Pink Noise LQR: How does Colored Noise affect the Optimal Policy in RL?

Jakob Hollenstein¹ Marko Zarić³ Samuele Tosatto^{1,2} Justus Piater^{1,2}

{jakob.hollenstein, samuele.tosatto, justus.piater}@uibk.ac.at

¹Department of Computer Science, University of Innsbruck, Austria

²Digital Science Center, University of Innsbruck, Austria

³Max Planck Institute for Intelligent Systems, Tübingen, Germany

Abstract

Colored noise, a class of temporally correlated noise processes, has shown promising results for improving exploration in deep reinforcement learning for both off-policy and on-policy algorithms. However, it is unclear how temporally correlated colored noise affects policy learning apart from changing exploration properties. In this paper, we investigate the implications of colored noise on on-policy deep reinforcement learning in a simplified setting, considering linear dynamics and a linear policy under quadratic costs. We derive a closed-form solution for the expected cost, revealing that colored noise affects both the expected cost and the optimal policy. Notably, the cost splits into a state-cost part equal to the unperturbed system's cost and a noise-cost term, affecting the policy, but independent of the initial state. While the cost changes depending on the noise, the expected trajectory remains independent of the noise color for a given linear policy. Far from the goal state, the state cost dominates, and the effect due to the noise is negligible: the policy approaches the optimal policy of the unperturbed system. Near the goal state, the noise cost dominates, changing the optimal policy.

1 Introduction

Deep reinforcement learning is an approximate dynamic programming technique to derive a policy (a controller) for a given environment, i.e., reward (=-cost) and dynamics. The policy is estimated based on trajectory samples gathered from the environment. To do so, the data collection, and thus the action selection, needs to be varied. This is typically done by randomly perturbing the action selection process, i.e., by action noise. Action noise can be applied additively to the deterministically selected action of a policy or by sampling from a stochastic policy. In continuous control settings, such as robotics, the system dynamics include integrative components: the action signal (e.g., force, torque, velocity), is integrated (velocity, position). This explains why temporally correlated action noise has been found to improve learning performance in reinforcement learning (Rückstieß et al., 2008; Raffin et al., 2021; Eberhard et al., 2023; Hollenstein et al., 2022; Chiappa et al., 2023; Hollenstein et al., 2024). In particular, the temporally correlated colored noise processes have shown promising results for continuous control (Eberhard et al., 2023; Pinneri et al., 2020; Hollenstein et al., 2024).

While empirically, these noise processes have shown improvements in learning performance, it is less clear how this noise affects the optimal policy. In this paper, we investigate this question in a simplified setting.

While continuous control reinforcement learning is able to deal with stochastic dynamics, i.e.,

$$s_{t+1} \sim p(\cdot|a_t, s_t)$$

in practice, often environments with deterministic dynamics are used, e.g., based on the MuJoCo simulator (Tassa et al., 2018; Todorov et al., 2012; Brockman et al., 2016).

To study the impact of colored noise on RL, we employ a simplified model, considering linear deterministic dynamics:

$$s_{t+1} = Gs_t + Ha_t$$

and a linear policy:

$$a_t = -Ks_t$$

Central to our investigation, we assume that the actions of the policy are perturbed by action noise drawn from \mathcal{C}_{β} , a colored noise process with noise color β : $\varepsilon_t \sim \mathcal{C}_{\beta}$:

$$a_t = -Ks_t + \varepsilon_t \tag{1}$$

Additionally, we assume the cost (= -reward) to be quadratic and the goal state to be s = 0. That is, we study the question of the impact of colored noise in the linear quadratic regulator (LQR) setting.

As is typical for reinforcement learning, we assume the initial state $s_0 \sim S$ to be sampled from a given initial state distribution. Furthermore, we limit the study to the episodic setting, limiting the length of the trajectories to T.

Since we are interested in understanding the effects of colored noise on the optimal policy in reinforcement learning, we are interested in the following questions:

- 1. How does colored noise affect the expected trajectory $\mathbb{E}[s_t]$, given an expected starting state \overline{s}_0 ?
- 2. How does the expected cost change when the noise color (β) is changed?
- 3. How does colored noise change the optimal policy?

Our contributions are:

- 1. We show that while the collected trajectories are affected by the noise, the *expected* trajectory remains unchanged regardless of the noise (Q1).
 - Context: On-policy RL results demonstrate improved performance with temporally correlated noise but do not directly address the noise's impact on data collection. This result applies to linear dynamics and extends to more complex environments if they are locally linear in the presence of noise.
- 2. We derive a closed-form solution for the expected cost.
 - Context: Sample-based estimation of cost in D-RL typically involves high variance, obscuring the effect of noise color, requiring a more reliable assessment of the impact of the noise.
- 3. We show how the cost is affected by the noise (Q2).
 - Context: While noise is primarily intended for exploration, and typically the unperturbed policy is evaluated, in on-policy RL, exploration and stochastic policy classes are more directly connected. Our result highlights how the change in data collection leads to a change in the cost, i.e., the optimization target. By dividing the cost into state-dependent terms, noise-dependent but policy-independent, and noise-policy dependent costs, we show how correlated noise can result in a different optimal policy.
- 4. We show how these different parts of the cost can affect the optimal policy K (Q3). Context: We empirically show that the impact of these noise terms varies, being generally small but more significant when the policy generates close to zero actions, e.g., when it aims to remain close to a goal.

1.1 Related Work

Exploration is critical for reinforcement learning. In continuous control deep reinforcement learning, the two most prominent noise types used for exploration are uncorrelated white noise (Haarnoja et al., 2019; Fujimoto et al., 2018; Abdolmaleki et al., 2018; Schulman et al., 2017) or temporally correlated Ornstein-Uhlenbeck (Uhlenbeck & Ornstein, 1930) noise, e.g., Lillicrap et al. (2016). The importance of temporally correlated noise has also been shown by methods that combine random aspects with deterministic state-to-action mappings (Raffin et al., 2021; Rückstieß et al., 2008; Chiappa et al., 2023). A further type of random exploration, more similar to white noise and Ornstein-Uhlenbeck noise is action noise exploration based on colored noise processes (Pinneri et al., 2020; Eberhard et al., 2023; Hollenstein et al., 2024).

