

Rank-1 Approximation of Inverse Fisher for Natural Policy Gradients in Deep Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

Natural gradients have long been studied in deep reinforcement learning due to their fast convergence properties and covariant weight updates. However, computing natural gradients requires inversion of the Fisher Information Matrix (FIM) at each iteration, which is computationally prohibitive in nature. In this paper, we present an efficient and scalable natural policy optimization technique that leverages a rank-1 approximation to full inverse-FIM. We theoretically show that under certain conditions, a rank-1 approximation to inverse-FIM converges faster than policy gradients and, under some conditions, enjoys the same sample complexity as stochastic policy gradient methods. We benchmark our method on a diverse set of environments and show that it achieves superior performance to standard actor-critic and trust-region baselines.

1 Introduction

Policy gradient methods (Sutton et al., 1999) have emerged as one of the most prominent candidates for modern reinforcement learning (RL). Specifically, they have seen widespread adoption in deep RL, where they incorporate differentiable policy parametrization and function approximators to help policy gradient methods scale up and learn robust controllers for high-dimensional state spaces (Arulkumaran et al., 2017; Wang et al., 2022). Despite this, a fundamental limitation of policy gradients is that RL objectives tend to be highly non-concave in nature. Consequently, algorithms based on standard gradient descent, such as REINFORCE (Williams, 1992), are prone to getting stuck at sub-optimal local minima when training high-dimensional policies. This is due to the gradient giving only the direction in which a policy’s parameters should be optimized, without indicating the size of the optimization step (Martens, 2020). Step size selection is, therefore, non-trivial. Small steps lead to stable but slow and sample-inefficient training. On the contrary, large steps can cause unstable updates. This is particularly dangerous in RL, as unstable policies collect poor data, making it difficult for the agent to recover and escape local optima (van Heeswijk, 2022).

One approach to reduce this sub-optimality and achieve stable training with large steps is to incorporate the geometry of the loss manifold into account and thus consider natural gradient descent methods (Amari, 1998). Specifically, when applied to policy gradient RL, where policies are modeled as probability distributions, the Fisher information matrix (FIM) is used to precondition the policy gradients, which gives rise to natural policy gradient (NPG) methods (Kakade, 2001; Bagnell & Schneider, 2003; Peters et al., 2005). NPG algorithms are of interest because they generally tend to converge faster and have superior performance compared to their standard policy gradient counterparts (Grondman et al., 2012). More recently, NPGs established the roots of some of the most popular RL algorithms used in practice, including Trust Region Policy Optimization (TRPO) (Schulman et al., 2015a) and Proximal Policy Optimization (PPO) (Schulman et al., 2017). Furthermore, the study of NPG methods is also motivated by their success in real-world applications (Richter et al., 2006; Su et al., 2019).

One of the main drawbacks of NPG algorithms is the high cost of computing and inverting the FIM. Storage is also an important concern, as for a parameterized policy, the FIM is a square matrix with dimensions equal to the number of parameters. Naturally, this issue becomes especially limiting when policies are parameterized by deep neural networks. In practice, NPG algorithms usually avoid computing the inverse Fisher directly

and instead rely on approximations. A simple option is to only consider the diagonal elements of the FIM (Becker & Lecun, 1989; LeCun et al., 2002), leading to efficient inversion and low storage requirements. A more sophisticated approach is K-FAC (Martens & Grosse, 2015), which uses Kronecker products to approximate the FIM efficiently. As an alternative to the true FIM, the empirical FIM can be computed from samples, which is especially convenient when the required underlying probability distributions are unknown (Schraudolph, 2002). Hessian-free methods give a less direct approach to using the FIM without ever explicitly computing or storing it (Martens, 2010). As an example, TRPO avoids computing the FIM, calculating the Fisher-vector product directly instead. Despite the extensive literature on the subject, choosing the right way of approximating natural gradients is usually not straightforward. This choice often involves a trade-off between computational efficiency and approximation accuracy.

The goal of this work is to bridge the theoretical convergence properties of NPG with a practical, fast, and computationally feasible approach to invert FIM for high-dimensional reinforcement learning tasks. In this work, we use the empirical FIM as an alternative to the FIM and incorporate the Sherman-Morrison formula for efficient inversions. We propose a rank-1 approximation to avoid explicit matrix operations, keeping the time complexity at $\mathcal{O}(d)$ for a policy parameterised by a vector θ of dimension d , without the need for inner loops like conjugate gradient, thus balancing computational cost and accuracy of curvature information.

The main contributions are listed below:

1. We provide a fast and efficient way to compute the inverse FIM, which can be computed in $\mathcal{O}(d)$ complexity rather than $\mathcal{O}(d^3)$ for naive inversion of FIM.
2. We theoretically show that under certain conditions, our NPG approximation using the empirical FIM and Sherman-Morrison formula enjoys global convergence when using log-linear function approximators.
3. We also prove that the error induced in the policy update due to this rank-1 approximation is bounded and reduces with the increase in iterations.
4. We propose a novel policy gradient algorithm using this update: SM-ActorCritic (**SMAC**). We show that SMAC achieves faster convergence in diverse OpenAI Gym tasks (classic control and MuJoCo environments (Todorov et al., 2012)), compared to vanilla policy gradient, Adam, and trust-region methods that use Conjugate Gradient solvers to approximate the direction of natural policy gradient.

2 Related Work

Gradient descent updates all the parameters of a model equally, meaning that weights are considered to be inside a Euclidean space. However, the loss function of the model actually induces a Riemannian manifold space. Therefore, the size of a gradient step does not reflect the true scale of change of the model’s loss. Consequently, standard gradient descent methods are highly sensitive to the choice of step size. To avoid this issue, Amari first framed gradient-based learning on a Riemannian manifold, introducing the natural gradient descent algorithm, which uses the update direction $F(\theta)^{-1}\nabla J(\theta)$, with the FIM $F(\theta)$ and objective $J(\theta)$ (Amari, 1998). Here, the FIM is a covariant measure of the amount of information contained in the parameters (Amari, 1998). Kakade extended this idea to reinforcement learning and coined the natural policy gradient (NPG) algorithm (Kakade, 2001). However, the computational cost associated with calculating the FIM, especially in models with a large number of parameters, can be prohibitively expensive. Therefore, practical implementations of natural gradients rely on approximations instead. We thus summarize the related work on these approximations as follows.

Kronecker-factored Approximate Curvature Block-diagonal approximations to the natural gradient, such as K-FAC (Martens & Grosse, 2015) and subsequent refinements (Zhang et al., 2023; Ren & Goldfarb, 2021; Bahamou et al., 2023; Benzing, 2022), have been widely adopted due to their computational efficiency. Related variants have also been developed for reinforcement learning (Wu et al., 2017).

Hessian-Free Methods Hessian-free optimization (Martens, 2010), another method for computing natural gradients, has seen widespread use in reinforcement learning, particularly due to its superior performance in TRPO (Schulman et al., 2015a). This approach employs the conjugate gradient (CG) method to solve the associated linear system, thereby avoiding explicitly storing FIM, keeping memory usage comparable to a single forward-backward pass. However, CG requires several additional forward-backward passes, which reduces computational efficiency. In addition, due to the limited precision of CG estimation, practitioners usually perform an extra KL-based line search to adapt the step size (Schulman et al., 2015a).

Empirical Fisher Information The empirical Fisher (EF) estimates the Fisher information matrix by replacing the exact expectation over gradient outer products with the outer product of the gradients computed from individual samples (Yang et al., 2022). However, the empirical Fisher’s simple estimate has been frequently criticized for deviating significantly from the true Fisher, and its limitations have been analyzed in many studies (Martens, 2020; Kunstner et al., 2019). Nevertheless, its very low computational cost has enabled a broad range of successful applications (Roux et al., 2007; Ren & Goldfarb, 2019; Wu et al., 2024). Compared with the extensive use of K-FAC and Hessian-free methods in RL, the EF approach remains under-explored and is seldom applied in this domain.

3 Preliminaries

Markov decision process (MDP) We consider an ~~episodic~~ MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ with a continuous or discrete state space \mathcal{S} , action space \mathcal{A} , transition kernel $P(s'|s, a)$, bounded reward function $r : \mathcal{S} \times \mathcal{A} \mapsto [-R, R]$, and discount factor $\gamma \in (0, 1)$. We also consider a stochastic policy $\pi_\theta(a|s)$ with parameters $\theta \in \mathbb{R}^d$ that induces a trajectory distribution $\pi_\theta(\tau_i)$ over $\tau_i = (s_0^i, a_0^i, s_1^i, a_1^i, \dots)$.

Performance objective and policy gradient Under the standard actor-critic framework, the RL objective is given by the expected discounted return of π_θ :

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right]. \quad (1)$$

To maximize Eq. 1, the advantage actor-critic (A2C) (Mnih et al., 2016) algorithm uses an advantage function $A(s, a) = Q(s, a) - V(s)$, defined through the critic’s state-value and action-value functions

$$V(s) = \mathbb{E}_{\tau \sim \pi_\theta} [G_t \mid s_t = s], \quad Q(s, a) = \mathbb{E}_{\tau \sim \pi_\theta} [G_t \mid s_t = s, a_t = a],$$

where $G_t = \sum_{l=0}^{\infty} \gamma^l r(s_{t+l}, a_{t+l})$. Following the policy gradient theorem, the A2C gradient is then given by

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{\infty} \nabla_\theta \log \pi_\theta(a_t | s_t) A(s_t, a_t) \right].$$

Fisher information matrix and Natural Policy Gradient For the exponential-family policy $\pi_\theta(\cdot | s)$, the natural Riemannian metric in parameter space ~~can be expressed in terms of state visitation distribution induced by the same policy $\pi_\theta(\cdot | s)$ as is the Fisher information matrix~~

$$\mathbf{F}_s(\theta) = \mathbb{E}_{a \sim \pi_\theta(\cdot | s)} [\nabla_\theta \log \pi_\theta(a | s) \nabla_\theta \log \pi_\theta(a | s)^\top]. \quad (2)$$

The invariant FIM $\mathbf{F}_\rho(\theta)$ can be obtained by averaging over the state visitation distribution leading ~~For a preconditioning matrix $\mathbf{F}_\rho(\theta)$, this leads~~ to the natural gradient update

$$\begin{aligned} \theta^{k+1} &= \theta^k + \eta \mathbf{F}_\rho(\theta^k)^{-1} \nabla_\theta J(\theta^k), \\ \mathbf{F}_\rho(\theta^k) &= \mathbb{E}_{s \sim d_\rho^{\pi_\theta}} [\mathbf{F}_s(\theta)], \end{aligned}$$

where η is the step size and $d_\rho^{\pi_\theta}$ represents the state visitation distribution under a policy π_θ and initial state distribution ρ :

$$d_\rho^{\pi_\theta}(s) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho} \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0, \pi_\theta) + \gamma \rho(s).$$

Empirical Fisher information Replacing the expectation in $\mathbf{F}(\theta)$ with a sample average over [one single trajectory](#) $\tau = \{(s_i, a_i)\}_{i=0}^{T-1}$ ~~a batch~~ $\{\tau_i\}_{i=1}^N$ gives the empirical Fisher approximation

$$\hat{\mathbf{F}}(\theta) = \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \nabla_{\theta} \log \pi_{\theta}(a_i | s_i)^{\top}.$$

Sherman-Morrison formula The Sherman-Morrison formula efficiently computes the inverse of a matrix after a “rank-1 update” has been applied to it, provided that the inverse of the original matrix was already known (Sherman & Morrison, 1950). Suppose $\mathbf{X} \in \mathbb{R}^{n \times n}$ is an **invertible square matrix** and $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ are **column vectors**. Then $\mathbf{X} + \mathbf{u}\mathbf{v}^{\top}$ is invertible **iff**

$$1 + \mathbf{v}^{\top} \mathbf{X}^{-1} \mathbf{u} \neq 0.$$

In this case, the Sherman-Morrison formula allows the computation of the inverse

$$(\mathbf{X} + \mathbf{u}\mathbf{v}^{\top})^{-1} = \mathbf{X}^{-1} - \frac{\mathbf{X}^{-1} \mathbf{u} \mathbf{v}^{\top} \mathbf{X}^{-1}}{1 + \mathbf{v}^{\top} \mathbf{X}^{-1} \mathbf{u}}. \quad (3)$$

If \mathbf{X}^{-1} is already known, this matrix inversion formula can bypass computing the inverse $(\mathbf{X} + \mathbf{u}\mathbf{v}^{\top})^{-1}$ directly, avoiding the cubic scaling with the dimension of \mathbf{X} . Instead, the update only requires three vector operations, which scale linearly with the dimension of \mathbf{X} . This update can use the full power of current GPU hardware platforms, since computing vector operations can be heavily parallelized. At the same time, it avoids the time-intensive matrix inversion processes. [Methods such as](#) ~~Specifically, methods like~~ Singular Value Decomposition (SVD), which decomposes a matrix into its singular values and vectors, and LU decomposition, which factors a matrix into a lower triangular matrix (L) and an upper triangular matrix (U), are much slower than the simple operations used in our approach.

For the rest of the paper, we abbreviate the actual FIM at θ^k , $\mathbf{F}(\theta^k)$ as \mathbf{F}^k and the estimate of a batch as $\hat{\mathbf{F}}^k$ for notational convenience.

