# Bounded Robustness in Reinforcement Learning
# via Lexicographic Objectives

**Daniel Jarne Ornia**                                    D.JARNEORNIA@TUDELFT.NL
*Delft University of Technology*

**Licio Romao**                                           LICIO.ROMAO@CS.OX.AC.UK
*University of Oxford*

**Lewis Hammond**                                         LEWIS.HAMMOND@CS.OX.AC.UK
*University of Oxford*

**Manuel Mazo Jr.**                                       M.MAZO@TUDELFT.NL
*Delft University of Technology*

**Alessandro Abate**                                      ALESSANDRO.ABATE@CS.OX.AC.UK
*University of Oxford*

## Abstract

Policy robustness in Reinforcement Learning may not be desirable at any cost: the alterations caused by robustness requirements from otherwise optimal policies should be explainable, quantifiable and formally verifiable. In this work we study how policies can be *maximally robust* to arbitrary observational noise by analysing how they are altered by this noise through a stochastic linear operator interpretation of the disturbances, and establish connections between robustness and properties of the noise kernel and of the underlying MDPs. Then, we construct sufficient conditions for policy robustness, and propose a robustness-inducing scheme, applicable to any policy gradient algorithm, that formally trades off expected policy utility for robustness through *lexicographic optimisation*, while preserving convergence and sub-optimality in the policy synthesis.

## 1. Introduction

Consider a dynamical system where we need to synthesise a controller (policy) through a model-free Reinfrocement Learning (Sutton and Barto, 2018) approach. When using a simulator for training we expect the deployment of the controller in the real system to be affected by different sources of noise, possibly not predictable or modelled (*e.g.* for networked components we may have sensor faults, communication delays, *etc*). In safety-critical systems, robustness (in terms of successfully controlling the system under disturbances) should preserve formal guarantees, and plenty of effort has been put on developing formal convergence guarantees on policy gradient algorithms (Agarwal et al., 2021; Bhandari and Russo, 2019). All these guarantees vanish under regularization and adversarial approaches, which are aimed to produce more robust policies. Therefore, for such applications one needs a scheme to regulate the robustness-utility trade-off in RL policies, that on the one hand preserves the formal guarantees of the original algorithms, and on the other attains sub-optimality conditions from the original problem. Additionally, if we do not know the structure of the disturbance (which holds in most applications), learning directly a policy for an arbitrarily disturbed environment will yield unexpected behaviours when deployed in the true system.

**Lexicographic Reinforcement Learning (LRL)**    Recently, lexicographic optimisation (Isermann, 1982; Rentmeesters et al., 1996) has been applied to the multi-objective RL setting (Skalse et al., 2022b). In an LRL setting some objectives may be more important than others, and so we may want to obtain policies that solve the multi-objective problem in a lexicographically prioritised way, *i.e.*, "find the policies that optimise objective $i$ (reasonably well), and from those the ones that optimise objective $i + 1$ (reasonably well), and so on".

**Previous Work**    In robustness against *model uncertainty*, the MDP may have noisy or uncertain reward signals or transition probabilities, as well as possible resulting *distributional shifts* in the training data (Heger, 1994; Xu and Mannor, 2006; Fu et al., 2018; Pattanaik et al., 2018; Pirotta et al., 2013; Abdullah et al., 2019), connecting to ideas on distributionally robust optimisation (Wiesemann et al., 2014; Van Parys et al., 2015). For *adversarial attacks or disturbances* on policies or action selection in RL agents (Gleave et al., 2020; Lin et al., 2017; Tessler et al., 2019; Pan et al., 2019; Tan et al., 2020; Klima et al., 2019; Liang et al., 2022), recently Gleave et al. (2020) propose to attack RL agents by swapping the policy for an adversarial one at given times. For a detailed review on Robust RL see Moos et al. (2022). Our work focuses in robustness versus *observational disturbances*, where agents observe a disturbed state measurement and use it as input for the policy (Kos and Song, 2017; Huang et al., 2017; Behzadan and Munir, 2017; Mandlekar et al., 2017; Zhang et al., 2020, 2021). Zhang et al. (2020) propose a *state-adversarial* MDP framework, and utilise adversarial regularising terms that can be added to different deep RL algorithms to make the resulting policies more robust to observational disturbances, and Zhang et al. (2021) study how LSTM increases robustness with optimal state-perturbing adversaries.

**Contributions**    Most existing work on RL with observational disturbances proposes modifying RL algorithms at the cost of *explainability* (in terms of sub-optimality bounds) and *verifiability*, since the induced changes in the new policies result in a loss of convergence guarantees. Our main contributions are summarised in the following points.

- We consider general unknown stochastic disturbances and formulate a quantitative definition of observational robustness that allows us to characterise the sets of robust policies for any MDP in the form of operator-invariant sets. We analyse how the structure of these sets depends on the MDP and noise kernel, and obtain an inclusion relation providing intuition into how we can search for robust policies more effectively.[1]

- We propose a meta-algorithm that can be applied to any existing policy gradient algorithm, Lexicographically Robust Policy Gradient (LRPG) that (1) Retains policy sub-optimality up to a specified tolerance while maximising robustness, (2) Formally controls the utility-robustness trade-off through this design tolerance, (3) Preserves formal guarantees.

Figure 1 represents a qualitative interpretation of the results in this work.

## 1.1. Preliminaries

---

1. There are strong connections between Sections 2-3 in this paper and the literature on planning for POMDPs (Spaan and Vlassis, 2004; Spaan, 2012) and MDP invariances (Ng et al., 1999; van der Pol et al., 2020; Skalse et al., 2022a), as well as recent work concerning robustness misspecification (Korkmaz, 2023).

**Notation** We use calligraphic letters $\mathcal{A}$ for collections of sets and $\Delta(\mathcal{A})$ as the space of probability measures over $\mathcal{A}$. For two probability distributions $P, P'$ defined on the same $\sigma$−algebra $\mathcal{F}$, $D_{TV}(P\|P') = \sup_{A\in\mathcal{F}}|P(A) - P'(A)|$ is the total variation distance. For two elements of a vector space we use $\langle\cdot,\cdot\rangle$ as the inner product. We use $\mathbf{1}_n$ as a column-vector of size $n$ that has all entries equal to 1. We say that an MDP is *ergodic* if for any policy the resulting Markov Chain (MC) is ergodic. We say that $S$ is a $n \times n$ row-stochastic matrix if $S_{ij} \geq 0$ and each
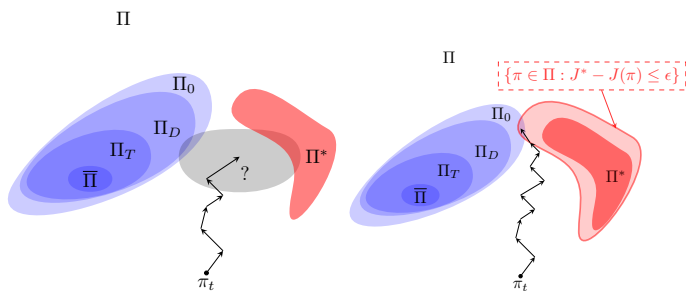


Figure 1: Qualitative representation LRPG (right), compared to usual robustness-inducing algorithms. The sets in blue are the robust policies to be defined in the coming sections. LRPG induces robustness while guaranteeing that the policies will deviate a bounded distance from the optimal.

row of $S$ sums 1. We assume all learning rates in this work $\alpha_t(x, u) \in [0, 1]$ ($\beta_t, \eta_t...$) satisfy the conditions $\sum_{t=1}^{\infty} \alpha_t(x, u) = \infty$ and $\sum_{t=1}^{\infty} \alpha_t(x, u)^2 < \infty$.

**Lexicographic Reinforcement Learning** Consider a parameterised policy $\pi_\theta$ with $\theta \in \Theta$, and two objective functions $K_1$ and $K_2$. PB-LRL uses a multi-timescale optimisation scheme to optimise $\theta$ faster for higher-priority objectives, iteratively updating the constraints induced by these priorities and encoding them via Lagrangian relaxation techniques (Bertsekas, 1997). Let $\theta' \in \text{argmax}_\theta K_1(\theta)$. Then, PB-LRL can be used to find parameters $\theta'' \in \{\text{argmax}_\theta K_2(\theta), \text{ s.t. } K_1(\theta) \geq K_1(\theta') - \epsilon\}$. This is done through the update:

$$\theta \leftarrow \text{proj}_\Theta\left[\theta + \nabla_\theta \hat{K}(\theta)\right], \quad \lambda \leftarrow \text{proj}_{\mathbb{R}_{\geq 0}}\left[\lambda + \eta_t(\hat{k}_1 - \epsilon_t - K_1(\theta))\right], \tag{1}$$

where $\hat{K}(\theta) := (\beta_t^1 + \lambda\beta_t^2) \cdot K_1(\theta) + \beta_t^2 \cdot K_2(\theta)$, $\lambda$ is a Langrange multiplier, $\beta_t^1, \beta_t^2, \eta_t$ are learning rates, and $\hat{k}_1$ is an estimate of $K_1(\theta')$. Typically, we set $\epsilon_t \to 0$, though we can use other tolerances too, *e.g.*, $\epsilon_t = 0.9 \cdot \hat{k}_1$. For more details see Skalse et al. (2022b).

## 2. Observationally Robust Reinforcement Learning

Robustness-inducing methods in model-free RL must address the following dilemma: How do we deal with uncertainty without an explicit mechanism to estimate such uncertainty during policy execution? Consider an example of an MDP where, at policy roll-out phase, there is a non-zero probability of measuring a "wrong" state. In such a scenario, measuring the wrong state can lead to executing unboundedly bad actions. This problem is represented by the following version of a noise-induced partially observable Markov Decision Process (Spaan, 2012).

**Definition 1** *An observationally-disturbed MDP (DOMDP) is (a POMDP) defined by the tuple $(X, U, P, R, T, \gamma)$ where $X$ is a finite set of states, $U$ is a set of actions, $P : U \times X \mapsto \Delta(X)$ is a probability measure of the transitions between states and $R : X \times U \times X \mapsto \mathbb{R}$ is a reward function. The map $T : X \mapsto \Delta(X)$ is a stochastic kernel induced by some unknown noise signal, such that $T(y \mid x)$ is the probability of measuring $y$ while the true state is $x$, and acts only on the state observations. At last $\gamma \in [0, 1]$ is a reward discount.*

A (memoryless) policy for the agent is a stochastic kernel $\pi : X \mapsto \Delta(U)$. For simplicity, we overload notation on $\pi$, denoting by $\pi(x,u)$ as the probability of taking action $u$ at state $x$. In a DOMDP[2] agents can measure the full state, but the measurement will be disturbed by some unknown random signal *in the policy deployment*. The difficulty of acting in such DOMDP is that agents will have to act based on disturbed states $\tilde{x} \sim T(\cdot \mid x)$. We then need to construct policies that will be as robust as possible against such noise *without the existance of a model to estimate, filter or reject disturbances*. The value function of a policy $\pi$ (*critic*), $V^\pi : X \mapsto \mathbb{R}$, is given by $V^\pi(x_0) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t R(x_t, \pi(x_t), x_{t+1})]$. The action-value function of $\pi$ (*Q-function*) is given by $Q^\pi(x,u) = \sum_{y \in X} P(x,u,y)(R(x,u,y) + \gamma V^\pi(y))$. We then define the objective function as $J(\pi) := \mathbb{E}_{x_0 \sim \mu_0}[V^\pi(x_0)]$ with $\mu_0$ being a distribution of initial states, and we use $J^* := \max_\pi J(\pi)$ and $\pi^*$ as the optimal policy, and $\Pi_\epsilon^* := \{\pi \in \Pi : J^* - J(\pi) \leq \epsilon\}$ is the set of $\epsilon$-optimal policies. If a policy is parameterised by $\theta \in \Theta$ we write $\pi_\theta$ and $J(\theta)$.

