

Evidence on the Regularisation Properties of Maximum-Entropy Reinforcement Learning

Rémy Hosseinkhan Boucher, Onofrio Semeraro, Lionel Mathelin,
LISN, CNRS, Université Paris-Saclay

What is the bias introduced by **entropy regularisation**? Are **complexity measures** linked to noise robustness?

Background

Partially Observable Markov Decision Process (POMDP)

$$X_{h+1} = F(X_h, U_h)$$

$$Y_h = G(X_h) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_Y^2 I_d)$$

F : state operator

G : observation operator

ϵ : Gaussian **noise**

Maximum-Entropy Objective

$$J^{\pi, \epsilon} = \mathbb{E}^{\pi, \epsilon} \left[\sum_{h=0}^H \gamma^h c(X_h, U_h) \right] - \alpha \mathbb{E}^{\pi, \epsilon} \left[\sum_{h=0}^H \gamma^h \mathcal{H}(\pi(\cdot | X_h)) \right]$$

α : entropy coefficient

\mathcal{H} : entropy

Notations

P_{ϵ}^{π} : trajectory probability with **observation noise**

$\mathbb{E}^{\pi, \epsilon}$: expectation under P_{ϵ}^{π}

$P^{\pi} = P_0^{\pi}$: no **noise** ($\epsilon \equiv 0$)

Excess Risk Metrics

$$\mathcal{R}^{\pi} = \mathbb{E}^{\pi, \epsilon} \left[\sum_{h=0}^H \gamma^h c(X_h, U_h) \right] - \mathbb{E}^{\pi} \left[\sum_{h=0}^H \gamma^h c(X_h, U_h) \right]$$

$$= J^{\pi, \epsilon} - J^{\pi}$$

$$\mathring{\mathcal{R}}^{\pi} = \frac{J^{\pi, \epsilon} - J^{\pi}}{J^{\pi}} = \frac{\mathcal{R}^{\pi}}{J^{\pi}}$$

Goal

Evaluating the Robustness $\mathring{\mathcal{R}}^{\pi}$ of Max-Entropy Policies under Observation Noise $\epsilon \sim \mathcal{N}(0, \sigma_Y^2 I_d)$

Noise-free training with PPO $\rightarrow P^{\pi} = P_0^{\pi}$: no **noise** ($\epsilon \equiv 0$)

5 coeff. α (x 10 seeds) \rightarrow 5 x 10 policies ($\pi_{\theta_{\alpha}^*}$)

Test $\pi_{\theta_{\alpha}^*}$ under different noise levels on Y

$\epsilon \nearrow \rightarrow J^{\pi^*, \epsilon} \nearrow$ (noise impacts perf)
 $\alpha > 0 \rightarrow \mathring{\mathcal{R}}^{\pi, \alpha} \searrow$ (robustness)

Find complexity measures $\mathcal{M}(\pi_{\theta})$ controlling the Excess Risk $\mathring{\mathcal{R}}^{\pi}$