Sherman-Morrison for computing the inverse Fisher The Sherman-Morrison formula can be used to incrementally build the inverse of the EF $(\hat{\mathbf{F}}^k)^{-1}$ at step k by reusing the previous estimate $(\hat{\mathbf{F}}^{k-1})^{-1}$ (Singh & Alistarh, 2020):

$$\begin{aligned} (\hat{\mathbf{F}}^k)^{-1} g^k &= \left[(\hat{\mathbf{F}}^{k-1} + \nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} \right]^{-1} g^k \underline{g^k (\hat{\mathbf{F}}^{k-1} + \nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top})^{-1} g^k} \\ &= (\hat{\mathbf{F}}^{k-1})^{-1} g^k - \frac{(\hat{\mathbf{F}}^{k-1})^{-1} \nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} (\hat{\mathbf{F}}^{k-1})^{-1} g^k}{1 + \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top} (\hat{\mathbf{F}}^{k-1})^{-1} \nabla_{\theta} \log \pi_{\theta}(a|s)} g^k. \end{aligned}$$

[Other approaches compute a fresh EF at each iteration](#) (Wu et al., 2024), [which is used in the natural gradient](#)

$$\Delta \theta = \eta (\lambda \mathbf{I} + \nabla_{\theta} \log \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s)^{\top})^{-1} \log \pi_{\theta}(a|s),$$

where $\lambda > 0$ is a fixed damping coefficient.

4 NPG Approximation

We adopt a matrix-free scheme for approximating the FIM to reduce memory and computational cost. At each step, we form a local empirical Fisher from the current batch, add a damping term λ , and apply the Sherman-Morrison formula to compute the inverse-vector product $\mathbf{F}(\theta)^{-1} g^k$ directly, [where \$g^k\$ is the policy gradient](#). This approach requires $\mathcal{O}(d)$ memory.

This formulation captures the local curvature of the loss landscape using only the current gradient direction. The damping term λ ensures numerical stability and prevents the update direction from being overly sensitive to individual gradients. Compared to the exact FIM, this method strikes a balance between computational efficiency and curvature-informed updates.

4.1 Matrix-Free Update

The Sherman-Morrison formula can be used to incrementally build the inverse of the EF $(\hat{\mathbf{F}}^k)^{-1}$ at step k by reusing the previous estimate $(\hat{\mathbf{F}}^{k-1})^{-1}$ (Singh & Alistarh, 2020):

Other approaches compute a fresh EF at each iteration (Wu et al., 2024), which is used in the natural gradient

The time and space complexity of explicitly computing and storing the FIM is $\mathcal{O}(d^2)$. However, the time and space complexity of explicitly calculating and storing the FIM is $\mathcal{O}(d^2)$. So, maintaining a previous estimate $(\hat{\mathbf{F}}^{k-1})^{-1}$ and applying $(\hat{\mathbf{F}}^k)^{-1}g$ can be infeasible for large-scale models, which is infeasible for large-scale models. In this work, we instead use the Sherman-Morrison formula to directly compute the product of the inverse empirical Fisher and the policy gradient.

We consider a one-sample empirical Fisher with damping, take $(s_k, a_k) \sim d^{\pi_{\theta^k}}(s)\pi_{\theta^k}(a | s)$, that is, one state-action pair visited when rolling out π_{θ^k} :

$$\hat{\mathbf{F}}^k = \lambda \mathbf{I} + \mathbf{F}^k(\theta) = \lambda \mathbf{I} + \nabla_{\theta} \log \pi_{\theta^k}(a_k | s_k) \nabla_{\theta} \log \pi_{\theta^k}(a_k | s_k)^{\top} \nabla_{\theta} \log \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top}.$$

where $\lambda > 0$ is a fixed damping coefficient. Applying the Sherman-Morrison formula then yields:

$$\begin{aligned} (\hat{\mathbf{F}}^k)^{-1} &= (\lambda \mathbf{I} + \nabla_{\theta} \log \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top})^{-1} \\ &= \frac{1}{\lambda} \mathbf{I} - \frac{\nabla_{\theta} \log \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top}}{\lambda^2 + \lambda \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top} \nabla_{\theta} \log \pi_{\theta}(a | s)} \\ &= \frac{1}{\lambda} \mathbf{I} - \frac{\mathbf{F}^k(\theta)}{\lambda^2 + \lambda \text{Tr}(\mathbf{F}^k(\theta))}. \end{aligned} \quad (4)$$

$$\begin{aligned} (\hat{\mathbf{F}}^k)^{-1} &= (\lambda \mathbf{I} + \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k) \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)^{\top})^{-1} \\ &= \frac{1}{\lambda} \mathbf{I} - \frac{\nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k) \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)^{\top}}{\lambda^2 + \lambda \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)^{\top} \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)} \\ &= \frac{1}{\lambda} \mathbf{I} - \frac{\mathbf{F}^k}{\lambda^2 + \lambda \text{Tr}(\mathbf{F}^k)}, \end{aligned} \quad (5)$$

Multiplying the inverse Fisher by the gradient gives the natural policy gradient update direction

$$\Delta \theta_k = \eta (\hat{\mathbf{F}}^k)^{-1} g^k = \eta \left[\frac{1}{\lambda} g^k - \frac{\nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k) \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)^{\top} \nabla_{\theta} \log \pi_{\theta}(a | s) \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top} g^k}{\lambda^2 + \lambda \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k)^{\top} \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k) \nabla_{\theta} \log \pi_{\theta}(a | s)^{\top} \nabla_{\theta} \log \pi_{\theta}(a | s)} \right], \quad (6)$$

where $g^k = \nabla_{\theta} \log \pi_{\theta^k}(a^k | s^k) \nabla_{\theta} \log \pi_{\theta}(a | s) A(s, a)$ is the one-sample policy gradient update direction, $A(s, a)$ is the advantage estimate and η is the step size.

This update involves only a matrix-vector operation and a vector outer product, making its time and memory complexity comparable to that of first-order methods, such as stochastic policy gradient.

4.2 Convergence Analysis

We take inspiration from the global convergence guarantees provided by Agarwal et al. (2021) and a follow-up work by Liu et al. (2022), which improved the global convergence results. To prove the global convergence, we take a standard Lipschitz policy assumption (Assumption 3) and average performance difference bounds (Lemma 4 from Liu et al. (2022)) and Lemma 2 take some previously proved bounds of the norm of exact policy gradient and the policy gradient defined till horizon H from Xu et al. (2020). We provide a theoretical

analysis showing that the Sherman–Morrison rank-1 approximation to the inverse Fisher achieves global convergence. Under strong convexity assumptions (already taken in Kakade (2001); Peters et al. (2005)) of the FIM, we show that the average performance between the Sherman-Morrison update and the optimal policy π^* can be bounded in terms of the variance of policy gradient estimator (σ^2), discount factor (γ), the error in the Advantage estimate induced by the neural approximator (ε_{bias}) and the Lipschitz constant of the policy gradient (L_J). We defer the proof to the appendix, and state the final theorem here.

Theorem 1. *If we take the stochastic Sherman-Morrison policy update in Equation 5 and take $\eta = \frac{1}{4L_J}$, $K = \mathcal{O}\left(\frac{1}{(1-\gamma)^2\varepsilon^2}\right)$, $N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right)$, and $H = \mathcal{O}\left(\log\left(\frac{1}{(1-\gamma)\varepsilon}\right)\right)$. Then, we have*

$$J(\pi^*) - \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[J(\theta^k)] \leq \frac{\sqrt{\varepsilon_{bias}}}{1-\gamma} + \varepsilon.$$

In total, stochastic PG samples $\tilde{\mathcal{O}}\left(\frac{\sigma^2}{(1-\gamma)^2\varepsilon^4}\right)$ trajectories, where ε_{bias} is the approximation error from the advantage function, and ε is an arbitrary precision level and $\tilde{\mathcal{O}}$ hides all the logarithmic factors.

Our proof relies on the fact that performance bound given by Lemma 4 (Liu et al., 2022) and we prove that this result holds with a sample complexity of $KN = \mathcal{O}\left(\frac{\sigma^2}{(1-\gamma)^2\varepsilon^4}\right)$, where N is the number of trajectories per update and K is the number of updates. This shows that to get an ε accurate policy, one needs $\mathcal{O}\left(\frac{1}{\varepsilon^2}\right)$ samples, and when the horizon H and number of samples N are large enough, then σ can be reduced to arbitrary precision. We would like to point out that the non-degeneracy of FIM [Assumption 2] and compatible function approximator [Assumption 4] are standard assumptions in stationary convergence of PG adapted by Agarwal et al. (2021) and can be realized by gaussian and softmax parametrized policies (See Section B.2 of Liu et al. (2022) and Section 8 of Ding et al. (2022) for further discussion). For compatible function approximation the value of ε_{bias} can be arbitrarily small for neural approximators (Wang et al., 2019) and is exactly zero for softmax parametrized policies (Ding et al., 2022). Detailed derivations can be found in Appendix 2.

Note that, the obtained bound stems from average regret to the global optimum and recently Yuan et al. (2022) have obtained $\tilde{\mathcal{O}}(\varepsilon^{-3})$ bound on the condition when the approximator error is negligible, or $\varepsilon_{bias} = 0$, but with a finite ε_{bias} the authors have recovered a bound of $\tilde{\mathcal{O}}(\varepsilon^{-4})$, similar to ours.

4.3 A Toy Problem

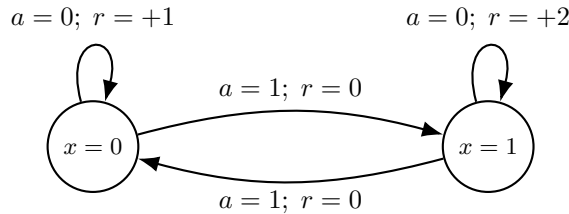


Figure 1: An example two-state MDP where each state x has two control actions, one which self-transition with reward +1 or +2) and another control that transitions the initial state to the other state with zero reward).

We adopt a classical reinforcement-learning MDP problem, been used by Kakade (2001) & Bagnell & Schneider (2003) with two states $x \in \{0, 1\}$ and two discrete actions $a \in \{0, 1\}$. The policy is parameterized by a vector $\theta = [\theta_0, \theta_1]^\top$, where each component governs an independent Bernoulli distribution over actions. Let the probability of taking action $a = 0$ in state x be

$$\pi_\theta(a = 0 \mid x = 0) = \frac{1}{1 + e^{\lambda_p \theta_0}}, \quad (7)$$

$$\pi_\theta(a = 0 \mid x = 1) = \frac{1}{1 + e^{\theta_1}} \quad (8)$$

where λ_p scales the first policy parameter and is used to test scaling invariance. Assume that the agent starts at state $x = 0$, then the objective function of the the MDP can be expressed as

$$J(\boldsymbol{\theta}) = \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0) + 2(1 - \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0)) \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 1). \quad (9)$$

Note that, for a converged policy the agent shall transit to state $x = 1$ and remain there and the final return shall be +2 which is the the action of self transitioning at state $x = 1$. So we start from writing the analytic policy gradient vector,

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} -\lambda_p \pi_0(1 - \pi_0)(1 - 2\pi_1) \\ -2(1 - \pi_0)\pi_1(1 - \pi_1) \end{bmatrix}. \quad (10)$$

and the FIM can be calculated as:

$$F(\boldsymbol{\theta}) = \begin{bmatrix} \lambda_p^2 \pi_0(1 - \pi_0) & 0 \\ 0 & \pi_1(1 - \pi_1) \end{bmatrix}. \quad (11)$$

Then NPG update direction can be evaluated as

$$\dot{\boldsymbol{\theta}}_{\text{NPG}} = F(\boldsymbol{\theta})^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} \frac{-(1 - 2\pi_1)}{\lambda_p} \\ -2(1 - \pi_0) \end{bmatrix}. \quad (12)$$

This direction eliminates the dependence on the arbitrary scaling λ_p , making it invariant under smooth reparameterizations of $\boldsymbol{\theta}$ as shown by Bagnell & Schneider (2003). Now the SM update direction in this setting can be evaluated as

$$\dot{\boldsymbol{\theta}}_{\text{SM}} = \frac{1}{\lambda} \left(\mathbf{I} - \frac{F(\boldsymbol{\theta})}{\lambda + \text{Tr}(F(\boldsymbol{\theta}))} \right) \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) \quad (13)$$

$$= \begin{bmatrix} \frac{1}{\lambda} (-\lambda_p \pi_0(1 - \pi_0)(1 - 2\pi_1)) - \frac{\lambda_p^3 \pi_0^2(1 - \pi_0)^2(1 - 2\pi_1)}{\lambda^2 + \lambda [\lambda_p^2 \pi_0(1 - \pi_0) + \pi_1(1 - \pi_1)]} \\ \frac{1}{\lambda} (-2(1 - \pi_0)\pi_1(1 - \pi_1)) - \frac{2 \pi_1^2(1 - \pi_1)^2(1 - \pi_0)}{\lambda^2 + \lambda [\lambda_p^2 \pi_0(1 - \pi_0) + \pi_1(1 - \pi_1)]} \end{bmatrix} \quad (14)$$

where λ is the damping factor. This type of approximation although looks highly nonlinear but we can carefully tune λ and make the update less sensitive to λ_p .

