
Discovered Policy Optimisation

Anonymous Authors¹

Abstract

The last decade has been revolutionary for *reinforcement learning* (RL) — it can now solve complex decision and control problems. Successful RL methods were handcrafted using mathematical derivations, intuition, and experimentation. This approach has a major shortcoming—it *results in specific solutions to the RL problem, rather than a protocol for discovering efficient and robust methods*. In contrast, the emerging field of *meta-learning* provides a toolkit for automatic *machine learning* method optimisation, potentially addressing this flaw. However, black-box approaches which attempt to discover RL algorithms with minimal prior structure have thus far not been successful. *Mirror Learning*, which includes RL algorithms, such as PPO, offers a potential framework. In this paper we explore the Mirror Learning space by meta-learning a “drift” function. We refer to the result as *Learnt Policy Optimisation* (LPO). By analysing LPO we gain original insights into policy optimisation which we use to formulate a novel, closed-form RL algorithm, *Discovered Policy Optimisation* (DPO). Our experiments in *Brax* environments confirm state-of-the-art performance of LPO and DPO, as well as their transfer to unseen settings.

1. Introduction

Recent advancements in deep learning have allowed reinforcement learning (RL) algorithms to successfully tackle large-scale problems. Early deep RL algorithms, such as A2C (Mnih et al., 2016) and DDPG (Lillicrap et al., 2015) follow the scheme of *generalized policy iteration* (GPI), where the RL agent alternates between policy evaluation and policy improvement phases (Sutton and Barto, 2018). Namely, they train a *critic* neural network that estimates

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

the value function, and an *actor* network which models the agent’s policy, and use these to estimate the update direction with policy gradients of the expected return. Another approach is that of *trust-region learning* (TRL) methods, which construct a surrogate objective that they optimise within a small region around the current policy (Schulman et al., 2015), aiming at guaranteeing stability of the policy iterates. Indeed, TRL-inspired algorithms, such as TRPO (Schulman et al., 2015) and PPO (Schulman et al., 2017), have been among the most widely used methods (Berner et al., 2019) and are known for their performance and stability. Nevertheless, although these research threads have delivered a handful of successful techniques, their design relies on concepts handcrafted by humans, rather than discovered in a *learning process*. As a possible consequence, these methods often suffer from various flaws, such as the brittleness to hyperparameter settings (Schulman et al., 2015; Haarnoja et al., 2018a), and a lack of optimality guarantees.

Unfortunately, the most promising alternative approach, *algorithm discovery*, thus far has been a tough nut to crack. First of all, techniques that enable optimisation of algorithm components in the *outer loop* process—meta RL (Schmidhuber, 1995; Finn et al., 2017)—are at an early stage of development, and are vastly more compute intensive than the RL problem in the *inner loop* (Xu et al., 2018). Second, to discover more complex algorithms, one needs to solve a more sophisticated, and computationally taxing meta-optimisation problem. As a result, most methods discover simple algorithms at best (Oh et al., 2020).

Recently, *Mirror Learning* (Kuba et al., 2022), a new theoretical framework, introduced an infinite space of provably correct algorithms, all of which share the same template. In a nutshell, a Mirror Learning algorithm is defined by four attributes, but in this work we focus on the *drift function*. A drift function guides the agent’s update, usually by penalising large changes. Any mirror learning algorithm provably achieves monotonic improvement of the return, and converges to an optimal policy (Kuba et al., 2022). Popular RL methods such as TRPO and PPO are instances of Mirror Learning.

In this paper, we use meta-learning to *discover* a new state-of-the-art (SOTA) RL algorithm within the Mirror Learning space. Our algorithm thus inherits theoretical conver-

gence guarantees by construction. Specifically we parameterise a drift function with a neural network, which we then meta-train using *evolution strategies* (Salimans et al., 2017, ES). The outcome of this meta-training is a specific mirror-learning algorithm which we name *Learnt Policy Optimisation* (LPO). However, we take a step further and learn from LPO about important RL concepts that it has discovered by visualising and analysing its optimisation objective.

Building upon these insights we propose a new, closed-form algorithm which we name —*Discovered Policy Optimisation* (DPO). We evaluate LPO and DPO in the *Brax* (Freeman et al., 2021) continuous control environments, where they obtain superior performance compared to PPO. Importantly, both LPO and DPO generalise to environments that were not used for training.

Related Work For an in-depth discussion on related work, we point the reader to Appendix D. Next, we continue by introducing the basic concepts and algorithms to understand our contributions.

2. Background

Please find details of the RL and Meta-RL problem formulations, as well as the Evolution Strategies frameworks for solving them, in Appendix E.

2.1. Mirror Learning

A mirror-learning agent (Kuba et al., 2022), in addition to value functions, has access to the following operators: the *drift function* $\mathfrak{D}_{\pi_k}(\pi|s)$ which, intuitively, evaluates the significance of change from policy π_k to π at state s ; the *neighbourhood operator* $\mathcal{N}(\pi_k)$ which forms a region around the policy π_k ; as well as sampling and drift distributions $\beta_{\pi_k}(s)$ and $\nu_{\pi_k}^\pi(s)$ over states. With these defined, a mirror-learning algorithm updates an agent’s policy by maximising the mirror objective as follows:

$$\pi_{k+1} = \arg \max_{\pi \in \mathcal{N}(\pi_k)} \left\{ \mathbb{E}_{s \sim \beta_{\pi_k}} [A_{\pi_k}(s, a)] - \mathbb{E}_{s \sim \nu_{\pi_k}^\pi} [\mathfrak{D}_{\pi_k}(\pi|s)] \right\}. \quad (1)$$

If, for all policies π and $\bar{\pi}$, the drift function satisfies the following conditions:

1. It is non-negative everywhere and zero at identity $\mathfrak{D}_{\pi_k}(\pi|s) \geq \mathfrak{D}_{\pi_k}(\pi_k|s) = 0$,
2. Its gradient with respect to π is zero at $\pi = \pi_k$,

then the Mirror Learning algorithm attains the monotonic improvement property, $\eta(\pi_{k+1}) \geq \eta(\pi_k)$, and converges

to the optimal return, $\eta(\pi_k) \rightarrow \eta(\pi^*)$, as $k \rightarrow \infty$ (Kuba et al., 2022). A mirror-learning agent can be implemented in practice by specifying functional forms of the drift function and neighbourhood operator, and parameterising the policy of the agent with a neural network, π_θ . As such, the agent approximates the objective in Equation (1) by sample averages, and maximises it with an optimisation method, like gradient ascent. It can be shown that PPO is a valid instance of Mirror Learning with drift function

$$\mathfrak{D}_{\pi_k}^{\text{PPO}}(\pi|s) \triangleq \mathbb{E}_{a \sim \pi_k} [f^{\text{PPO}}], \quad (2)$$

$$f^{\text{PPO}} = \text{ReLU} \left([r(\pi) - \text{clip}(r(\pi), 1 \pm \epsilon)] A_{\pi_k}(s, a) \right).$$

Although in its formulation, PPO puts no explicit constraints on its update size (Schulman et al., 2017), as its maximisation oracle (see Equation (1)) is N steps of gradient ascent, with learning rate α and gradient clipping threshold c , it implicitly employs a neighbourhood of an Euclidean ball of radius $N\alpha c$ around θ_k .

