Noether Networks: Meta-Learning Useful Conserved Quantities

Ferran Alet^{*1}, Dylan Doblar^{*1}, Allan Zhou², Josh Tenenbaum¹, Kenji Kawaguchi³, Chelsea Finn²; ¹MIT, ²Stanford, ³NUS; *Equal contribution

Noether Networks meta-learn inductive biases in the form of useful conserved quantities that improve predictions when optimized inside the prediction function.

Motivation and background

Successful **inductive biases** often exploit **symmetries**, which are often unknown and hard to discover, or difficult to encode.

We learn **conserved quantities**, inspired by **Noether's Theorem**:

For every **continuous symmetry** property of a dynamical system, there is a corresponding **conservation law**.

Tailoring framework [1]: impose inductive biases at **prediction time**

- **Optimize** unsupervised *tailoring loss* **inside** prediction function
- **Fine-tune** the model to the **particular** query
- **Drawback**: tailoring losses must be **hand-coded**

Prediction and training with neural Noether loss

 $\mathcal{L}_{\text{Noether}}(x_0, \tilde{x}_{1:T}; g_{\phi}) = \sum_{t=1}^T |g_{\phi}(x_0) - g_{\phi}(\tilde{x}_t)|^2$

1: **procedure** PREDICTSEQUENCE($x_0; \theta, \phi$)

- $x_t \leftarrow f_{\theta(x_0;\phi)}(x_{t-1}) \ \forall t \in \{1,\ldots,T\}$ J.
- return $\hat{x}_{1:T}$ 6:
- 7: procedure TRAIN
- $\phi \leftarrow$ randomly initialized embedding weights 8:
- $\theta \leftarrow$ randomly initialized predictive model weights 9:
- 10: while not done do
- Sample batch $x_{0:T}^{(0)}, \ldots, x_{0:T}^{(N)} \sim \mathcal{D}_{\text{train}}$ 11:
- for $0 \le n \le N$ do 12:

13:
$$\hat{x}_{1:T}^{(n)} \leftarrow \text{PREDICTSEQUENCE}(x_0^{(n)}; \theta, \phi)$$

14:
$$\phi \leftarrow \phi - \lambda_{\text{emb}} \nabla_{\phi} \sum_{n=0}^{N} \mathcal{L}_{\text{task}}(\hat{x}_{1:T}^{(n)}, x_{1:T}^{(n)})$$

15:
$$\theta \leftarrow \theta - \lambda_{\text{out}} \nabla_{\theta} \sum_{n=0}^{N} \mathcal{L}_{\text{task}}(\hat{x}_{1:T}^{(n)}, x_{1:T}^{(n)})$$

16: **return** ϕ, θ



Learning from raw pixels in real-world video prediction



Optimizing Noether conservation loss at prediction time

Physics 101 [5]: colliding objects • Grad-CAM heatmaps: embeddings attend to relevant pixels (right) • Improves performance, esp. with long horizon (below, left and center) • Generalizes beyond a single test-time inner step (below, right)



Theoretical advantages: enforcing conservations

Each conservation law decreases ξ by 1, lowering the train-test gap.

Noether Nets recover known conservation laws

Noether loss is parameterized as symbolic formula in a simple DSL.

Method		Description		RMSE		Method	D	escription	RMSE	
Vanilla N Noether	MLP Nets	N/A $p^2-2.99{ m cc}$	$\operatorname{s}(q)$	$0.0563 \\ 0.0423$		Vanilla MLF Noether Nets	s q^2	N/A $+ 1.002 p^2$	$.0174 \\ .0165$	
True \mathcal{H} [o	oracle]	$p^2 - 3.00 { m cc}$	$\operatorname{s}(q)$	0.0422	-> -	True \mathcal{H} [oracl	e] q^2	$+1.000 p^2$.0166	
Ideal pendulum						Ideal spring				
MSE between coordinates Change in HNN-conserved quantity										
0.	020	Baseline Noether	e NN Network	J.M	0.	0 min	2~	\sim		
0.	015	Hand-co	oded loss nian NN	r	-0.		- Ground	d truth		
0.	010		1	MA:	-0. -0.	3	Noethe Oracle	er Network loss onian NN		
0.	005	M	way	-	-0.	4				
0.000										
0 10 Prediction horizon						0 10 Prediction horizon				
Real (dissipative) pendulum system										
		()					, , , , , , , , , ,			
proving predictions with controlled dynamics										
Baseline	Ground Truth	Noether Network	1e-3	Tes	t MS	SE	1.00	Test SS	Moether Network	
J	J	6	3				0.99		– SVG (more steps) – SVG	
•			2			SSIM	0.97 0.96 0.95			
1	-	-					0.94			
-	-	-	0 5	5 10 Predicti	1 on h	5 20 25 orizon	0	5 10 1 Prediction I	15 20 25	
Natod work										



Related work

[1] Tailoring: encoding inductive biases by optimizing losses optimized at prediction time; Alet et al. [2] Meta-learning symmetries by reparameterization; Zhou et al. [3] Hamiltonian neural networks; Greydanus et al. [4] Stochastic video generation with a learned prior; Denton and Fergus. (Used as base video prediction model) [5] Physics 101: Learning physical object properties from unlabeled videos; Wu et al. [6] AI Poincaré: Machine learning conservation laws from trajectories; Liu and Tegmark.



Theorem 1. Let $\rho \in \mathbb{N}^+$. Then, for any $\delta > 0$, with probability at least $1 - \delta$ over an iid draw of n examples $((x_i, y_i))_{i=1}^n$, the following holds for all $\theta \in \Theta$:

 $G(\theta) \le C\sqrt{\frac{\xi \ln(\max(\sqrt{\xi}, 1)) + \xi \ln(2R(\zeta^{1-1/\rho})\sqrt{n}) + \ln(1/\delta)}{2\pi}} + \mathbb{1}\{\xi \ge 1\}\sqrt{\frac{\zeta^{2/\rho}}{\pi}}.$ (1)

where $\xi = m$ if $g_{\phi}(f_{\theta}(x)) = g_{\phi}(x)$ for any $x \in \mathbb{R}^d$ and $\theta \in \Theta$, and $\xi = d$ otherwise.