As shown in Fig.2, the plot in the left panel update the trajectories in the log-odds space for all damping settings $\lambda_p \in \{0.75, 1.0, 2.0\}$, Compared with the PG updates, the SM updates follow a direction much closer to that of NPG and maintain a step size more consistent. And PG updates deviate significantly as λ changes, showing its sensitivity to parameter scaling. The plot in the right panel shows the returns over iterations, where SM update achieves faster and more stable improvement than PG and approaches the same final performance as NPG. Therefore, the SM update preserves the geometric characteristics of the NPG while being more computational efficient in high-dimensional cases.

4.4 SM-ActorCritic

We plug our approximation into a standard A2C algorithm. Let $\pi_{\boldsymbol{\theta}}$ be the policy (actor) and V_{ϕ} the value network (critic), with parameters $\boldsymbol{\theta}$ and ϕ , respectively. At each iteration, we collect T transition tuples and compute the generalized advantage estimator (Schulman et al., 2015b, GAE) with a fixed $\lambda_{\text{GAE}} \in [0, 1]$:

$$\delta_t = r_t + \gamma V_{\phi}(s_{t+1}) - V_{\phi}(s_t), \quad (15)$$

$$A_t = \sum_{l=0}^{T-1-t} (\gamma \lambda_{\text{GAE}})^l \delta_{t+l}. \quad (16)$$

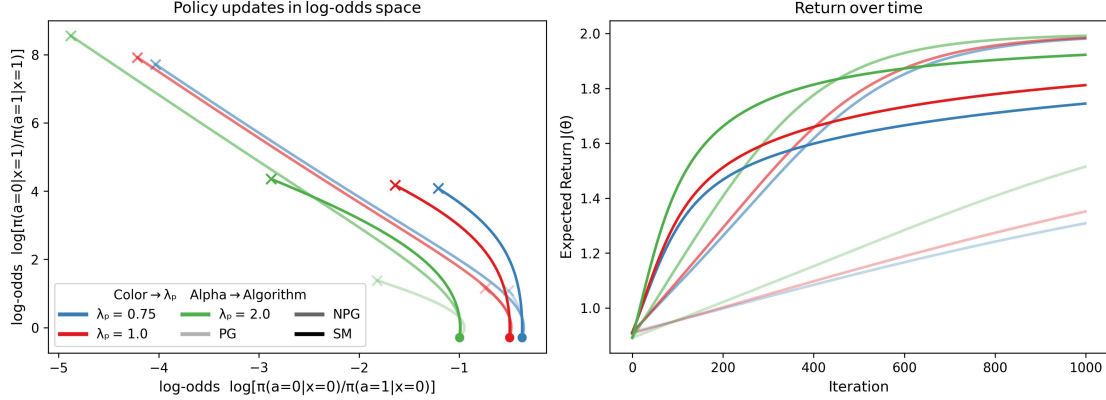


Figure 2: Results on toy MDP problem: Left side plot represents the log-odds plot for state $x = 0$ in x-axis and state $x = 1$ in y-axis and the right side plot shows the returns over iteration for PG, NPG and SM updates for different $\lambda_p \in \{0.75, 1.0, 2.0\}$ coefficient in policy parametrization. Note that, the least intense lines represent PG direction, medium is NPG and maximum deep lines represent SM update direction.

The critic is then trained by minimizing the mean squared error over the batch:

$$\mathcal{L}_{\text{critic}}(\phi) = \frac{1}{T} \sum_{t=0}^{T-1} (R_t - V_\phi(s_t))^2,$$

using the Adam optimizer with a learning rate of α , where $R_t = A_t + V_\phi(s_t)$ is the Monte-Carlo return target. The actor is updated by a single natural policy gradient step obtained from the matrix-free Fisher inverse. The full SMAC update is given in Algorithm 1.

Algorithm 1 SM-ActorCritic

- 1: **input:** steps per update T , step size η , damping λ
 - 2: initialise θ, ϕ
 - 3: **for** iteration $k = 0, 1, \dots$ **do**
 - 4: rollout T steps, store (s^k, a^k, r^k)
 - 5: compute advantages A_k with GAE
 - 6: **for all** samples k **do**
 - 7: $g^k \leftarrow \nabla_{\theta} \log \pi_{\theta}(a|s) A(s, a)$
 - 8: $\theta^{k+1} \leftarrow \theta^k + \eta \hat{\mathbf{F}}^k{}^{-1} g^k$ with equation 6
 - 9: **end for**
 - 10: update critic: $\phi^{k+1} \leftarrow \phi^k - \alpha \nabla_{\phi} \sum_t (R_t^k - V_{\phi}(s_t))^2$
 - 11: **end for**
-

5 Experimental Setup

To evaluate the performance of our proposed approach, SMAC, we provide experiments in several environments and compare against four baselines. We begin with three simple control environments, before moving on to six MuJoCo tasks (Todorov et al., 2012). For all experiments, we used the Gymnasium framework (Towers et al., 2024).

We implement SMAC on top of the A2C algorithm. In all environments, we evaluate SMAC against three baselines. To measure the benefits of using our NPG approximator, we consider an algorithm that is similar

| | AC-SGD | AC-Adam | AC-CG | AC-KFAC | SMAC (Ours) |
|--------------|-------------------|------------------------------------|------------------------------------|--------------------|-------------------------------------|
| Acrobot | -86.2 ± 4.7 | -88.4 ± 5.7 | -93.3 ± 11.6 | -170.7 ± 16.5 | -94.1 ± 17.0 |
| Cartpole | 962.3 ± 38.5 | 989.5 ± 14.1 | 977.9 ± 34.4 | 331.4 ± 142.3 | 973.4 ± 49.5 |
| Pendulum | -2116 ± 678 | -5731 ± 272 | -190 ± 23 | -1161 ± 15 | -2226 ± 1138 |
| Half Cheetah | -1672 ± 10.6 | -646 ± 9.7 | 1766 ± 864 | 513 ± 50.3 | 2936 ± 946 |
| Hopper | 180.5 ± 9.1 | 303.8 ± 23.9 | 172.7 ± 10.9 | 305.0 ± 33.6 | 331.2 ± 20.9 |
| Swimmer | 12.4 ± 3.0 | 17.3 ± 0.9 | 23.8 ± 1.7 | 24.2 ± 2.5 | 28.5 ± 1.4 |
| Walker | 213.4 ± 36.9 | 79.3 ± 58.7 | 229.8 ± 38.3 | 83.7 ± 72.9 | 204.1 ± 58.0 |
| Humanoid | 4539 ± 171 | 3105 ± 345 | 3175 ± 1445 | 260.8 ± 11.15 | 4625 ± 274 |
| Pusher | -433.6 ± 41.6 | -568.3 ± 16.4 | -441.0 ± 35.6 | -1086.3 ± 16.2 | -408.2 ± 31.5 |

Table 1: Average returns per episode, measured at the end of training, for each task in both classic control and MuJoCo environments. All results are averaged over 5 seeds (\pm standard deviations).

to SMAC in all aspects, except that it uses first-order updates instead of natural gradients. This is equivalent to A2C with a simple stochastic gradient descent (SGD) optimizer, which we denote as AC-SGD. To ensure a fair comparison to vanilla A2C, we also compare against the AC-Adam baseline, where A2C is trained using the more complex Adam optimizer. We also test SMAC against other NPG approximators. The AC-CG baseline uses the conjugate gradient algorithm to directly compute the NPG during the A2C update, without ever computing or storing the Fisher information matrix itself, similarly to TRPO (Schulman et al., 2015a). Finally, the second NPG approximator uses the Kronecker-factored approximate curvature (Martens & Grosse, 2015, KFAC). Our implementation, AC-KFAC, follows the application of KFAC to actor-critic algorithms as proposed by Wu et al. (2017), but omits the trust region.

SMAC and all the baselines we compare against compute advantages using the GAE. This is because GAE handles the bias-variance trade-off better than the simpler advantage estimator used by A2C, and has been shown to be effective in high-dimensional continuous control. To better handle continuous control problems, SMAC and all the baselines we compare against compute advantages using the GAE (Schulman et al., 2015b).

For each of the five algorithms in each environment, we train five models using different random seeds. The results reported throughout training are averaged over these seeds, with standard deviations included. We evaluate agent performance using (undiscounted) returns, averaged across episodes. Besides the maximum average returns achieved, we are also interested in convergence speed.

Furthermore, we also track the average log-probabilities of the actions taken. We treat these as a measure of the agents’ confidence in their actions. We therefore expect a well-performing policy to become more deterministic, achieving both higher log-probabilities and higher returns. On the contrary, we can also diagnose overconfident policies, where an increase in log-probabilities does not correlate to an increase in returns. The agent may therefore prematurely converge to an overly deterministic sub-optimal policy, which is likely to hinder exploration and thus future improvements. We report the log-probabilities measured in the supplementary material. Furthermore, we discuss the average log-probabilities of the actions taken as a measure of the agents’ confidence in the actions they take. We report these in detail in the supplementary material.

In all environments, episodes have a horizon of 1000 timesteps. To increase the readability of the plots, we group timesteps into 100 bins with no overlap and report the average of each bin. We apply additional smoothing by computing an exponentially weighted mean over the bins, with a smoothing factor of 0.1. The average performances and standard deviations we report are then computed from these binned and smoothed results. However, for the sake of completeness, we also plot the raw, unbinned, and unsmoothed results in the supplementary material.

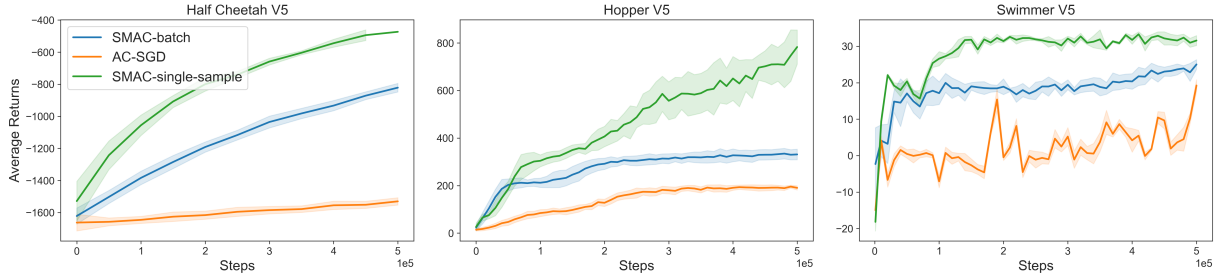


Figure 3: Effect of Batch Size.

5.1 Practical MethodAblation: Effect of Batch Size

The matrix-free updates operate on a single trajectory sample. Although the per-step cost is only $\mathcal{O}(d)$, when applying this per-sample scheme to a full batch of size b , the update must be executed b times, so the effective time complexity for the batch becomes $\mathcal{O}(bd)$, so the computational cost can be high to update a whole batch. ~~many more updates are required to process a fixed number of transitions, so the computational cost can be very high.~~

To investigate this trade-off, we evaluate an alternative that uses the average gradient of log-probability over all sampled state-action transitions in the batch:

$$\bar{\ell} = \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} \log \pi_{\theta}(a_i | s_i).$$

Where B denotes the total number of state-action pairs collected in the batch. The regularized empirical Fisher then becomes

$$\tilde{\mathbf{F}} \hat{\mathbf{F}}^k = \lambda \mathbf{I} + \bar{\ell}^k (\bar{\ell}^k)^{\top},$$

This is a rank-one moment approximation that assumes vectors in batch have small angular variance. The bias can be written explicitly:

$$\mathbb{E}[\tilde{\mathbf{F}} - \hat{\mathbf{F}}] = \frac{1}{B} \text{Cov}(\ell) \succeq 0$$

so $\tilde{\mathbf{F}}$ underestimates $\hat{\mathbf{F}}$ by a positive semidefinite term that shrinks as gradients concentrate or as B grows. and the update direction of SMAC is given by

$$\theta^{k+1} \underline{\Delta \theta^k} = \theta^k - \eta \left[\frac{1}{\lambda} g^k - \frac{\bar{\ell}^k (\bar{\ell}^k)^{\top} g^k}{\lambda^2 + \lambda (\bar{\ell}^k)^{\top} \bar{\ell}^k} \right].$$

Results. Fig. 3 compares the single-sample and batch-mean variants on Mujoco HalfCheetah and Hopper. Using $B = 1000$ reduces the overall optimization time by roughly 80% while achieving returns that remain close to the single-sample baseline and clearly outperform the standard A2C optimized with SGD. Therefore, to further balance computational resources and performance, we implement the version with a batch size of 1000 in all subsequent experiments.

5.2 Effect of damping coefficient λ

We use the EF $\hat{\mathbf{F}}^k = \lambda \mathbf{I} + \ell^k (\ell^k)^\top$ with $\lambda > 0$. After applying the Sherman-Morrison formula, its inverse has the closed form

$$(\lambda \mathbf{I} + \ell^k (\ell^k)^\top)^{-1} = \frac{1}{\lambda} \left(\mathbf{I} - \frac{\ell^k (\ell^k)^\top}{\lambda + (\ell^k)^\top \ell^k} \right).$$

It scales directions unevenly, components perpendicular to ℓ are reduced by a factor $\frac{1}{\lambda}$, while the component along ℓ is reduced by $\frac{1}{\lambda + (\ell^k)^\top \ell^k}$.