2.2. Evolution Strategies

Please see Appendix E.1

3. Methods

Our overall approach is to *meta-learn* a drift function to perform policy optimisation over a fixed episode length K . Hence, our meta-objective is the expected final return

$$F(\phi) = \mathbb{E}[\eta(\pi_K)|\phi].$$

In Subsection 3.1 we describe the parameterisation of our learnt drift function, while in Subsection 3.2 we provide a detailed description of our meta-learning approach.

3.1. Drift Function Network

The drift function that we learn takes form $\mathfrak{D}_{\pi_k}(\pi|s) = \mathbb{E}_{a \sim \pi_k} [f_\phi(\mathbf{x})|s]$, where $f_\phi(\mathbf{x})$ as a fully-connected neural network parameterised by ϕ . Our drift network is a function of the probability ratio between a candidate and the old policy, $r = \pi(a|s)/\pi_k(a|s)$, and of the advantage $A = A_{\pi_k}(s, a)$ (which we assume to be normalised across each batch). To ease learning complicated mappings, we include non-linear transformations of these arguments, ultimately forming the following input:

$$\mathbf{x}_{r,A} = [(1-r), (1-r)^2, (1-r)A, (1-r)^2A, \log(r), \log(r)^2, \log(r)A, \log(r)^2A].$$

In order to guarantee that the neural network is a valid drift function, it suffices to impose that $f_\phi(\mathbf{x}_{r,A}) = 0$ and $\nabla_r f_\phi(\mathbf{x}_{r,A}) = \mathbf{0}$ whenever $r = 1$, and $f_\phi(\mathbf{x}_{r,A}) \geq 0$ everywhere. As in our model, $\mathbf{x}_{r,A} = 0$ whenever $r =$

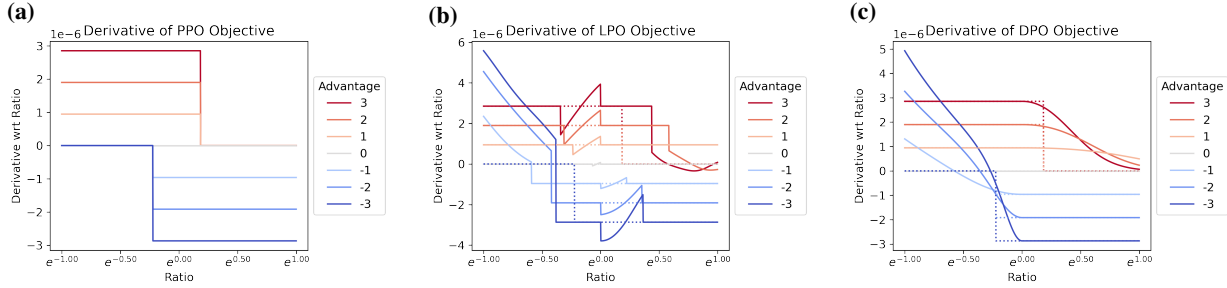


Figure 1: Objective visualisation: Comparing the derivatives at fixed advantage values across policy ratios. (a) - PPO, (b) - LPO, and (c) - DPO. See Appendix B for heat maps of the ratio derivatives.

1, the former condition is guaranteed by excluding bias terms from the network architecture. To meet the latter two conditions, we apply the ReLU activation at the last layer with a slight shift, $x \mapsto \text{ReLU}(x - \xi)$, where $\xi = 10^{-6}$.

3.2. Meta-Training the Drift Function Network

The meta-objective we optimise is the performance of the learner policy at the end of training: $F(\phi) = \mathbb{E}[\eta(\pi_{\theta_K}) | \text{alg}_\phi]$, where θ_K is the K^{th} (last) iterate of the RL training under the mirror-learning algorithm alg_ϕ . The expectation is taken over the randomness of the initial parameter θ_0 and stochasticity of the environment. We solve this problem using evolution strategies. At each generation (outer loop iteration), we sample a batch of perturbations of ϕ , initialise the policy parameters θ_0 , and then train the policy under alg_ϕ , using the drift function’s parameter ϕ , for K iterations. At the end of the inner-loop training, we estimate the return of the final policy ϕ_{θ_K} , and use it to estimate the gradient of $\nabla_\phi F(\phi)$ as in Equation (4).

We meta-learn the drift function and evaluate policies trained by it in the *Brax* (Freeman et al., 2021) physics simulator environments. We implement our method on top of the *Brax* version of PPO, which provides a Mirror Learning-friendly code template, keeping the policy architecture and training hyperparameters unchanged. For meta-training we use both *evosax* (Lange, 2022) and the *Learned_optimization* (Metz et al., 2022) libraries.

4. Empirical Studies

We consider two different meta-training setups. First, as described in Subsection E.3 we attempt to learn a drift function completely from scratch to investigate how similar it is to existing algorithms like PPO. Second, in Subsection 4.1 we ask whether we can learn a drift function that successfully generalises to multiple environments, if it is initialised near PPO. See Appendix E.3 for drift learning from scratch.

4.1. Learning with the PPO Initialisation

In this setting, f_ϕ is a small neural network, with a single hidden layer with 128 neurons, with bias terms removed, and tanh activation function. We add PPO to the output of the last hidden layer before passing it to the shifted ReLU,

$$f_\phi(\mathbf{x}_{r,A}) = \text{ReLU}\left(\tilde{f}_\phi(\mathbf{x}_{r,A}) - \xi + \text{ReLU}\left([r - \text{clip}(r, 1 \pm \epsilon)] A\right)\right),$$

where $\tilde{f}_\phi(\mathbf{x}_{r,A})$ is the output of the last hidden layer of the drift network. As such, the resulting drift function is similar to that of PPO at initialisation.