**Assumption 1** *For any DOMDP and policy $\pi$, the resulting MC is irreducible and aperiodic.*

We now formalise a notion of *observational robustness*. Firstly, due to the presence of the stochastic kernel $T$, the policy we are applying is altered as we are applying a collection of actions in a possibly wrong state. Then, $\langle \pi, T \rangle(x,u) := \sum_{y \in X} T(y \mid x)\pi(y,u)$, where $\langle \pi, T \rangle : X \mapsto \Delta(U)$ is the *disturbed* policy, which averages the current policy given the error induced by the presence of the stochastic kernel. Notice that $\langle \cdot, T \rangle(x) : \Pi \mapsto \Delta(U)$ is an averaging operator yielding the alteration of the policy due to noise. We define the *robustness regret*[3]: $\rho(\pi, T) := J(\pi) - J(\langle \pi, T \rangle)$.

**Definition 2 (Policy Robustness)** *A policy $\pi$ is $\kappa$-robust against a stochastic kernel $T$ if $\rho(\pi, T) \leq \kappa$. If $\pi$ is $0$-robust it is maximally robust. The sets of $\kappa$-robust policies are $\Pi_\kappa := \{\pi \in \Pi : \rho(\pi, T) \leq \kappa\}$, with $\Pi_0$ being the maximally robust policies.*

One can motivate the characterisation and models above from a control perspective, where policies use as input discretised state measurements with possible sensor measurement errors. Formally ensuring robustness properties when learning RL policies will, in general, force the resulting policies to deviate from optimality in the undisturbed MDP. We propose then the following problem.

**Problem 1** *Consider a DOMDP model as per Definition 1 and let $\epsilon$ be a non-negative tolerance level. Our goal is to find amongst all $\epsilon$-optimal policies those that minimize the robustness level $\kappa$:*

$$\text{minimize } \kappa \ \ s.t. \pi \in \Pi_\epsilon^\star \cap \Pi_\kappa.$$

Note that this is formulated as general as possible with respect to the robustness of the policies: We would like to find a policy that, trading off $\epsilon$ in terms of cumulative rewards, observes the same discounted rewards when disturbed by $T$.

## 3. Characterisation of Robust Policies

An important question to be addressed before trying to synthesise robust policies is what these robust policies look like, and how they are related to DOMDP properties. A policy $\pi$ is said to be

---

2. Definition 1 is a generalised form of the State-Adversarial MDP used by Zhang et al. (2020): the adversarial case is a particular form of DOMDP where $T$ assigns probability 1 to one adversarial state.

3. The robustness regret satisfies $\rho(\pi^*, T) \geq 0$ for all kernels $T$, and it allows us to directly compare the robustness regret with the utility regret of the policy.

constant if $\pi(x) = \pi(y)$ for all $x, y \in X$, and the collection of all constant policies is denoted by $\bar{\Pi}$. A policy is called a fixed point of $\langle \cdot, T \rangle$ if $\pi(x) = \langle \pi, T \rangle(x)$ for all $x \in X$. The collection of all fixed points is $\Pi_T$. Observe furthermore that $\Pi_T$ *only depends on the kernel $T$ and the set*[4] $X$. Let us assume we have a policy iteration algorithm that employs an action-value function $Q^\pi$ and policy $\pi$. The advantage function for $\pi$ is defined as $A^\pi(x, u) := Q^\pi(x, u) - V^\pi(x)$. We can similarly define the *noise disadvantage* of policy $\pi$ as:

$$D^\pi(x, T) := V^\pi(x) - \mathbb{E}_{u \sim \langle \pi, T \rangle(x)}[Q^\pi(x, u)], \tag{2}$$

which measures the difference of applying at state $x$ an action according to the policy $\pi$ with that of playing an action according to $\langle \pi, T \rangle$ and then continuing playing an action according to $\pi$. Our intuition says that if it happens to be the case that $D^\pi(x, T) = 0$ for all states in the DOMDP, then such a policy is maximally robust. And this is indeed the case, as shown in the next proposition.

**Proposition 3** *Consider a DOMDP as in Definition 1 and the robustness notion as in Definition 2. If a policy $\pi$ is such that $D^\pi(x, T) = 0$ for all $x \in X$, then $\pi$ is maximally robust, i.e., let $\Pi_D := \{\pi \in \Pi : \mu_\pi(x) D^\pi(x, T) = 0 \, \forall \, x \in X\}$, then we have that $\Pi_D \subseteq \Pi_0$.*
**Proof** *We want to show that $D^\pi(x, T) = 0 \implies \rho(\pi, T) = 0$. Taking $D^\pi(x, T) = 0$ one has a policy that produces an disadvantage of zero when noise kernel $T$ is applied. Then, $\forall \, x \in X$,*

$$D^\pi(x, T) = 0 \implies \mathbb{E}_{u \sim \langle \pi, T \rangle(x)}[Q^\pi(x, u)] = V^\pi(x). \tag{3}$$

*Now define the value of the disturbed policy as $V^{\langle \pi, T \rangle}(x) = \mathbb{E}_{\substack{u \sim \langle \pi, T \rangle(x), \\ y \sim P(\cdot|x,u)}} \left[ r(x, u, y) + \gamma V^{\langle \pi, T \rangle}(y) \right].$*
*We will now show that $V^\pi(x) = V^{\langle \pi, T \rangle}(x)$, for all $x \in X$. Observe, from (3) using $V^\pi(x) = \mathbb{E}_{u \sim \langle \pi, T \rangle(x)}[Q^\pi(x, u)]$, we have $\forall x \in X$:*

$$\begin{aligned} V^\pi(x) - V^{\langle \pi, T \rangle}(x) &= \mathbb{E}_{u \sim \langle \pi, T \rangle(x)}[Q^\pi(x, u)] - \mathbb{E}_{\substack{u \sim \langle \pi, T \rangle(x) \\ y \sim P(\cdot|x,u)}} \left[ r(x, u, y) + \gamma V^{\langle \pi, T \rangle}(y) \right] = \\ &= \mathbb{E}_{\substack{u \sim \langle \pi, T \rangle(x) \\ y \sim P(\cdot|x,u)}} \left[ \gamma V^\pi(y) - \gamma V^{\langle \pi, T \rangle}(y) \right] = \gamma \mathbb{E}_{y \sim P(\cdot|x,u)} \left[ V^\pi(y) - V^{\langle \pi, T \rangle}(y) \right]. \end{aligned} \tag{4}$$

*Now, taking the sup norm at both sides of (4) we get*

$$\|V^\pi(x) - V^{\langle \pi, T \rangle}(x)\|_\infty = \gamma \left\| \mathbb{E}_{y \sim P(\cdot|x,u)} \left[ V^\pi(y) - V^{\langle \pi, T \rangle}(y) \right] \right\|_\infty. \tag{5}$$

*Since the norm on the right hand side of (5) is over $y \in X$ and $\gamma < 1$, it follows that $\|V^\pi(x) - V^{\langle \pi, T \rangle}(x)\|_\infty = 0$. Finally, $\|V^\pi(x) - V^{\langle \pi, T \rangle}(x)\|_\infty = 0 \implies V^\pi(x) - V^{\langle \pi, T \rangle}(x) = 0 \, \forall x \in X$, and $V^\pi(x) - V^{\langle \pi, T \rangle}(x) = 0 \, \forall \, x \in X \implies J(\pi) = J(\langle \pi, T \rangle) \implies \rho(\pi, T) = 0$.* ∎

So far we have shown that both the set of fixed points $\bar{\bar{\Pi}}$ and the set of policies for which the disadvantage function is equal to zero $\Pi_D$ are contained in the set of maximally robust policies. We now show how the defined robust policy sets can be linked in a single result through the following policy inclusions.

---

4. There is a (natural) bijection between the set of constant policies and the space $\Delta(U)$. The set of fixed points of the operator $\langle \cdot, T \rangle$ also has an algebraic characterisation in terms of the null space of the operator $\mathrm{Id}(\cdot) - \langle \cdot, T \rangle$. We are not exploiting the later characterisation in this paper.

**Theorem 4 (Policy Inclusions)** *For a DOMDP with noise kernel $T$, consider the sets $\overline{\Pi}, \Pi_T, \Pi_D$ and $\Pi_0$. Then, the following inclusion relation holds: $\overline{\Pi} \subseteq \Pi_T \subseteq \Pi_D \subseteq \Pi_0$. Additionally, the sets $\overline{\Pi}, \Pi_T$ are* convex *for all MDPs and kernels $T$, but $\Pi_D, \Pi_0$ may not be.*

**Proof** If a policy $\pi \in \Pi$ is a fixed point of the operator $\langle \cdot, T \rangle$, then $\rho(\pi, T) = J(\pi) - J(\langle \pi, T \rangle) = J(\pi) - J(\pi) = 0 \implies \pi \in \Pi_0$. Therefore, $\Pi_T \subseteq \Pi_0$. Now, the space of stochastic kernels $\mathcal{K} : X \mapsto \Delta(X)$ is equivalent to the space of row-stochastic $|X| \times |X|$ matrices, therefore one can write $T(y \mid x) \equiv T_{xy}$ as the $xy$−th entry of the matrix $T$. Then, the representation of a constant policy as an $X \times U$ matrix can be written as $\overline{\pi} = \mathbf{1}_{|X|} v^\top$, where $\mathbf{1}_{|X|}$ where $v \in \Delta(U)$ is any probability distribution over the action space. Observe that, applying the operator $\langle \pi, T \rangle$ to a constant policy yields $\langle \overline{\pi}, T \rangle = T\mathbf{1}_{|X|} v^\top$. By the Perron-Frobenius Theorem (Horn and Johnson, 2012), since $T$ is row-stochastic it has at least one eigenvalue $\text{eig}(T) = 1$, and this admits a (strictly positive) eigenvector $T\mathbf{1}_{|X|} = \mathbf{1}_{|X|}$. Therefore, $\langle \overline{\pi}, T \rangle = T\mathbf{1}_{|X|} v^\top = \mathbf{1}_{|X|} v^\top = \overline{\pi} \implies \overline{\Pi} \subseteq \Pi_T$. Combining this result with Proposition 3, we simply need to show that $\Pi_T \subset \Pi_D$. Take $\pi$ to be a fixed point of $\langle \pi, T \rangle$. Then $\langle \pi, T \rangle = \pi$, and from the definition in (2):

$$D^\pi(x, T) = V^\pi(x) - \mathbb{E}_{u \sim \langle \pi, T \rangle(x, \cdot)}[Q^\pi(x, u)] = V^\pi(x) - \mathbb{E}_{u \sim \pi(x, \cdot)}[Q^\pi(x, u)] = 0.$$

Therefore, $\pi \in \Pi_D$, which completes the sequence of inclusions. Convexity of $\overline{\Pi}, \Pi_T$ follows from considering the convex hulls of two constant or fixed point policies. ∎

Let us reflect on the inclusion relations of Theorem 4. The inclusions are in general not strict, and in fact the geometry of the sets (as well as whether some of the relations are in fact equalities) is highly dependent on the reward function, and in particular on the complexity (from an information-theoretic perspective) of the reward function. As an intuition, less complex reward functions (more uniform) will make the inclusions above expand to the entire policy set, and more complex reward functions will make the relations collapse to equalities.

