • DeepSeaTreasure: A classic MORL benchmark that requires exploration in a gridworld to  
 260 find treasures with different values and depths.
- 261 • FruitTree: A multi-objective variant of the classic gridworld problem, where an agent has to  
 262 collect different types of fruits with varying rewards and penalties.
- 263 • HalCheetah, Hopper, MountainCar: Three continuous control tasks that require controlling  
 264 different agents for task solving and minimizing energy usage.

265 More details about the environments are gathered in Appendix C.

266 **Baselines.** We compare DPMORL with four state-of-the-art baselines in the context of distributional  
 267 preferences of MORL: Optimistic Linear Support (OLS) [34, 12]; Prediction-Guided Multi-Objective  
 268 Reinforcement Learning (PGMORL) [4]; Generalized Policy Improvement with Linear Support  
 269 (GPI-LS) [3]; Generalized Policy Improvement with Prioritized Dyna (GPI-PD) [3].

270 **Training Details.** For all methods, including DPMORL and the baselines, each policy was trained  
 271 for  $1 \times 10^7$  steps. We learn a set of  $N = 20$  policies for DPMORL and all of the baselines to ensure  
 272 a fair comparison. Finally, we use the learned  $N = 20$  policies in each method for evaluations.

273 **Implementation Details.** We use PPO algorithms implemented in Stable Baselines 3 [35] in  
 274 Algorithm 1. We use a normalization technique by linearly mapping the return into scale  $[0, 1]^K$   
 275 without modifying the optimal policies. Detailed implementations are provided in Appendix B.

276 **Evaluation Metrics.** To thoroughly evaluate the performance of DPMORL and compare it with  
 277 the baseline methods, we employed four distinct metrics. These comprise two existing MORL  
 278 metrics, HyperVolume and Expected Utility, in addition to two novel metrics developed specifically  
 279 for assessing distributional preference, i.e., Constraint Satisfaction and Variance Objective. The latter  
 280 two are designed to underscore the optimality of multivariate distributions associated with the learned  
 281 policies. In terms of Constraint Satisfaction, we randomly generate  $M = 100$  constraints for the  
 282 policy set produced. The Constraint Satisfaction metric is then computed by considering the highest  
 283 probability within the policy set that satisfies each individual constraint. The Variance Objective  
 284 metric, on the other hand, involves generating  $M = 100$  random linear weights. These weights  
 285 are applied to both the expected returns and the standard deviations of returns in each dimension.  
 286 This objective encourages attaining greater expected returns while simultaneously reducing variance,  
 287 thereby catering to dynamic preferences. Further details about the implementation of these evaluation  
 288 metrics are provided in Appendix B.

### 289 5.3 Results

290 Since PGMORL works on continuous action spaces, here we omit the results of PGMORL on  
 291 environments with discrete action space, DeepSeaTreasure and FruitTree. More results are gathered  
 292 in Appendix C.

Table 1: Experimental results of each method on each standard MORL metric on all five environments. “EU” stands for Expected Utility, and “HV” stands for Hypervolume.

Environment Metric	Hopper		HalfCheetah		MountainCar		DeepSeaTreasure		FruitTree	
	EU	HV	EU	HV	EU	HV	EU	HV	EU	HV
GPI-PD	2374.44	5033458.34	412.62	1083227.57	-55.44	7367.59	5.93	9.75	<b>6.98</b>	0.14
GPI-LS	1398.00	1446705.10	98.31	2839051.60	-37.37	8052.58	5.04	9.36	3.84	1.49
OLS	175.07	3298.13	-580.38	1420270.13	-470.00	2.45	4.64	<b>10.28</b>	6.27	7.49
PGMORL	300.19	32621.25	111.30	383681.25	-429.48	94.45	-	-	-	-
DPMORL (Ours)	<b>3492.93</b>	<b>12154967.99</b>	<b>1189.68</b>	<b>8593769.62</b>	<b>-29.89</b>	<b>8663.80</b>	<b>6.70</b>	8.68	6.89	<b>16.39</b>

Table 2: Experimental results of each method on Distributional metric on all five environments. “Constraints” stands for Constraints Satisfaction, and “Var” stands for Variance Objective.