Surprisingly, we have found that meta-training in a **single** environment is sufficient to generate drift functions whose abilities transfer to unseen tasks. Moreover, we found that the learnt drifts generally display similar characteristics. For readability, we chose the drift function that was trained on *Ant*, whose induced algorithm we refer to as *Learnt Policy Optimisation* (LPO).

The results on Appendix A Figure 2 show that LPO, trained only on *Ant*, outperforms PPO in unseen environments. Furthermore, the *Brax* PPO implementation uses different hyperparameters, such as the number of update epochs and the total number of timesteps, for each of the tasks. This means that LPO, which was trained on *Ant* with hyperparameters associated to it, is robust not only against new environments, but also against new hyperparameters. We visualise the derivative of the LPO loss in Figure 5, which enables us to derive an analytical version of it in Section 5.

5. Analysis of LPO

In this section, we analyse the two key features that are consistently learnt and contribute most to LPO’s performance. We then interpret their effect on *policy entropy*, and the *update asymmetry* discovered by LPO, through which it differs largely from PPO (see Appendix B Figure 5 for visualisation). This analysis often refers to the heatmaps of the learned objectives, which can be found in Appendix B.

Rollback for negative advantage. In the bottom-left quadrant of the heat map, which corresponds to $A < 0$ and $r < 1$, we observe that the ratio derivative of the LPO objective is positive in a large region, roughly corresponding to $r < 1 - \epsilon$. This implies that actions that fall into this quadrant, although seemingly not appealing, are encouraged to be taken by the agent, which can be interpreted as a form of *rollback* (Wang et al., 2020). Hence, LPO learns to decrease r down to $1 - \epsilon$, but unlike PPO, encourages r to stay precisely around that value. By doing so, LPO prevents the agent from giving up on actions that appear poor at the moment, and encourages it to keep exploring them at a moderate frequency.

Cautious optimism for positive advantage. The upper-right quadrant corresponds to $A > 0$ and $r > 1$, which is induced by actions that seem the most appealing to update to. Nevertheless, LPO is cautious in doing so, gradually decreasing the pace of its update towards them, and eventually abstaining from chasing the most extreme advantage values—these may come from critic errors. We want to highlight that this view of LPO on positive advantages is different than that of PPO, which simply removes any incentive from updating towards actions with $r > 1 + \epsilon$, and thus can be viewed as *optimistic* relative to PPO.

Implicit entropy maximisation. Together, these two central features of LPO encourage the agent to spread its policy probability mass moderately over all actions, thus leading to larger entropy and allowing for richer exploration. Thus, LPO *has implicitly discovered entropy maximisation*, which we demonstrate in Appendix C Figure 7.

We provide more analysis in Appendix F.

6. Discovered Policy Optimisation: A New RL Algorithm Inspired By LPO

In this section, advancing concepts that LPO has discovered, we introduce a novel algorithm—Discovered Policy Optimization (DPO).

6.1. The Discovered Drift Function Model

Combining the key features identified in Section 5, we construct a closed-form model of LPO that can easily be implemented with just a few lines of code on top of an existing PPO implementation. We name the algorithm *Discovered Policy Optimisation* (DPO) because we have not derived it—it was instead discovered in the meta-learning process. DPO is a mirror-learning algorithm, with a drift function that takes different functional forms, depending on the sign of advantage A , as dictated by the update asymmetry principle from the previous section. Specifically, we have found

that the (parameter-free) drift function $f(r, A)$ given by

$$\begin{cases} \text{ReLU}((r - 1)A - \alpha \tanh((r - 1)A/\alpha)) & A \geq 0 \\ \text{ReLU}(\log(r)A - \beta \tanh(\log(r)A/\beta)) & A < 0 \end{cases}$$

faithfully reproduces the key features of LPO (cautious optimism and rollback) for appropriate constants $\alpha = 2$ and $\beta = 0.6$ (see Appendix E.2 for verification of the drift conditions). We visualise DPO in Figure 6 and note that even the “crossing-over” of gradient slices of LPO on Figure 5 is faithfully reproduced.

6.2. Results

We compare DPO to PPO and LPO on all Brax environments with provided hyperparameters in Appendix A Figure 2. LPO and DPO significantly outperform PPO on most environments. We use the PPO implementation provided by Brax, which we enhanced with advantage normalisation as we observed it to improve performance across the majority of the environments. Our methods also use this implementation technique. While evaluating DPO, similarly to LPO, we do not re-tune *any* hyperparameters that were originally selected for PPO in Brax. The results on Figure 2 reveal that DPO matches the performance of LPO and outperforms PPO on the evaluated environments, *despite being a two-line analytic model of LPO based on two key features*. This enables RL practitioners to implement DPO as easily as PPO with a performance on par with our best learnt drift function.

7. Conclusion

In this paper, we performed *algorithm discovery* by restricting our meta-learning to the space of valid *Mirror Learning* algorithms. Specifically, we optimised a drift function parameterised by a neural network, which we trained with *Evolution Strategies*. We consider this work to be an example of the new, promising paradigm of RL algorithm discovery. Namely, our strategy was to develop a high-performing RL algorithm by combining theoretical insights with large-scale computational techniques. As a result of the training, we obtained a theoretically sound method that we named *Learnt Policy Optimisation* (LPO), which outperforms a state-of-the-art baseline (PPO) in unseen environments, and with unseen hyperparameter settings. After analysing the learned features discovered by LPO, we introduced *Discovered Policy Optimisation* (DPO)—a closed-form approximation to LPO. Our experimental results show that DPO matches LPO in performance and robustness to hyperparameters. In the future, we plan to expand the variety of inputs to the learnt drift function, as well as to meta-learn other attributes of mirror learning. We expect these advancements to provide more insights into policy optimisation, ultimately resulting in more robust and better performing RL algorithms.