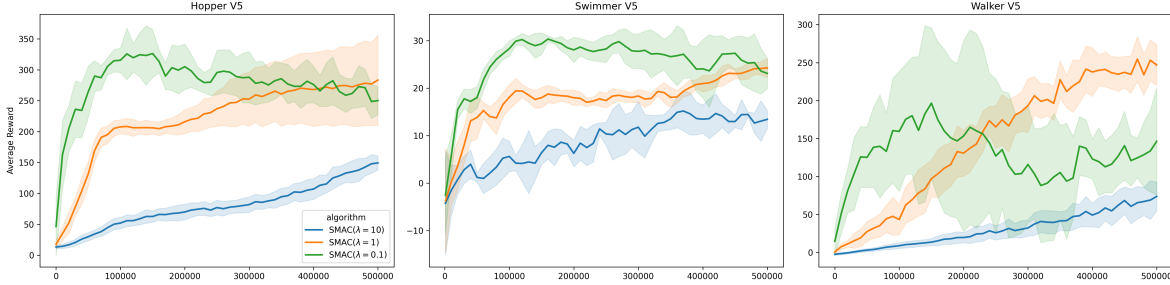


Figure 4: Effect of different λ

Results. We evaluate SMAC on three Mujoco tasks, compare three damping coefficient value $\lambda \in \{0.1, 1, 10\}$. Overall, the results agree with our analysis: too large λ slows learning, too small λ can make the update noisy, and a moderate λ gives the best stability-speed tradeoff on average. The choices in practice are more complex, for example, the choice of λ is coupled with the learning rate because both scale the step, so a balance is needed. A feasible path is to use a small damping coefficient to keep the EF term close to its undamped form, and then apply a normalization so the update magnitude is preserved.

5.3 Effect of previous estimate EF

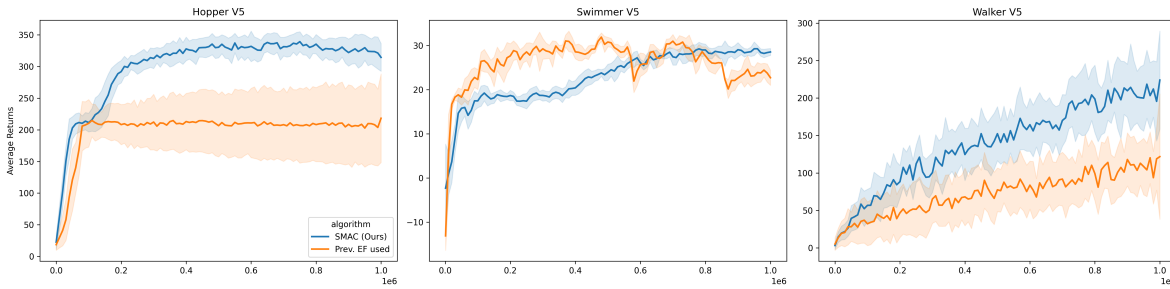


Figure 5: Previous EF used vs. Ours

Using the previous EF not only increases the memory and time complexity to $\mathcal{O}(d^2)$, but introduces a lag bias: after each update the policy and state distribution change, so the old EF reflects the old distribution and can mis-scale directions when curvature changes quickly. In this experiments we used a mixed estimator $\mathbf{F}^k = \beta \mathbf{F}^{k-1} + (1 - \beta) \ell^k (\ell^k)^\top$ with a fixed $\beta = 0.5$ to limit reliance on the old EF. Results show that on Swimmer it gives faster early improvement but remains less stable, while on Hopper and Walker it learns more slowly with larger variance. Therefore, if the KL change is strictly constrained, using the previous estimate EF is still feasible, and the choice of β should depend on the rate of policy change.

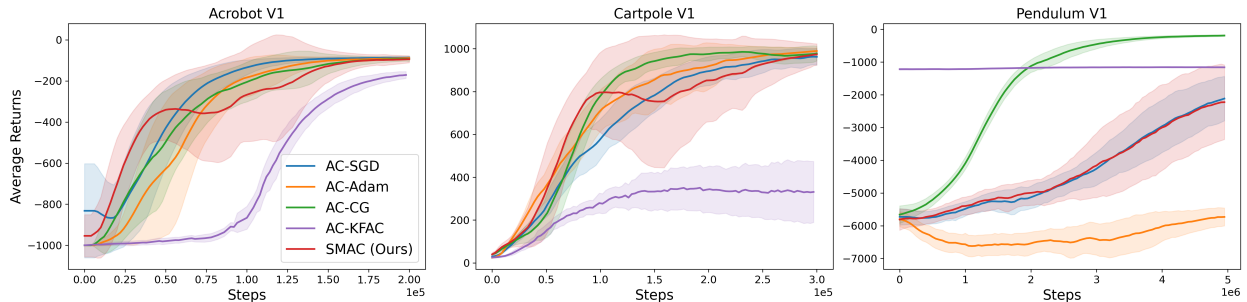


Figure 6: Results on three classic control tasks. We report the average return per episode. All results are first grouped into bins and smoothed, and then averaged over 5 random seeds, with shaded areas showing standard deviation.

6 Results

We provide the final performance achieved by training each of the five algorithms on each of the nine tasks. The results in Table 1 show that our method, SMAC, can achieve competitive results in two of the three simple control tasks. More importantly, the benefits of SMAC appear to increase in more complex environments, such as MuJoCo, where SMAC outperforms the other baselines in terms of final returns in all but one task. A more in-depth analysis of these results, with a special focus on convergence rate, is given in this section.

6.1 Classic Control

In the classic control environment, our method, SMAC, reaches a high initial average return faster than the baselines on two of the tasks, as shown in Fig. 6 (see Fig. 8 for raw results). This, however, comes at the cost of additional variance across seeds. The overall performance of SMAC is also competitive in those two tasks.

In *Acrobot*, SMAC only requires 40000 timesteps to pass the average return threshold of -400 . This is faster than AC-Adam, AC-SGD, AC-CG, and AC-KFAC, which require 70%, 35%, 50%, and 235% more timesteps, respectively. This accelerated learning is, however, followed by unstable training until approximately 125000 timesteps. Despite this, SMAC stabilizes and converges, reaching a final performance of -94.1 ± 6.3 , which is similar to three of the baselines, and higher than AC-KFAC.

A similar behavior can be observed in *Cartpole*. SMAC agents achieve 75% of the maximum return possible in 87000 timesteps, slightly faster than AC-CG, which needs approximately 10% more timesteps. AC-Adam and AC-SGD are even slower, requiring 28% and 62% more timesteps, respectively, to reach the same threshold. At the same time, the average performance of AC-KFAC never improves above 40% of the optimal policy’s return. SMAC becomes more unstable after this initial period, but stabilizes and reaches a competitive performance of 973.4 ± 29.8 by the end of training.

In contrast, in the last environment, *Pendulum*, our method fails to outperform AC-CG in terms of convergence speed and overall performance. AC-CG requires only 1.5 million timesteps to match the final performance of SMAC, indicating that our proposed approximation of the FIM can struggle to improve learning in some environments. Moreover, SMAC is also outperformed by AC-KFAC, which quickly finds a good policy but then stops learning. However, SMAC still achieves a final return of -2226 ± 859 , outperforming AC-Adam and matching AC-SGD, with the former also displaying much slower and unstable learning.

As proposed in Sec. 5, the validity of our results is also confirmed by the action log-probabilities recorded (Fig. 10). In *Acrobot* and *Cartpole*, the two tasks in which SMAC performs competitively, the log-probabilities converge faster than those of the baseline methods. Together with the increase in performance achieved by

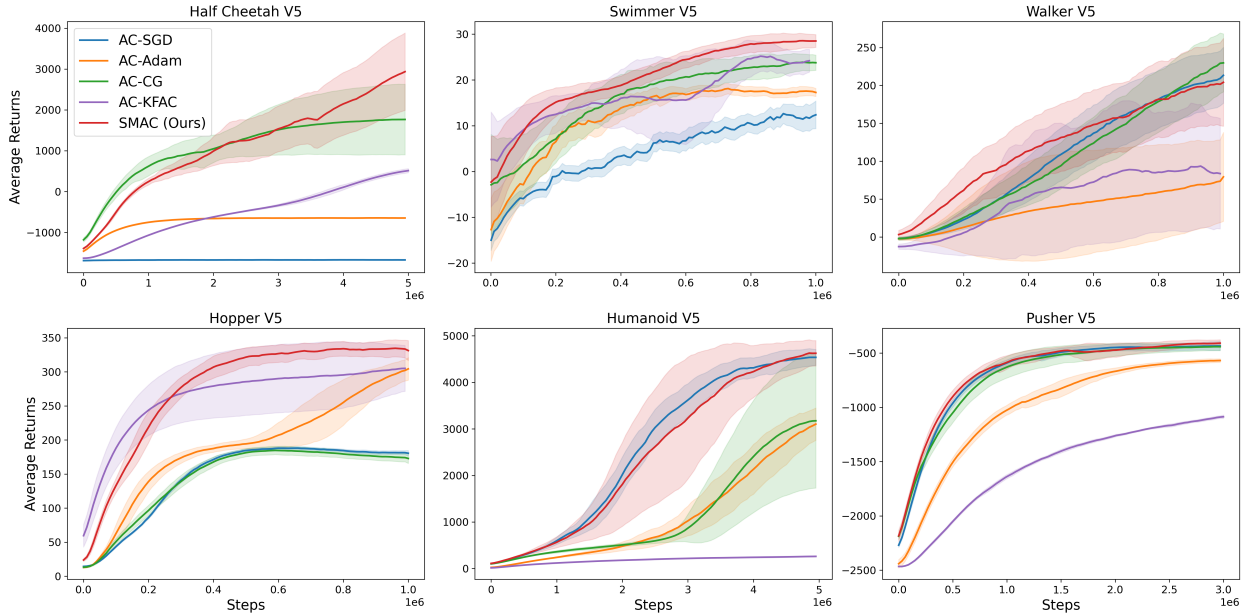


Figure 7: Results on six MuJoCo tasks. We report the average return per episode. All results are first grouped into bins and smoothed, and then averaged over 5 random seeds, with shaded areas showing standard deviation.

SMAC in these tasks, this may indicate that our agents can rapidly become confident in their policies when taking actions that yield increasingly higher returns. At the same time, in *Pendulum*, the log-probabilities of SMAC match those of AC-SGD and are bounded above by AC-CG. This could indicate that our method becomes overconfident in its actions too early, finally leading to a sub-optimal policy.

6.2 MuJoCo

In the more complex MuJoCo environment, we have found that throughout training, our method, SMAC, consistently outperforms or matches the baseline algorithms in terms of performance in 5 out of 6 tasks, as shown in Fig. 7 (see Fig. 9 for raw results) (Fig. 7). In the sixth task, SMAC still achieves comparable results. Furthermore, SMAC exhibits a noticeably higher sample efficiency in three tasks.

SMAC shows improvements over existing methods in *Half Cheetah*, where it achieves an average return of 2936 ± 946 , surpassing the next-best algorithm, AC-CG, by over 1100 points. While it learns more slowly for the first 2 million timesteps, SMAC ultimately matches the final performance of AC-CG after only 3.4 million timesteps, which represents just 68% of the training budget. It then continues to improve, while AC-CG plateaus. This superior performance comes at the cost of lower stability, with both SMAC and AC-CG being much more unstable than AC-KFAC, and the simpler AC-Adam and AC-SGD. The largest improvement is in *Half Cheetah*, where SMAC achieves an average return of 863.8 ± 38.7 by the end of training. The next-best algorithm, AC-CG, achieves only 444.6 ± 12.1 , being surpassed by our method by over 1300 points. Importantly, SMAC also learns three times faster, matching the final performance of AC-CG after only 1.7 million timesteps. Additionally, note that while SMAC has a higher variance across seeds mid-training, this reduces by the end of learning.

Similar trends can be observed in *Hopper*, where SMAC achieves the final performance of AC-Adam using almost three times fewer timesteps. Our method achieves a final average return of 331.2 ± 29.7 , outperforming the second-best performances of AC-Adam and AC-KFAC by approximately 30 points. The stability of SMAC appears to be comparable to that of AC-Adam, higher than that of AC-KFAC, but lower than that of the other two baselines.

Likewise, in *Swimmer*, SMAC consistently outperforms the second-best baselines, with a final average return of 28.5 ± 0.8 , which is approximately 5 points higher than those of AC-CG and AC-KFAC. Moreover, it learns the fastest, matching the peak performance of AC-CG almost twice as fast. The stability of SMAC is also greater than that of AC-SGD, AC-CG, and AC-KFAC, and competitive with that of AC-Adam.

In *Walker*, our method learns faster than the baselines in the first 60% of the timesteps. While its final performance of 204.1 ± 87.9 is lower than those of AC-CG and AC-SGD, which achieve 229.8 ± 23.9 and 213.4 ± 42.9 , respectively, SMAC reaches the threshold of 115 (half of the final return of AC-CG) in only 410000 timesteps, while AC-CG itself needs 39% more. Additionally, our approach outperforms AC-Adam and AC-KFAC in both average return and stability.