Research in the LQR setting has investigated noise based exploration strategies, establishing regret bounds for learning unknown dynamics. Simchowitz & Foster (2020) and Simchowitz et al. (2020) show that "naive" exploration, e.g., white-noise injection, is rate-optimal for online LQR when aiming to identify the system and recover the noise-free optimal controller. These studies focus on minimizing regret with process or observation noise, often assuming i.i.d. Gaussian noise, rather than on how temporally correlated noise affects the optimal policy for known dynamics. Classical LQG control also treats noisy dynamics and observations—and even temporally correlated noise (Kwong, 1987; Escobedo-Trujillo & Garrido-Meléndez, 2021)—but typically does not model additive action noise directly in the control law.

In this work, we fill that gap by studying

$$a_t = -Ks_t + \varepsilon_t,$$

where ε_t is drawn from a colored noise process generated in the frequency domain—an approach particularly relevant to deep RL implementations (Eberhard et al., 2023; Hollenstein et al., 2024; Pinneri et al., 2020).

2 Background

Colored noise processes are a class of temporally correlated noise processes that are parameterized by the noise color β . This class includes temporally uncorrelated noise (white noise, $\beta=0$) and temporally correlated red noise ($\beta=2$), which is exhibited by, e.g., Brownian motion. Sequences of different noise colors are illustrated in the time domain in Figure 1. The noise color β describes how the expected power spectral density (PSD) behaves, i.e., the power components scale with $\frac{1}{f\beta}$. This is illustrated in Figure 2.

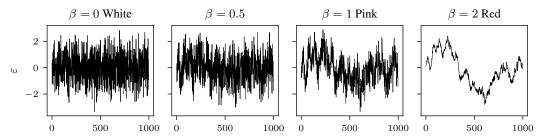


Figure 1: Colored noise processes generate temporally correlated noise with varying degrees of temporal correlation depending on the noise color β . From left to right: the temporal correlation increases with increasing β , from $\beta = 0$ (temporally uncorrelated white noise) to $\beta = 2$ (highly temporally correlated red noise).

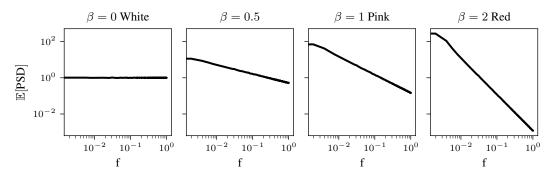


Figure 2: Colored noise is defined by the slope of the expected power spectral density: the power components scale with $\frac{1}{f\beta}$

Generating colored noise Following the work by Eberhard et al. (2023); Hollenstein et al. (2024) on colored noise in deep reinforcement learning, we generate colored noise in the frequency domain and apply the Inverse Real (Fast) Fourier (IRFFT) to retrieve a sequence of perturbation ε_t , i.e., noise generation follows the algorithm proposed by Timmer & König (1995). This means that for a specific rollout the frequency components Φ are sampled once $\Phi \sim \mathcal{C}_{\beta}$, and remain fixed for the entire episode. ε_t can then be computed by the inverse Fourier transform, which amounts to a weighted summation of the components of Φ . This summation can be expressed as the inner product between Φ^{\top} and the time dependent Fourier coefficients f_t :

$$\varepsilon_t = \Phi^\top \cdot f_t \tag{2}$$

The perturbation ε_t can be interpreted as a discrete-time signal at time index $t \in \{0, \dots, M-1\}$ 1}, derived using the inverse real Fourier transform from $N = \lfloor M/2 \rfloor + 1$ frequency components. For simplicity in the colored noise generation and inverse Fourier transform, we assume both M and N to be even valued. The derivations in both cases are analogous. For details on the noise generation process see Appendix A.1.

$$\varepsilon_t = \frac{1}{N-1} \sum_{k=0}^{N-1} \left[c_k \varphi_{2k} \cos(-k \frac{t}{M} \cdot 2\pi) + c_k \varphi_{2k+1} \sin(-k \frac{t}{M} \cdot 2\pi) \right]$$
(3)

where $\varphi_{2k}, \varphi_{2k+1}$ denote the real, respectively imaginary part of the frequency domain Fourier coefficients and c_k denotes a scaling factor.

$$c_k = \begin{cases} 0 & \text{if } k \in \{0, N-1\} \\ 1 & \end{cases}$$

The sum can be rewritten as the dot product

$$\varepsilon_t = \Phi^{\top} f_t$$

where Φ and f_t are defined as:

$$\Phi = \begin{bmatrix} \varphi_1 \\ \vdots \\ \varphi_{2N} \end{bmatrix} \quad f_t = \begin{bmatrix} c_0 \cos(-0 \cdot 2\pi \frac{i}{M}) \\ c_0 \sin(-0 \cdot 2\pi \frac{i}{M}) \\ \vdots \\ c_{N-1} \cos(-(N-1) \cdot 2\pi \frac{i}{M}) \\ c_{N-1} \sin(-(N-1) \cdot 2\pi \frac{i}{M}) \end{bmatrix}$$
(4)

The components φ_i of Φ are independently sampled, depending on the noise color β , and the sequence length M:

1: **procedure** $C(M, \beta)$

 $N \leftarrow \frac{M}{2} + 1$ $f \leftarrow \{\frac{1}{M}, \frac{1}{M}, \dots, \frac{i}{M}, \dots, \frac{N-1}{M}\}$ $\sigma^2 \leftarrow \{\dots, f_i^{-\beta}, \dots\}$

▷ Calculate scales

5:
$$\frac{1}{c^2} \leftarrow \left(\frac{2}{M}\right)^2 \cdot \sum w_i^2 | w_i \in \{\sigma_1, \dots, \sigma_{N-2}, \frac{\sigma_{N-1}}{2}\}$$

$$\mathcal{N}(0, c \cdot \sigma_0 \cdot \sqrt{2})$$

$$\mathcal{N}(0, c \cdot \sigma_0 \cdot 0)$$

$$\mathcal{N}(0, c \cdot \sigma_1)$$

$$\mathcal{N}(0, c \cdot \sigma_1)$$

$$\vdots$$

$$\mathcal{N}(0, c \cdot \sigma_{N-1})$$

$$\mathcal{N}(0, c \cdot \sigma_{N-1})$$

7: return Φ

8: end procedure

3 Q1: Expected Trajectory remains unchanged

The addition of action noise changes the distribution of sampled trajectories. The distribution also changes when the noise color is varied. This is illustrated in Figure 3. Empirically this figure also shows, that despite the widely different distribution of trajectories, the expected trajectory remains unchanged. In this section we derive the expected trajectory.