References

- Ferran Alet, Martin F. Schneider, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. Meta-learning curiosity algorithms. *CoRR*, abs/2003.05325, 2020. URL <https://arxiv.org/abs/2003.05325>.
- Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemyslaw Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- John D Co-Reyes, Yingjie Miao, Daiyi Peng, Esteban Real, Sergey Levine, Quoc V Le, Honglak Lee, and Aleksandra Faust. Evolving reinforcement learning algorithms. *arXiv preprint arXiv:2101.03958*, 2021.
- Yan Duan, John Schulman, Xi Chen, Peter L. Bartlett, Ilya Sutskever, and Pieter Abbeel. RL²: Fast reinforcement learning via slow reinforcement learning. *CoRR*, abs/1611.02779, 2016. URL <http://arxiv.org/abs/1611.02779>.
- Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Vlad Mnih, Tom Ward, Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *International Conference on Machine Learning*, pages 1407–1416. PMLR, 2018.
- Xidong Feng, Oliver Slumbers, Ziyu Wan, Bo Liu, Stephen McAleer, Ying Wen, Jun Wang, and Yaodong Yang. Neural auto-curricula in two-player zero-sum games. *Advances in Neural Information Processing Systems*, 34, 2021.
- Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017.
- C Daniel Freeman, Erik Frey, Anton Raichuk, Sertan Girgin, Igor Mordatch, and Olivier Bachem. Brax—a differentiable physics engine for large scale rigid body simulation. *arXiv preprint arXiv:2106.13281*, 2021.
- Scott Fujimoto, Herke Hoof, and David Meger. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pages 1587–1596. PMLR, 2018.
- Juan Jose Garau-Luis, Yingjie Miao, John D Co-Reyes, Aaron Parisi, Jie Tan, Esteban Real, and Aleksandra Faust. Multi-objective evolution for generalizable policy gradient algorithms. *arXiv preprint arXiv:2204.04292*, 2022.
- Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International Conference on Machine Learning*, pages 1861–1870. PMLR, 2018a.
- Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018b.
- Rein Houthoofd, Yuhua Chen, Phillip Isola, Bradley Stadie, Filip Wolski, OpenAI Jonathan Ho, and Pieter Abbeel. Evolved policy gradients. *Advances in Neural Information Processing Systems*, 31, 2018.
- Chloe Ching-Yun Hsu, Celestine Mendler-Düner, and Moritz Hardt. Revisiting design choices in proximal policy optimization. *arXiv preprint arXiv:2009.10897*, 2020.
- Louis Kirsch, Sjoerd van Steenkiste, and Jürgen Schmidhuber. Improving generalization in meta reinforcement learning using learned objectives. *arXiv preprint arXiv:1910.04098*, 2019.
- Jakub Grudzien Kuba, Christian Schroeder de Witt, and Jakob Foerster. Mirror learning: A unifying framework of policy optimisation. *arXiv preprint arXiv:2201.02373*, 2022.
- Robert Tjarko Lange. evosax: Jax-based evolution strategies, 2022. URL <http://github.com/RobertTLange/evosax>.
- Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- Luke Metz, C Daniel Freeman, Samuel S Schoenholz, and Tal Kachman. Gradients are not all you need. *arXiv preprint arXiv:2111.05803*, 2021.
- Luke Metz, C Daniel Freeman, James Harrison, Niru Maheswaranathan, and Jascha Sohl-Dickstein. Practical tradeoffs between memory, compute, and performance in learned optimizers. *arXiv preprint arXiv:2203.11860*, 2022.
- Volodymyr Mnih, Adria Puigdomenech Badia, Mehdi Mirza, Alex Graves, Timothy Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR, 2016.

- 275 Junhyuk Oh, Matteo Hessel, Wojciech M Czarnecki, Zhong-
 276 wen Xu, Hado P van Hasselt, Satinder Singh, and David
 277 Silver. Discovering reinforcement learning algorithms.
 278 *Advances in Neural Information Processing Systems*, 33:
 279 1060–1070, 2020.
- 280 Art B Owen. Monte carlo theory, methods and examples
 281 (book draft), 2014.
- 282 Ingo Rechenberg. Evolutionsstrategie–optimierung tech-
 283 nischer systeme nach prinzipien der biologischen evolu-
 284 tion. 1973.
- 285 Tim Salimans, Jonathan Ho, Xi Chen, Szymon Sidor,
 286 and Ilya Sutskever. Evolution strategies as a scalable
 287 alternative to reinforcement learning. *arXiv preprint*
 288 *arXiv:1703.03864*, 2017.
- 289 Jürgen Schmidhuber. On learning how to learn learning
 290 strategies. 1995.
- 291 John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan,
 292 and Philipp Moritz. Trust region policy optimization.
 293 In *International conference on machine learning*, pages
 294 1889–1897. PMLR, 2015.
- 295 John Schulman, F. Wolski, Prafulla Dhariwal, Alec Rad-
 296 ford, and Oleg Klimov. Proximal policy optimization
 297 algorithms. *ArXiv*, abs/1707.06347, 2017.
- 298 David Silver, Guy Lever, Nicolas Heess, Thomas Degris,
 299 Daan Wierstra, and Martin Riedmiller. Deterministic
 300 policy gradient algorithms. In *International conference*
 301 *on machine learning*, pages 387–395. PMLR, 2014.
- 302 Richard S Sutton and Andrew G Barto. *Reinforcement*
 303 *learning: An introduction*. 2018.
- 304 Yuhui Wang, Hao He, and Xiaoyang Tan. Truly proximal
 305 policy optimization. In *Uncertainty in Artificial Intelli-*
 306 *gence*, pages 113–122. PMLR, 2020.
- 307 Paul J Werbos. Backpropagation through time: what it
 308 does and how to do it. *Proceedings of the IEEE*, 78(10):
 309 1550–1560, 1990.
- 310 Yuhuai Wu, Mengye Ren, Renjie Liao, and Roger Grosse.
 311 Understanding short-horizon bias in stochastic meta-
 312 optimization. *arXiv preprint arXiv:1803.02021*, 2018.
- 313 Zhongwen Xu, Hado P van Hasselt, and David Silver. Meta-
 314 gradient reinforcement learning. *Advances in neural in-*
 315 *formation processing systems*, 31, 2018.
- 316 Tom Zahavy, Zhongwen Xu, Vivek Veeriah, Matteo Hessel,
 317 Junhyuk Oh, Hado van Hasselt, David Silver, and
 318 Satinder Singh. A self-tuning actor-critic algorithm. In
 319 *Proceedings of the 34th International Conference on*
 320 *Neural Information Processing Systems*, NIPS’20, Red
 321 Hook, NY, USA, 2020. Curran Associates Inc. ISBN
 322 9781713829546.
- 323 Tingting Zhao, Hirotaka Hachiya, Gang Niu, and Masashi
 324 Sugiyama. Analysis and improvement of policy gradient
 325 estimation. In *NIPS*, pages 262–270. Citeseer, 2011.