**Corollary 5** *For any* ergodic *DOMDP there exist reward functions $\overline{R}$ and $\underline{R}$ such that the resulting DOMDP satisfies A) $\Pi_D = \Pi_0 = \Pi$ (any policy is max. robust) if $R = \overline{R}$, and B) $\Pi_T = \Pi_D = \Pi_0$ (only fixed point policies are maximally robust) if $R = \underline{R}$.*

**Proof** [Corollary 5] For statement A) let $\overline{R}(\cdot, \cdot, \cdot) = c$ for some constant $c \in \mathbb{R}$. Then, $J(\pi) = \mathbb{E}_{x_0 \sim \mu_0}[\sum_t \gamma^t \overline{r}_t \mid \pi] = \frac{c\gamma}{1-\gamma}$, which does not depend on the policy $\pi$. For any noise kernel $T$ and policy $\pi$, $J(\pi) - J\langle \pi, T \rangle = 0 \implies \pi \in \Pi_0$. For statement B assume $\exists \pi \in \Pi_0 : \pi \notin \Pi_T$. Then, $\exists x^* \in X$ and $u^* \in U$ such that $\pi(x^*, u^*) \neq \langle \pi, T \rangle(x^*, u^*)$. Let $\underline{R}(x, u, x') := c$ if $x = x^*$ and $u = u^*$, 0 otherwise. Then, $\mathbb{E}[R(x, \pi(x), x'] < \mathbb{E}[R(x, \langle \pi, T \rangle(x), x']$ and since the MDP is ergodic $x$ is visited infinitely often and $J(\pi) - J(\langle \pi, T \rangle) > 0 \implies \pi \notin \Pi_0$, which contradicts the assumption. Therefore, $\Pi_0 \setminus \Pi_T = \emptyset \implies \Pi_0 = \Pi_T$. ∎

We can now summarise the insights from Theorem 4 and Corollary 5 in the following conclusions: (1) The set $\overline{\Pi}$ is maximally robust, convex and *independent of the DOMDP*, (2) The set $\Pi_T$ is maximally robust, convex, includes $\overline{\Pi}$, and its properties *only depend* on $T$, (3) The set $\Pi_D$ includes $\Pi_T$ and is maximally robust, but its properties *depend on the DOMDP*.

## 4. Robustness through Lexicographic Objectives

To be able to apply LRL results to our robustness problem we need to first cast robustness as a valid objective to be maximised, and then show that a stochastic gradient descent approach would indeed find a global maximum of the objective, therefore yielding a maximally robust policy. [5]

---

**Algorithm 1** LRPG

  **input** Simulator, $\tilde{T}$, $\epsilon$
  initialise $\theta$, critic (if using), $\lambda$, $\{\beta_t^1, \beta_t^2, \eta\}$
  set $t = 0$, $x_t \sim \mu_0$
  **while** $t < \text{max\_iterations}$ **do**
    perform $u_t \sim \pi_\theta(x_t)$
    observe $r_t$, $x_{t+1}$, sample $y \sim \tilde{T}(\cdot \mid x)$
    **if** $\hat{K}_1(\theta)$ not converged **then**
      $\hat{k}_1 \leftarrow \hat{K}_1(\theta)$
    **end if**
    update critic (if using)
    update $\theta$ using (8) and $\lambda$ using (1)
  **end while**
  **output** $\theta$

---

**Proposed approach** Following the framework presented in previous sections, we propose the following approach to obtain lexicographic robustness. In the introduction, we emphasised that the motivation for this work comes partially from the fact that we may not know $T$ in reality, or have a way to estimate it. However, the theoretical results until now depend on $T$. Our proposed solution to this lies in the results of Theorem 4. We can use a *design* generator $\tilde{T}$ to perturb the policy during training such that $\tilde{T}$ has the *smallest possible fixed point set* (i.e. the constant policy set, $\tilde{T}$ satisfies $\Pi_{\tilde{T}} = \overline{\Pi}$), and any algorithm that drives the policy towards the set of fixed points of $\tilde{T}$ *will also drive the policy towards fixed points of $T$*: from Theorem 4, $\Pi_{\tilde{T}} \subseteq \Pi_T$.

### 4.1. Lexicographically Robust Policy Gradient

Consider then the objective to be minimised:

$$K_{\tilde{T}}(\theta) = \frac{1}{2} \sum_{x \in X} \mu_{\pi_\theta}(x) \sum_{u \in U} \left( \pi_\theta(x, u) - \langle \pi_\theta, \tilde{T} \rangle(x, u) \right)^2, \tag{6}$$

Notice that optimising (6) projects the current policy onto the set of fixed points of the operator $\langle \cdot, \tilde{T} \rangle$, and due to Assumption 1, which requires $\mu_{\pi_\theta}(x) > 0$ for all $x \in X$, the optimal solution is equal to zero if and only if there exists a value of the parameter $\theta$ for which the corresponding $\pi_\theta$ is a fixed point of $\langle \cdot, \tilde{T} \rangle$. We present now the proposed LRPG meta-algorithm to achieve lexicographic robustness for any policy gradient algorithm at choice. From Skalse et al. (2022b), the convergence of PB-LRL algorithms is guaranteed as long as the original policy gradient algorithm for each single objective converges.

**Assumption 2** *The policy is updated through an algorithm (e.g. A2C, PPO...) such that $\theta_{t+1} \leftarrow \text{proj}_\Theta \left[ \theta_t + \alpha_t \nabla_{\theta_t} \hat{K}_1 \right]$ converges* a.s. *to a (local or global) optimum $\theta^*$.*

**Theorem 6** *Consider a DOMDP as in Definition 1 and let $\pi_\theta$ be a parameterised policy. Take a design kernel $\tilde{T} \in \{T : \Pi_T = \overline{\Pi}\}$. Consider the following modified gradient for objective $K_{\tilde{T}}(\theta)(x)$ and sampled point $y \sim \tilde{T}(\cdot \mid x)$:*

$$\nabla_\theta \hat{K}_{\tilde{T}}'(\theta) = \mathbb{E}_{x \sim \mu_{\pi_\theta}} \Big[ \sum_{u \in U} (\pi_\theta(x, u) - \pi_\theta(y, u)) \nabla_\theta \pi_\theta(x, u) \Big]. \tag{7}$$

---

5. The advantage of using LRL is that we can formally bound the trade-off between *robustness and optimality* through $\epsilon$, determinining how far we allow our resulting policy to be from an optimal policy in favour of it being more robust.

*Given an $\epsilon > 0$, if Assumptions 1 and 2 hold, then the following iteration (LRPG):*

$$\theta \leftarrow \text{proj}_{\Theta} \left[ \theta + (\beta_t^1 + \lambda \beta_t^2) \cdot \nabla_\theta \hat{K}_1(\theta) + \beta_t^2 \nabla_\theta \hat{K}_{\tilde{T}}'(\theta) \right] \tag{8}$$

*converges a.s. to parameters $\theta^\epsilon$ that satisfy $\theta^\epsilon \in \text{argmin}_{\theta \in \Theta'} K_{\tilde{T}}(\theta)$ such that $K_1^* \geq K_1(\theta^\epsilon) - \epsilon$, where $\Theta' = \Theta$ if $\theta^*$ is globally optimal and a compact local neighbourhood of $\theta^*$ otherwise.*

**Proof** To apply LRL results, we need to show that both gradient descent schemes converge (separately) to local or global maxima. Let us first show that $\theta_{t+1} = \text{proj}_{\Theta} \left[ \theta_t - \alpha_t \nabla_\theta \hat{K}_{\tilde{T}}'(\theta_t) \right]$ converges *a.s.* to parameters $\tilde{\theta}$ satisfying $K_{\tilde{T}} = 0$. We prove this making use of fixed point iterations with non-expansive operators (specifically, Theorem 4, section 10.3 in Borkar (2008)). First, observe that for a tabular representation, $\pi_\theta(x, u) = \theta_{xu}$, and $\nabla_\theta \pi_\theta(x, u)$ is a vector of zeros, with value 1 for the position $\theta_{xu}$. We can then write the SGD in terms of the policy for each state $x$, considering $\pi(x) \equiv (\theta_{xu_1}, \theta_{xu_2}, ..., \theta_{xu_k})^T$. Let $y \sim \tilde{T}(\cdot \mid x)$. Then:

$$\pi_{t+1}(x) = \pi_t(x) - \alpha_t \big( \pi_t(x) - \pi_t(y) \big) = \pi_t(x) - \alpha_t \Big( \pi_t(x) - \langle \pi_t, \tilde{T} \rangle(x) - \big( \pi_t(y) - \langle \pi_t, \tilde{T} \rangle(x) \big) \Big).$$

We now need to verify that the necessary conditions for applying Theorem 4, section 10.3 in Borkar (2008) hold. First, making use of the property $\|\tilde{T}\|_\infty = 1$ for any row-stochastic matrix $\tilde{T}$, for any two policies $\pi_1, \pi_2 \in \Pi$:

$$\|\langle \pi_1, \tilde{T} \rangle - \langle \pi_2, \tilde{T} \rangle\|_\infty = \|\tilde{T}\pi_1 - \tilde{T}\pi_2\|_\infty = \|\tilde{T}(\pi_1 - \pi_2)\|_\infty \leq \|\tilde{T}\|_\infty \|\pi_1 - \pi_2\|_\infty = \|\pi_1 - \pi_2\|_\infty.$$

Therefore, the operator $\langle \cdot, \tilde{T} \rangle$ is non-expansive with respect to the sup-norm. For the final condition:

$$\mathbb{E}_{y \sim \tilde{T}(\cdot|x)} \left[ \pi_t(y) - \langle \pi_t, \tilde{T} \rangle(x) \mid \pi_t, \tilde{T} \right] = \sum_{y \in X} \tilde{T}(y \mid x) \pi_t(y) - \langle \pi_t, \tilde{T} \rangle(x) = 0.$$

Therefore, the difference $\pi_t(y) - \langle \pi_t, \tilde{T} \rangle(x)$ is a martingale difference for all $x$. One can then apply Theorem 4, sec. 10.3 (Borkar, 2008) to conclude that $\pi_t(x) \to \tilde{\pi}(x)$ almost surely. Finally from Assumption 1, for any policy all states $x \in X$ are visited infinitely often, therefore $\pi_t(x) \to \tilde{\pi}(x) \forall x \in X \implies \pi_t \to \tilde{\pi}$ and $\tilde{\pi}$ satisfies $\langle \tilde{\pi}, \tilde{T} \rangle = \tilde{\pi}$, and $K_{\tilde{T}}(\tilde{\pi}) = 0$.

Now, from Assumption 2, the iteration $\theta \leftarrow \text{proj}_{\Theta} \left[ \theta + \alpha_t \nabla_\theta \hat{K}_1 \right]$ converges *a.s.* to a (local or global) optimum $\theta^*$. Then, both objectives are invex Ben-Israel and Mond (1986b) (either locally or globally), and any linear combination of them will also be invex (again, locally or globally). Finally, we can directly apply the results from Skalse et al. (2022b), and

$$\theta \leftarrow \text{proj}_{\Theta} \left[ \theta + (\beta_t^1 + \lambda \beta_t^2) \cdot \nabla_\theta \hat{K}_1(\theta) + \beta_t^2 \nabla_\theta \hat{K}_{\tilde{T}}'(\theta) \right]$$

converges *a.s.* to parameters $\theta^\epsilon$ that satisfy $\theta^\epsilon \in \text{argmin}_{\theta \in \Theta'} K_{\tilde{T}}(\theta)$ such that $K_1^* \geq K_1(\theta^\epsilon) - \epsilon$, where $\Theta' = \Theta$ if $\theta^*$ is globally optimal and a compact local neighbourhood of $\theta^*$ otherwise. ∎

**Remark 7** *Observe that (7) is not the true gradient of (6), and $\theta^\epsilon \in \text{argmin}_{\theta \in \Theta'} K_{\tilde{T}}(\theta)$ if there exists a (local) minimum of $K_{\tilde{T}}$ in $\Theta^\epsilon := \{\theta : K_1^* \geq K_1(\theta) - \epsilon\}$. However, from Theorem 6 we know that the (pseudo) gradient descent scheme converges to a global minimum in the tabular case, therefore $\langle \nabla_\theta \hat{K}_{\tilde{T}}'(\theta), \nabla_\theta \hat{K}_{\tilde{T}}(\theta) \rangle < 0$ (Borkar, 2008), and gradient-like descent schemes will converge to (local or) global minimisers, which motivates the choice of this gradient approximation.*

8

We reflect again on Figure 1. The main idea behind LRPG is that by formally expanding the set of acceptable policies with respect to $K_1$, we may find robust policies more effectively while guaranteeing a minimum performance in terms of expected rewards. This addresses directly the premise behind Problem 1. In LRPG the first objective is still to minimise the distance $J^* - J(\pi)$ up to some tolerance. Then, from the policies that satisfy this constraint, we want to steer the learning algorithm towards a maximally robust policy, and we can do so without knowing $T$.