Parameterised Policy π_{θ}

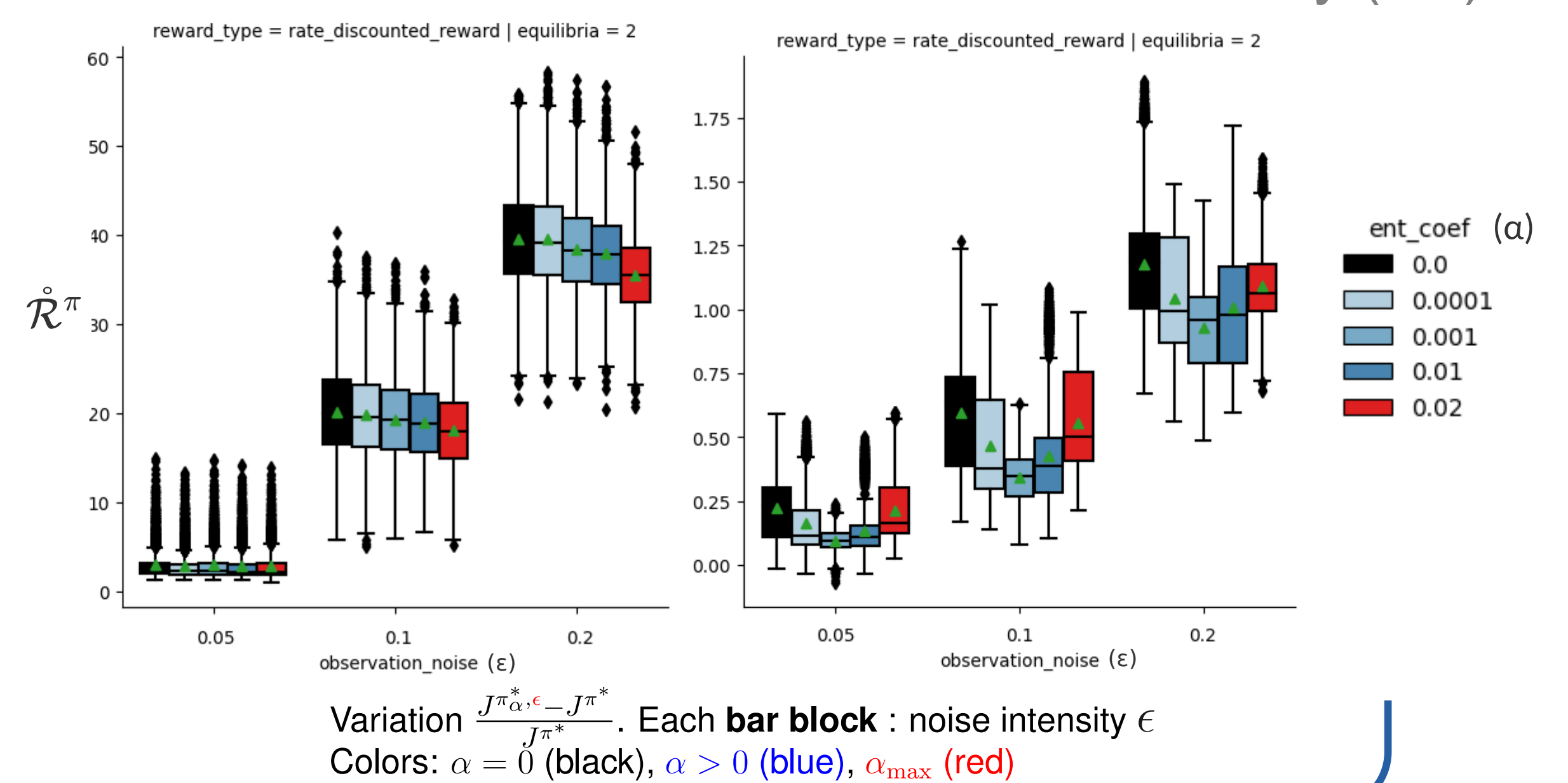
Parameter Space $\theta \in \Theta$

Complexity measures quantify model **complexity**

Regularisation \rightarrow Low complexity

Complexity measure: $\mathcal{M} : \Theta \rightarrow \mathbb{R}_+$

RL Environments: Lorenz - Kuramoto-Sivashinsky (KS)



Question

Which **complexity measures** indicate **noise robustness**? Why do **high entropy** policies learn **better final solutions**?

Contribution

Introduction of complexity measures from Statistical Learning

Norm-based Complexity Measures

$$\pi_{\theta}(\cdot | X_k) \sim \mathcal{N}(\mu_{\theta}(X_k), \theta_{\sigma_{\pi}} I)$$

$$\text{If } \mu_{\theta}(x) = (\sigma_l \circ \sigma_{l-1} \circ \dots \circ \sigma_1)(x), \text{ Lips}(\mu_{\theta}) \leq \prod_{i=1}^l \text{Lips}(\sigma_i) = \prod_{i=1}^l \|\theta_i\|$$

Regularity - Lipschitz Continuity \rightarrow Lipschitz bound (norm product)?

Sharpness-based Complexity Measures

Curvature - Hessian \rightarrow Fisher Information?

Objective landscape - Robustness \rightarrow Flat minima?

$$\bullet \mathcal{M}(\pi_{\theta}, \mathcal{D}) = \text{Tr}(\mathcal{I}(\theta_{\mu})) = \text{Tr}(-\mathbb{E}^{X \sim \rho^{\pi}, U \sim \pi_{\theta}(\cdot | X)} [\nabla_{\theta_{\mu}}^2 \log \pi_{\theta}(U | X)])$$