For the *Humanoid* environment, SMAC achieves the top final performance of 4625 ± 277 average return, followed closely by AC-SGD with 4539 ± 110 . It greatly outperforms AC-Adam and AC-CG in terms of both final returns and final stability, with only AC-SGD and AC-KFAC being more stable. However, the performance of AC-KFAC is much lower than that of any of the other algorithms. While AC-SGD is slightly faster than our method, SMAC still learns much faster than AC-Adam and AC-CG, matching their final performance after using just 59% of the total allocated training budget.

In *Pusher*, the performance of SMAC closely follows that of AC-SGD, slightly outperforming the latter with a final average return of -408.2 ± 32.2 . Our method is also more sample efficient, passing the -441 return threshold (final performance of AC-CG) 21% faster than AC-CG and 9% faster than AC-SGD. While all agents stabilize by the end of training, only AC-Adam and AC-KFAC end up being more stable than SMAC, with AC-SGD having the highest standard deviation.

Finally, we turn our attention to the average log-probabilities of the actions taken during training (Fig. 11), as proposed in Sec. 5. The increase in log-probabilities shown by SMAC is stable and correlates with the increase in average returns given in Fig. 7. This may indicate that the policy quickly becomes confident in the actions it chooses, without becoming overconfident. Note that the log-probabilities of our method converge the fastest in 24 out of 6 tasks. While the log-probabilities of AC-CG converge first in *Humanoid* and *Pusher*, this does not lead to better performance. On the contrary, in the former task, AC-CG agents become overconfident in their sub-optimal actions. While in *Half Cheetah* the log-probabilities of AC-CG converge first, SMAC ultimately achieves a higher confidence in its actions, which matches its superior final performance. A similar trend appears in *Walker*, with AC-KFAC converging the fastest, but to a sub-optimal policy. The slow convergence of the other two agents’ log-probabilities in most tasks may indicate uncertainty in their actions, even when their performance matches that of SMAC.

7 Conclusion

In this work, we introduced a simple yet effective rank-1 approximation to the Natural Policy Gradient in the Actor-Critic framework. By using a regularized Empirical Fisher matrix and the Sherman–Morrison formula, we enable scalable and efficient natural policy gradient updates with only $\mathcal{O}(d)$ complexity. We theoretically prove the convergence properties of this approximation and integrate it into the standard Actor-Critic algorithm to demonstrate its practical utility. Experimental results on classic control and MuJoCo tasks show that our method achieves both fast and stable convergence in most environments, outperforms methods such as Advantage Actor-Critic and conjugate gradient in TRPO in several environments, offering a promising alternative for efficient NPG approximations.

8 Broader Impact

Reinforcement learning is now used in many domains, from fine-tuning language models to controlling systems such as power grids or data centres, and this broad use can bring both benefits and risks. It can improve efficiency, safety, and energy use, but it can also amplify existing biases or support decision systems that have harmful social or economic effects. Our work focuses only on the optimisation aspect: we study how to solve a given policy optimisation problem faster, so that a fixed objective is reached in fewer gradient steps. If the training data for language models or the environment model for control tasks reflects real-world

dynamics and values accurately, a more efficient optimiser can produce better policies with less computation and a smaller training carbon footprint. At the same time, any negative societal impact would mainly arise from biased, incomplete, or misaligned data, objectives, and reward designs, rather than from the particular optimisation method used.

References

- Alekh Agarwal, Sham M. Kakade, Jason D. Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76, 2021. URL <http://jmlr.org/papers/v22/19-736.html>.
- Shun-ichi Amari. Natural gradient works efficiently in learning. *Neural Comput.*, 10(2):251–276, 1998. doi: 10.1162/089976698300017746. URL <https://doi.org/10.1162/089976698300017746>.
- Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. Deep reinforcement learning: A brief survey. *IEEE signal processing magazine*, 34(6):26–38, 2017.
- J. Andrew Bagnell and Jeff Schneider. Covariant policy search. In *Proceedings of the 18th International Joint Conference on Artificial Intelligence, IJCAI’03*, pp. 1019–1024, San Francisco, CA, USA, 2003. Morgan Kaufmann Publishers Inc.
- Achraf Bahamou, Donald Goldfarb, and Yi Ren. A mini-block fisher method for deep neural networks. In Francisco J. R. Ruiz, Jennifer G. Dy, and Jan-Willem van de Meent (eds.), *International Conference on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia, Spain*, volume 206 of *Proceedings of Machine Learning Research*, pp. 9191–9220. PMLR, 2023. URL <https://proceedings.mlr.press/v206/bahamou23a.html>.
- Jonathan Baxter and Peter L. Bartlett. Infinite-horizon policy-gradient estimation. *J. Artif. Intell. Res.*, 15: 319–350, 2001. doi: 10.1613/JAIR.806. URL <https://doi.org/10.1613/jair.806>.
- Suzanna Becker and Yann Lecun. Improving the convergence of back-propagation learning with second-order methods. In *Proceedings of the 1988 Connectionist Models Summer School*, 01 1989.
- Frederik Benzing. Gradient descent on neurons and its link to approximate second-order optimization. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp. 1817–1853. PMLR, 2022. URL <https://proceedings.mlr.press/v162/benzing22a.html>.
- Yuhao Ding, Junzi Zhang, and Javad Lavaei. On the global optimum convergence of momentum-based policy gradient. In *International Conference on Artificial Intelligence and Statistics*, pp. 1910–1934. PMLR, 2022.
- Ivo Grondman, Lucian Busoniu, Gabriel AD Lopes, and Robert Babuska. A survey of actor-critic reinforcement learning: Standard and natural policy gradients. *IEEE Transactions on Systems, Man, and Cybernetics, part C (applications and reviews)*, 42(6):1291–1307, 2012.
- Sham M. Kakade. A natural policy gradient. In Thomas G. Dietterich, Suzanna Becker, and Zoubin Ghahramani (eds.), *Advances in Neural Information Processing Systems 14 [Neural Information Processing Systems: Natural and Synthetic, NIPS 2001, December 3-8, 2001, Vancouver, British Columbia, Canada]*, pp. 1531–1538. MIT Press, 2001. URL <https://proceedings.neurips.cc/paper/2001/hash/4b86abe48d358ecf194c56c69108433e-Abstract.html>.
- Frederik Kunstner, Philipp Hennig, and Lukas Balles. Limitations of the empirical fisher approximation for natural gradient descent. In Hanna M. Wallach, Hugo Larochelle, Alina Beygelzimer, Florence d’Alché-Buc, Emily B. Fox, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pp. 4158–4169, 2019. URL <https://proceedings.neurips.cc/paper/2019/hash/46a558d97954d0692411c861cf78ef79-Abstract.html>.

- Yann LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pp. 9–50. Springer, 2002.
- Yanli Liu, Kaiqing Zhang, Tamer Basar, and Wotao Yin. An improved analysis of (variance-reduced) policy gradient and natural policy gradient methods. *CoRR*, abs/2211.07937, 2022. doi: 10.48550/ARXIV.2211.07937. URL <https://doi.org/10.48550/arXiv.2211.07937>.
- James Martens. Deep learning via hessian-free optimization. In Johannes Fürnkranz and Thorsten Joachims (eds.), *Proceedings of the 27th International Conference on Machine Learning (ICML-10), June 21-24, 2010, Haifa, Israel*, pp. 735–742. Omnipress, 2010. URL <https://icml.cc/Conferences/2010/papers/458.pdf>.
- James Martens. New insights and perspectives on the natural gradient method. *J. Mach. Learn. Res.*, 21: 146:1–146:76, 2020. URL <https://jmlr.org/papers/v21/17-678.html>.
- James Martens and Roger B. Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 2408–2417. JMLR.org, 2015. URL <http://proceedings.mlr.press/v37/martens15.html>.
- 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*, pp. 1928–1937. PmLR, 2016.
- Jan Peters, Sethu Vijayakumar, and Stefan Schaal. Natural actor-critic. In *European Conference on Machine Learning*, pp. 280–291. Springer, 2005.
- Matteo Pirodda, Marcello Restelli, and Luca Bascetta. Adaptive step-size for policy gradient methods. In Christopher J. C. Burges, Léon Bottou, Zoubin Ghahramani, and Kilian Q. Weinberger (eds.), *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States*, pp. 1394–1402, 2013. URL <https://proceedings.neurips.cc/paper/2013/hash/f64eac11f2cd8f0efa196f8ad173178e-Abstract.html>.
- Yi Ren and Donald Goldfarb. Efficient subsampled gauss-newton and natural gradient methods for training neural networks. *CoRR*, abs/1906.02353, 2019. URL <http://arxiv.org/abs/1906.02353>.
- Yi Ren and Donald Goldfarb. Tensor normal training for deep learning models. In Marc’Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 26040–26052, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/dae3312c4c6c7000a37ecfb7b0aeb0e4-Abstract.html>.
- Silvia Richter, Douglas Aberdeen, and Jin Yu. Natural actor-critic for road traffic optimisation. *Advances in neural information processing systems*, 19, 2006.
- Nicolas Le Roux, Pierre-Antoine Manzagol, and Yoshua Bengio. Topmoumoute online natural gradient algorithm. In John C. Platt, Daphne Koller, Yoram Singer, and Sam T. Roweis (eds.), *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pp. 849–856. Curran Associates, Inc., 2007. URL <https://proceedings.neurips.cc/paper/2007/hash/9f61408e3afb633e50cdf1b20de6f466-Abstract.html>.
- Nicol N Schraudolph. Fast curvature matrix-vector products for second-order gradient descent. *Neural computation*, 14(7):1723–1738, 2002.

- John Schulman, Sergey Levine, Pieter Abbeel, Michael I. Jordan, and Philipp Moritz. Trust region policy optimization. In Francis R. Bach and David M. Blei (eds.), *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 1889–1897. JMLR.org, 2015a. URL <http://proceedings.mlr.press/v37/schulman15.html>.
- John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional continuous control using generalized advantage estimation. *arXiv preprint arXiv:1506.02438*, 2015b.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Jack Sherman and Winifred J. Morrison. Adjustment of an inverse matrix corresponding to a change in one element of a given matrix. *Annals of Mathematical Statistics*, 21(1):124–127, 1950. doi: 10.1214/aoms/1177729893.
- Sidak Pal Singh and Dan Alistarh. Woodfisher: Efficient second-order approximation for neural network compression. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL <https://proceedings.neurips.cc/paper/2020/hash/d1ff1ec86b62cd5f3903ff19c3a326b2-Abstract.html>.
- Zhiqiang Su, Meng Zhou, Fangfang Han, Yiwu Zhu, Dalei Song, and Tingting Guo. Attitude control of underwater glider combined reinforcement learning with active disturbance rejection control. *Journal of Marine Science and Technology*, 24(3):686–704, 2019.
- Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In S. Solla, T. Leen, and K. Müller (eds.), *Advances in Neural Information Processing Systems*, volume 12. MIT Press, 1999. URL https://proceedings.neurips.cc/paper_files/paper/1999/file/464d828b85b0bed98e80ade0a5c43b0f-Paper.pdf.
- Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 5026–5033. IEEE, 2012. doi: 10.1109/IROS.2012.6386109.
- Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
- Wouter JA van Heeswijk. Natural policy gradients in reinforcement learning explained. *arXiv preprint arXiv:2209.01820*, 2022.
- Lingxiao Wang, Qi Cai, Zhuoran Yang, and Zhaoran Wang. Neural policy gradient methods: Global optimality and rates of convergence. *arXiv preprint arXiv:1909.01150*, 2019.
- Xu Wang, Sen Wang, Xingxing Liang, Dawei Zhao, Jincai Huang, Xin Xu, Bin Dai, and Qiguang Miao. Deep reinforcement learning: A survey. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):5064–5078, 2022.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3):229–256, 1992.
- Xiaodong Wu, Wenyi Yu, Chao Zhang, and Philip C. Woodland. An improved empirical fisher approximation for natural gradient descent. In Amir Globersons, Lester Mackey, Danielle Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural Information Processing Systems 38: Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024. URL http://papers.nips.cc/paper_files/paper/2024/hash/f23098fa0cfcdef0e743b134d380eeb9-Abstract-Conference.html.

- Yuhuai Wu, Elman Mansimov, Roger B. Grosse, Shun Liao, and Jimmy Ba. Scalable trust-region method for deep reinforcement learning using kronecker-factored approximation. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp. 5279–5288, 2017. URL <https://proceedings.neurips.cc/paper/2017/hash/361440528766bbaaaa1901845cf4152b-Abstract.html>.
- Pan Xu, Felicia Gao, and Quanquan Gu. An improved convergence analysis of stochastic variance-reduced policy gradient. In Amir Globerson and Ricardo Silva (eds.), *Proceedings of the Thirty-Fifth Conference on Uncertainty in Artificial Intelligence, UAI 2019, Tel Aviv, Israel, July 22-25, 2019*, volume 115 of *Proceedings of Machine Learning Research*, pp. 541–551. AUAI Press, 2019. URL <http://proceedings.mlr.press/v115/xu20a.html>.
- Pan Xu, Felicia Gao, and Quanquan Gu. Sample efficient policy gradient methods with recursive variance reduction. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJlxIJBFDr>.
- Minghan Yang, Dong Xu, Zaiwen Wen, Mengyun Chen, and Pengxiang Xu. Sketch-based empirical natural gradient methods for deep learning. *J. Sci. Comput.*, 92(3):94, 2022. doi: 10.1007/S10915-022-01911-X. URL <https://doi.org/10.1007/s10915-022-01911-x>.
- Rui Yuan, Robert M. Gower, and Alessandro Lazaric. A general sample complexity analysis of vanilla policy gradient. In Gustau Camps-Valls, Francisco J. R. Ruiz, and Isabel Valera (eds.), *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pp. 3332–3380. PMLR, 28–30 Mar 2022. URL <https://proceedings.mlr.press/v151/yuan22a.html>.
- Lin Zhang, Shaohuai Shi, and Bo Li. Eva: Practical second-order optimization with kronecker-vectorized approximation. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL https://openreview.net/forum?id=_Mic8V96Voy.
- Tingting Zhao, Hirotaka Hachiya, Gang Niu, and Masashi Sugiyama. Analysis and improvement of policy gradient estimation. *Neural Networks*, 26:118–129, 2012. doi: 10.1016/J.NEUNET.2011.09.005. URL <https://doi.org/10.1016/j.neunet.2011.09.005>.