For a given Φ , the action noise is fixed for the whole duration of a trajectory. The policy can be included in the dynamical system to make it autonomous, i.e., the system evolution only depends on the initial state. This means that the trajectory, or more precisely, the state s_t , can be expressed as a sum capturing the recursion starting at the given state s_0 (Appendix B):

$$s_1 = Gs_0 + Ha_0 \tag{5}$$

$$s_1 = Gs_0 - HKs_0 + H\Phi^{\top} f_0 \tag{6}$$

$$s_2 = Gs_1 + Ha_1 \tag{7}$$

$$\vdots (8)$$

$$s_t = (G - HK)^t s_0 + \sum_{i=1}^t (G - HK)^{t-i} H \Phi^\top f_{i-1}$$
(9)

Assuming that the initial state is randomly chosen, $s_0 \sim S_0$ and the expected value exists $\mathbb{E}[s_0] = \overline{s_0}$, the expected value for $\mathbb{E}[s_t]$ can be computed:

$$\mathbb{E}[s_t] = \mathbb{E}\left[(G - HK)^t s_0 + \sum_{i=1}^t (G - HK)^{t-i} H \Phi^\top f_{i-1} \right]$$

Because $\mathbb{E}[\Phi] = \mathbf{0}$, by definition of the chosen colored noise generation process, using the linearity property of the expectation, this simplifies to

$$\mathbb{E}\left[s_{t}\right] = \mathbb{E}\left[\left(G - HK\right)^{t} s_{0}\right] = \left(G - HK\right)^{t} \mathbb{E}\left[s_{0}\right]$$

That is: the *expected trajectory* under colored noise remains unchanged.

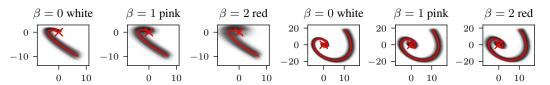


Figure 3: Effect of different noise colors on sampled trajectories. (first three plots) Double Integrator environment, (last three plots) Randomly generated test environment. While the distribution of trajectories (black) changes depending on the noise color, the expected trajectory (red) remains unchanged, reaching the same goal state (marked X) reliably.

4 Q2: Effect of Colored Noise on Cost

In the previous section, we demonstrated that the noise color alters the distribution of sampled trajectories, but the expected trajectory stays the same. This raises the question of how the noise color affects the cost, particularly the expected cost.

We assume quadratic costs:

$$J = \sum_{t=0}^{T} s_t^{T} Q s_t + a_t^{T} R a_t = J_Q + J_R$$
 (10)

Because of the presence of action noise, e.g., Equation (1), we are interested in the expected cost $\overline{J} = \mathbb{E}[J]$, which we derive (see Appendix C for details) as follows:

$$\mathbb{E}\left[J\right] = \mathbb{E}\left[J_Q\right] + \mathbb{E}\left[J_R\right] = \mathbb{E}\left[\sum_{t=0}^{\top} s_t^{\top} Q s_t\right] + \mathbb{E}\left[\sum_{t=0}^{\top} a_t^{\top} R a_t\right] = \tag{11}$$

$$\sum_{t=0}^{T} \left(\overline{s}_0^{\mathsf{T}} S_t \overline{s}_0 + \operatorname{tr} \left(S_t \operatorname{Cov}[s_0] \right) \right)$$
 (12)

$$+f_t^{\top} \mathbb{E} \left[\Phi R \Phi^{\top} \right] f_t + \tag{13}$$

$$+ \sum_{i=1}^{t} \sum_{j=1}^{t} f_{j-1}^{\top} \mathbb{E} \left[\Phi(W_{i,j,t} + B_{i,j,t}) \Phi^{\top} \right] f_{i-1} + \sum_{i=1}^{t} f_{i-1}^{\top} \mathbb{E} \left[\Phi Y_{t,i} \Phi^{\top} \right] f_{t}$$
(14)

where
$$C^t = (G - HK)^t$$

$$B_{i,j,t} := H^\top C^{t-j^\top} K^\top RKC^{t-i} H$$

$$W_{i,j,t} = H^\top C^{t-j^\top} QC^{t-i} H$$

$$Y_{t,i} := H^\top (C^{t-i})^\top K^\top RKC^t$$

$$S_t := C^{t^\top} QC^t + (C^t)^\top K^\top RKC^t$$

Note that the cost splits into a state dependent term equal to the cost of the noise-free system (12) and three noise terms, one independent of the policy K (13), and two dependent on the policy K (14). Interestingly, all three of the noise terms are independent of the expected initial state \bar{s}_0 . These additional terms show analytically that the noise influences the cost. This is illustrated empirically in Figure 4.

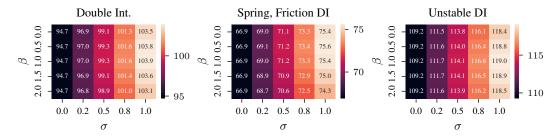


Figure 4: Effect of different noise colors β and noise scales $\sigma \cdot \Phi$ on the expected cost. Generally, the cost increases with a larger noise scale σ . Whether the cost is higher for a specific noise color β depends on the dynamics. The experiments are performed with a horizon T=30.

5 Q3: Optimal Policy K under Colored Noise

In the previous section we showed that the noise color affects the expected cost and that the expected cost splits into a state dependent cost term and a noise dependent cost term.

The split of the expected cost \overline{J} Equation (11), into the state dependent cost

$$J_{s_0} = \sum_{t=0}^{\top} \left(\overline{s}_0^{\top} S_t \overline{s}_0 + \operatorname{tr} \left(S_t \operatorname{Cov}[s_0] \right) \right)$$

and noise dependent cost (keeping only terms dependent on K)

$$J_{\varepsilon} = \sum_{t=0}^{\top} \left(\sum_{i=1}^{t} \sum_{j=1}^{t} f_{j-1}^{\top} \mathbb{E} \left[\Phi(W_{i,j,t} + B_{i,j,t}) \Phi^{\top} \right] f_{i-1} + \sum_{i=1}^{t} f_{i-1}^{\top} \mathbb{E} \left[\Phi Y_{t,i} \Phi^{\top} \right] f_{t} \right)$$

shows a dependency of the optimal policy, i.e., a change in optimal K in the presence of noise.