## 5. Considerations on Noise Generators

A natural question emerging is how to choose $\tilde{T}$, and how the choice influences the resulting policy robustness towards any other true $T$. In general, for any arbitrary policy utility landscape in a given MDP, there is no way of bounding the distance of the resulting policies for two different noise kernels $T_1, T_2$. However, *the optimality of the policy* remains bounded: Through LRPG guarantees we know that, for both cases, the utility of the resulting policy will be at most $\epsilon$ far from the optimal.

**Corollary 8** *Take $T$ to be any arbitrary noise kernel, and $\tilde{T}$ to satisfy $\tilde{T} \in \{T : \Pi_T = \overline{\Pi}\}$. Let $\pi$ be a policy resulting from a LRPG algorithm. Assume that $\min_{\pi' \in \Pi_{\tilde{T}}} D_{TV}(\pi \| \pi') = a$ for some $a < 1$. Then, it holds for any $T$ that $\min_{\pi' \in \Pi_T} D_{TV}(\pi \| \pi') \leq a$.*

**Proof** The proof follows by the inclusion results in Theorem 4. If $\Pi_{\tilde{T}} = \overline{\Pi}$, then $\Pi_{\tilde{T}} \subseteq \Pi_T$ for any other $T$. Then, the distance from $\pi$ to the set $\Pi_T$ is at most the distance to $\Pi_{\tilde{T}}$. ∎

That is, when using LRPG to obtain a robust policy $\pi$, the resulting policy is at most $a$ far from the set of fixed points (and therefore a maximally robust policy) with respect to the true $T$. This is the key argument behind our choices for $\tilde{T}$: A priori, the most sensible choice is a kernel that has no other fixed point than the set of constant policies. This fixed point condition is satisfied in the discrete state case for any $\tilde{T}$ that induces an irreducible Markov Chain, and in continuous state for any $\tilde{T}$ that satisfies a reachability condition (*i.e.* for any $x_0 \in X$, there exists a finite time for which the probability of reaching any ball $B \subset X$ of radius $r > 0$ through a sequence $x_{t+1} = T(x_t)$ is measurable). This holds for (additive) uniform or Gaussian disturbances.

## 6. Experiments

We verify the theoretical results of LRPG in a series of experiments on discrete state/action safety-related environments (Chevalier-Boisvert et al., 2018) (for extended experiments in continuous control tasks, hyperparameters *etc.* see extended version). We use A2C (Sutton and Barto, 2018) (LR-A2C) and PPO (Schulman et al., 2017) (LR-PPO) for our implementations of LRPG. In all cases, the lexicographic tolerance was set to $\epsilon = 0.99\hat{k}_1$ to deviate as little as possible from the primary objective. We compare against the baseline algorithms and against SA-PPO (Zhang et al., 2020) which is among the most effective (adversarial) robust RL approaches in literature. We trained 10 independent agents for each algorithm, and reported scores of the median agent (as in Zhang et al. (2020)) for 50 roll-outs. To simulate $\tilde{T}$ we disturb $x$ as $\tilde{x} = x + \xi$ for (1) a uniform bounded noise signal $\xi \sim \mathcal{U}_{[-b,b]}$ ($\tilde{T}^u$) and (2) and a Gaussian noise ($\tilde{T}^g$) such that $\xi \sim \mathcal{N}(0, 0.5)$. We test the resulting policies against a noiseless environment ($\emptyset$), a kernel $T_1 = \tilde{T}^u$, a kernel $T_2 = \tilde{T}^g$ and against two different state-adversarial noise configurations ($T_{adv}^2$) as proposed by Zhang et al. (2021) to evaluate how effective LRPG is at rejecting adversarial disturbances.

| | *PPO on MiniGrid Environments* | | | | | *A2C on MiniGrid Environments* | | | |
|---|---|---|---|---|---|---|---|---|---|
| Noise | PPO | $LR_{PPO}(K^u_T)$ | $LR_{PPO}(K^g_T)$ | SA-PPO | ‖ | A2C | $LR_{A2C}(K^u_T)$ | $LR_{A2C}(K^g_T)$ | $LR_{A2C}(K_D)$ |
| *LavaGap* | | | | | | | | | |
| $\emptyset$ | **0.95±0.003** | **0.95±0.075** | **0.95±0.101** | 0.94±0.068 | | **0.94±0.004** | **0.94±0.005** | **0.94±0.003** | **0.94±0.006** |
| $T_1$ | 0.80±0.041 | **0.95±0.078** | 0.93±0.124 | 0.88±0.064 | | 0.83±0.061 | **0.93±0.019** | 0.89±0.032 | 0.91±0.088 |
| $T_2$ | 0.92±0.015 | **0.95±0.052** | **0.95±0.094** | 0.93±0.050 | | 0.89±0.029 | **0.94±0.008** | 0.93±0.011 | 0.93±0.021 |
| $T^2_{adv}$ | 0.01±0.051 | 0.71±0.251 | 0.21±0.357 | **0.87±0.116** | | 0.27±0.119 | **0.79±0.069** | 0.68±0.127 | 0.56±0.249 |
| *LavaCrossing* | | | | | | | | | |
| $\emptyset$ | **0.95±0.023** | 0.93±0.050 | 0.93±0.018 | 0.88±0.091 | | 0.91±0.024 | 0.91±0.063 | 0.90±0.017 | **0.92±0.034** |
| $T_1$ | 0.50±0.110 | **0.92±0.053** | 0.89±0.029 | 0.64±0.109 | | 0.66±0.071 | **0.78±0.111** | 0.72±0.073 | 0.76±0.098 |
| $T_2$ | 0.84±0.061 | **0.92±0.050** | **0.92±0.021** | 0.85±0.094 | | 0.78±0.054 | 0.83±0.105 | 0.86±0.029 | **0.87±0.063** |
| $T^2_{adv}$ | 0.0±0.004 | 0.50±0.171 | 0.38±0.020 | **0.82±0.072** | | 0.06±0.056 | 0.04±0.030 | 0.01±0.008 | **0.09±0.060** |
| *DynamicObstacles* | | | | | | | | | |
| $\emptyset$ | **0.91±0.002** | **0.91±0.008** | **0.91±0.007** | **0.91±0.131** | | **0.91±0.011** | 0.88±0.020 | 0.89±0.009 | **0.91±0.013** |
| $T_1$ | 0.23±0.201 | **0.77±0.102** | 0.61±0.119 | 0.45±0.188 | | 0.27±0.104 | 0.43±0.108 | 0.45±0.162 | **0.56±0.270** |
| $T_2$ | 0.50±0.117 | **0.75±0.075** | 0.70±0.072 | 0.68±0.490 | | 0.45±0.086 | 0.53±0.109 | 0.52±0.161 | **0.67±0.203** |
| $T^2_{adv}$ | -0.49±0.312 | 0.51±0.234 | 0.33±0.202 | **0.55±0.170** | | -0.54±0.209 | -0.21±0.192 | -0.53±0.261 | **-0.51±0.260** |

Table 1: Reward values gained by LRPG and baselines on discrete control tasks.

**Robustness Results** We use objectives as defined in (6). Additionally, we aim to test the hypothesis: If we have an estimator for the critic $Q^\pi$ we can obtain robustness without inducing regularity in the policy using $D^\pi$, yielding a larger policy subspace to steer towards, and hopefully achieving policies closer to optimal. For this, we consider the objective $K_D(\theta)(x) := \frac{1}{2}\|D^{\pi_\theta}(x, T)\|_2^2$ by modifying A2C to retain a Q critic. We investigate the impact of LRPG PPO and A2C for discrete action-space problems on Gymnasium (Brockman et al., 2016). *Minigrid-LavaGap* (fully observable), *Minigrid-LavaCrossing* (partially observable) are safe exploration tasks where the agent needs to navigate an environment with cliff-like regions. *Minigrid-DynamicObstacles* (stochastic, partially observable) is a dynamic obstacle-avoidance environment. See Table 1.

## 7. Discussion

**Experiments** We applied LRPG on PPO and A2C (and SAC algorithms), for a set of discrete and continuous control environments. These environments are particularly sensitive to robustness problems; the rewards are sparse, and applying a sub-optimal action at any step of the trajectory often leads to terminal states with zero (or negative) reward. LRPG successfully induces lower robustness regrets in the tested scenarios, and the use of $K_D$ as an objective (even though we did not prove the convergence of a gradient based method with such objective) yields a better compromise between robustness and rewards. When compared to recent observational robustness methods, LRPG obtains similar robustness results while *preserving the original guarantees of the chosen algorithm.*

**Shortcomings and Contributions** The motivation for LRPG comes from situations where, when deploying a model-free controller in a dynamical system, we do not have a way of estimating the noise generation and we *are required to retain convergence guarantees of the algorithms used.* Although LRPG is a useful approach for learning policies in control problems where the noise sources are unknown, questions emerge on whether there are more effective methods of incorporating robustness into RL policies when guarantees are not needed. Specifically, since a completely model-free approach does not allow for simple alternative solutions such as filtering or disturbance rejection, there are reasons to believe it could be outperformed by model-based (or model learning) approaches. However, we argue that in completely model-free settings, LRPG provides a rational strategy to robustify RL agents.

# References

Mohammed Amin Abdullah, Hang Ren, Haitham Bou Ammar, Vladimir Milenkovic, Rui Luo, Mingtian Zhang, and Jun Wang. Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*, 2019.

Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *J. Mach. Learn. Res.*, 22 (98):1–76, 2021.

Vahid Behzadan and Arslan Munir. Vulnerability of deep reinforcement learning to policy induction attacks. In *International Conference on Machine Learning and Data Mining in Pattern Recognition*, pages 262–275. Springer, 2017.

A. Ben-Israel and B. Mond. What is invexity? *The Journal of the Australian Mathematical Society. Series B. Applied Mathematics*, 28(1):1–9, 1986a.

Adi Ben-Israel and Bertram Mond. What is invexity? *The ANZIAM Journal*, 28(1):1–9, 1986b.

Dimitri Bertsekas. *Nonlinear Programming*. Athena Scientific, 1999.

Dimitri P Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3): 334–334, 1997.

Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.

Vivek S. Borkar. *Stochastic Approximation*. Hindustan Book Agency, 2008.

V.S. Borkar and K. Soumyanatha. An analog scheme for fixed point computation. i. theory. *IEEE Transactions on Circuits and Systems I: Fundamental Theory and Applications*, 44(4):351–355, 1997. doi: 10.1109/81.563625.

Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.

Maxime Chevalier-Boisvert, Lucas Willems, and Suman Pal. Minimalistic gridworld environment for openai gym. https://github.com/maximecb/gym-minigrid, 2018.