Environment Metric	Hopper		HalfCheetah		MountainCar		DeepSeaTreasure		FruitTree	
	Constraints	Var	Constraint	Var	Constraint	Var	Constraint	Var	Constraint	Var
GPI-PD	0.47	979.74	0.64	83.49	<b>1.00</b>	-31.15	0.85	<b>2.59</b>	0.65	<b>3.42</b>
GPI-LS	0.25	607.47	0.60	51.67	<b>1.00</b>	-22.73	0.85	1.99	0.33	1.35
OLS	0.05	75.48	0.43	-311.27	0.05	-244.21	0.80	1.86	0.50	2.98
PGMORL	0.05	98.47	0.50	45.48	0.11	-249.69	-	-	-	-
DPMORL (Ours)	<b>0.76</b>	<b>1645.89</b>	<b>0.82</b>	<b>431.26</b>	<b>1.00</b>	<b>-16.32</b>	<b>0.90</b>	2.54	<b>0.67</b>	3.21

293 **Generation of Utility Functions.** In accordance with the methodology detailed in Section 4.2.1,  
 294 we employ non-decreasing neural networks in conjunction with diversity-focused objectives to  
 295 generate a diverse assortment of utility functions. The generated functions are visually represented in  
 296 Figure 5. The outcomes clearly indicate that optimizing diversity-based objective functions allows  
 297 for generating a broad range of non-linear utility functions, thereby encompassing an expansive array  
 298 of preferences with respect to returns. Subsequently, we utilize the first  $N = 20$  utility functions  
 299 depicted in Figure 5 to train an equivalent number of policies under DPMORL.

300 **Standard MORL Metrics.** The results under standard MORL metrics are shown in Table 1. Focusing  
 301 on the Expected Utility, the performance of DPMORL is the highest in the Hopper, HalfCheetah,  
 302 MountainCar, and DeepSeaTreasure environments (4/5), indicating that our method outperforms  
 303 others in terms of expected utility. While for FruitTree environment, DPMORL also obtains consistent  
 304 performance with the best contender, GPI-LS. In terms of the HyperVolume, DPMORL also performs  
 305 the best in the Hopper, HalfCheetah, MountainCar, and FruitTree environments (4/5). The results show  
 306 that DPMORL can yield better performance across both metrics and most environments (Hopper,  
 307 HalfCheetah, MountainCar, and DeepSeaTreasure), meanwhile delivering robust and consistent  
 308 results on the other one (FruitTree).

309 **Distributional MORL Metrics.** The results under distributional metrics are shown in Table 2.  
 310 DPMORL outperforms other methods in most environments on both Constraint Satisfaction (5/5)  
 311 and Variance Objective (3/5), indicating its strong ability to handle the distributional multi-objective  
 312 reinforcement learning problem. For the rest environments, DPMORL also obtained comparable  
 313 results compared with the best contender, GPI-LS. The results show that the policies learned by  
 314 DPMORL have a higher probability of satisfying randomly generated constraints, and can better  
 315 balance the trade-off between expectations and variances.

316 **Return Distributions of Learned Policies by DPMORL.** We provide visualizations of the multi-  
 317 dimensional return distributions of each policy learned by DPMORL on all five environments. On  
 318 Hopper, HalfCheetah, DeepSeaTreasure and FruitTree, DPMORL learns a diverse set of policies  
 319 with different distributional properties in returns. We also visualize the learning process of different  
 320 policies. We provide these results in Appendix C.

## 321 6 Conclusion

322 In this work, we initialized the study of distributional MORL, specifically when preferences over  
 323 different objectives and their return distributions are uncertain. We introduced the concept of  
 324 Distributional Pareto-Optimal (DPO) policies with rigorous theoretical analysis, which extend the  
 325 notion of Pareto optimality in MORL to include preferences over the entire distribution of returns, not  
 326 just their expected values. To obtain such desirable policies, we proposed a new algorithm, DPMORL,  
 327 designed to learn DPO policies with non-linear utility functions. DPMORL allows for expressing  
 328 a wide range of distributional preferences, providing a flexible and expressive approach to MORL.  
 329 Experiment results showed that DPMORL consistently outperformed existing MORL methods in  
 330 terms of optimizing policies for multivariate expected and distributional preferences.