The effect of σ and β on K is illustrated empirically in Figure 5. Here, the Double Integrator is studied with initial state $s_0^{\top} = [0.5 \quad 0.5]$, horizon T = 32, the optimal policy K, $(K \in \mathbb{R}^2)$ is found numerically from the closed form solution of the expected cost. With larger scale σ , or change in β the components of the gain matrix policy K change.

However, if the noise is scaled down, $\sigma \cdot \Phi$ for $\sigma \ll 1$, J_{s_0} dominates the combined cost and the optimal policy K approaches the optimal policy of the unperturbed system. On the flip side, when J_{s_0} has little influence, J_{ε} dominates the combined cost, causing a shift in the optimal policy K to counteract the noise effect. The influence of J_{s_0} is small when the system is close to the goal state ($s_{\text{goal}} = \mathbf{0}$). This indicates that the cost is likely to be dominated by the state cost J_{s_0} at the beginning of the trajectory, shifting to J_{ε} towards the end, i.e., close to the goal.

This has several interesting implications:

- J_{ε} is independent of s_0 and will thus not converge to zero. For an infinite horizon $\lim_{T\to\infty} J$ might diverge and average or discounted cost formulations need to be investigated.
- Hollenstein et al. (2022) suggests reducing the influence of the noise over the course of the training process, i.e., scheduling σ in $\sigma \cdot \Phi$ to shrink over the training process. This would reduce the influence of J_{ε} and recover the optimal policy of the unperturbed system.
- Close to the goal state $s=\mathbf{0}$, the unperturbed policy would not take any action $a=-Ks=\mathbf{0}$. In this case, J_{ε} would dominate over J_{s_0} . However, in practical applications, the policy will either have to take actions, suggesting $s\neq\mathbf{0}$, or, if the system is required to stay close to the goal state, the action noise scale needs to be small to prevent the system from deviating too far from $\mathbf{0}$. Both of these factors would lead to $J_{s_0}\gg J_{\varepsilon}$, suggesting that the optimal policy K approaches the solution of the unperturbed system.

This shift from the state dependent cost J_{s_0} to the noise dependent cost J_{ε} is illustrated in Figure 6 for the Double Integrator, $s_0 = \begin{bmatrix} 10 \\ 10 \end{bmatrix}$, for K the infinite horizon LQR solution is

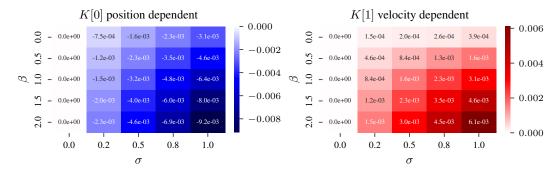


Figure 5: The optimal policy K, is affected by the noise color β and the noise scale σ . We numerically derive the optimal K for the Double Integrator for $\overline{s}_0^\top = \begin{bmatrix} 0 & 0 \end{bmatrix}$ and horizon T = 30. The plots show the difference to the optimal K for the action-noise free setting for both dimensions of K (state, velocity) separately. The color gradients show that the growing discrepancy to the noise-free policy is driven by the increase in noise scale and the change in noise color (i.e., temporal correlation of the noise).

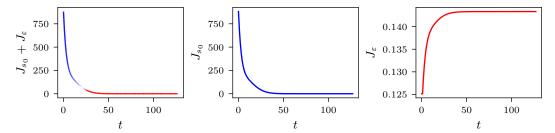


Figure 6: Visualization of the contributions to the cost for each timestep t: The closed-form solution for the expected cost consists of a state dependent term $J_{s_0}(t)$ (\blacksquare) and a noise dependent term $J_{\varepsilon}(t)$ (\blacksquare) At the beginning of the trajectory, the cost is dominated by $J_{s_0}(t)$, close to the goal state, the cost is dominated by the cost incurred by the noise $J_{\varepsilon}(t)$. The noise cost $J_{\varepsilon}(t)$ appears to reach an equilibrium. Overall J_{s_0} , which is independent of the action noise, dominates.

used, and T = 120. In this example, the total cost of the trajectory is determined mostly by the state dependent cost accounting for 99.7% of the total cost. Here, the influence of the noise on the policy would be marginal.

6 Conclusion

In this paper, we investigated the effect of colored action noise on the optimal policy in a simplified LQR setting. We found that the expected trajectory for a given policy remains unchanged in the presence of colored noise but that the expected cost changes. Associated with this change in cost is a change in the optimal policy. The change in cost is due to an additional cost term compared to the cost of the unperturbed system, which is independent of the starting state and instead depends on the noise color, system dynamics, and policy matrix. This effect is relevant close to the goal state, but has little impact further away from the goal. This suggests that while colored noise can change the optimal policy, this change is likely to be small in practice.

Acknowledgments

This research was partially funded by the Austrian Science Fund (FWF): I 5755-N (ELSA), and by the Autonomous Province of Bolzano-Bozen - South Tyrol under Funding Agreement 10/2024, Abstractron.

References

Abdolmaleki, A., Springenberg, J. T., Tassa, Y., Munos, R., Heess, N., and Riedmiller, M. A. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*, 2018. https://openreview.net/forum?id=S1ANxQWOb.

Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. OpenAI Gym. arXiv:1606.01540 [cs], June 2016. http://arxiv.org/abs/1606.01540.

Chiappa, A. S., Vargas, A. M., Huang, A. Z., and Mathis, A. Latent Exploration for Reinforcement Learning. [object Object], 2023. doi: 10.48550/ARXIV.2305.20065. https://arxiv.org/abs/2305.20065.

Eberhard, O., Hollenstein, J., Pinneri, C., and Martius, G. Pink Noise Is All You Need: Colored Noise Exploration in Deep Reinforcement Learning. In *International Conference on Learning Representations*, February 2023. https://openreview.net/forum?id=hQ9V5QN27eS.

