

Differentiable Implicit Layers

A. Look¹, S. Doneva², M. Kandemir¹, R. Gemulla², J. Peters³

¹Bosch Center for Artificial Intelligence, ²University Mannheim, ³TU Darmstadt



1 Motivation

Implicit functions are used in a wide range of domains, e.g. physics, numerics, or math. Learning of implicitly defined functions is not feasible with the standard deep learning practice. We introduce the framework of unconstrained and non-convex **Differentiable Implicit Layers (DIL)** as a plug-and-play extension for neural networks that enables efficient learning of such implicitly defined problems.

Aim: Scalable general purpose implicit layers.

2 Implicit Layers

An implicit layer obtains an output \mathbf{y} as an argmin-solution [1]

$$\mathbf{y} := \underset{\mathbf{u}}{\operatorname{argmin}} f(\mathbf{u}; \mathbf{x}, \boldsymbol{\theta}). \quad (1)$$

- Output $\mathbf{y} \in \mathbb{R}^{D_y}$
- Score function $f: \mathbb{R}^{D_y+D_x+D_\theta} \rightarrow \mathbb{R}$
- Learnable weights $\boldsymbol{\theta} \in \mathbb{R}^{D_\theta}$
- Optional input $\mathbf{x} \in \mathbb{R}^{D_x}$

During training the solution \mathbf{y} is evaluated on a **scalar loss function** $\mathcal{L}(\mathbf{y})$. Backpropagation $\partial \mathcal{L} / \partial \boldsymbol{\theta}$ via **Implicit Function Theorem (IFT)**

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = -\frac{\partial \mathcal{L}}{\partial \mathbf{y}} \left(\frac{\partial^2 f}{\partial \mathbf{y}^2} \right)^{-1} \left(\frac{\partial^2 f}{\partial \boldsymbol{\theta} \partial \mathbf{y}} \right). \quad (2)$$

Problem: Costly inverse Hessian $\mathbf{H}^{-1} = (\partial^2 f / \partial \mathbf{y}^2)^{-1}$.

3 The proposed Framework

Instead of explicitly estimating the inverse Hessian, rather estimate the vector-inverse Hessian product

$$\mathbf{g}^T = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \left(\frac{\partial^2 f}{\partial \mathbf{y}^2} \right)^{-1}. \quad (3)$$

Obtain \mathbf{g} as the solution to the linear system of equations

$$\mathbf{H} \left(\overbrace{\mathbf{H}^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{y}}}^{\mathbf{g}} \right) = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}. \quad (4)$$

- Hessian \mathbf{H} is positive semi-definite
- LSE can be solved via the **Conjugate Gradient** method [2] without evaluating the Hessian

Algorithm 1 Forward/ Backward Evaluation for a DIL

```

Input: Score Function  $f(\cdot; \mathbf{x}, \boldsymbol{\theta})$ , Parameters  $\boldsymbol{\theta}$ , solver( $\cdot$ )
function Forward( $\mathbf{x}$ )
     $\mathbf{y} = \text{solver}(f(\cdot; \mathbf{x}, \boldsymbol{\theta}))$  ▷ Optional Input  $\mathbf{x}$ 
    return  $\mathbf{y}$  ▷ Solve for argmin-solution  $\mathbf{y}$  with user defined solver
function Backward( $\mathbf{y}, \partial \mathcal{L} / \partial \mathbf{y}$ )
    Set  $\mathbf{g}_1 = \partial f / \partial \mathbf{y}, \mathbf{g}_2 = \partial \mathcal{L} / \partial \mathbf{y}$  ▷ Via IFT
     $\mathbf{g} = \text{VJP\_CG}(\mathbf{y}, \mathbf{g}_1, \mathbf{g}_2)$  ▷ Score, Loss function gradient at optimal solution
    return  $d\mathcal{L}/d\mathbf{x}, d\mathcal{L}/d\boldsymbol{\theta}$  ▷ vector-inv. Hessian product  $\mathbf{g}^T = \frac{\partial \mathcal{L}}{\partial \mathbf{y}} \left( \frac{\partial^2 f}{\partial \mathbf{y}^2} \right)^{-1}$ 
     $d\mathcal{L}/d\mathbf{x} = -\text{grad}(\mathbf{g}_1, \mathbf{x}, \text{grad\_outputs} = \mathbf{g})^T$  ▷ Return  $-\mathbf{g}^T \frac{\partial^2 f}{\partial \mathbf{x} \partial \mathbf{y}}$  in  $\frac{d\mathcal{L}}{d\mathbf{x}}$ 
     $d\mathcal{L}/d\boldsymbol{\theta} = -\text{grad}(\mathbf{g}_1, \boldsymbol{\theta}, \text{grad\_outputs} = \mathbf{g})^T$  ▷ Return  $-\mathbf{g}^T \frac{\partial^2 f}{\partial \boldsymbol{\theta} \partial \mathbf{y}}$  in  $\frac{d\mathcal{L}}{d\boldsymbol{\theta}}$ 
function VJP_CG( $\mathbf{y}, \partial f / \partial \mathbf{y}, \partial \mathcal{L} / \partial \mathbf{y}$ ) ▷ CG with efficient grad calls
    Init  $\mathbf{x}_0$ 
    Set  $\mathbf{r}_0 = \partial \mathcal{L} / \partial \mathbf{y} - (\text{grad}(\partial f / \partial \mathbf{y}, \mathbf{y}, \text{grad\_outputs} = \mathbf{x}_0))^T + \epsilon \mathbf{x}_0$ 
    Set  $\mathbf{p}_0 = \mathbf{r}_0, k = 0$ 
    while  $\|\mathbf{r}_k\| > \text{tol}$  do
         $\mathbf{A}\mathbf{p}_k = \text{grad}(\partial f / \partial \mathbf{y}, \mathbf{y}, \text{grad\_outputs} = \mathbf{p}_k)^T + \epsilon \mathbf{p}_k$  ▷ VJP = JVPT
         $\alpha_k = (\mathbf{r}_k^T \mathbf{r}_k) / (\mathbf{p}_k^T \mathbf{A}\mathbf{p}_k)$ 
         $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha \mathbf{p}_k$ 
         $\mathbf{r}_{k+1} = \mathbf{r}_k - \alpha \mathbf{A}\mathbf{p}_k$ 
         $\beta_k = (\mathbf{r}_{k+1}^T \mathbf{r}_{k+1}) / (\mathbf{r}_k^T \mathbf{r}_k)$ 
         $\mathbf{p}_{k+1} = \mathbf{r}_{k+1} + \beta_k \mathbf{p}_k$ 
         $k = k + 1$ 
    return  $\mathbf{x}_{k+1}$  ▷ Solution  $\mathbf{g}$  in  $\mathbf{H}\mathbf{g} = \frac{\partial \mathcal{L}}{\partial \mathbf{y}}$ 
    
```