Justin Fu, Katie Luo, and Sergey Levine. Learning robust rewards with adverserial inverse reinforcement learning. In *International Conference on Learning Representations*, 2018.

Adam Gleave, Michael Dennis, Cody Wild, Neel Kant, Sergey Levine, and Stuart Russell. Adversarial policies: Attacking deep reinforcement learning. In *International Conference on Learning Representations*, 2020.

Morgan A Hanson. On sufficiency of the kuhn-tucker conditions. *Journal of Mathematical Analysis and Applications*, 80(2):545–550, 1981.

Matthias Heger. Consideration of risk in reinforcement learning. In *Machine Learning Proceedings 1994*, pages 105–111. Elsevier, 1994.

Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.

Sandy Huang, Nicolas Papernot, Ian Goodfellow, Yan Duan, and Pieter Abbeel. Adversarial attacks on neural network policies. *arXiv preprint arXiv:1702.02284*, 2017.

H Isermann. Linear lexicographic optimization. *Operations-Research-Spektrum*, 4(4):223–228, 1982.

Richard Klima, Daan Bloembergen, Michael Kaisers, and Karl Tuyls. Robust temporal difference learning for critical domains. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '19, page 350–358, Richland, SC, 2019. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450363099.

Ezgi Korkmaz. Adversarial robust deep reinforcement learning requires redefining robustness. *arXiv preprint arXiv:2301.07487*, 2023.

Jernej Kos and Dawn Song. Delving into adversarial attacks on deep policies. *arXiv preprint arXiv:1705.06452*, 2017.

Yongyuan Liang, Yanchao Sun, Ruijie Zheng, and Furong Huang. Efficient adversarial training without attacking: Worst-case-aware robust reinforcement learning. *Advances in Neural Information Processing Systems*, 35:22547–22561, 2022.

Yen-Chen Lin, Zhang-Wei Hong, Yuan-Hong Liao, Meng-Li Shih, Ming-Yu Liu, and Min Sun. Tactics of adversarial attack on deep reinforcement learning agents. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 3756–3762, 2017.

Ajay Mandlekar, Yuke Zhu, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Adversarially robust policy learning: Active construction of physically-plausible perturbations. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3932–3939. IEEE, 2017.

Janosch Moos, Kay Hansel, Hany Abdulsamad, Svenja Stark, Debora Clever, and Jan Peters. Robust reinforcement learning: A review of foundations and recent advances. *Machine Learning and Knowledge Extraction*, 4(1):276–315, 2022.

Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations: Theory and application to reward shaping. In *Proc. of the Sixteenth International Conference on Machine Learning, 1999*, 1999.

Xinlei Pan, Daniel Seita, Yang Gao, and John Canny. Risk averse robust adversarial reinforcement learning. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 8522–8528. IEEE, 2019.

Santiago Paternain, Luiz F. O. Chamon, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained reinforcement learning has zero duality gap. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 7553–7563, 2019.

Anay Pattanaik, Zhenyi Tang, Shuijing Liu, Gautham Bommannan, and Girish Chowdhary. Robust deep reinforcement learning with adversarial attacks. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '18, page 2040–2042, Richland, SC, 2018. International Foundation for Autonomous Agents and Multiagent Systems.

Matteo Pirotta, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe policy iteration. In *International Conference on Machine Learning*, pages 307–315. PMLR, 2013.

Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of Machine Learning Research*, 22(268):1–8, 2021. URL http://jmlr.org/papers/v22/20-1364.html.

Mark J Rentmeesters, Wei K Tsai, and Kwei-Jay Lin. A theory of lexicographic multi-criteria optimization. In *Proceedings of ICECCS'96: 2nd IEEE International Conference on Engineering of Complex Computer Systems (held jointly with 6th CSESAW and 4th IEEE RTAW)*, pages 76–79. IEEE, 1996.

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

Joar Skalse, Matthew Farrugia-Roberts, Stuart Russell, Alessandro Abate, and Adam Gleave. Invariance in policy optimisation and partial identifiability in reward learning. *arXiv preprint arXiv:2203.07475*, 2022a.

Joar Skalse, Lewis Hammond, Charlie Griffin, and Alessandro Abate. Lexicographic multi-objective reinforcement learning. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*, pages 3430–3436, jul 2022b. doi: 10.24963/ijcai.2022/476.

Morton Slater. Lagrange multipliers revisited. Cowles Commission Discussion Paper No. 403, 1950.

Matthijs TJ Spaan. Partially observable markov decision processes. In *Reinforcement Learning*, pages 387–414. Springer, 2012.

Matthijs TJ Spaan and N Vlassis. A point-based pomdp algorithm for robot planning. In *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA'04. 2004*, volume 3, pages 2399–2404. IEEE, 2004.

Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.

Kai Liang Tan, Yasaman Esfandiari, Xian Yeow Lee, Soumik Sarkar, et al. Robustifying reinforcement learning agents via action space adversarial training. In *2020 American control conference (ACC)*, pages 3959–3964. IEEE, 2020.

Chen Tessler, Yonathan Efroni, and Shie Mannor. Action robust reinforcement learning and applications in continuous control. In *International Conference on Machine Learning*, pages 6215–6224. PMLR, 2019.

Elise van der Pol, Thomas Kipf, Frans A. Oliehoek, and Max Welling. Plannable approximations to mdp homomorphisms: Equivariance under actions. In *Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems*, AAMAS '20, page 1431–1439, Richland, SC, 2020. International Foundation for Autonomous Agents and Multiagent Systems. ISBN 9781450375184.

Bart PG Van Parys, Daniel Kuhn, Paul J Goulart, and Manfred Morari. Distributionally robust control of constrained stochastic systems. *IEEE Transactions on Automatic Control*, 61(2):430–442, 2015.

Wolfram Wiesemann, Daniel Kuhn, and Melvyn Sim. Distributionally robust convex optimization. *Operations Research*, 62(6):1358–1376, 2014.

Huan Xu and Shie Mannor. The robustness-performance tradeoff in markov decision processes. *Advances in Neural Information Processing Systems*, 19, 2006.

Huan Zhang, Hongge Chen, Chaowei Xiao, Bo Li, Mingyan Liu, Duane Boning, and Cho-Jui Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations. *Advances in Neural Information Processing Systems*, 33:21024–21037, 2020.

Huan Zhang, Hongge Chen, Duane Boning, and Cho-Jui Hsieh. Robust reinforcement learning on state observations with learned optimal adversary. In *International Conference on Learning Representation (ICLR)*, 2021.

Shangtong Zhang. Modularized implementation of deep rl algorithms in pytorch. https://github.com/ShangtongZhang/DeepRL, 2018.

## Appendix A.  Examples and Further Considerations

We provide here two examples to show how we can obtain limit scenarios $\Pi_0 = \Pi$ (any policy is maximally robust) or $\Pi_0 = \Pi_T$ (Example 1), and how for some MDPs the third inclusion in Theorem 4 is strict (Example 2).

**Example 1**  Consider the simple MDP in Figure 2. First, consider the reward function $R_1(x_1, \cdot, \cdot) = 10$, $R_1(x_2, \cdot, \cdot) = 0$. This produces a "dummy" MDP where all policies have the same reward sum. Then, $\forall T, \pi, V^{\langle \pi, T \rangle} = V^\pi$, and therefore we have $\Pi_D = \Pi_0 = \Pi$.
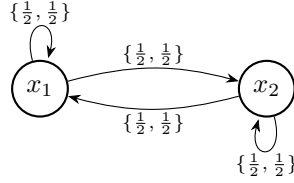


Figure 2: Example MDP. Values in brackets represent $\{P(\cdot, u_1, \cdot), P(\cdot, u_2, \cdot)\}$.

Now, consider the reward function $R_2(x_1, u_1, \cdot) = 10$, $R_2(\cdot, \cdot, \cdot) = 0$ elsewhere. Take a non-constant policy $\pi$, *i.e.*, $\pi(x_1) \neq \pi(x_2)$. In the example DOMDP (assuming the initial state is drawn uniformly from $X_0 = \{x_1, x_2\}$) one can show that at any time in the trajectory, there is a stationary probability $\Pr\{x_t = x_1\} = \frac{1}{2}$. Let us abuse notation and write $\pi(x_i) = (\ \pi(x_i, u_1) \quad \pi(x_i, u_2)\ )^\top$ and $R(x_i) = (\ R(x_i, u_1, \cdot) \quad R(x_i, u_2, \cdot)\ )^\top$. For the given reward structure we have $R(x_2) = (\ 0 \quad 0\ )^\top$, and therefore:

$$J(\pi) = E_{x_0 \sim \mu_0}\left[ \sum_{t=0}^\infty \gamma^t R_t \right] = \frac{1}{2}\langle R(x_1), \pi(x_1) \rangle \frac{\gamma}{1-\gamma}. \tag{9}$$

Since the transitions of the MDP are independent of the actions, following the same principle as in (9): $J\langle \pi, T \rangle = \frac{1}{2}\langle R(x_1), \langle \cdot, T\rangle(\pi)(x_1)\rangle \frac{\gamma}{1-\gamma}$. For any noise map $\langle \cdot, T \rangle \neq \mathrm{Id}$, for the two-state policy it holds that $\pi \notin \Pi_T \implies \langle \pi, T \rangle \neq \pi$. Therefore $\langle \pi, T \rangle(x_1) \neq \pi(x_1)$ and:

$$J(\pi) - J(\langle \pi, T \rangle) = \frac{5\gamma}{1-\gamma} \cdot \big( \pi(x_1, 1) - \langle \pi, T \rangle(x_1, 1) \big) \neq 0,$$

which implies that $\pi \notin \Pi_0$.

**Example 2**  Consider the same MDP in Figure 2 with reward function $R(x_1, u_1, \cdot) = R(x_2, u_1, \cdot) = 10$, and a reward of zero for all other transitions. Take a policy $\pi(x_1) = (1\ 0)$, $\pi(x_2) = (0\ 1)$. The policy yields a reward of 10 in state $x_1$ and a reward of 0 in state $x_2$. Again we assume the initial state is drawn uniformly from $X_0 = \{x_1, x_2\}$. Then, observe:

$$J(\pi) = E_{x_0 \sim \mu_0}\left[ \sum_{t=0}^\infty \gamma^t R_t \right] = \frac{1}{2}\langle R(x_1), \pi(x_1) \rangle \frac{\gamma}{1-\gamma} =$$
$$= \frac{1}{2}\frac{10\gamma}{1-\gamma} = \frac{5\gamma}{1-\gamma}.$$

Define now noise map $T(\cdot \mid x_1) = (\frac{1}{2} \frac{1}{2})$ and $T(\cdot \mid x_2) = (\frac{1}{2} \frac{1}{2})$. Observe this noise map yields a policy with non-zero disadvantage, $D^\pi(x_1, T) = \frac{5\gamma}{1-\gamma} - \left(\frac{5\gamma}{1-\gamma} - 2.5\right) = 2.5$ and similarly $D^\pi(x_2, T) = -2.5$, therefore $\pi \notin \Pi_D$. However, the policy *is maximally robust*:

$$J(\langle \pi, T \rangle) = \frac{1}{2}\langle R(x_1), \langle \pi, T \rangle(x_1) \rangle \frac{\gamma}{1-\gamma} +$$
$$+\frac{1}{2}\langle R(x_2), \langle \pi, T \rangle(x_2) \rangle \frac{\gamma}{1-\gamma} = \frac{1}{2}\frac{\gamma}{1-\gamma}(5+5) = \frac{5\gamma}{1-\gamma}. \tag{10}$$

Therefore, $\pi \in \Pi_0$.

