



34th Conference on Neural Information Processing Systems NeurIPS 2020

Ode to an ODE

Krzysztof Choromanski, Jared Quincy Davis, Valerii Likhosherstov, Xingyou Song, Jean-Jacques Slotine, Jacob Varley, Honglak Lee, Adrian Weller, Vikas Sindhwani



Neural ODEs:

• Continuous variants of standard ResNet networks:

$$\frac{d\mathbf{x}(t)}{dt} = f(\mathbf{x}_t, t, \theta) \qquad \mathbf{x}_t = \mathbf{x}_{t_0} + \int_{t_0}^t f(\mathbf{x}_s, s, \theta) ds \quad (1)$$

- Emulate deep discrete neural networks with **compact** number of parameters.
- Parameters of the Neural ODEs encapsulated in the mapping $\theta(t)$. How to design it ?
- As every deep neural network system, suffer from exploding/vanishing gradients which makes training challenging. Can we robustify Neural ODEs ?



Ode to an ODE System:

• IDEA: Design $\theta(t)$, so that when integrated, Neural ODE emulates deep ResNet with orthogonal connection matrices.



• This leads to the matrix-flow on the **orthogonal group** and effectively: to a **nested system of flows**, where the orthogonal flow encoding $\theta(t)$ determines main flow. How to design learnable orthogonal flows and why are they good ?



• b_{ψ} can be modeled by a neural network producing skew-symmetric matrices or via parameterized isospectral flows (e.g. double-bracket flows)

Lemma 4.1 (ODEtoODES for gradient stabilization). Consider a Neural ODE on time interval [0, T] and given by Formula 2. Let $\mathcal{L} = \mathcal{L}(\mathbf{x}_T)$ be a differentiable loss function. The following holds for any $t \in [0, 1]$, where e = 2.71828... is Euler constant:

$$\frac{1}{e} \| \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} \|_2 \le \| \frac{\partial \mathcal{L}}{\partial \mathbf{x}_t} \|_2 \le e \| \frac{\partial \mathcal{L}}{\partial \mathbf{x}_T} \|_2.$$

Ode to an ODE System in Action:







The Alan Turing Institute



Thank you for your Attention !

https://arxiv.org/abs/2006.11421