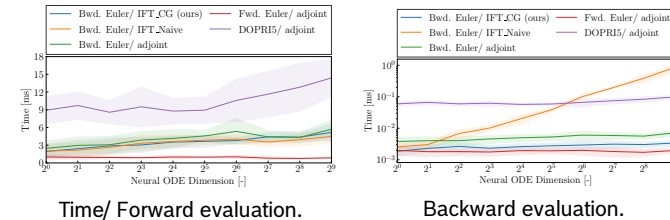
4 Applications

4.1 Implicit Neural Ordinary Differential Equations

Solve a Neural Ordinary Differential Equation [3] with the backward Euler method [4] and obtain

$$d\mathbf{x} = \mathbf{h}(\mathbf{x}; \boldsymbol{\theta}) dt \xrightarrow[\text{Euler}]{\text{Backward}} \mathbf{x}_{t+1} = \mathbf{x}_t + \mathbf{h}(\mathbf{x}_{t+1}; \boldsymbol{\theta}) \Delta t. \quad (5)$$

Runtime Profiles. Forward/ backward evaluation times for random NODE initializations (Layers: 2, Hidden size: 30).



- Backward Euler faster than default NODE solver DOPRI5 [3]
- Backpropagation via IFT faster than adjoint method [3]

Predictive Performance. Average MSE on extrapolation tasks.

Model	Van der Pol	Spiral Data	CMU Walking
DOPRI5 _{adj.}	0.89 ± 0.15	0.13 ± 0.01	15.92 ± 2.10
Fwd. Euler _{adj.}	0.68 ± 0.07	0.20 ± 0.01	12.17 ± 1.39
Bwd. Euler _{adj.}	0.67 ± 0.12	0.09 ± 0.01	13.68 ± 2.02
Bwd. Euler _{IFT, Naive}	0.38 ± 0.05	0.09 ± 0.00	11.57 ± 1.79
Bwd. Euler _{IFT, CG (ours)}	0.40 ± 0.06	0.09 ± 0.01	11.43 ± 1.27

- Backward Euler outperforms other solvers
- Backpropagation via CG on par with naive IFT application

4.2 Differentiable Path Planning

Observe current state of the system $\mathbf{x}_{obs.}$ and plan optimal trajectory on a limited horizon H [5] by minimizing the total cost $\sum c$. Afterwards, execute first control \mathbf{u}_0 and move time step forwards.

$$\underset{\mathbf{u}_{0:H}}{\operatorname{argmin}} \sum_{t=0}^H c(\mathbf{x}_t, \mathbf{u}_t; \boldsymbol{\theta}_c), \text{ s.t. } \mathbf{x}_{t+1} = \mathbf{h}(\mathbf{x}_t, \mathbf{u}_t; \boldsymbol{\theta}_h), \mathbf{x}_0 = \mathbf{x}_{obs.}, \quad (6)$$

Imitation Learning from Observations. Learn jointly linear cost c and neural dynamical model h by observing N expert trajectories $\mathbf{x}_{1:T}^N$ with horizon T on cartpole swing-up task. Train on dataset with low variance at initial state \mathbf{x}_0 .

Model	Low Variance $\mathbf{x}_0 \sim \mathcal{N}(0, 0.04\mathbf{I})$	High Variance $\mathbf{x}_0 \sim \mathcal{N}(0, 0.08\mathbf{I})$
Expert	9.2 ± 0.0	(9.3 ± 0.1)
Behavioral Cloning	14.6 ± 0.6	18.4 ± 1.4
MPC _{IFT} (ours)	13.7 ± 0.4	14.7 ± 1.9

- Our method outperforms behavioral cloning and comes with improved generalization capabilities.

References

- [1] S. Gould, R. Hartley, and D. Campbell. Deep Declarative Networks: A New Hope. *arXiv*, abs/1909.04866, 2019.
- [2] J. R. Shewchuk. An introduction to the conjugate gradient method without the agonizing pain. Technical report, 1994.
- [3] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud. Neural Ordinary Differential Equations. In *NeurIPS*. 2018.
- [4] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*. Springer, 1993.
- [5] M. Diehl. Numerical optimal control, 2011.