Escobedo-Trujillo, B. and Garrido-Meléndez, J. Stochastic LQR optimal control with white and colored noise: Dynamic programming technique. Revista Mexicana de Ingeniería

- $\label{eq:quimica} Quimica,~20(2):1111-1127,~2021.~~ http://rmiq.org/iqfvp/Numbers/V20/No2/Sim2353.~~pdf.$
- Fujimoto, S., van Hoof, H., and Meger, D. Addressing Function Approximation Error in Actor-Critic Methods. In *International Conference on Machine Learning*, October 2018. http://arxiv.org/abs/1802.09477.
- Haarnoja, T., Zhou, A., Hartikainen, K., Tucker, G., Ha, S., Tan, J., Kumar, V., Zhu, H., Gupta, A., Abbeel, P., and Levine, S. Soft Actor-Critic Algorithms and Applications. arXiv:1812.05905 [cs, stat], January 2019. http://arxiv.org/abs/1812.05905.
- Hollenstein, J., Auddy, S., Saveriano, M., Renaudo, E., and Piater, J. Action Noise in Off-Policy Deep Reinforcement Learning: Impact on Exploration and Performance. *Transactions on Machine Learning Research*, November 2022. ISSN 2835-8856. https://openreview.net/forum?id=NljBlZ6hmG&referrer=%5BAuthor% 20Console%5D(%2Fgroup%3Fid%3DTMLR%2FAuthors%23your-submissions).
- Hollenstein, J., Martius, G., and Piater, J. Colored noise in PPO: Improved exploration and performance through correlated action sampling. In *Conference of the Association for the Advancement of Artificial Intelligence*, February 2024.
- Kwong, R. On the LQG problem with correlated noise and its relation to minimum variance control. In 26th IEEE Conference on Decision and Control, pp. 763-767, Los Angeles, California, USA, December 1987. IEEE. doi: 10.1109/CDC.1987.272492. http://ieeexplore.ieee.org/document/4049369/.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016. http://arxiv.org/abs/1509.02971.
- Pinneri, C., Sawant, S., Blaes, S., Achterhold, J., Stueckler, J., Rolínek, M., and Martius, G. Sample-efficient cross-entropy method for real-time planning. In Kober, J., Ramos, F., and Tomlin, C. J. (eds.), 4th Conference on Robot Learning, CoRL 2020, 16-18 November 2020, Virtual Event / Cambridge, MA, USA, volume 155 of Proceedings of Machine Learning Research, pp. 1049-1065. PMLR, 2020. https://proceedings.mlr.press/v155/pinneri21a.html.
- Raffin, A., Kober, J., and Stulp, F. Smooth exploration for robotic reinforcement learning. In Faust, A., Hsu, D., and Neumann, G. (eds.), *Conference on Robot Learning*, 2021. https://proceedings.mlr.press/v164/raffin22a.html.
- Rückstieß, T., Felder, M., and Schmidhuber, J. State-Dependent Exploration for Policy Gradient Methods. In Daelemans, W., Goethals, B., and Morik, K. (eds.), *Machine Learning and Knowledge Discovery in Databases*, volume 5212, pp. 234–249. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. ISBN 978-3-540-87480-5 978-3-540-87481-2. doi: 10.1007/978-3-540-87481-2_16. http://link.springer.com/10.1007/978-3-540-87481-2_16.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal Policy Optimization Algorithms. *CoRR*, abs/1707.06347, 2017.
- Simchowitz, M. and Foster, D. Naive Exploration is Optimal for Online LQR. In *Proceedings of the 37th International Conference on Machine Learning*, pp. 8937–8948. PMLR, November 2020. https://proceedings.mlr.press/v119/simchowitz20a.html.
- Simchowitz, M., Singh, K., and Hazan, E. Improper Learning for Non-Stochastic Control. In *Proceedings of Thirty Third Conference on Learning Theory*, pp. 3320-3436. PMLR, July 2020. https://proceedings.mlr.press/v125/simchowitz20a.html.
- Tassa, Y., Doron, Y., Muldal, A., Erez, T., Li, Y., Casas, D. d. L., Budden, D., Abdolmaleki, A., Merel, J., Lefrancq, A., Lillicrap, T., and Riedmiller, M. DeepMind Control Suite. arXiv:1801.00690 [cs], January 2018. http://arxiv.org/abs/1801.00690.
- Timmer, J. and König, M. On generating power law noise. *Astron. Astrophys*, 300:707–710, 1995.

- Todorov, E., Erez, T., and Tassa, Y. MuJoCo: A physics engine for model-based control. In 2012 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 5026–5033, October 2012. doi: 10.1109/IROS.2012.6386109.
- Uhlenbeck, G. E. and Ornstein, L. S. On the theory of the Brownian motion. *Physical review*, 36(5):823, 1930.

A Colored Noise

A.1 Generating Colored Noise

```
1: procedure GenerateColoredNoise(N, \beta)
                      L \leftarrow \lfloor N/2 \rfloor
                   f \leftarrow \left\{\frac{1}{N}, \frac{1}{N}, \dots, \frac{i}{N}, \dots, \frac{L}{N}\right\}
s \leftarrow \left\{\dots, f_i^{-\beta/2}, \dots\right\}
w_L \leftarrow \begin{cases} s_L, & \text{if } L \text{ is odd} \\ s_L/2, & \text{otherwise} \end{cases}
                                                                                                                                                                            \triangleright Frequencies of components 0 \dots L
                                                                                                                                                                                                                                   ▷ Calculate scales
                     \sigma \leftarrow \frac{2}{N} \cdot \sqrt{\sum w_i^2}
\alpha = \{\dots, \alpha_i, \dots\} : \alpha_i \sim \mathcal{N}(0, s_i)
\beta = \{\dots, \beta_i, \dots\} : \beta_i \sim \mathcal{N}(0, s_i)
   7:
                                                                                                                                                                                                                                                       ▶ Real part
                                                                                                                                                                                                                                     ▶ Imaginary part
                      \alpha_0 \sim \mathcal{N}(0, s_0 \cdot \sqrt{2})
10:
11:
                    \alpha_L \sim \begin{cases} \mathcal{N}(0, s_0 \cdot \sqrt{2}), & \text{if } odd \\ \mathcal{N}(0, s_0), & \text{otherwise} \end{cases}
12:
                   \beta_L \sim \mathcal{N}(0, s_0) \cdot \begin{cases} 0, & \text{if } odd \\ 1, & \text{otherwise} \end{cases}
\gamma \leftarrow \{\dots, \gamma_i, \dots\} : \gamma_i = \alpha_i + \mathbf{i}\beta_i
\tau_{\varepsilon} = \mathcal{F}^{-1}[\gamma] \cdot 1/\sigma
13:
14:
15:
16:
                                                                                                                                                                        \triangleright Return noise sequence of length N
                      return \tau_{\varepsilon}
17: end procedure
```