## Appendix B. Theoretical Results

### B.1. Auxiliary Results

**Theorem 9 (Stochastic Approximation with Non-Expansive Operator)** *Let $\{\xi_t\}$ be a random sequence with $\xi_t \in \mathbb{R}^n$ defined by the iteration:*

$$\xi_{t+1} = \xi_t + \alpha_t(F(\xi_t) - \xi_t + M_{t+1}),$$

*where:*

1. *The step sizes $\alpha_t$ satisfy standard learning rate assumptions.*

2. *$F : \mathbb{R}^n \mapsto \mathbb{R}^n$ is a $\|\cdot\|_\infty$ non-expansive map. That is, for any $\xi_1, \xi_2 \in \mathbb{R}^n$, $\|F(\xi_1) - F(\xi_2)\|_\infty \leq \|\xi_1 - \xi_2\|_\infty$.*

3. *$\{M_t\}$ is a martingale difference sequence with respect to the increasing family of $\sigma-$fields $\mathcal{F}_t := \sigma(\xi_0, M_0, \xi_1, M_1, ..., \xi_t, M_t)$.*

*Then, the sequence $\xi_t \to \xi^*$ almost surely where $\xi^*$ is a fixed point such that $F(\xi^*) = \xi^*$.*

**Proof** See Borkar and Soumyanatha (1997). ∎

**Theorem 10 (PB-LRL Convergence with 2 objectives.(Skalse et al., 2022b))** *Let $\mathcal{M}$ be a multi-objective MDP with objectives $K_i$, $i \in \{1, 2\}$. Assume a policy $\pi$ is twice differentiable in parameters $\theta$, and if using a critic $V_i$ assume it is continuously differentiable on parameters $w_i$. Choose a tolerance $\epsilon$, and suppose that if PB-LRL is run for $T$ steps, there exists some limit point $w_i \to w_i^*(\theta)$ when $\theta$ is held fixed. If for both objectives there exists a gradient descent scheme such that $\lim_{T \to \infty} \mathbb{E}_t[\theta_t] \in \Theta_i^\epsilon$ then combining the objectives as in (1) yields $\lim_{T \to \infty} \mathbb{E}_t[\theta_t] \in \{\text{argmax}_\theta K_2(\theta), \text{s.t. } K_1(\theta) \geq K_1(\theta') - \epsilon\}$.*

**Proof** [Proof Sketch] We refer the interested reader to Skalse et al. (2022b) for a full proof, and here attempt to provide the intuition behind the result in the form of a proof sketch.

Let us begin by briefly recalling the general problem statement: we wish to take a multi-objective MDP $\mathcal{M}$ with $m$ objectives, and obtain a lexicographically optimal policy (one that optimises the first objective, and then subject to this optimises the second objective, and so on). More precisely, for a policy $\pi$ parameterised by $\theta$, we say that $\pi$ is (globally) *lexicographically $\epsilon$-optimal* if $\theta \in \Theta_m^\epsilon$, where $\Theta_0^\epsilon = \Theta$ is the set of all policies in $\mathcal{M}$, $\Theta_{i+1}^\epsilon := \{\theta \in \Theta_i^\epsilon \mid \max_{\theta' \in \Theta_i^\epsilon} K_i(\theta') - K_i(\theta) \leq \epsilon_i\}$, and $\mathbb{R}^{m-1} \ni \epsilon \succcurlyeq 0.$[6]

---

6. The proof in Skalse et al. (2022b) also considers *local* lexicographic optima, though for the sake of simplicity, we do not do so here.

The basic idea behind policy-based lexicographic reinforcement learning (PB-LRL) is to use a multi-timescale approach to first optimise $\theta$ using $K_1$, then at a slower timescale optimise $\theta$ using $K_2$ while adding the condition that the loss with respect to $K_1$ remains bounded by its current value, and so on. This sequence of constrained optimisations problems can be solved using a Lagrangian relaxation (Bertsekas, 1999), either in series or – via a judicious choice of learning rates – simultaneously, by exploiting a separation in timescales (Borkar, 2008). In the simultaneous case, the parameters of the critic $w_i$ (if using an actor-critic algorithm, if not this part of the argument may be safely ignored) for each objective are updated on the fastest timescale, then the parameters $\theta$, and finally (i.e., most slowly) the Lagrange multipliers for each of the remaining constraints.

The proof proceeds via induction on the number of objectives, using a standard stochastic approximation argument (Borkar, 2008). In particular, due to the learning rates chosen, we may consider those more slowly updated parameters fixed for the purposes of analysing the convergence of the more quickly updated parameters. In the base case where $m = 1$, we have (by assumption) that $\lim_{T \to \infty} \mathbb{E}_t[\theta] \in \Theta_1^\epsilon$. This is simply the standard (non-lexicographic) RL setting. Before continuing to the inductive step, Skalse et al. (2022b) observe that because gradient descent on $K_1$ converges to globally optimal stationary point when $m = 1$ then $K_1$ must be globally *invex* (where the opposite implication is also true) (Ben-Israel and Mond, 1986a).[7]

The reason this observation is useful is that because each of the objectives $K_i$ shares the same functional form, they are all invex, and furthermore, invexity is conserved under linear combinations and the addition of scalars, meaning that the Lagrangian formed in the relaxation of each constrained optimisation problem is also invex. As a result, if we assume that $\lim_{T \to \infty} \mathbb{E}_t[\theta] \in \Theta_i^\epsilon$ as our inductive hypothesis, then the stationary point of the Lagrangian for optimising objective $K_{i+1}$ is a global optimum, given the constraints that it does not worsen performance on $K_1, \ldots, K_i$. Via Slater's condition (Slater, 1950) and standard saddle-point arguments (Bertsekas, 1999; Paternain et al., 2019), we therefore have that $\lim_{T \to \infty} \mathbb{E}_t[\theta] \in \Theta_{i+1}^\epsilon$, completing the inductive step, and thus the overall inductive argument.

This concludes the proof that $\lim_{T \to \infty} \mathbb{E}_t[\theta] \in \Theta_m^\epsilon$. We refer the reader to Skalse et al. (2022b) for a discussion of the error $\epsilon$, but intuitively it corresponds to a combination of the representational power of $\theta$, the critic parameters $w_i$ (if used), and the duality gap due to the Lagrangian relaxation (Paternain et al., 2019). In cases where the representational power of the various parameters is sufficiently high, then it can be shown that $\epsilon = 0$. ∎

### B.2. On Adversarial Disturbances and other Noise Kernels

A problem that remains open after this work is what constitutes an appropriate choice of $\tilde{T}$, and what can we expect by restricting a particular class of $\tilde{T}$. We first discuss adversarial examples, and then general considerations on $\tilde{T}$ versus $T$.

**Adversarial Noise**  As mentioned in the introduction, much of the previous work focuses on adversarial disturbances. We did not directly address this in the results of this work since our motivation lies in the scenarios where the disturbance is not adversarial and is unknown. However, following the results of Section 3, we are able to reason about adversarial disturbances. Consider

---

7. A differentiable function $f : \mathbb{R}^n \to \mathbb{R}$ is (globally) invex if and only if there exists a function $g : \mathbb{R}^n \times \mathbb{R}^n \to \mathbb{R}^n$ such that $f(x_1) - f(x_2) \geq g(x_1, x_2)^\top \nabla f(x_2)$ for all $x_1, x_2 \in \mathbb{R}^n$ (Hanson, 1981).

JARNE ORNIA ROMAO HAMMOND MAZO JR. ABATE

an adversarial map $T_{adv}$ to be

$$\langle \pi, T_{adv} \rangle (x) = \pi(y), \quad y \in \text{argmax}_{y \in X_{ad}(x)} \, d\big(\pi(x), \pi(y)\big),$$

with $X_{ad}(x) \subseteq X$ being a set of admissible disturbance states for $x$, and $d(\cdot, \cdot)$ is a distance measure between distributions (*e.g.* 2-norm).

**Proposition 11** *Constant policies are a fixed point of $T_{adv}$, and are the only fixed points if for all pairs $x_0, x_k$ there exists a sequence $\{x_0, ..., x_k\} \subseteq X$ such that $x_i \in X_{ad}(x_i)$.*

**Proof** First, it is straight-forward that if $\overline{\pi} \in \overline{\Pi} \implies \langle \overline{\pi}, T_{adv} \rangle (x) = \overline{\pi}(x)$. To show they are the only fixed points, assume that there is a non-constant policy $\pi'$ that is a fixed point of $T_{ad}$. Then, there exists $x, z$ such that $\pi'(x) \neq \pi'(z)$. However, by assumption, we can construct a sequence $\{x, ..., z\} \subseteq X$ that connects $x$ and $z$ and every state in the sequence is in the admissible set of the previous one. Assume without loss of generality that this sequence is $\{x, y, z\}$. Then, if $\pi'$ is a fixed point, $\langle \pi', T_{adv} \rangle (x) = \pi'(x)$, $\langle \pi', T_{adv} \rangle (y) = \pi'(y)$ and $\langle \pi', T_{adv} \rangle (z) = \pi'(z)$. However, $\pi'(x) \neq \pi'(z)$, so either $\pi'(x) \neq \pi'(y) \implies d(\pi'(x), \pi'(y)) \neq 0$ or $\pi'(y) \neq \pi'(z) \implies d(\pi'(y), \pi'(z)) \neq 0$, therefore $\pi'$ cannot be a fixed point of $T_{adv}$. ∎

The main difference between an adversarial operator and the random noise considered throughout this work is that $T_{adv}$ is *not a linear operator*, and additionally, it is time varying (since the policy is being modified at every time step of the PG algorithm). Therefore, including it as a LRPG objective would invalidate the assumptions required for LRPG to retain formal guarantees of the original PG algorithm used, and it is not guaranteed that the resulting policy gradient algorithm would converge.

## Appendix C. Experiment Methodology

We use in the experiments well-tested implementations of A2C, PPO and SAC from Stable Baselines 3 (Raffin et al., 2021) to include the computation of the lexicographic parameters in (1). All experiments were run on an Ubuntu 18.04 system, with a 16 core CPU and a graphic card Nvidia GeForce 3060.

**LRPG Parameters.** The LRL parameters are initialised in all cases as $\beta_0^1 = 2$, $\beta_0^2 = 1$, $\lambda = 0$ and $\eta = 0.001$. The LRL tolerance is set to $\epsilon_t = 0.99\hat{k}_1$ to ensure we never deviate too much from the original objective, since the environments have very sparse rewards. We use a first order approximation to compute the LRL weights from the original LMORL implementation.

