



SHINE: SHaring the INverse Estimate from the forward pass for bi-level optimization and implicit models

ICLR 2022 Spotlight

Zaccharie Ramzi

CEA & Inria Parietal

Joint work with: F. Mannel, S. Bai, P. Ciuciu, J.-L. Starck, T. Moreau

Deep Equilibrium networks

Deep Equilibrium networks (DEQs) (Bai, Kolter, et al., 2019) are a type of implicit model. The output is the solution to a fixed-point equation.

$$h_{\theta}(x) = z^*, \text{ where } z^* = f_{\theta}(z^*, x)$$

Deep Equilibrium networks

Deep Equilibrium networks (DEQs) (Bai, Kolter, et al., 2019) are a type of implicit model. The output is the solution to a fixed-point equation.

$$h_{\theta}(x) = z^*, \text{ where } z^* = f_{\theta}(z^*, x)$$

This approximates an infinite depth network:

$$z_n = f_{\theta}(z_{n-1}), \quad \forall n \rightarrow \infty$$

In practice, we work with root finding algorithms using $g_{\theta} = id - f_{\theta}$.

quasi-Newton methods

For the forward pass' root finding problem

$$g(z^*) = 0$$

Newton Methods: $z_{n+1} = z_n - J_g(z_n)^{-1}g(z_n)$

quasi-Newton methods

For the forward pass' root finding problem

$$g(z^*) = 0$$

Newton Methods: $z_{n+1} = z_n - J_g(z_n)^{-1}g(z_n)$

Idea: replace the costly Jacobian inverse $J_g(z_n)^{-1}$ with a qN matrix B_n^{-1} .

quasi-Newton Methods: $z_{n+1} = z_n - B_n^{-1}g(z_n)$.

DEQ's backward pass

DEQs gradient computation using Implicit Function Theorem:

$$\frac{\partial \mathcal{L}}{\partial \theta} \Big|_{z^*} = \nabla_z \mathcal{L}(z^*)^\top J_{g_\theta}(z^*)^{-1} \frac{\partial g_\theta}{\partial \theta} \Big|_{z^*},$$

we need to invert a huge matrix $J_{g_\theta}(z^*)$ in a certain direction $\nabla_z \mathcal{L}(z^*)$.

In practice this is done using an iterative algorithm.

The limits of DEQs

DEQs achieve excellent results in NLP (Natural Language Processing) and CV (Computer Vision) tasks, but they are slow to train.

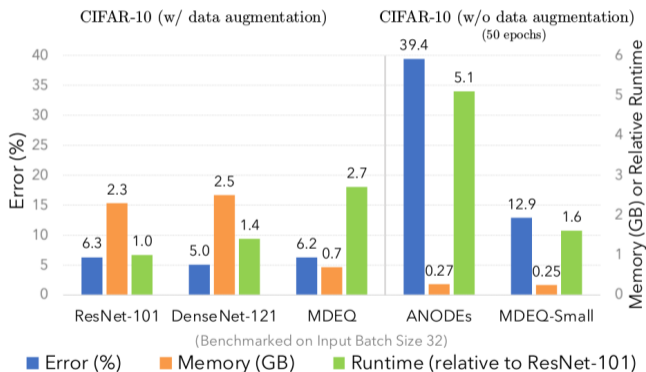


Figure: Performance, memory and training speed of DEQs. (Bai, Koltun, et al., 2020)

Can we avoid the Jacobian inversion?

We introduced **SHINE: SHaring the INverse Estimate**.

$$B^{-1} \approx J_{g\theta}(z^*)^{-1}$$

Can we avoid the Jacobian inversion?

We introduced **SHINE: SHaring the INverse Estimate**.

$$B^{-1} \approx J_{g\theta}(z^*)^{-1}$$

True Jacobian inverse

Can we avoid the Jacobian inversion?

We introduced **SHINE: SHaring the INverse Estimate**.

$$\underbrace{B^{-1}}_{\text{quasi-Newton matrix}} \approx \underbrace{J_{g\theta}(z^*)^{-1}}_{\text{True Jacobian inverse}}$$

Can we avoid the Jacobian inversion?

We introduced **SHINE: SHaring the INverse Estimate**.

$$\underbrace{B^{-1}}_{\text{quasi-Newton matrix}} \approx \underbrace{J_{g_\theta}(z^*)^{-1}}_{\text{True Jacobian inverse}}$$

Properties of B :

- It is computed when solving $z^* - f_\theta(z^*, x) = 0$ using a quasi-Newton method.

Can we avoid the Jacobian inversion?

We introduced **SHINE: SHaring the INverse Estimate**.

$$\underbrace{B^{-1}}_{\text{quasi-Newton matrix}} \approx \underbrace{J_{g_\theta}(z^*)^{-1}}_{\text{True Jacobian inverse}}$$

Properties of B :

- It is computed when solving $z^* - f_\theta(z^*, x) = 0$ using a quasi-Newton method.
- It is easily invertible using the Sherman-Morrison formula, because low-rank.

SHINE direction convergence

Theorem (Convergence of SHINE to the Hypergradient using ULI)

Under the Uniform Linear Independence (ULI) assumption and some additional smoothness and convexity assumptions, for a given parameter θ , (z_n) converges q -superlinearly to z^ and*

$$\lim_{n \rightarrow \infty} \nabla_z \mathcal{L}(z_n)^\top \mathbf{B}_n^{-1} \frac{\partial g_\theta}{\partial \theta} \Big|_{z_n} = \frac{\partial \mathcal{L}}{\partial \theta} \Big|_{z^*}.$$

Computer vision results

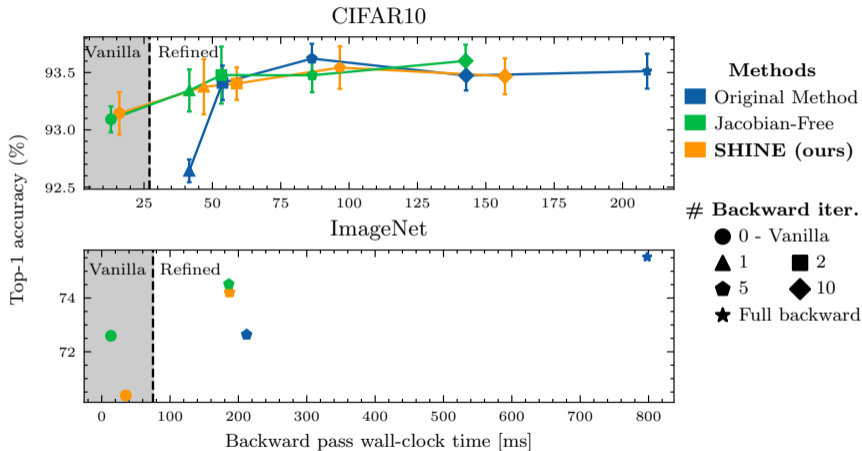


Figure: MDEQs (Bai, Koltun, et al., 2020) with SHINE.

Come check our poster #6363

SHINE accelerates the training of **Deep Equilibrium Networks**.

Come check our poster #6363

SHINE accelerates the training of **Deep Equilibrium Networks**.

Come chat with us at our poster to see:

- How we can obtain better theoretical guarantees by modifying the forward pass.
- Our application to Bi-Level optimization.