B Derivation of s_t

$$a_{t} = -Ks_{t} + \Phi^{\top} f_{t}$$

$$s_{1} = Gs_{0} + Ha_{0}$$

$$s_{1} = Gs_{0} - HKs_{0} + H\Phi^{\top} f_{0}$$

$$s_{1} = (G - HK)s_{0} + H\Phi^{\top} f_{0}$$

$$s_{2} = Gs_{1} + Ha_{1}$$

$$s_{2} = Gs_{1} - HKs_{1} + H\Phi^{\top} f_{1}$$

$$s_{2} = (G - HK)s_{1} + H\Phi^{\top} f_{1}$$

$$s_{2} = (G - HK)(G - HK)s_{0} + H\Phi^{\top} f_{0}) + H\Phi^{\top} f_{1}$$

$$s_{2} = (G - HK)^{2}s_{0} + (G - HK)(H\Phi^{\top} f_{0}) + H\Phi^{\top} f_{1}$$

$$s_{3} = Gs_{2} + Ha_{2}$$

$$s_{3} = (G - HK)(G - HK)s_{2} + H\Phi^{\top} f_{2}$$

$$s_{3} = (G - HK)((G - HK)^{2}s_{0} + (G - HK)(H\Phi^{\top} f_{0}) + H\Phi^{\top} f_{1}) + H\Phi^{\top} f_{2}$$

$$s_{3} = (G - HK)^{3}s_{0} + (G - HK)^{2}(H\Phi^{\top} f_{0}) + (G - HK)H\Phi^{\top} f_{1} + H\Phi^{\top} f_{2}$$

$$\vdots$$

$$s_{t} = (G - HK)^{t}s_{0} + \sum_{i=1}^{t} (G - HK)^{t-i}\Phi^{\top} f_{i-1}$$

C Derivation of closed-form expected cost $\mathbb{E}[J]$

From Equations (1) and (9) and let C = (G - HK). For a given trajectory length T, the noise sample ε_t is generated by the Fourier series at time t. This amounts to a weighted sum of the frequency component random variables: $\varepsilon_t = \Phi^{\top} f_t$.

$$\mathbb{E}\left[J\right] = \mathbb{E}\left[J_Q\right] + \mathbb{E}\left[J_R\right] = \mathbb{E}\left[\sum_{t=0}^T s_t^\top Q s_t\right] + \mathbb{E}\left[\sum_{t=0}^T a_t^\top R a_t\right] = \sum_{t=0}^T \mathbb{E}\left[s_t^\top Q s_t\right] + \sum_{t=0}^T \mathbb{E}\left[a_t^\top R a_t\right] = \sum_{t=0}^T \mathbb{E}\left[s_t^\top Q s_t\right] + \mathbb{E}\left[a_t^\top R a_t\right]$$

We derive $\mathbb{E}\left[s_t^{\top}Qs_t\right]$ and $\mathbb{E}\left[a_t^{\top}Ra_t\right]$ in separate subsections and combine the results afterwards into the final closed-form expected cost solution.

C.1 Derivation of $\mathbb{E}\left[s_t^{\top}Qs_t\right]$

$$\mathbb{E}\left[s_{t}^{\top}Qs_{t}\right] = \mathbb{E}\left[\left(C^{t}s_{0} + \sum_{i=1}^{t}C^{t-i}H\Phi^{\top}f_{i-1}\right)^{\top}Q\left(C^{t}s_{0} + \sum_{i=1}^{t}C^{t-i}H\Phi^{\top}f_{i-1}\right)\right] =$$

$$= \mathbb{E}\left[s_{0}^{\top}C^{t}^{\top}QC^{t}s_{0}\right] + \mathbb{E}\left[2s_{0}^{\top}C^{t}^{\top}Q\sum_{i=1}^{t}C^{t-i}H\Phi^{\top}f_{i-1}\right] +$$

$$+ \mathbb{E}\left[\sum_{j=1}^{t}f_{j-1}^{\top}\Phi H^{\top}C^{t-j}^{\top}Q\sum_{i=1}^{t}C^{t-i}H\Phi^{\top}f_{i-1}\right] =$$

$$= \mathbb{E}\left[s_{0}^{\top}\right]C^{t}^{\top}QC^{t}\mathbb{E}\left[s_{0}\right] + \text{tr}(C^{t}^{\top}QC^{t}\text{Cov}[s_{0}]) + 2\mathbb{E}\left[s_{0}^{\top}\right]C^{t}^{\top}Q\sum_{i=1}^{t}C^{t-i}\mathbb{E}\left[\Phi^{\top}\right]f_{i-1} +$$

$$+ \sum_{j=1}^{t}\sum_{i=1}^{t}f_{j-1}^{\top}H^{\top}\mathbb{E}\left[\Phi C^{t-j}^{\top}QC^{t-i}H\Phi^{\top}\right]f_{i-1} =$$

$$= \overline{s_{0}}^{\top}C^{t}^{\top}QC^{t}\overline{s_{0}} + \text{tr}(C^{t}^{\top}QC^{t}\text{Cov}[s_{0}]) + \sum_{j=1}^{t}\sum_{i=1}^{t}f_{j-1}^{\top}\mathbb{E}\left[\Phi H^{\top}C^{t-j}^{\top}QC^{t-i}H\Phi^{\top}\right]f_{i-1} =$$

$$= \overline{s_{0}}^{\top}C^{t}^{\top}QC^{t}\overline{s_{0}} + \text{tr}(C^{t}^{\top}QC^{t}\text{Cov}[s_{0}]) + \sum_{j=1}^{t}\sum_{i=1}^{t}f_{j-1}^{\top}\mathbb{E}\left[\Phi W_{i,j,t}\Phi^{\top}\right]f_{i-1}$$