A Appendix

A.1 Notations.

Let $\gamma \in (0, 1)$ be the discount, $r : \mathcal{S} \times \mathcal{A} \rightarrow [-R, R]$ is the reward function, and $H \in \mathbb{N} \cup \{\infty\}$ the horizon. Define the truncated objective

$$J^H(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[\sum_{t=0}^{H-1} \gamma^t r(s_t, a_t) \right], \quad J(\theta) = J^\infty(\theta). \quad (17)$$

For each trajectory $\tau^H = (s_0, a_0, \dots, s_{H-1}, a_{H-1})$ sampled under π_θ , let $g(\tau^H \mid \theta)$ denote a single-trajectory unbiased estimator of the truncated policy gradient:

$$\mathbb{E}[g(\tau^H \mid \theta)] = \nabla J^H(\theta). \quad (18)$$

$$g(\tau^H \mid \theta) = \sum_{t=0}^{H-1} \gamma^t G_t \nabla_\theta \log \pi_\theta(a_t \mid s_t) \quad (19)$$

A.2 Assumptions

Assumption 1. *Let $g(\tau^H \mid \theta)$ be an unbiased single-trajectory estimator of $\nabla J^H(\theta)$. Assume there exists $\sigma > 0$ such that for all θ ,*

$$\mathbb{E}[\|g(\tau^H \mid \theta) - \mathbb{E}[g(\tau^H \mid \theta)]\|^2] \leq \sigma^2,$$

Assumption 2. *The Fisher Information matrix induced by the policy π_θ is Positive Definite and satisfies*

$$\mathbf{F}(\theta) \mathbf{F}_\rho(\theta) = \mathbb{E}_{s \sim d_\rho^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta} [\nabla_\theta \log \pi_\theta \cdot \nabla_\theta \log \pi_\theta^T] \succcurlyeq \mu_F \mathbf{I} \quad (20)$$

Assumption 3. *The policy gradient estimates are bounded by a constant $G \geq 0$ and the change in policy gradient at different timestep are bounded by the change in policy parametrization weights. So,*

$$\|\nabla_\theta \log \pi_\theta\| \leq G \quad (21)$$

$$\|\nabla_\theta \log \pi_{\theta_1} - \nabla_\theta \log \pi_{\theta_2}\| \leq M \|\theta_1 - \theta_2\| \quad (22)$$

Assumption 4. Transferrable Compatible Function Approximator: *When we approximate the Advantage function using the policy π_θ it induces a transfer error defined as ϵ_{bias} which is zero for softmax parametrization and is very small for dense neural policy class.*

$$L_{\nu^*}(w_\star^\theta; \theta) = \mathbb{E}_{(s,a) \sim \nu^*} \left[(A^{\pi_\theta}(s, a) - (1 - \gamma)(w_\star^\theta)^\top \nabla_\theta \log \pi_\theta(a \mid s))^2 \right] \leq \epsilon_{bias} \quad (23)$$

where $\nu^*(s, a) = d_\rho^{\pi^*}(s) \cdot \pi^*(s, a)$ is the state-action distribution induced by an optimal policy π^* and $w_\star^\theta = \operatorname{argmin}_w L_{\nu_\rho^{\pi_\theta}}(w; \theta)$ is obtained via the full natural policy gradient direction at the parametrization θ .

Remark 1. Regarding $\mathbf{F}(\theta)$, we know from Assumptions 1 and 2 that

$$\mu_F I_d \preccurlyeq \mathbf{F}(\theta) \preccurlyeq G^2 I_d \text{ for any } \theta \in \mathbf{R}^d.$$

Remark 2. Now we can get an upper bound to the second term of Equation 5 using Remark 1 as:

$$\left\| \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \operatorname{Tr}(\hat{\mathbf{F}}^k)} \right\| \leq \frac{G^2}{\lambda^2 + \lambda \mu_F}.$$

Remark 3. Assumption 1 treats σ^2 as a constant upper bound on the variance of the truncated policy gradient estimator, although in general $\text{Var}(g(\tau^H | \theta))$ may depend on the horizon H . More precisely, for each finite H we write

$$\text{Var}(g(\tau^H | \theta)) = \sigma_H^2.$$

For likelihood-ratio estimators such as REINFORCE, existing analyses show that σ_H^2 can grow with H , but remains finite under smooth Gaussian policies with bounded rewards and features (Zhao et al., 2012)(Pirodda et al., 2013)(Baxter & Bartlett, 2001). In our proof we use a finite truncated horizon and define

$$\sigma^2 := \sup_{1 \leq H \leq H_{\max}} \sigma_H^2 < \infty.$$

This form of bounded-variance assumption is standard in recent convergence analyses of policy gradient methods (Liu et al., 2022) (Xu et al., 2020)(Xu et al., 2019).

A.3 Helper Lemmas

We will establish some helper lemma to prove the global convergence of Sherman-Morrison policy update.

Lemma 1. (Liu et al., 2022) In the stochastic PG update, ~~by choosing $\eta = \frac{1}{4L_J}$, $K = \frac{32L_J(J^{H,*} - J^H(\theta_0))}{\varepsilon}$, and $N = \frac{6\sigma^2}{\varepsilon}$~~ , we have

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2] \leq \frac{\frac{J^{H*} - J^H(\theta_0)}{K} + (\frac{\eta}{2} + L_J \eta^2) \frac{\sigma^2}{N}}{\frac{\eta}{2} - L_J \eta^2} \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2] \leq \varepsilon. \quad (24)$$

In total, stochastic PG samples $\mathcal{O}\left(\frac{\sigma^2}{(1-\gamma)^2 \varepsilon^2}\right)$ trajectories.

Lemma 2. (Liu et al., 2022)~~(Xu et al., 2020)~~ Let $J(\theta), J^H(\theta)$ are L_J -smooth, where $L_J = \frac{MR}{(1-\gamma)^2} + \frac{2G^2R}{(1-\gamma)^3}$ and

$$\|\nabla J^H(\theta) - \nabla J(\theta)\| \leq GR \left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right) \gamma^H$$

Lemma 3. (Liu et al., 2022)~~Let $g^k = \frac{1}{N} \sum_{i=1}^N g(\tau_i^H | \theta^k)$ be the policy gradient estimate for samples until horizon H . Then, we have~~

Lemma 4. (Liu et al., 2022) Let $\mathbf{w}_*^k = \underline{F^{-1} F_\rho^{-1}(\theta^k)} \nabla J(\theta^k)$ be the exact NPG update direction at θ^k . Then, we have

$$\begin{aligned} J(\pi^*) - \frac{1}{K} \sum_{k=0}^{K-1} J(\theta^k) &\leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \frac{1}{\eta K} \mathbb{E}_{s \sim d_{\rho^*}} [KL(\pi^*(\cdot|s) || \pi_{\theta^0}(\cdot|s))] \\ &\quad + \frac{G}{K} \sum_{k=0}^{K-1} \|\mathbf{w}^k - \mathbf{w}_*^k\| + \frac{M\eta}{2K} \sum_{k=0}^{K-1} \|\mathbf{w}^k\|^2. \end{aligned} \quad (25)$$

A.4 Global Convergence of SM Update

Let us take \mathbf{w}^k as the update direction of Sherman-Morrison. To this end, we need to upper bound $\frac{1}{K} \sum_{k=0}^{K-1} \|\mathbf{w}^k - \mathbf{w}_*^k\|$, $\frac{1}{K} \sum_{k=0}^{K-1} \|\mathbf{w}^k\|^2$, and $\frac{1}{K} \mathbb{E}_{s \sim d_{\rho^*}} [KL(\pi^*(\cdot|s) || \pi_{\theta^0}(\cdot|s))]$, where $\mathbf{w}_*^k = F^{-1}(\theta^k) \nabla J(\theta^k)$ is the exact NPG update direction at θ^k .

- Bounding $\frac{1}{K} \sum_{k=0}^{K-1} \|\mathbf{w}^k - \mathbf{w}_\star^k\|$.

We know from Jensen's inequality [that](#) $\left(\mathbb{E}[\|\mathbf{w}_t^{j+1} - \mathbf{w}_{t,\star}^{j+1}\|]\right)^2 \leq \mathbb{E}[\|\mathbf{w}_t^{j+1} - \mathbf{w}_{t,\star}^{j+1}\|^2]$ [so](#), [and](#)

$$\begin{aligned}
& \left(\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}_\star^k\|] \right)^2 \\
& \leq \frac{1}{K} \sum_{k=0}^{K-1} (\mathbb{E}[\|\mathbf{w}^k - \mathbf{w}_\star^k\|])^2 \\
& \leq \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}_\star^k\|^2] \\
& \leq \frac{2}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \nabla J(\theta^k)\|^2] + \frac{2}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J(\theta^k) - \mathbf{w}_\star^k\|^2]
\end{aligned} \tag{26}$$

Let

$$\mathbf{g}^k = \frac{1}{N} \sum_{i=1}^N g(\tau_i^H | \theta^k),$$

be the PG direction and

$$\mathbf{w}^k = \frac{1}{\lambda} \mathbf{g}^k - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} \mathbf{g}^k$$

be the weight update from the Sherman-Morrison update then we have from [Lemma 2](#) and [Remark 2](#) [Lemma 3](#) and [Assumption 3](#) that

$$\begin{aligned}
& \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \nabla J(\theta^k)\|^2] \\
& = \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \frac{1}{\lambda} \mathbf{g}^k - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} \mathbf{g}^k - \nabla J(\theta^k) \right\|^2 \right] \\
& = \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \frac{1}{\lambda} (\mathbf{g}^k - \nabla J^H(\theta^k)) - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} (\mathbf{g}^k - \nabla J^H(\theta^k)) - (\nabla J(\theta^k) - \nabla J^H(\theta^k)) \right. \right. \\
& \quad \left. \left. + \frac{1}{\lambda} \nabla J^H(\theta^k) - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} \nabla J^H(\theta^k) - \nabla J^H(\theta^k) \right\|^2 \right] \\
& \leq \frac{\sigma^2}{N} \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) + \frac{42}{N} G^2 R^2 \left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right)^2 \gamma^{2H} \\
& \quad + \frac{42}{N} \left(1 + \frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2]
\end{aligned} \tag{27}$$

Furthermore, Assumption 2 tells us that

$$\begin{aligned}
& \mathbb{E}[\|\nabla J(\theta^k) - \mathbf{w}_\star^k\|^2] \\
& = \mathbb{E}[\|\nabla J(\theta^k) - F^{-1}(\theta^k) \nabla J(\theta^k)\|^2] \\
& \leq \left(1 + \frac{1}{\mu_F} \right)^2 \mathbb{E}[\|\nabla J(\theta^k)\|^2] \\
& \leq \left(1 + \frac{1}{\mu_F} \right)^2 \left(2\mathbb{E}[\|\nabla J^H(\theta^k)\|^2] + 2G^2 R^2 \left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right)^2 \gamma^{2H} \right).
\end{aligned} \tag{28}$$

Combining equation 27 and equation 28 with equation 26 gives

$$\begin{aligned}
& \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}_\star^k\|] \\
& \leq \left(4 \frac{\sigma^2}{N} \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) + 4G^2 R^2 \left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right)^2 \gamma^{2H} \left(1 + \left(1 + \frac{1}{\mu_F} \right)^2 \right) \right. \\
& \quad \left. + 4 \left(1 + \left(1 + \frac{1}{\mu_F} \right)^2 + \frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2] \right)^{0.5}
\end{aligned} \tag{29}$$

And recall from equation 24 that

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2] \leq \frac{\frac{J^{H,\star} - J^H(\theta_0)}{K} + (\frac{\eta}{2} + L_J \eta^2) \frac{\sigma^2}{N}}{\frac{\eta}{2} - L_J \eta^2}.$$

Let us take $\eta = \frac{1}{4L_J}$. Then we get

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\nabla J^H(\theta^k)\|^2] \leq \frac{16L_J(J^{H,\star} - J^H(\theta_0))}{K} + \frac{3\sigma^2}{N}$$

In addition, let H , N , and K satisfy

$$\begin{aligned}
& \frac{1}{3} \left(\frac{\varepsilon}{3G} \right)^2 \geq 4G^2 R^2 \left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right)^2 \gamma^{2H} \left(1 + \left(1 + \frac{1}{\mu_F} \right)^2 \right)^2 \\
& N \geq \frac{4\sigma^2 \left(3 + \frac{4}{\lambda^2} + \frac{4G^2}{\lambda^2 + \lambda \mu_F} + 3 \left(1 + \frac{1}{\mu_F} \right)^2 \right)}{\frac{1}{3} \left(\frac{\varepsilon}{3G} \right)^2}, \\
& K \geq \frac{64L_J(J^{H,\star} - J^H(\theta_0)) \left(\left(1 + \frac{1}{\mu_F} \right)^2 + 1 + \frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right)}{\frac{1}{3} \left(\frac{\varepsilon}{3G} \right)^2}.
\end{aligned} \tag{30}$$

Then, we have

$$\frac{G}{K} \sum_{k=0}^{K-1} \mathbb{E}[\|\mathbf{w}^k - \mathbf{w}_\star^k\|] \leq \frac{\varepsilon}{3}. \tag{31}$$

- Bounding $\frac{1}{K} \sum_{k=0}^{K-1} \|\mathbf{w}^k\|^2$.