A. All Environment Results

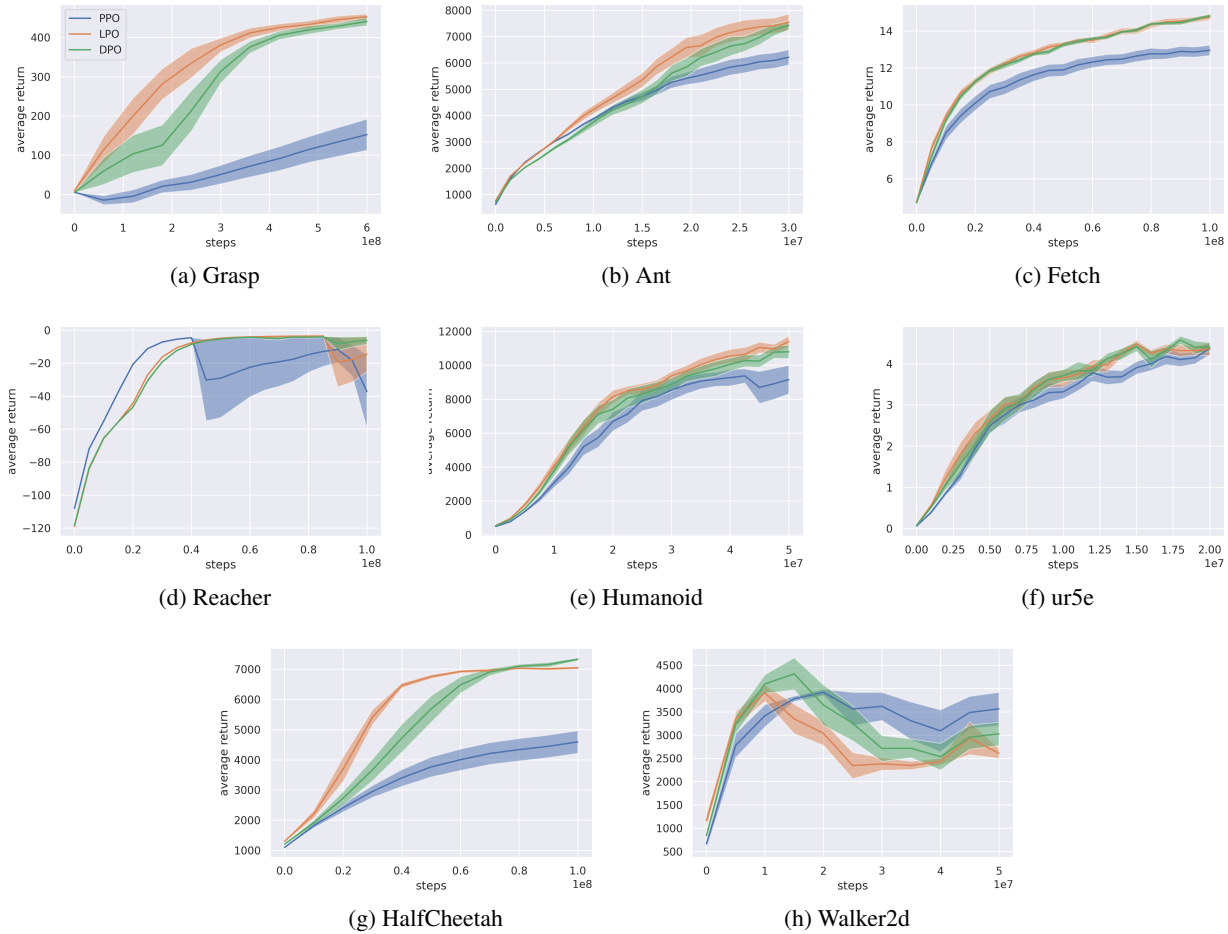


Figure 2: Performance comparison between PPO (blue), LPO (orange), and DPO (green) in Brax environments. The curves represent mean evaluation return across 10 random seeds, with error bars showing standard error. Both LPO and DPO beat PPO across most environments even though they were only meta-trained on Ant, and make use of no hyper parameter tuning. Furthermore, both LPO and DPO are more consistent across runs.

B. Objective Visualisations

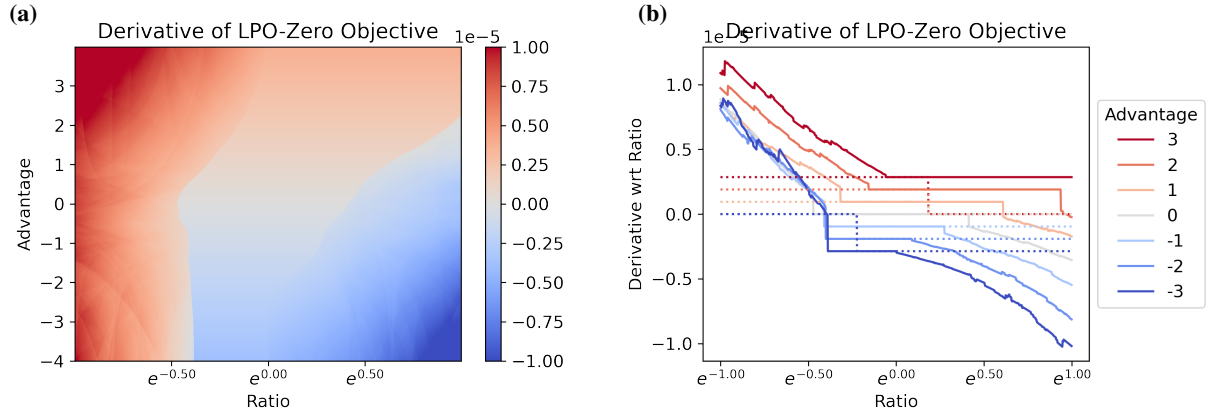


Figure 3: Visualisation of the LPO-Zero (learning from scratch) objective: (a) is the heat map of the ratio derivative of the LPO-Zero objective, and (b) shows its slices for fixed advantage values. The algorithm encourages updates towards actions with positive values of the ratio derivative.