331 **References**

- 332 [1] Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of multi-  
333 objective sequential decision-making. *Journal of Artificial Intelligence Research*, 48:67–113,  
334 2013.
- 335 [2] Conor F. Hayes, Roxana Radulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane,  
336 Mathieu Reymond, Timothy Verstraeten, Luisa M. Zintgraf, Richard Dazeley, Fredrik Heintz,  
337 Enda Howley, Athirai A. Irissappane, Patrick Mannion, Ann Nowé, Gabriel de Oliveira Ramos,  
338 Marcello Restelli, Peter Vamplew, and Diederik M. Roijers. A practical guide to multi-objective  
339 reinforcement learning and planning. *Auton. Agents Multi Agent Syst.*, 36(1):26, 2022.
- 340 [3] Lucas Nunes Alegre, Ana L. C. Bazzan, Diederik M. Roijers, Ann Nowé, and Bruno C. da Silva.  
341 Sample-efficient multi-objective learning via generalized policy improvement prioritization.  
342 *CoRR*, abs/2301.07784, 2023.
- 343 [4] Jie Xu, Yunsheng Tian, Pingchuan Ma, Daniela Rus, Shinjiro Sueda, and Wojciech Matusik.  
344 Prediction-guided multi-objective reinforcement learning for continuous robot control. In  
345 *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18*  
346 *July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages  
347 10607–10616. PMLR, 2020.
- 348 [5] Anirban Majumdar and Claire J Tomlin. Risk-sensitive sequential decision making: Uncertainty-  
349 based performance metrics and risk-sensitive control. *Foundations and Trends® in Systems and*  
350 *Control*, 4(4):265–365, 2017.
- 351 [6] Shiao Hong Lim and Ilyas Malik. Distributional reinforcement learning for risk-sensitive  
352 policies. In *Advances in Neural Information Processing Systems*, 2022.
- 353 [7] Aivar Sootla, Alexander I Cowen-Rivers, Taher Jafferjee, Ziyang Wang, David H Mguni, Jun  
354 Wang, and Haitham Ammar. Sauté rl: Almost surely safe reinforcement learning using state  
355 augmentation. In *International Conference on Machine Learning*, pages 20423–20443. PMLR,  
356 2022.
- 357 [8] Mathieu Reymond, Conor F Hayes, Denis Steckelmacher, Diederik M Roijers, and Ann Nowé.  
358 Actor-critic multi-objective reinforcement learning for non-linear utility functions. *Autonomous*  
359 *Agents and Multi-Agent Systems*, 37(2):23, 2023.
- 360 [9] Haim Levy. Stochastic dominance and expected utility: Survey and analysis. *Management*  
361 *science*, 38(4):555–593, 1992.
- 362 [10] Moshe Shaked and J George Shanthikumar. *Stochastic orders*. Springer Science Business  
363 Media, 2007.
- 364 [11] Patrick Mannion, Jim Duggan, and Enda Howley. Dynamic treatment regimes in mobile health.  
365 In *2016 International Joint Conference on Neural Networks (IJCNN)*, pages 526–533. IEEE,  
366 2016.
- 367 [12] Lucas Nunes Alegre, Ana L. C. Bazzan, and Bruno C. da Silva. Optimistic linear support  
368 and successor features as a basis for optimal policy transfer. In Kamalika Chaudhuri, Stefanie  
369 Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato, editors, *International*  
370 *Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*,  
371 volume 162 of *Proceedings of Machine Learning Research*, pages 394–413. PMLR, 2022.
- 372 [13] Axel Abels, Diederik M. Roijers, Tom Lenaerts, Ann Nowé, and Denis Steckelmacher. Dynamic  
373 weights in multi-objective deep reinforcement learning. In Kamalika Chaudhuri and Ruslan  
374 Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning,*  
375 *ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of*  
376 *Machine Learning Research*, pages 11–20. PMLR, 2019.
- 377 [14] Matthieu Geist, Julien Pérolat, Mathieu Laurière, Romuald Elie, Sarah Perrin, Olivier Bachem,  
378 Rémi Munos, and Olivier Pietquin. Concave utility reinforcement learning: The mean-field game  
379 viewpoint. In Piotr Faliszewski, Viviana Mascardi, Catherine Pelachaud, and Matthew E. Taylor,