We have from equation 27 and equation 24 that

$$\begin{aligned}
& \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \|\mathbf{w}^k\|^2 \\
&= \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \frac{1}{\lambda} \mathbf{g}^k - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} \mathbf{g}^k \right\|^2 \right] \\
&= \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} \left[\left\| \frac{1}{\lambda} (\mathbf{g}^k - \nabla J^H(\theta^k)) - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} (\mathbf{g}^k - \nabla J^H(\theta^k)) + \frac{1}{\lambda} \nabla J^H(\theta^k) \right. \right. \\
&\quad \left. \left. - \frac{\hat{\mathbf{F}}^k}{\lambda^2 + \lambda \text{Tr}(\hat{\mathbf{F}}^k)} \nabla J^H(\theta^k) \right\|^2 \right] \\
&\leq 2 \frac{\sigma^2}{N} \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) + 2 \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) \frac{1}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\nabla J^H(\theta^k)\|^2] \\
&\leq \frac{168}{N} \frac{\sigma^2}{\lambda^2} \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) + \frac{3216 L_J J^{H*} - J^H(\theta_0)}{K} \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right).
\end{aligned}$$

Where to get to the last inequality we have taken $\eta = \frac{1}{4L_J}$ and applied Lemma 3.

$$\begin{aligned}
N &\geq \frac{168 \eta \sigma^2 \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right)}{\frac{\varepsilon}{6}}, \\
K &\geq \frac{3216 L_J \eta (J^{H*} - J^H(\theta_0)) \left(\frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right)}{\frac{\varepsilon}{6}},
\end{aligned} \tag{32}$$

we arrive at

$$\frac{\eta}{K} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathbf{w}^k\|^2] \leq \frac{\varepsilon}{3}. \tag{33}$$

- Bounding $\frac{1}{K} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} [\text{KL}(\pi^*(\cdot|s) || \pi_{\theta^0}(\cdot|s))]$.

By taking

$$K \geq \frac{3 \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} [\text{KL}(\pi^*(\cdot|s) || \pi_{\theta^0}(\cdot|s))]}{\eta \varepsilon} \tag{34}$$

we have

$$\frac{1}{\eta K} \mathbb{E}_{s \sim d_{\rho}^{\pi^*}} [\text{KL}(\pi^*(\cdot|s) || \pi_{\theta^0}(\cdot|s))] \leq \frac{\varepsilon}{3}, \tag{35}$$

From the Equation 30

$$\left(\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \right)^2 \gamma^{2H} \leq C_0 \varepsilon^2, \tag{36}$$

for some constant $C_0 > 0$ that does not depend on H or ε .

For all $H \geq 0$

$$\frac{H+1}{1-\gamma} + \frac{\gamma}{(1-\gamma)^2} \leq \frac{H+1}{1-\gamma} + \frac{1}{(1-\gamma)^2} \leq \frac{C_1(H+1)}{(1-\gamma)^2},$$

for a constant $C_1 \geq 1$ that depends only on γ . Substituting this into Equation 36 gives

$$(H+1)^2 \gamma^{2H} \leq C_2 (1-\gamma)^4 \varepsilon^2, \quad (37)$$

for constant $C_2 > 0$ independent of H and ε .

Since the prefactor grows polynomially in H while γ^{2H} decays geometrically,

$$\gamma^{2H} \leq c_1 \varepsilon^2,$$

for some constant $c_1 > 0$ independent of H . Taking logarithm on both side yields

$$H = \mathcal{O}(\log((1-\gamma)^{-1} \varepsilon^{-1})). \quad (38)$$

The lower bounds on N in equation 30 and equation 32 have the form

$$N \geq c_2 \sigma^2 \left(1 + \left(1 + \frac{1}{\mu_F}\right)^2 + \frac{1}{\lambda^2} + \frac{G^2}{\lambda^2 + \lambda \mu_F} \right) \varepsilon^{-2},$$

for a constant c_2 that absorbs numerical factors. Since the expression in parentheses is independent of ε , we have

$$N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right). \quad (39)$$

The lower bounds on K collected from equation 30, equation 32, and equation 34 can be written as

$$K \geq c_3 \left(L_J(J^{H,\star} - J^H(\theta_0)) + \mathbb{E}_{s \sim d_p^\pi} [\text{KL}(\pi^\star(\cdot|s) || \pi_{\theta^0}(\cdot|s))] \right) \varepsilon^{-2},$$

for a constant c_3 depending only on $(1-\gamma)^{-1}$ and the model parameters. Since both terms in parentheses scale as $\mathcal{O}((1-\gamma)^{-2})$. We have

$$K = \mathcal{O}\left(\frac{1}{(1-\gamma)^2 \varepsilon^2}\right). \quad (40)$$

In summary, we require N , K and H to satisfy equation 30, equation 32, and equation 34, which leads to

$$N = \mathcal{O}\left(\frac{\sigma^2}{\varepsilon^2}\right), \quad K = \mathcal{O}\left(\frac{1}{(1-\gamma)^2 \varepsilon^2}\right), \quad H = \mathcal{O}(\log((1-\gamma)^{-1} \varepsilon^{-1})).$$

By combining equation 31, equation 33, equation 35 and equation 25, we can conclude that

$$J(\pi^\star) - \frac{1}{K} \sum_{k=0}^{K-1} J(\theta^k) \leq \frac{\sqrt{\varepsilon_{\text{bias}}}}{1-\gamma} + \varepsilon.$$

In total, stochastic Sherman-Morrison Policy gradient requires to sample $KN = \mathcal{O}\left(\frac{\sigma^2}{(1-\gamma)^2 \varepsilon^4}\right)$ trajectories.

A.5 Detailed Derivation of the Toy Problem

The policy is parameterized by a vector $\boldsymbol{\theta} = [\theta_0, \theta_1]^\top$, with independent Bernoulli distribution:

$$\pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0) = \frac{1}{1 + e^{\lambda_p \theta_0}}, \quad (41)$$

$$\pi_{\boldsymbol{\theta}}(a = 0 \mid x = 1) = \frac{1}{1 + e^{\theta_1}}, \quad (42)$$

where λ_p scales the first policy parameter and is used to test scaling invariance. If we assume that the agent starts at the state $x = 0$ then the objective (expected return) in this simplified setting is

$$J(\boldsymbol{\theta}) = \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0) + 2(1 - \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0)) \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 1). \quad (43)$$

The policy gradient theorem states:

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{x,a} [\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|x) Q(x, a)]. \quad (44)$$

For the Bernoulli parameterization and noting that the policy function resembles sigmoid activation function the derivative of the log-policy can be simplified as

$$\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a|x) = \begin{cases} -\lambda_p(1 - \pi_0) & \text{if } x = 0, a = 0 \\ \lambda_p \pi_0 & \text{if } x = 0, a = 1 \\ -(1 - \pi_1) & \text{if } x = 1, a = 0 \\ \pi_1 & \text{if } x = 1, a = 1 \end{cases} \quad (45)$$

where $\pi_0 = \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0)$ and $\pi_1 = \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 1)$. Using the rewards in the environment, one can derive the analytic policy gradient vector

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \begin{bmatrix} -\lambda_p \pi_0 (1 - \pi_0) (1 - 2\pi_1) \\ -2(1 - \pi_0) \pi_1 (1 - \pi_1) \end{bmatrix}. \quad (46)$$

The geometry of the policy manifold is captured by the FIM and as the states are independent of each other we can write the per state FIM as,

$$F_{ii} = \mathbb{E}_{a \sim \pi(\cdot|x=i)} [\nabla_{\theta_i} \log \pi_{\boldsymbol{\theta}}(a|x=i) \nabla_{\theta_i} \log \pi_{\boldsymbol{\theta}}(a|x=i)^\top] \quad (47)$$

So For State ($x = 0$): The Score functions for $a \in \{0, 1\}$ can be expressed as:

$$\frac{\partial}{\partial \theta_0} \log \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 0) = -\lambda_p(1 - \pi_0), \quad \frac{\partial}{\partial \theta_0} \log \pi_{\boldsymbol{\theta}}(a = 1 \mid x = 0) = \lambda_p \pi_0 \quad (48)$$

Hence,

$$F_{00} = \pi_0 [-\lambda_p(1 - \pi_0)]^2 + (1 - \pi_0) [\lambda_p \pi_0]^2 \quad (49)$$

$$= \lambda_p^2 [\pi_0(1 - \pi_0)^2 + (1 - \pi_0)\pi_0^2] \quad (50)$$

$$= \lambda_p^2 \pi_0(1 - \pi_0) [(1 - \pi_0) + \pi_0] \quad (51)$$

$$= \lambda_p^2 \pi_0(1 - \pi_0). \quad (52)$$

Similarly for state ($x = 1$): The Score functions for $a \in \{0, 1\}$ can be expressed as:

$$\frac{\partial}{\partial \theta_1} \log \pi_{\boldsymbol{\theta}}(a = 0 \mid x = 1) = -(1 - \pi_1), \quad \frac{\partial}{\partial \theta_1} \log \pi_{\boldsymbol{\theta}}(a = 1 \mid x = 1) = \pi_1 \quad (53)$$

Hence,

$$F_{11} = \pi_1[-(1 - \pi_1)]^2 + (1 - \pi_1)[\pi_1]^2 \quad (54)$$

$$= \pi_1(1 - \pi_1)^2 + (1 - \pi_1)\pi_1^2 \quad (55)$$

$$= \pi_1(1 - \pi_1)[(1 - \pi_1) + \pi_1] \quad (56)$$

$$= \pi_1(1 - \pi_1) \quad (57)$$

Owing to independence between states, $F(\theta)$ is diagonal:

$$F(\theta) = \begin{bmatrix} \lambda_p^2 \pi_0(1 - \pi_0) & 0 \\ 0 & \pi_1(1 - \pi_1) \end{bmatrix}. \quad (58)$$

Then the NPG update direction is obtained by preconditioning the gradient with the inverse Fisher matrix,

$$\dot{\theta}_{\text{NPG}} = F(\theta)^{-1} \nabla_{\theta} J(\theta) = \begin{bmatrix} \frac{-(1 - 2\pi_1)}{\lambda_p} \\ -2(1 - \pi_0) \end{bmatrix}. \quad (59)$$

To reduce computation while preserving approximate invariance, the SM approach uses a low-rank inverse metric \tilde{F}^{-1} defined in terms of damping factor λ and exact FIM $F(\theta)$ capturing dominant curvature direction,

$$\tilde{F}^{-1} = \frac{1}{\lambda} \left(\mathbf{I} - \frac{F(\theta)}{\lambda + \text{Tr}(F(\theta))} \right) \quad (60)$$

$$(61)$$

Hence the resulting SM update direction can be written as:

$$\dot{\theta}_{\text{SM}} = \tilde{F}^{-1} \nabla_{\theta} J(\theta) \quad (62)$$

$$= \frac{1}{\lambda} \left(\mathbf{I} - \frac{F(\theta)}{\lambda + \text{Tr}(F(\theta))} \right) \nabla_{\theta} J(\theta) \quad (63)$$

$$= \frac{1}{\lambda} \left(\nabla_{\theta} J(\theta) - \frac{F(\theta) \nabla_{\theta} J(\theta)}{\lambda + \text{Tr}(F(\theta))} \right) \quad (64)$$

$$= \begin{bmatrix} \frac{1}{\lambda} (-\lambda_p \pi_0(1 - \pi_0)(1 - 2\pi_1)) - \frac{\lambda_p^3 \pi_0^2(1 - \pi_0)^2(1 - 2\pi_1)}{\lambda^2 + \lambda [\lambda_p^2 \pi_0(1 - \pi_0) + \pi_1(1 - \pi_1)]} \\ \frac{1}{\lambda} (-2(1 - \pi_0)\pi_1(1 - \pi_1)) - \frac{2\pi_1^2(1 - \pi_1)^2(1 - \pi_0)}{\lambda^2 + \lambda [\lambda_p^2 \pi_0(1 - \pi_0) + \pi_1(1 - \pi_1)]} \end{bmatrix} \quad (65)$$

A.6 Additional Results

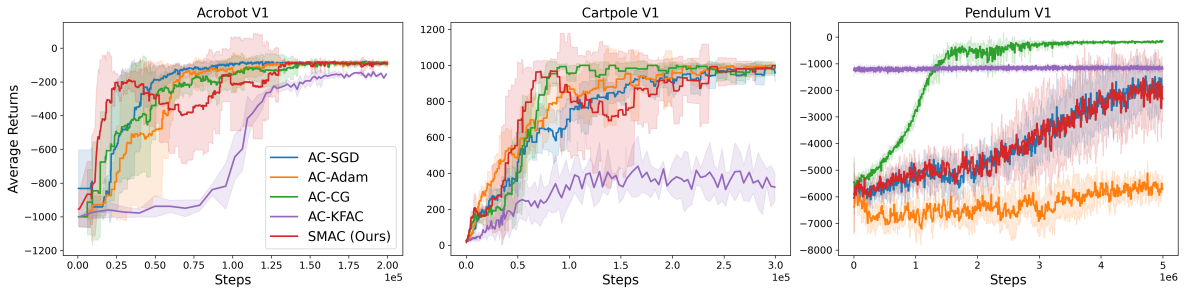


Figure 8: Raw results on three classic control tasks. These correspond to the results reported in Fig. 6, but without binning or smoothing.