### C.1. Discrete Control

The discrete control environments used can be seen in Figure 3. Since all the environments use a pixel representation of the observation, we use a shared representation for the value function and policy, where the first component is a convolutional network, implemented as in Zhang (2018). The hyper-parameters of the neural representations are presented in Table 2.
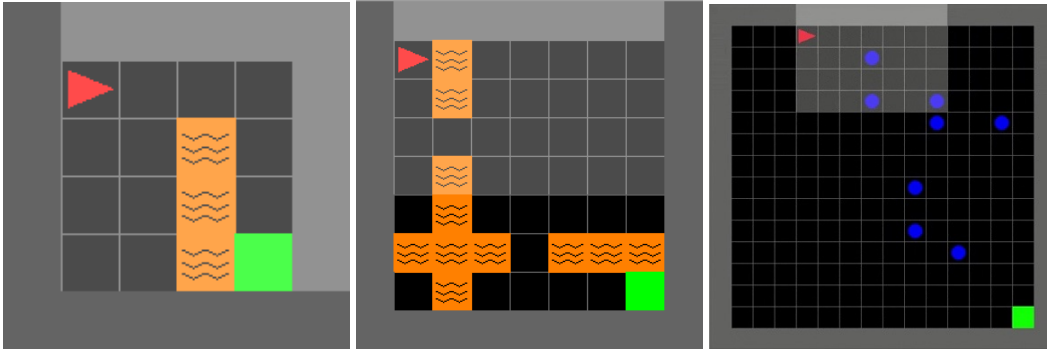
Figure 3: Screenshots of the environments used, from left: LavaGap, LavaCrossing and DynamicObstables.

| Layer | Output | Func. |
|-------|--------|-------|
| Conv1 | 16 | ReLu |
| Conv2 | 32 | ReLu |
| Conv3 | 64 | ReLu |

Table 2: Shared Observation Layers

The actor and critic layers, for both algorithms, are a fully connected layer with $64$ features as input and the corresponding output. We used in all cases an Adam optimiser. We optimised the parameters for each (vanilla) algorithm through a quick parameter search, and apply the same parameters for the Lexicographically Robust versions.

| | LavaGap | LavaCrossing | DynObs |
|-------|---------|--------------|--------|
| Parallel Envs | 16 | 16 | 16 |
| Steps | $2 \cdot 10^6$ | $2 \cdot 10^6$ | $8 \times 10^6$ |
| $\gamma$ | 0.99 | 0.99 | 0.98 |
| $\alpha$ | 0.00176 | 0.00176 | 0.00181 |
| $\epsilon$(Adam) | $10^{-8}$ | $10^{-8}$ | $10^{-8}$ |
| Grad. Clip | 0.9 | 0.9 | 0.5 |
| Gae | 0.95 | 0.95 | 0.95 |
| Rollout | 64 | 64 | 64 |
| E. Coeff | 0.01 | 0.014 | 0.011 |
| V. Coeff | 0.05 | 0.05 | 0.88 |

Table 3: A2C Parameters

|  | LavaGap | LavaCrossing | DynObs |
|---|---|---|---|
| Parallel Envs | 8 | 8 | 8 |
| Steps | $6 \cdot 10^6$ | $2 \cdot 10^6$ | $8 \times 10^5$ |
| $\gamma$ | 0.95 | 0.99 | 0.97 |
| $\alpha$ | 0.001 | 0.001 | 0.001 |
| $\epsilon$(Adam) | $10^{-8}$ | $10^{-8}$ | $10^{-8}$ |
| Grad. Clip | 1 | 1 | 0.1 |
| Ratio Clip | 0.2 | 0.2 | 0.2 |
| Gae | 0.95 | 0.95 | 0.95 |
| Rollout | 256 | 512 | 256 |
| Epochs | 10 | 10 | 10 |
| E. Coeff | 0 | 0.1 | 0.01 |

Table 4: PPO Parameters

For the implementation of the LRPG versions of the algorithms, in all cases we allow the algorithm to iterate for $1/3$ of the total steps before starting to compute the robustness objectives. In other words, we use $\hat{K}(\theta) = K_1(\theta)$ until $t = \frac{1}{3} \max\_steps$, and from this point we resume the lexicographic robustness computation as described in Algorithm 1. This is due to the structure of the environments simulated. The rewards (and in particular the positive rewards) are very sparse in the environments considered. Therefore, when computing the policy gradient steps, the loss for the primary objective is practically zero until the environment is successfully solved at least once. If we implement the combined lexicographic loss from the first time step, many times the algorithm would converge to a (constant) policy without exploring for enough steps, leading to convergence towards a maximally robust policy that does not solve the environment.

**Noise Kernels.** We consider two types of noise; a normal distributed noise $\tilde{T}^g$ and a uniform distributed noise $\tilde{T}^u$. For the environments LavaGap and DynamicObstacles, the kernel $\tilde{T}^u$ produces a disturbed state $\tilde{x} = x + \xi$ where $\|\xi\|_\infty \leq 2$, and for LavaCrossing $\|\xi\|_\infty \leq 1.5$. The normal distributed noise is in all cases $\mathcal{N}(0, 0.5)$. The maximum norm of the noise is quite large, but this is due to the structure of the observations in these environments. The pixel values are encoded as integers $0 - 9$, where each integer represents a different feature in the environment (empty space, doors, lava, obstacle, goal...). Therefore, any noise $\|\xi\|_\infty \leq 0.5$ would most likely not be enough to *confuse* the agent. On the other hand, too large noise signals are unrealistic and produce pathological environments. All the policies are then tested against two "true" noise kernels, $T_1 = \tilde{T}^u$ and $T_2 = \tilde{T}^g$. The main reason for this is to test both the scenarios where we assume a *wrong* noise kernel, and the case where we are training the agents with the correct kernel.

**Comparison with SA-PPO.** One of the baselines included is the State-Adversarial PPO algorithm proposed in Zhang et al. (2020). The implementation includes an extra parameter that multiplies the regularisation objective, $k_{ppo}$. Since we were not able to find indications on the best parameter for discrete action environments, we implemented $k_{ppo} \in \{0.1, 1, 2\}$ and picked the best result for each entry in Table 1. Larger values seemed to de-stabilise the learning in some cases. The rest of the parameters are kept as in the vanilla PPO implementation.

| | PPO on MiniGrid Environments | | | | A2C on MiniGrid Environments | | | |
|---|---|---|---|---|---|---|---|---|
| Noise | Vanilla | $\text{LR}_{\text{PPO}}(K_T^u)$ | $\text{LR}_{\text{PPO}}(K_T^g)$ | SA-PPO ‖ | Vanilla | $\text{LR}_{\text{A2C}}(K_T^u)$ | $\text{LR}_{\text{A2C}}(K_T^g)$ | $\text{LR}_{\text{A2C}}(K_D)$ |
| *LavaGap* | | | | | | | | |
| ∅ | **0.95±0.003** | **0.95±0.075** | **0.95±0.101** | 0.94±0.068 | **0.94±0.004** | **0.94±0.005** | **0.94±0.003** | **0.94±0.006** |
| $T_1$ | 0.80±0.041 | **0.95±0.078** | 0.93±0.124 | 0.88±0.064 | 0.83±0.061 | **0.93±0.019** | 0.89±0.032 | 0.91±0.088 |
| $T_2$ | 0.92±0.015 | **0.95±0.052** | **0.95±0.094** | 0.93±0.050 | 0.89±0.029 | **0.94±0.008** | 0.93±0.011 | 0.93±0.021 |
| $T_{adv}^{0.5}$ | 0.56±0.194 | **0.93±0.101** | 0.91±0.076 | 0.90±0.123 | 0.92±0.034 | **0.94±0.003** | **0.94±0.007** | 0.93±0.015 |
| $T_{adv}^1$ | 0.20±0.243 | **0.90±0.124** | 0.68±0.190 | **0.90±0.135** | 0.75±0.123 | **0.94±0.006** | 0.92±0.038 | 0.88±0.084 |
| $T_{adv}^2$ | 0.01±0.051 | 0.71±0.251 | 0.21±0.357 | **0.87±0.116** | 0.27±0.119 | **0.79±0.069** | 0.68±0.127 | 0.56±0.249 |
| *LavaCrossing* | | | | | | | | |
| ∅ | **0.95±0.023** | 0.93±0.050 | 0.93±0.018 | 0.88±0.091 | 0.91±0.024 | 0.91±0.063 | 0.90±0.017 | **0.92±0.034** |
| $T_1$ | 0.50±0.110 | **0.92±0.053** | 0.89±0.029 | 0.64±0.109 | 0.66±0.071 | **0.78±0.111** | 0.72±0.073 | 0.76±0.098 |
| $T_2$ | 0.84±0.061 | **0.92±0.050** | **0.92±0.021** | 0.85±0.094 | 0.78±0.054 | 0.83±0.105 | 0.86±0.029 | **0.87±0.063** |
| $T_{adv}^{0.5}$ | 0.29±0.098 | **0.91±0.081** | **0.91±0.054** | 0.87±0.045 | 0.56±0.039 | 0.51±0.089 | 0.43±0.041 | **0.68±0.126** |
| $T_{adv}^1$ | 0.03±0.022 | 0.83±0.122 | 0.86±0.132 | **0.87±0.059** | 0.27±0.158 | 0.25±0.118 | 0.17±0.067 | **0.43±0.060** |
| $T_{adv}^2$ | 0.0±0.004 | 0.50±0.171 | 0.38±0.020 | **0.82±0.072** | 0.06±0.056 | 0.04±0.030 | 0.01±0.008 | **0.09±0.060** |
| *DynamicObstacles* | | | | | | | | |
| ∅ | **0.91±0.002** | **0.91±0.008** | **0.91±0.007** | **0.91±0.131** | **0.91±0.011** | 0.88±0.020 | 0.89±0.009 | **0.91±0.013** |
| $T_1$ | 0.23±0.201 | **0.77±0.102** | 0.61±0.119 | 0.45±0.188 | 0.27±0.104 | 0.43±0.108 | 0.45±0.162 | **0.56±0.270** |
| $T_2$ | 0.50±0.117 | **0.75±0.075** | 0.70±0.072 | 0.68±0.490 | 0.45±0.086 | 0.53±0.109 | 0.52±0.161 | **0.67±0.203** |
| $T_{adv}^{0.5}$ | 0.74±0.230 | 0.89±0.118 | 0.85±0.061 | **0.90±0.142** | 0.46±0.214 | 0.55±0.197 | 0.51±0.371 | **0.62±0.249** |
| $T_{adv}^1$ | 0.26±0.269 | 0.79±0.157 | 0.68±0.144 | **0.84±0.150** | 0.19±0.284 | **0.35±0.197** | 0.23±0.370 | 0.10±0.379 |
| $T_{adv}^2$ | -0.49±0.312 | 0.51±0.234 | 0.33±0.202 | **0.55±0.170** | -0.54±0.209 | -0.21±0.192 | -0.53±0.261 | **-0.51±0.260** |

Table 5: Extended Reward Results.

### C.1.1. EXTENDED RESULTS: ADVERSARIAL DISTURBANCES

Even though we do not use an adversarial attacker or disturbance in our reasoning through this work, we implemented a policy-based state-adversarial noise disturbance to test the benchmark algorithms against, and evaluate how well each of the methods reacts to such adversarial disturbances.

**Adversarial Disturbance**  We implement a bounded policy-based adversarial attack, where at each state $x$ we maximise for the KL divergence between the disturbed and undisturbed state, such that the adversarial operator is:

$$T_{adv}^{\varepsilon}(y \mid x) = 1 \implies y \in \underset{\tilde{x}}{\arg\max} \, D_{KL}(\pi(x), \pi(\tilde{x}))$$
$$s.t. \; \|x - \tilde{x}\|_2 \leq \varepsilon.$$

The optimisation problem is solved at every point by using a Stochastic Gradient Langevin Dynamics (SGLD) optimiser. The results are presented in Table 5.

This type of adversarial attack with SGLD optimiser was proposed in Zhang et al. (2020). As one can see, the adversarial disturbance is quite successful at severely lowering the obtained rewards in all scenarios. Additionally, as expected SA-PPO was the most effective at minimizing the disturbance effect (as it is trained with adversarial disturbances), although LRPG produces reasonably robust policies against this type of disturbances as well. At last, A2C appears to be much more sensitive to adversarial disturbances than PPO, indicating that the policies produced by PPO are by default more robust than A2C.