C.2 Derivation of $\mathbb{E}\left[a_t^{\top} R a_t\right]$

Calculating the expected action cost for time step a_t results in three separate action cost terms:

$$\begin{split} \mathbb{E}\left[a_t^{\top}Ra_t\right] = & \mathbb{E}\left[(-Ks_t + \Phi^{\top}f_t)^{\top}R(-Ks_t + \Phi^{\top}f_t)\right] = \\ = & \mathbb{E}\left[f_t^{\top}\Phi R\Phi^{\top}f_t\right] + \mathbb{E}\left[s_t^{\top}K^{\top}RKs_t\right] - 2\mathbb{E}\left[s_t^{\top}K^{\top}R\Phi^{\top}f_t\right] \end{split}$$

First action cost term:

$$\mathbb{E}\left[f_t^\top \Phi R \Phi^\top f_t\right] = f_t^\top \mathbb{E}\left[\Phi R \Phi^\top\right] f_t$$

Second action cost term:

$$\mathbb{E}\left[s_{t}^{\top}K^{\top}RKs_{t}\right] = \mathbb{E}\left[s_{0}^{\top}(C^{t})^{\top}K^{\top}RKC^{t}s_{0}\right] + \operatorname{tr}\left((C^{t})^{\top}K^{\top}RKC^{t}\operatorname{Cov}[s_{0}]\right) + \\ + \mathbb{E}\left[2s_{0}^{\top}C^{t}^{\top}K^{\top}R\sum_{i=1}^{t}KC^{t-i}H\Phi^{\top}f_{i-1}\right] + \\ + \mathbb{E}\left[\sum_{i=1}^{t}\sum_{j=1}^{t}\left(C^{j-i}H\Phi^{\top}f_{j-1}\right)^{\top}K^{\top}RKC^{t-i}H\Phi^{\top}f_{i-1}\right] = \\ = \mathbb{E}\left[s_{0}^{\top}\right]\left(C^{t}\right)^{\top}K^{\top}RKC^{t}\mathbb{E}\left[s_{0}\right] + \operatorname{tr}\left((C^{t})^{\top}K^{\top}RKC^{t}\operatorname{Cov}[s_{0}]\right) + \\ + \mathbb{E}\left[\sum_{i=1}^{t}\sum_{j=1}^{t}f_{j-1}^{\top}\Phi H^{\top}C^{j-i}^{\top}K^{\top}RKC^{t-i}H\Phi^{\top}f_{i-1}\right] = \\ = \overline{s}_{0}^{\top}(C^{t})^{\top}K^{\top}RKC^{t}\overline{s}_{0} + \operatorname{tr}\left((C^{t})^{\top}K^{\top}RKC^{t}\operatorname{Cov}[s_{0}]\right) + \\ + \sum_{i=1}^{t}\sum_{j=1}^{t}f_{j-1}^{\top}\mathbb{E}\left[\Phi H^{\top}C^{j-i}^{\top}K^{\top}RKC^{t-i}H\Phi^{\top}\right]f_{i-1} = \\ [\operatorname{substitute}:B_{i,j,t}:=H^{\top}C^{j-i}^{\top}K^{\top}RKC^{t-i}H] \\ = \overline{s}_{0}^{\top}(C^{t})^{\top}K^{\top}RKC^{t}\overline{s}_{0} + \operatorname{tr}\left((C^{t})^{\top}K^{\top}RKC^{t}\operatorname{Cov}[s_{0}]\right) + \sum_{i=1}^{t}\sum_{j=1}^{t}f_{j-1}^{\top}\mathbb{E}\left[\Phi B_{i,j,t}\Phi^{\top}\right]f_{i-1}$$

Third action cost term:

$$\mathbb{E}\left[s_t^{\top} K^{\top} R \Phi^{\top} f_t\right] = \mathbb{E}\left[\left(C^t s_0 + \sum_{i=1}^t C^{t-i} H \Phi^{\top}\right)^{\top} K^{\top} R \Phi^{\top} f_t\right] =$$

$$= \mathbb{E}\left[s_0^{\top} (C^t)^{\top} K^{\top} R \Phi^{\top} f_t\right] + \mathbb{E}\left[\sum_{i=1}^t f_{i-1}^{\top} \Phi H^{\top} (C^{t-i})^{\top} K^{\top} R \Phi^{\top} f_t\right] =$$

$$\left[\mathbb{E}\left[s_0^{\top} (C^t)^{\top} K^{\top} R \Phi^{\top} f_t\right] = 0 \quad \& \quad \text{substitute } Y_{t,i} := H^{\top} (C^{t-i})^{\top} K^{\top} R\right]$$

$$= \sum_{i=1}^t f_{i-1}^{\top} \mathbb{E}\left[\Phi Y_{t,i} \Phi^{\top}\right] f_t$$

C.3 Closed-form expected cost $\mathbb{E}[J]$

When combining the derived parts, we group them based on

- (15) **Initial state dependency**: We use the linear property of the quadratic form and trace operator to merge the state cost and action cost parts with initial state dependency.
- (16) **Noise dependency**: The action cost term that only depends on the noise.
- (17) **Noise** + **Policy dependency**: State and action cost have a quadratic double sum term, which we combine to one (linearity of expectation). This term and the linear action cost term are both noise and policy dependent but have no dependents on the initial state.

$$\mathbb{E}\left[J\right] = \sum_{t=0}^{T} \mathbb{E}\left[s_{t}^{\top}Qs_{t}\right] + \sum_{t=0}^{T} \mathbb{E}\left[a_{t}^{\top}Ra_{t}\right] = \sum_{t=0}^{T} \left(S_{t}^{\top}QC^{t} + \left(C^{t}\right)^{\top}K^{\top}RKC^{t}\right) \cdot \left(C^{t}^{\top}QC^{t} + \left(C^{t}\right)^{\top}K^{\top}RKC^{t}\right) \cdot \left(S_{0}\right)\right)$$
(15)
$$+ f_{t}^{\top}\mathbb{E}\left[\Phi R\Phi^{\top}\right] f_{t} +$$
(16)
$$+ \sum_{i=1}^{t} \sum_{j=1}^{t} f_{j-1}^{\top}\mathbb{E}\left[\Phi(W_{i,j,t} + B_{i,j,t})\Phi^{\top}\right] f_{i-1} + \sum_{i=1}^{t} f_{i-1}^{\top}\mathbb{E}\left[\Phi Y_{t,i}\Phi^{\top}\right] f_{t}\right)$$
(17)
$$\text{where } C^{t} = (G - HK)^{t} \quad B_{i,j,t} := H^{\top}C^{t-j^{\top}}K^{\top}RKC^{t-i}H$$

$$W_{i,j,t} = H^{\top}C^{t-j^{\top}}QC^{t-i}H \quad Y_{t,i} := H^{\top}(C^{t-i})^{\top}K^{\top}R$$