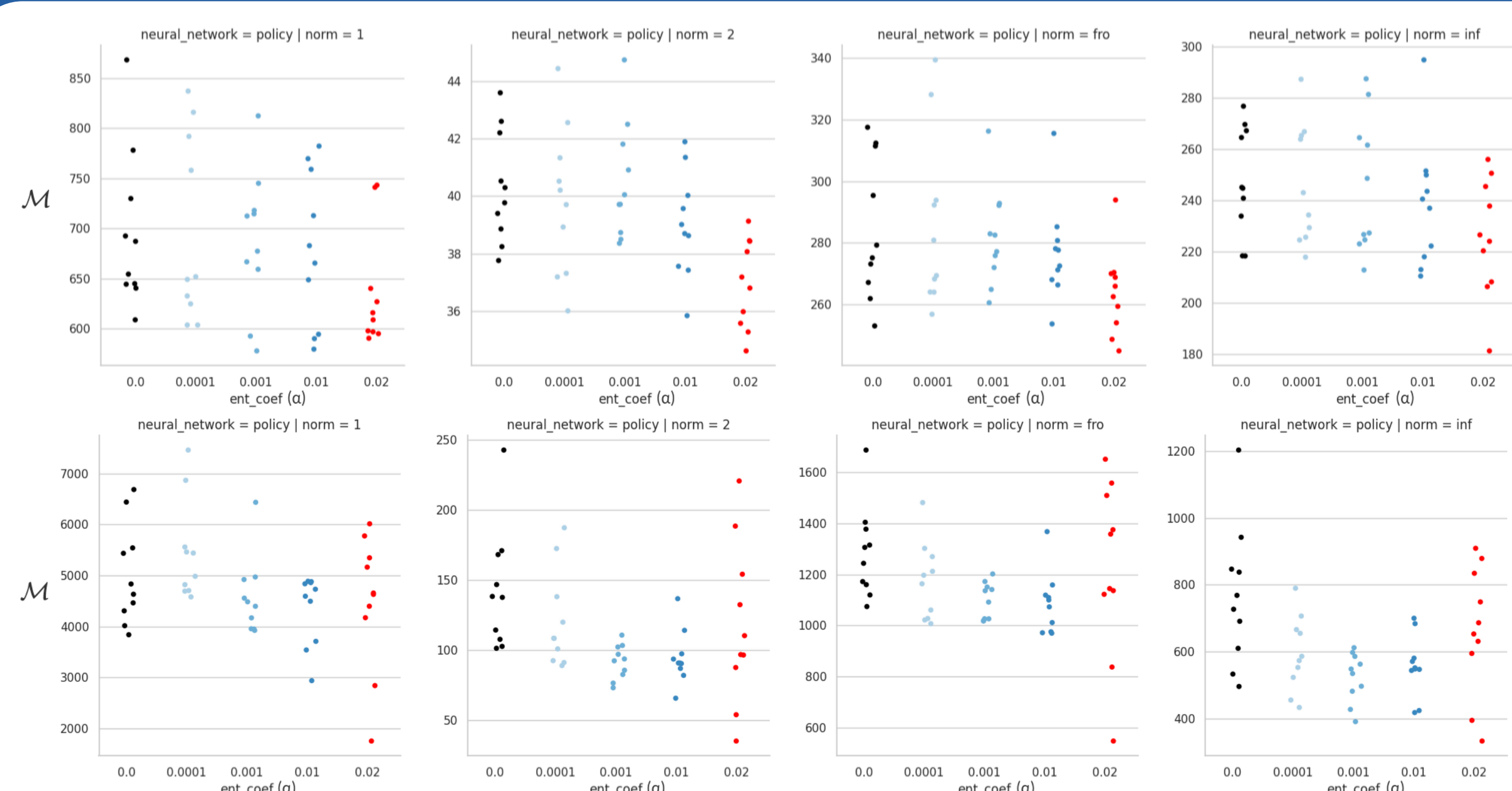
$$\bullet \mathcal{M}(\pi_{\theta}, \mathcal{D}) = \|\theta_{\mu}\|_p$$

$$\bullet \mathcal{M}(\pi_{\theta}, \mathcal{D}) = \prod_{i=1}^l \|\theta_{\mu}^i\|_p \text{ where } \theta_{\mu}^i \text{ is the } i^{\text{th}} \text{ layer of the network } \mu_{\theta_{\mu}}$$

$$\nabla_{\theta}^2 J^{\pi_{\theta}} = \mathbb{E}^{\pi_{\theta}} \left[\sum_{h,i,j=0}^H c(X_h, U_h) \left(\nabla_{\theta} \log \pi_{\theta}(U_i | X_i) \nabla_{\theta} \log \pi_{\theta}(U_j | X_j)^T + \nabla_{\theta}^2 \log \pi_{\theta}(U_i | X_i) \right) \right]$$