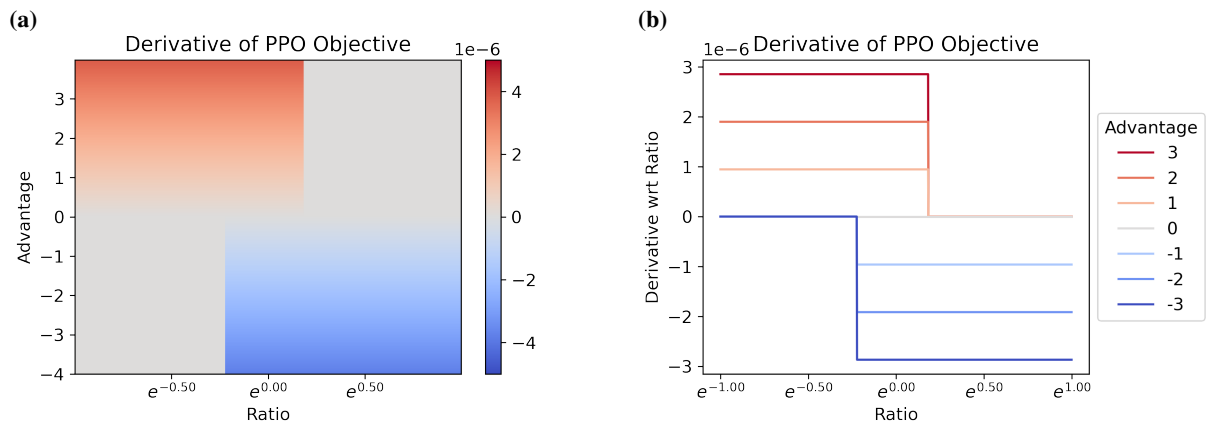


Figure 4: PPO objective visualisation: (a) is the heat map of the ratio derivative of the PPO objective, and (b) shows its slices for fixed advantage values. The algorithm encourages updates towards actions with positive values of the ratio derivative.

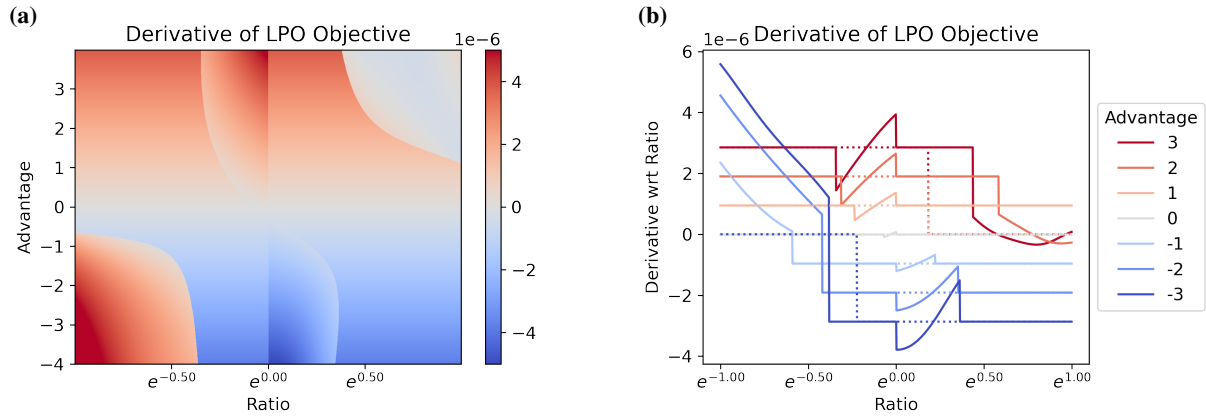


Figure 5: Visualisation of the LPO objective: (a) is the heat is the heat map of the ratio derivative of the LPO objective, and (b) shows its slices for fixed advantage values. The algorithm encourages updates towards actions with positive values of the ratio derivative.

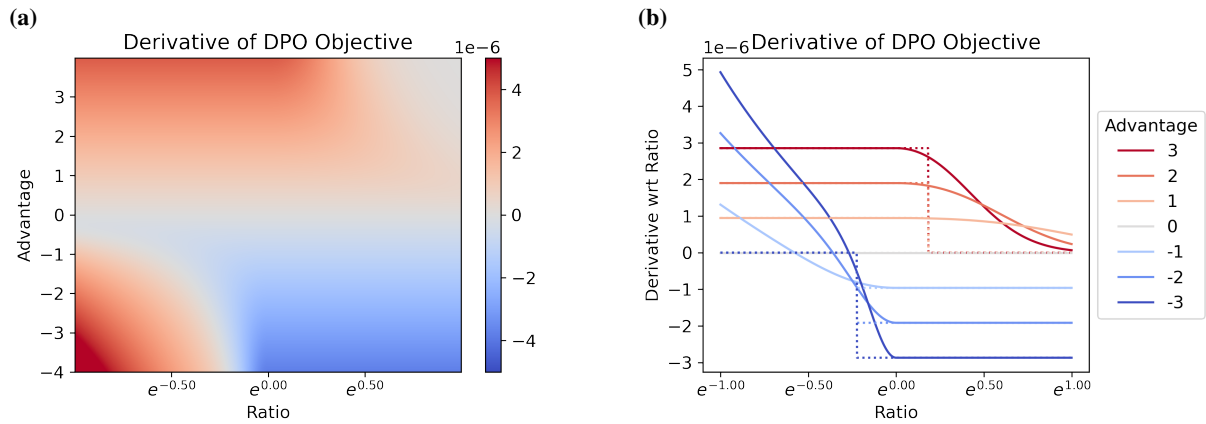


Figure 6: Visualisation of the DPO objective: (a) is the heat is the heat map of the ratio derivative of the DPO objective, and (b) shows its slices for fixed advantage values. Positive values of the derivative encourage updates towards the action.

C. Entropy Visualisation

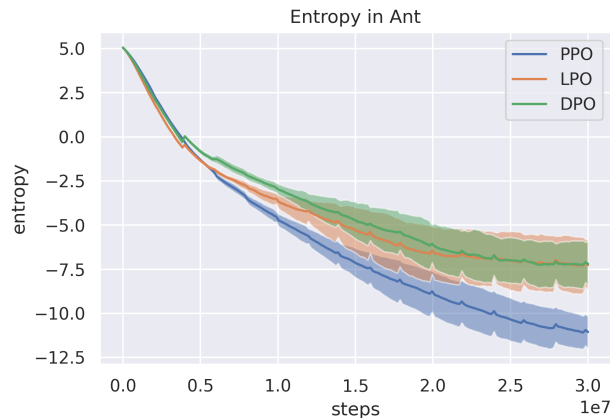


Figure 7: Entropy comparison, throughout training on Ant, between PPO (blue), LPO (orange), and DPO (green, see Section 6) across 10 seeds. Error bars denote standard error. While entropy of all methods decrease throughout training, the entropy of policy learned by both LPO and DPO remain significantly higher than that of PPO.