### C.2. Continuous Control

The continuous control environments simulated are MountainCar, LunarLander and BipedalWalker. The policies used are in all cases MLP policies with ReLU gates and a $(64, 64)$ feature extractor plus a fully connected layer to output the values and actions unless stated otherwise. The hyperparameters can be found in tables C.2 and 8. The implementation is based on Stable Baselines 3
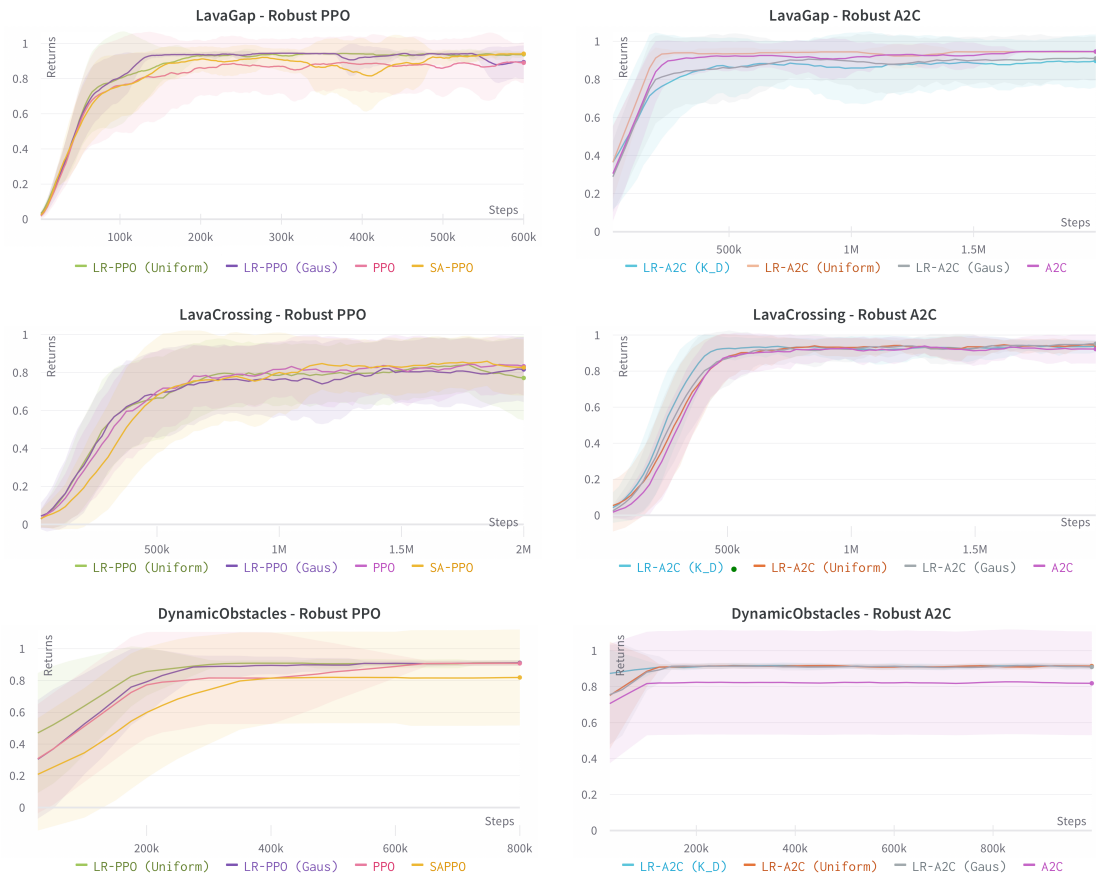
Figure 4: Learning Plots for Discrete Control Environments.

| Noise | *PPO on Continuous Environments* | | | | | *SAC on Continuous Environments* | | |
| | Vanilla | LR$_{\text{PPO}}$ ($K^u_T$) | LR$_{\text{PPO}}$ ($K^g_T$) | SA-PPO | | Vanilla | LR$_{\text{SAC}}$ ($K^u_T$) | LR$_{\text{SAC}}$ ($K^g_T$) |
|---|---|---|---|---|---|---|---|---|
| *MountainCar* | | | | | | | | |
| $\emptyset$ | **94.77±0.26** | 93.17±0.89 | 94.66±1.61 | 88.69±3.93 | | 93.52±0.05 | **94.43±0.19** | 93.84±0.05 |
| $T_1$ | 88.67±1.41 | 91.46±1.22 | **94.91±1.35** | 88.41±3.99 | | 1.89±65.31 | 71.81±13.04 | **76.90±7.11** |
| $T_2$ | 92.22±1.11 | 92.40±1.28 | **94.76±1.42** | 89.32±3.79 | | -27.82±73.10 | **72.93±8.57** | 69.41±13.03 |
| *LunarLander* | | | | | | | | |
| $\emptyset$ | 267.99±38.04 | **269.76±22.93** | 243.08±37.03 | 220.18±98.78 | | 268.96±51.52 | 275.17±14.04 | **282.24±15.95** |
| $T_1$ | 156.09±22.87 | **280.91±20.34** | 182.80±49.26 | 164.53±45.48 | | 128.18±17.73 | **187.64±76.30** | 153.81±33.16 |
| $T_2$ | 158.02±46.57 | **276.76±16.20** | 212.62±37.56 | 221.84±73.61 | | 140.92±20.61 | **187.82±25.27** | 158.18±28.60 |
| *BipedalWalker* | | | | | | | | |
| $\emptyset$ | 265.39±82.36 | 261.39±83.19 | **276.66±44.85** | 251.60±103.08 | | 236.39±157.03 | 302.56±70.79 | **313.56±52.17** |
| $T_1$ | 174.15±170.30 | 253.56±72.66 | 220.28±118.61 | **264.69±61.63** | | 203.93±167.83 | 241.45±124.54 | **241.60±139.93** |
| $T_2$ | 135.16±182.30 | 243.27±89.86 | **265.37±80.60** | 255.21±90.61 | | 84.10±198.12 | 198.20±151.64 | **229.75±166.87** |

Table 6: Reward values gained by LRPG and baselines on continuous control tasks.

| | MountainCarContinuous | LunarLanderContinuous | BipedalWalker-v3 |
|---|---|---|---|
| Parallel Envs | 1 | 16 | 32 |
| Steps | $2 \times 10^4$ | $1 \times 10^6$ | $5 \times 10^6$ |
| $\gamma$ | 0.9999 | 0.999 | 0.999 |
| $\alpha$ | $3 \times 10^{-4}$ | $3 \times 10^{-4}$ | $3 \times 10^{-4}$ |
| Grad. Clip | 5 | 0.5 | 0.5 |
| Ratio Clip | 0.2 | 0.2 | 0.18 |
| Gae | 0.9 | 0.98 | 0.95 |
| Epochs | 10 | 4 | 10 |
| E. Coeff | 0.00429 | 0.01 | 0 |

Table 7: PPO Parameters for Continuous Control

(Raffin et al., 2021) tuned algorithms. **Noise Kernels.** We consider again two types of noise; a normal distributed noise $\tilde{T}^g$ and a uniform distributed noise $\tilde{T}^u$. In all cases, algorithms are implemented with a state observation normalizer. That is, assimptotically all states will be observed to be in the set $(-1, 1)$. For this reason, the uniform noise is bounded at lower values than for the discrete control environments. For BipedalWalker $\|\xi\|_\infty \leq 0.05$ and for Lunarlander and MountainCar $\|\xi\|_\infty \leq 0.1$. Larger values were shown to destabilize learning.

| | MountainCarContinuous | LunarLanderContinuous | BipedalWalker-v3 |
|---|---|---|---|
| Steps | $5 \times 10^4$ | $5 \times 10^5$ | $5 \times 10^5$ |
| $\gamma$ | 0.9999 | 0.99 | 0.98 |
| $\alpha$ | $3 \times 10^{-4}$ | $7.3 \times 10^{-4}$ | $7.3 \times 10^{-4}$ |
| $\tau$ | 0.01 | 0.01 | 0.01 |
| Train Freq. | 32 | 1 | 64 |
| Grad. Steps | 32 | 1 | 64 |
| MLP Arch | (64,64) | (400,300) | (400,300) |

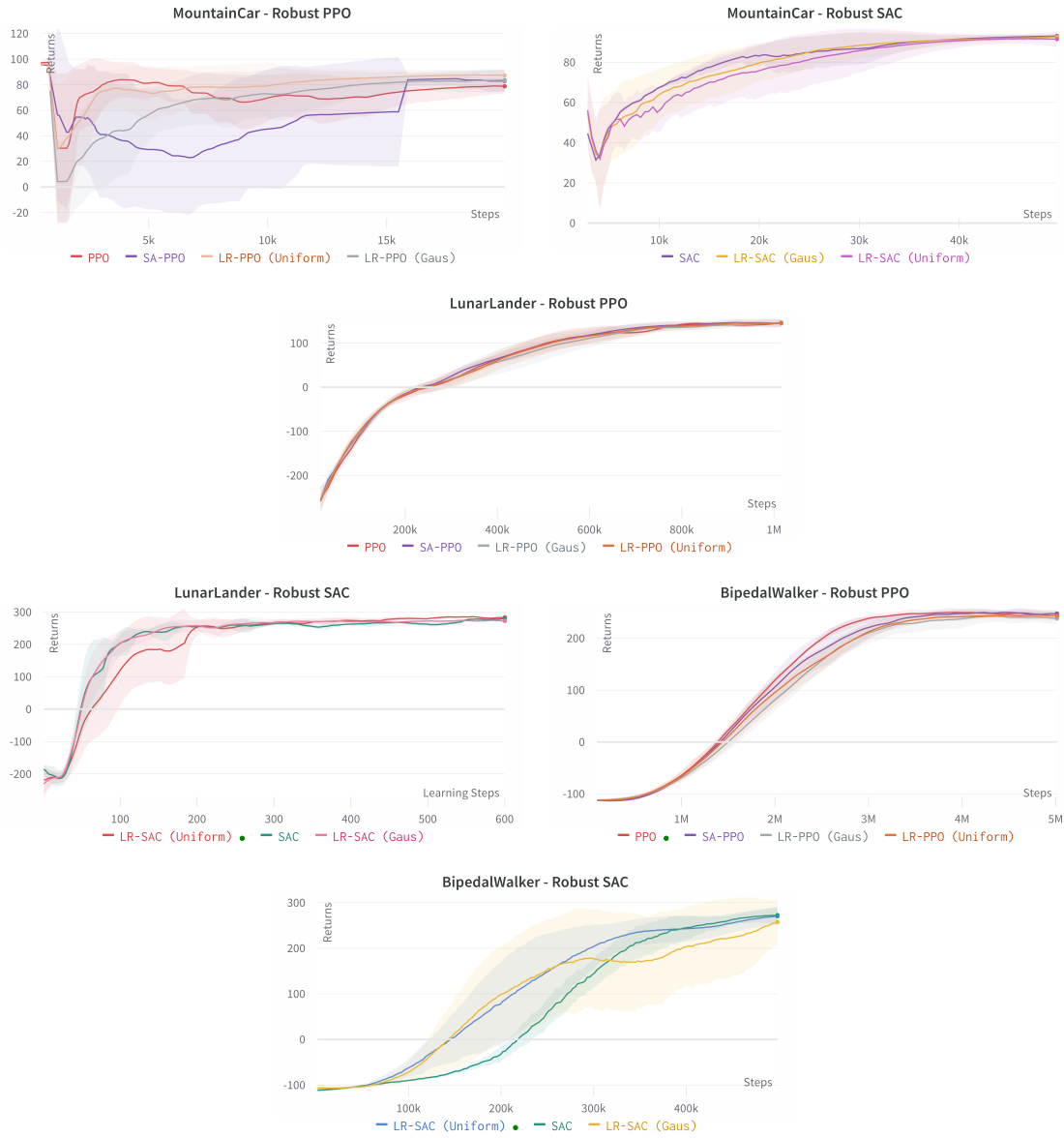Table 8: SAC Parameters for Continuous Control

Figure 5: Learning Plots for Continuous Control Environments.

**Learning processes**   In general, learning was not severlely affected by the LRPG scheme. However, it was shown to induce a larger variance in the trajectories observed, as seen in LunarLander with LR-SAC and BipedalWalker with LR-SAC.