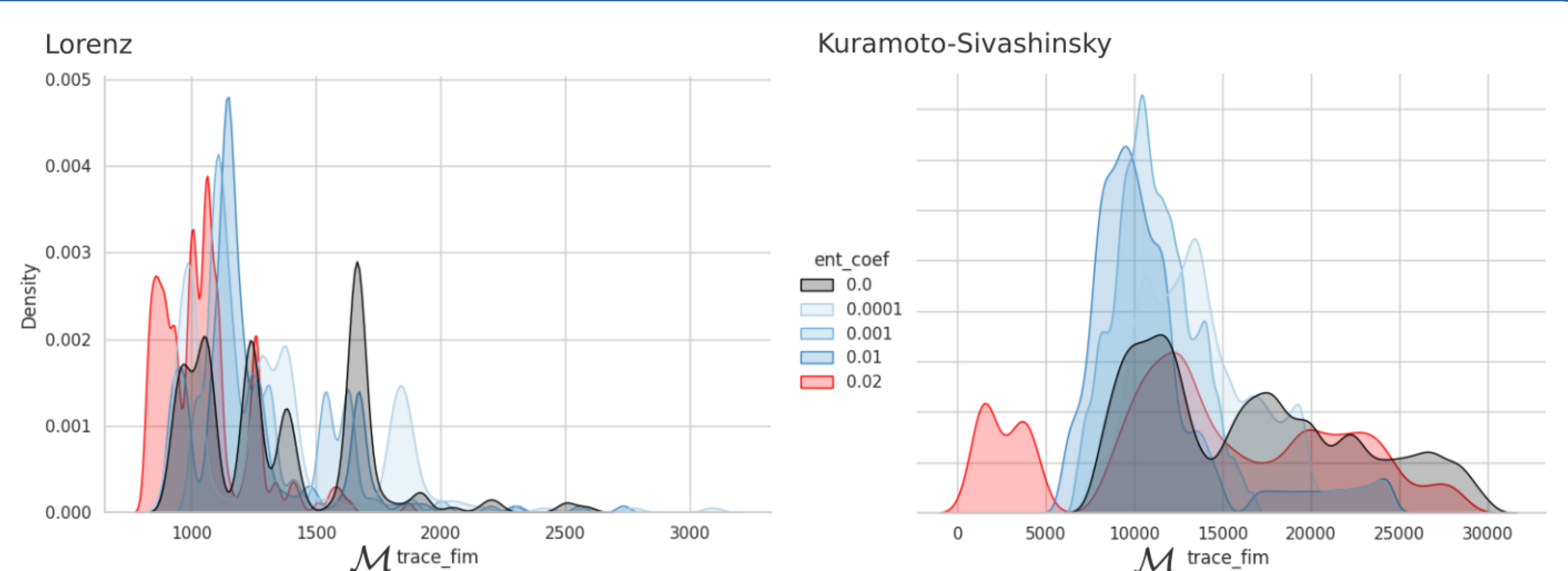
$$\mathcal{I}(\theta) = -\mathbb{E}^{X \sim \rho, U \sim \pi_{\theta}(\cdot | X)} [\nabla_{\theta}^2 \log \pi_{\theta}(U | X)]$$

Results



$\mathcal{M}(\pi_{\theta}^{\alpha}) = \prod_{i=1}^l \|\theta_{\mu}^i\|$
Colors: $\alpha = 0$, $\alpha > 0$, α_{\max}
Top: Lorenz, Bottom: KS

✓ Low $\mathcal{M}(\pi_{\theta_{\mu}}^{\alpha})$ corresponds to low $\mathring{\mathcal{R}}^{\pi}$



$$\mathcal{M}(\pi_{\theta}^{\alpha}) = -\text{Tr}(\mathbb{E}^{X \sim \rho^{\pi_{\theta}}, U \sim \pi_{\theta}(\cdot | X)} [\nabla_{\theta_{\mu}}^2 \log \pi_{\theta}(U | X)])$$

Colors: $\alpha = 0$, $\alpha > 0$, α_{\max}

✓ Measure $\mathcal{M}(\pi_{\theta_{\mu}}^{\alpha})$ distribution with **fat-right tail (extremely large value)** for $\mathring{\mathcal{R}}^{\pi}$ large

✓ Complexity measures can explain noise robustness