D. Related Work

In the last few years, significant effort has been directed at developing effective RL algorithms through both algorithmic and implementational advances. Fujimoto et al. (2018, TD3) combine DDPG policy training with estimates of pessimistic Bellman targets from a separate critic. Hsu et al. (2020) stabilise the, previously unsuccessful (Schulman et al., 2017), KL-penalised version of PPO and improve its robustness through novel policy design choices. Haarnoja et al. (2018b) introduce a mechanism that automatically adjusts the temperature parameter of the *entropy bonus* in SAC. However, none of these hand-crafted efforts succeeds in fully mitigating common RL pathologies, such as sensitivity to hyperparameter choices and lack of domain generalisation (Duan et al., 2016). This motivates radically expanding the RL algorithm search space through automated means.

Indeed, recent work explores applications of meta-RL techniques in guiding RL algorithm discovery and design. RL² equips a learning agent with a recurrent neural network conditioning on transitions between tasks, and adapts the agent’s behaviour to the current environment (Duan et al., 2016). Similarly, a MAML agent learns policy meta-parameters which can adapt to any task with a few steps of gradient descent (Finn et al., 2017). Neither RL² nor MAML, however, go beyond improving the robustness of ultimately hand-crafted algorithms. To overcome this, Xu et al. (2018, FRODO) introduce an actor-critic method that adjusts its hyperparameters online using *meta-gradients* that are updated with every few inner iterations. Similarly, STAC (Zahavy et al., 2020) uses implementation techniques from IMPALA (Espeholt et al., 2018) and auxiliary loss-guided meta-parameter tuning to further improve on FRODO.

Such advances have inspired extending meta-gradient RL techniques to more ambitious objectives, including the discovery of algorithms *ab initio*. Notably, Oh et al. (2020) succeeded in meta-learning an RL algorithm, LPG, that can solve simple tasks efficiently without explicitly relying on concepts such as value functions and policy gradients. Similarly, Evolved Policy Gradients (Houthoofd et al., 2018, EPG) meta-trains a policy loss network function with Evolution Strategies (Salimans et al., 2017, ES). Although EPG surpasses PPO in average performance, it suffers from much larger variance (Houthoofd et al., 2018). MetaGenRL (Kirsch et al., 2019) instead meta-learns the loss function for deterministic policies which are inherently less affected by estimators’ variance (Silver et al., 2014). MetaGenRL, however, fails to improve upon DDPG (Lillicrap et al., 2015) in terms of performance, despite building up on it. Neither EPG nor MetaGenRL have resulted in the discovery of novel analytical RL algorithms, perhaps due to the limited interpretability of the loss functions learnt. Lastly, Co-Reyes et al. (2021), Garau-Luis et al. (2022) and Alet et al. (2020) discover and improve standard RL conventions by evolving, symbolically, algorithms represented as graphs, which leads to improved performance in simple tasks. However, none of those trained-from-scratch methods inherit correctness guarantees, limiting our certainty of the generality of their abilities. In contrast our method, LPO, is meta-developed in a Mirror Learning space (Kuba et al., 2022), where every algorithm is guaranteed convergence to an optimal policy. As a result to this construction, meta-training of LPO is easier than that of

methods that learn “from scratch”, and achieves great performance across environments. Furthermore, thanks to the clear meta-structure of Mirror Learning, LPO is interpretable, and lets us discover new learning strategies. This lets us introduce DPO—an efficient algorithm with a closed-form formulation that exploits the discovered learning concepts.

E. Additional Background

E.1. Reinforcement Learning

Formulation We formulate the reinforcement learning (RL) problem as a *Markov decision process* (MDP) (Sutton and Barto, 2018) represented by a tuple $\langle \mathcal{S}, \mathcal{A}, r, P, \gamma, d \rangle$ which defines the experience of a learning agent as follows: at time step $t \in \mathbb{N}$, the agent is at state $s_t \in \mathcal{S}$ (where $s_0 \sim d$) and takes an action $a_t \in \mathcal{A}$ according to its stochastic policy $\pi(\cdot | s_t)$, which is a member of the policy space Π . The environment then emits the reward $r(s_t, a_t)$ and transits to the next state s_{t+1} drawn from the transition function, $s_{t+1} \sim P(\cdot | s_t, a_t)$. The agent aims to maximise the expected value of the total discounted return

$$\eta(\pi) \triangleq \mathbb{E}_\pi [R^\gamma] = \mathbb{E}_{s_0 \sim d, a_t \sim \pi, s_t \sim P} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right].$$

The agent guides its learning process with value functions that evaluate the expected return conditioned on states or state-action pairs

$$\begin{aligned} V_\pi(s) &\triangleq \mathbb{E}[R^\gamma | \pi, s_0 = s] \\ Q_\pi(s, a) &\triangleq \mathbb{E}[R^\gamma | \pi, s_0 = s, a_0 = a] \end{aligned}$$

respectively. The function that the agent is concerned about most is the *advantage function*, which computes relative values of actions at different states,

$$A_\pi(s, a) \triangleq Q_\pi(s, a) - V_\pi(s).$$

Policy Optimisation In fact, by updating its policy simply to maximise the advantage function at every state, the agent is guaranteed to improve its policy, $\eta(\pi_{\text{new}}) \geq \eta(\pi_{\text{old}})$ (Sutton and Barto, 2018). This fact, although requiring a maximisation operation that is intractable in large state-space settings tackled by deep RL (where the policy π_θ is parameterised by weights θ of a neural network), has inspired a range of algorithms that perform it approximately. For example, A2C (Mnih et al., 2016) updates the policy by a step of policy gradient (PG) ascent

$$\theta_{k+1} = \theta_k + \frac{\alpha}{B} \sum_{b=1}^B A_{\pi_{\theta_k}}(s_b, a_b) \nabla_\theta \log \pi_{\theta_k}(a_b | s_b)$$

estimated from a batch of B transitions with $\alpha \in (0, 1)$. Nevertheless, such simple adoptions of *generalized policy iteration* (Sutton and Barto, 2018, GPI) suffer from large variance and instability (Zhao et al., 2011; Silver et al., 2014; Schulman et al., 2017). Hence, methods that constrain (either explicitly or implicitly) the policy update size are preferred (Schulman et al., 2015). Among the most popular, as well as successful ones, is *Proximal Policy Optimization* (Schulman et al., 2017, PPO), inspired by *trust region learning* (Schulman et al., 2015), which updates its policy by maximising the PPO-clip objective

$$\begin{aligned} \pi_{k+1} &= \arg \max_{\pi \in \Pi} \mathbb{E}_{s \sim \rho_{\pi_k}, a \sim \pi_k} [L^{\text{PPO}}], \\ L^{\text{PPO}} &= \min \left(r(\pi) A_{\pi_k}(s, a), \text{clip}(r(\pi), 1 \pm \epsilon) A_{\pi_k}(s, a) \right) \end{aligned} \tag{3}$$

where $r(\pi) \triangleq \pi(a | s) / \pi_k(a | s)$ is the policy ratio and $\text{clip}(\cdot, 1 \pm \epsilon)$ clips the input inside the $[1 - \epsilon, 1 + \epsilon]$ interval if necessary. In deep RL, the maximisation oracle in Equation (3) is implemented through a few steps of gradient ascent on policy parameters.