- 380 editors, *21st International Conference on Autonomous Agents and Multiagent Systems, AAMAS*  
381 *2022, Auckland, New Zealand, May 9-13, 2022*, pages 489–497. International Foundation for  
382 Autonomous Agents and Multiagent Systems (IFAAMAS), 2022.
- 383 [15] Mridul Agarwal, Vaneet Aggarwal, and Tian Lan. Multi-objective reinforcement learning  
384 with non-linear scalarization. In Piotr Faliszewski, Viviana Mascardi, Catherine Pelachaud,  
385 and Matthew E. Taylor, editors, *21st International Conference on Autonomous Agents and*  
386 *Multiagent Systems, AAMAS 2022, Auckland, New Zealand, May 9-13, 2022*, pages 9–17.  
387 International Foundation for Autonomous Agents and Multiagent Systems (IFAAMAS), 2022.
- 388 [16] Mridul Agarwal, Vaneet Aggarwal, and Tian Lan. Multi-objective reinforcement learning with  
389 non-linear scalarization. In *Proceedings of the 21st International Conference on Autonomous*  
390 *Agents and Multiagent Systems*, pages 9–17, 2022.
- 391 [17] Marc G Bellemare, Will Dabney, and R’emi Munos. A distributional perspective on reinforce-  
392 ment learning. In *Proceedings of the 34th International Conference on Machine Learning-*  
393 *Volume 70*, pages 449–458. JMLR. org, 2017.
- 394 [18] Will Dabney, Mark Rowland, Marc G Bellemare, and Remi Munos. Distributional reinforcement  
395 learning with quantile regression. In *Thirty-Second AAAI Conference on Artificial Intelligence*,  
396 2018.
- 397 [19] Gabriel Barth-Maron, Matthew W. Hoffman, David Budden, Will Dabney, Dan Horgan, Dhruva  
398 TB, Alistair Muldal, Nicolas Heess, and Timothy P. Lillicrap. Distributed distributional deter-  
399 ministic policy gradients. In *6th International Conference on Learning Representations,*  
400 *ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.*  
401 OpenReview.net, 2018.
- 402 [20] Vivek S Borkar. Risk-sensitive reinforcement learning. *SIAM Journal on Control and Optimiza-*  
403 *tion*, 40(3):818–834, 2002.
- 404 [21] Aviv Tamar, Yonathan Glassner, and Shie Mannor. Optimality and approximation in online and  
405 risk-averse markov decision processes. *The Journal of Machine Learning Research*, 16(1):2553–  
406 2592, 2015.
- 407 [22] L A Prashanth and Shalabh Bhatnagar. Risk-sensitive reinforcement learning. *arXiv preprint*  
408 *arXiv:1604.00147*, 2016.
- 409 [23] Javier Garc’ia and Fernando Fern’andez. Safe exploration in markov decision processes. *arXiv*  
410 *preprint arXiv:1205.1111*, 2012.
- 411 [24] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization.  
412 *arXiv preprint arXiv:1705.10528*, 2017.
- 413 [25] Yinlam Chow, Ofir Nachum, Edgar A Duenez-Guzman, and Mohammad Ghavamzadeh.  
414 Lyapunov-based safe policy optimization for continuous control. *arXiv preprint*  
415 *arXiv:1805.08711*, 2018.
- 416 [26] Zsolt Gabor and Zoltan Koles. Multi-objective reinforcement learning: A comprehensive  
417 overview. *Proceedings of the IEEE-INNS-ENNS International Joint Conference on Neural*  
418 *Networks*, 3:2030–2035, 1998.
- 419 [27] Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Multi-objective deep reinforcement  
420 learning. *arXiv preprint arXiv:1610.02707*, 2016.
- 421 [28] Brandon Amos, Lei Xu, and J Zico Kolter. Input convex neural networks. In *International*  
422 *Conference on Machine Learning*, pages 146–155. PMLR, 2017.
- 423 [29] Tabish Rashid, Mikayel Samvelyan, Christian Schroeder De Witt, Gregory Farquhar, Jakob  
424 Foerster, and Shimon Whiteson. Monotonic value function factorisation for deep multi-agent  
425 reinforcement learning. *The Journal of Machine Learning Research*, 21(1):7234–7284, 2020.
- 426 [30] Charles Dugas, Yoshua Bengio, François Bélisle, Claude Nadeau, and René Garcia. Incorporat-  
427 ing functional knowledge in neural networks. *Journal of Machine Learning Research*, 10(6),  
428 2009.

- 429 [31] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal  
430 policy optimization algorithms. *CoRR*, abs/1707.06347, 2017.
- 431 [32] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy  
432 maximum entropy deep reinforcement learning with a stochastic actor. In Jennifer G. Dy  
433 and Andreas Krause, editors, *Proceedings of the 35th International Conference on Machine  
434 Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*, volume 80 of  
435 *Proceedings of Machine Learning Research*, pages 1856–1865. PMLR, 2018.
- 436 [33] Lucas N. Alegre, Florian Felten, El-Ghazali Talbi, Grégoire Danoy, Ann Nowé, Ana L. C.  
437 Bazzan, and Bruno C. da Silva. MO-Gym: A library of multi-objective reinforcement learn-  
438 ing environments. In *Proceedings of the 34th Benelux Conference on Artificial Intelligence  
439 BNAIC/Benelearn 2022*, 2022.
- 440 [34] Diederik M Roijers. Multi-objective decision-theoretic planning. *AI Matters*, 2(4):11–12, 2016.
- 441 [35] Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah  
442 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of  
443 Machine Learning Research*, 22(268):1–8, 2021.