To make this a closed-form solution we have to evaluate all expectations in the formula which all can be evaluated in the following way analogously:

$$\mathbb{E}\left[\Phi Z\Phi^{\top}\right] = \\
\mathbb{E}\left[\begin{bmatrix}\varphi_{1,1} & \cdots & \varphi_{M,1} \\ \vdots & \vdots & \vdots \\ \varphi_{1,N} & \cdots & \varphi_{M,N}\end{bmatrix} Z\begin{bmatrix}\varphi_{1,1} & \cdots & \varphi_{M,1} \\ \vdots & \vdots & \vdots \\ \varphi_{1,N} & \cdots & \varphi_{M,N}\end{bmatrix}\right] = \\
\mathbb{E}\left[\begin{bmatrix}\sum_{k,l} Z_{k,l} \varphi_{k,0} \varphi_{l,0} & \cdots & \sum_{k,l} Z_{k,l} \varphi_{k,0} \varphi_{l,N} \\ \vdots & \vdots & \vdots \\ \sum_{k,l} Z_{k,l} \varphi_{k,N} \varphi_{l,0} & \cdots & \sum_{k,l} Z_{k,l} \mathbb{E}\left[\varphi_{k,0} \varphi_{l,N}\right]\right] = \\
\begin{bmatrix}\sum_{k,l} Z_{k,l} \mathbb{E}\left[\varphi_{k,0} \varphi_{l,0}\right] & \cdots & \sum_{k,l} Z_{k,l} \mathbb{E}\left[\varphi_{k,0} \varphi_{l,N}\right] \\ \vdots & \vdots & \vdots \\ \sum_{k,l} Z_{k,l} \mathbb{E}\left[\varphi_{k,N} \varphi_{l,0}\right] & \cdots & \sum_{k,l} Z_{k,l} \mathbb{E}\left[\varphi_{k,N} \varphi_{l,N}\right]\right] = \\
\begin{bmatrix}\sum_{k} Z_{k,k} \mathbb{E}\left[\varphi_{k,0} \varphi_{k,0}\right] & \cdots & \sum_{k} Z_{k,k} \mathbb{E}\left[\varphi_{k,0} \varphi_{k,N}\right] \\ \vdots & \vdots & \vdots \\ \sum_{k} Z_{k,k} \mathbb{E}\left[\varphi_{k,N} \varphi_{k,0}\right] & \cdots & \sum_{k} Z_{k,k} \mathbb{E}\left[\varphi_{k,N} \varphi_{k,N}\right]\right] \\
&= \mathbb{E}\left[\varphi_{k,i} \varphi_{k,j}\right] = 0 \quad \text{for} \quad i \neq j \\
\begin{bmatrix}\sum_{k} Z_{k,k} \mathbb{Var}[\varphi_{k,0}] & \cdots & \sum_{k} Z_{k,k} \mathbb{E}\left[\varphi_{k,N} \varphi_{k,N}\right]\right] \\
&= \frac{\sum_{k} Z_{k,k} \operatorname{Var}[\varphi_{k,0}]}{\sum_{k} Z_{k,k} \operatorname{Var}[\varphi_{k,N}]} = \frac{\sum_{k} Z_{k,k} \operatorname{Var}[\varphi_{k,N}]}{\sum_{k} Z_{k,k} \operatorname{Var}[\varphi_{k,N}]} = \frac{\operatorname{diag}_{V \to M}\left(\operatorname{diag}_{M \to V}(Z) \cdot \operatorname{Var}[\Phi]\right)\right)}{\sum_{k} Z_{k,k} \operatorname{Var}[\Phi](18)}$$

By definition of the colored noise generation process $Var[\Phi]$ is known which results in the following closed form solution:

$$\mathbb{E}\left[J\right] = \sum_{t=0}^{T} \mathbb{E}\left[s_{t}^{\top}Qs_{t}\right] + \sum_{t=0}^{T} \mathbb{E}\left[a_{t}^{\top}Ra_{t}\right] =$$

$$\sum_{t=0}^{T} \left(\overline{s}_{0}^{\top}S_{t}\overline{s}_{0} + \operatorname{tr}\left(S_{t}\operatorname{Cov}[s_{0}]\right) + \right.$$

$$\left. + f_{t}^{\top}\operatorname{diag}_{V \to M}\left(\operatorname{diag}_{M \to V}(R) \cdot \operatorname{Var}[\Phi]\right)\right)f_{t} +$$

$$\left. + \sum_{i=1}^{t} \sum_{j=1}^{t} f_{j-1}^{\top}\operatorname{diag}_{V \to M}\left(\operatorname{diag}_{M \to V}(W_{i,j,t} + B_{i,j,t}) \cdot \operatorname{Var}[\Phi]\right)\right)f_{i-1} +$$

$$\left. + \sum_{i=1}^{t} f_{i-1}^{\top}\operatorname{diag}_{V \to M}\left(\operatorname{diag}_{M \to V}(Y_{t,i}) \cdot \operatorname{Var}[\Phi]\right)\right)f_{t} \right)$$
where $C^{t} := (G - HK)^{t}$

$$B_{i,j,t} := H^{\top}C^{t-j^{\top}}K^{\top}RKC^{t-i}H$$

$$W_{i,j,t} := H^{\top}C^{t-j^{\top}}QC^{t-i}H$$

$$Y_{t,i} := H^{\top}(C^{t-i})^{\top}K^{\top}R$$

$$S_{t} := C^{t^{\top}}QC^{t} + (C^{t})^{\top}K^{\top}RKC^{t}$$