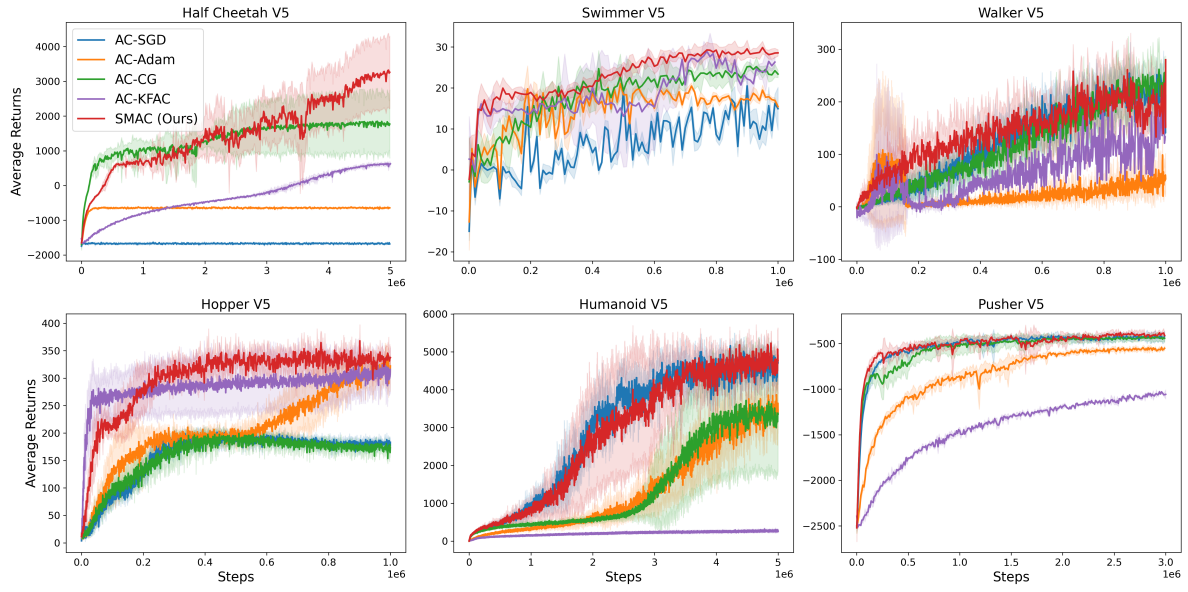


Figure 9: Results on six MuJoCo tasks. These correspond to the results reported in Fig. 7, but without binning or smoothing.

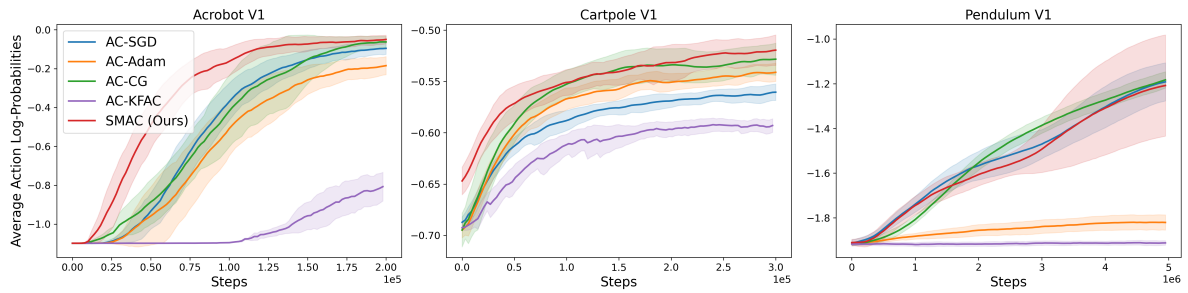


Figure 10: Average log-probabilities of the actions taken on three classic control tasks. All results are averaged over 5 random seeds, with shaded areas showing standard deviation.

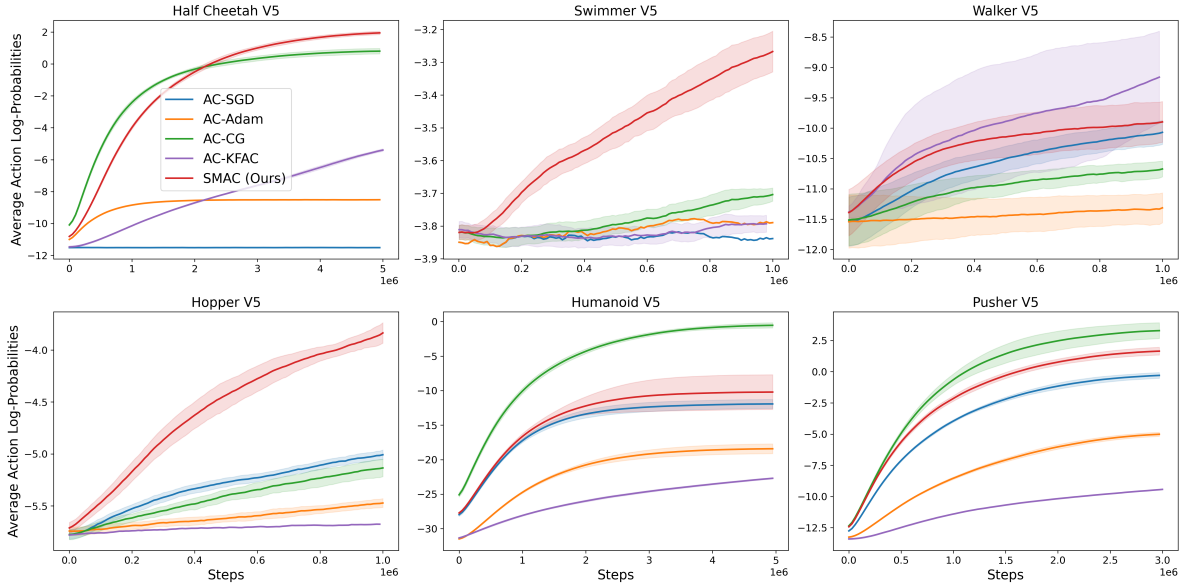


Figure 11: Average log-probabilities of the actions taken on six MuJoCo tasks. All results are averaged over 5 random seeds, with shaded areas showing standard deviation.

A.7 Hyper-parameters

We consider the following hyper-parameters:

- η (actor): step size for the policy-network update.
- α (critic): learning rate for the value network, optimized with Adam.
- T : number of environment steps collected before each parameter update.
- γ : reward discount factor.
- λ_{GAE} : trace-decay parameter used by Generalized Advantage Estimation when computing the advantage values A_t .
- λ : damping factor used in the matrix-free Sherman-Morrison approximation of the empirical Fisher.

We tune five of these on each environment. We omit tuning T , as preliminary runs have shown that performance is not highly correlated with its value. We select values for α (critic) from the interval $[1 \times 10^{-5}, 5 \times 10^{-3}]$, for γ from the set $\{0.9, 0.95, 0.99\}$, and for λ_{GAE} from the set $\{0.8, 0.85, 0.9, 0.95, 0.99\}$. For the rest of the hyper-parameters, we select values according to Tables 2 and 5. The goal is to find hyper-parameter configurations that maximize the return averaged over the last 100 epochs, for each environment and algorithm.

We find that learning rates have the highest impact on performance. On the contrary, $\gamma = 0.99$ and $\lambda_{\text{GAE}} = 0.9$ perform well in all environments. For each environment and algorithm, we select the hyper-parameter configuration that performs best, and then use it to generate the results in Sec. 6. These configurations are shown in Tables 3 and 4 (and Table 5 for AC-KFAC). For hyper-parameters specific to AC-KFAC, we found a damping of 1×10^{-3} , factor decay of 0.99, KL clipping value of 1×10^{-3} , and an update frequency of 1 iteration to perform well in all environments except the ones specified in Table 5.

| Algorithm | η (actor) | λ |
|-----------|--|---|
| SMAC | $[5 \times 10^{-5}, 5 \times 10^{-1}]$ | $\{10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1\}$ |
| AC-Adam | $[5 \times 10^{-6}, 1 \times 10^{-2}]$ | - |
| AC-SGD | $[5 \times 10^{-5}, 5 \times 10^{-1}]$ | - |
| AC-CG | $[5 \times 10^{-3}, 9 \times 10^{-1}]$ | - |
| AC-KFAC | $[1 \times 10^{-5}, 1 \times 10^{-3}]$ | - |

Table 2: Hyper-parameter values evaluated during tuning. Additionally, for each algorithm, we select values for α (critic) from $[1 \times 10^{-5}, 5 \times 10^{-3}]$, for γ from $\{0.9, 0.95, 0.99\}$, and for λ_{GAE} from $\{0.8, 0.85, 0.9, 0.95, 0.99\}$. See Table 5 for additional details on tuning hyperparameters specific to AC-KFAC.

| Environment | Algorithm | η (actor) | α (critic) | T | γ | λ_{GAE} | λ |
|-------------|-----------|--------------------|--------------------|------|----------|------------------------|-----------|
| Acrobot | SMAC | 5×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | 0.1 |
| | AC-Adam | 6×10^{-4} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 2×10^{-1} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 6×10^{-1} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| Cartpole | SMAC | 5×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | 0.1 |
| | AC-Adam | 7×10^{-5} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 7×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 8×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| Pendulum | SMAC | 6×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | 0.1 |
| | AC-Adam | 7×10^{-4} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 5×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 3×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | | | | | | |

Table 3: Training hyper-parameters for the classic control environments.

| Environment | Algorithm | η (actor) | α (critic) | T | γ | λ_{GAE} | λ |
|--------------|-----------|----------------------|--------------------|------|----------|------------------------|-----------|
| Half Cheetah | SMAC | 8×10^{-4} | 6×10^{-4} | 1000 | 0.99 | 0.9 | 0.01 |
| | AC-Adam | 5×10^{-3} | 6×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-1} | 9×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 9×10^{-1} | 8×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| Hopper | SMAC | 1×10^{-2} | 1×10^{-4} | 1000 | 0.99 | 0.9 | 1.0 |
| | AC-Adam | 1×10^{-4} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-3} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 1×10^{-2} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | | | | | | |
| Swimmer | SMAC | 3×10^{-2} | 1×10^{-4} | 1000 | 0.99 | 0.9 | 1.0 |
| | AC-Adam | 1×10^{-4} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-3} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 1×10^{-2} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-2} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| Walker | SMAC | 3×10^{-2} | 1×10^{-5} | 1000 | 0.99 | 0.9 | 1.0 |
| | AC-Adam | 1×10^{-5} | 1×10^{-5} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-4} | 1×10^{-5} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 1×10^{-2} | 1×10^{-5} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-3} | 1×10^{-4} | 1000 | 0.99 | 0.9 | - |
| Humanoid | SMAC | 1×10^{-4} | 1×10^{-3} | 1000 | 0.99 | 0.9 | 0.01 |
| | AC-Adam | 3×10^{-4} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 9×10^{-1} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | | | | | | |
| Pusher | SMAC | 1.5×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | 0.01 |
| | AC-Adam | 2.5×10^{-3} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-SGD | 1×10^{-1} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-CG | 9×10^{-1} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |
| | AC-KFAC | 1×10^{-2} | 1×10^{-3} | 1000 | 0.99 | 0.9 | - |

Table 4: Training hyper-parameters for the MuJoCo environments.

| Hyper-parameter | Tuning values | Environment | Final value |
|------------------------|--|-------------|----------------------|
| Damping | $[5 \times 10^{-5}, 5 \times 10^{-2}]$ | Pendulum | 1×10^{-2} |
| | | Humanoid | 1.5×10^{-2} |
| Decay factor | $\{0.85, 0.9, 0.95, 0.99, 1.0\}$ | Pendulum | 0.85 |
| | | Humanoid | 0.9 |
| KL clip | $[5 \times 10^{-5}, 3 \times 10^{-2}]$ | Pendulum | 2.5×10^{-2} |
| | | Humanoid | 4×10^{-3} |
| Update frequency | $\{1, 2, 4, 8, 16, 32, 64\}$ | Pendulum | 16 |
| | | Humanoid | 64 |

Table 5: Hyper-parameter values evaluated when tuning AC-KFAC on each environment. We include the values we select for the Pendulum and Humanoid environments. All other environments use 1×10^{-3} damping, 0.99 factor decay, 1×10^{-3} KL clip, and 1 update frequency. These tuning values and final values complement Tables 2, 3, and 4.