Meta-RL The above approaches to policy optimisation rely on human-possessed knowledge, and thus are limited by humans’ understanding of the problem. The goal of *meta-RL* is to optimise the algorithm of the learning agent. Formally, suppose that an RL algorithm alg_ϕ , parameterised by ϕ , trains an agent for K iterations. Meta-RL aims to find the meta-parameter $\phi = \phi^*$ such that the expected return of the output policy, $\mathbb{E}[\eta(\pi_K) | \text{alg}_\phi]$, is highest.

Evolution Strategies Evolution Strategies [ES] (Rechenberg, 1973; Salimans et al., 2017) is a backpropagation-free approach to optimisation of stochastic functions. At their core lies the following identity, which holds for any continuously differentiable function F of ϕ , and any positive scalar σ :

$$\nabla_{\phi} \mathbb{E}_{\epsilon \sim \mathcal{N}}[F(\phi + \sigma\epsilon)] = \frac{1}{\sigma} \mathbb{E}_{\epsilon \sim \mathcal{N}}[F(\phi + \sigma\epsilon)\epsilon], \quad (4)$$

where $\mathcal{N} \triangleq \mathcal{N}(0, I)$ denotes the standard multivariate normal distribution. By taking the limit $\sigma \rightarrow 0$, the gradient on the left-hand side recovers the gradient of $\nabla_{\phi} F(\phi)$. These facts inspire an approach of optimising F with respect to ϕ without estimating gradients with backpropagation—for a random sample $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}$, the vector $\frac{1}{n\sigma} \sum_{i=1}^n F(\phi + \sigma\epsilon_i)\epsilon_i$ is an unbiased gradient estimate. To reduce variance of this estimator, antithetic sampling is commonly used (Owen, 2014). In the context of Meta-RL, where ϕ is the meta-parameter of an RL algorithm alg_{ϕ} , the role of $F(\phi)$ is played by the average return after the training, $F(\phi) = \mathbb{E}[\eta(\pi_K)|\phi]$. As oppose to the meta-gradient approaches described in Section D, ES does not require backpropagation of the gradient through the whole training episode—a cumbersome procedure which, often approximated by the truncated backpropagation, introduces bias (Verbos, 1990; Wu et al., 2018; Oh et al., 2020; Feng et al., 2021; Metz et al., 2021).

E.2. DPO Drift Proof

The DPO drift function $f(r, A)$ is given by

$$\begin{cases} \text{ReLU}((r-1)A - \alpha \tanh((r-1)A/\alpha)) & A \geq 0 \\ \text{ReLU}(\log(r)A - \beta \tanh(\log(r)A/\beta)) & A < 0. \end{cases}$$

The first condition for a valid drift is that f be non-negative everywhere, which trivially holds since $\text{ReLU}(x) \geq 0$ for all $x \in \mathbb{R}$.

The second condition is that f be zero at $\pi = \pi_{\text{old}}$. Now $r = \pi/\pi_{\text{old}} = 1$ implies $r-1 = 0$ and $\log r = 0$, which combined with $\tanh(0) = 0$ imply that $f = 0$ as required.

The final condition is that the gradient of f with respect to π be zero at $\pi = \pi_{\text{old}}$. This is equivalent to having zero gradient with respect to $r = \pi/\pi_{\text{old}}$ at $r = 1$ since the gradients are equal up to a constant. Now writing

$$\begin{aligned} f^+ &= (r-1)A - \alpha \tanh((r-1)A/\alpha) \\ f^- &= \log(r)A - \beta \tanh(\log(r)A/\beta) \end{aligned}$$

for $A \geq 0$ and $A < 0$ respectively, we have

$$\begin{aligned} \frac{\partial f^+}{\partial r} &= A - A \cosh^{-2}((r-1)A/\alpha) \\ \frac{\partial f^-}{\partial r} &= \frac{A}{r} - \frac{A}{r} \cosh^{-2}(\log(r)A/\beta) \end{aligned}$$

which both evaluate to 0 at $r = 1$, since $\cosh(0) = 1$. This implies for $A \geq 0$ that

$$\frac{\partial f}{\partial r} = \frac{\partial \text{ReLU}(f^+)}{\partial r} = \begin{cases} \frac{\partial f^+}{\partial r} & \text{if } f^+ \geq 0 \\ 0 & \text{if } f^+ < 0 \end{cases} = 0$$

at $r = 1$ and for $A < 0$ that

$$\frac{\partial f}{\partial r} = \frac{\partial \text{ReLU}(f^-)}{\partial r} = \begin{cases} \frac{\partial f^-}{\partial r} & \text{if } f^- \geq 0 \\ 0 & \text{if } f^- < 0 \end{cases} = 0$$

at $r = 1$. Taken together we conclude, for all A , that f has zero gradient at $r = 1$.

E.3. Learning drift functions from scratch

In this setting, f_ϕ is a neural network with two hidden layers of size 256 and a ReLU activation function. We meta-train it across 5 Brax environments. We name the resulting algorithm LPO-Zero, and visualise it in Figure 3.

Interestingly, LPO-Zero appears to have learnt a few PPO-like features, as can be observed on Figure 3. For example, it appears to have learnt to clip the update incentive at a specific ratio threshold, much like PPO; however, it only does so for negative advantages. Nevertheless, LPO-Zero largely underperforms with respect to LPO, and possibly requires much more training to catch up.

F. More Analysis of LPO

Update asymmetry. LPO learns *asymmetric* features that respect a natural asymmetry of behaviour change in RL: increasing r for positive advantage A may encourage exploration of a newly-found action *or* strengthen a dominant action, whereas decreasing r for negative A will always discourage exploration of that action *and* strengthen a dominant action. In this context, the two discussed features of LPO make it completely unlike PPO, which clips the update incentives symmetrically around the origin.

Secondary Features. LPO, but not LPO-Zero, appears to consistently learn objectives with gradient spikes around $r = 1$ in the upper left and lower right quadrants. Nevertheless, adding them to our analytic model of LPO did not improve performance. We speculate, therefore, that these spikes are mostly artifacts of the network parameterisation.