

On Training Implicit Models

Zhengyang Geng^{1,2}, Xin-Yu Zhang², Shaojie Bai³, Yisen Wang², Zhouchen Lin^{2,4}
¹Zhejiang Lab, ²Peking University, ³Carnegie Mellon University, ⁴Pazhou Lab

arXiv: <https://arxiv.org/abs/2111.05177>
 Github: https://github.com/Gsunshine/phantom_grad
 Email: zhengyanggeng@gmail.com



Background

- Implicitly-defined neural networks have achieved competitive performances compared with explicit models.
- Implicit models treat the evolution of hidden states as certain dynamics, *e.g.*, fixed-point equations or ordinary differential equations (ODEs);
- The forward passes are formulated as black-box solvers of the underlying dynamics, and the backward passes are performed via implicit differentiation.
- In this work, we argue that a carefully designed inexact gradient, named phantom gradient, is sufficient to efficiently and effectively train implicit models.

Implicit Differentiation

We adopt the formulation of DEQ models [1].

- Input projection module \mathcal{M} : $\mathbf{u} = \mathcal{M}(\mathbf{x})$, where \mathbf{x} is the input data;
- Equilibrium module \mathcal{F} and the equilibrium state \mathbf{h}^* given by

$$\mathbf{h}^* = \mathcal{F}(\mathbf{h}^*, \mathbf{z}),$$

where \mathbf{z} is a union of the module's input \mathbf{u} and parameters $\boldsymbol{\theta}$;

- Post-processing module \mathcal{G} : $\hat{\mathbf{y}} = \mathcal{G}(\mathbf{h}^*)$, where $\hat{\mathbf{y}}$ is the predicted label of \mathbf{x} ;
- Loss function \mathcal{L} and the training objective, *i.e.*, the expected loss:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{P}} [\mathcal{L}(\hat{\mathbf{y}}(\mathbf{x}; \boldsymbol{\theta}), \mathbf{y})],$$

where \mathbf{y} is the true label of \mathbf{x} .

- Using Implicit Differentiation, the gradient of \mathbf{h}^* *w.r.t.* \mathbf{z} is given by

$$\frac{\partial \mathbf{h}^*}{\partial \mathbf{z}} = \frac{\partial \mathcal{F}}{\partial \mathbf{z}} \bigg|_{\mathbf{h}^*} \left(\mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \bigg|_{\mathbf{h}^*} \right)^{-1}.$$

The gradient of \mathcal{L} *w.r.t.* \mathbf{z} is thus given by

$$\frac{\partial \mathcal{L}}{\partial \mathbf{z}} = \frac{\partial \mathbf{h}^*}{\partial \mathbf{z}} \frac{\partial \mathcal{L}}{\partial \mathbf{h}} = \frac{\partial \mathcal{F}}{\partial \mathbf{z}} \bigg|_{\mathbf{h}^*} \left(\mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \bigg|_{\mathbf{h}^*} \right)^{-1} \frac{\partial \mathcal{L}}{\partial \mathbf{h}^*}.$$

Future Perspectives

- The phantom gradient may come with a structured bias in comparison with the exact one; how to eliminate the structured bias?
- The UPG and its precision in the training process suggest developing an adaptive gradient solver.
- (Aggressive) The loss landscape and the training strategy are the two sides of the same coin; how to study their interaction in training implicit models?

Phantom Gradient

- **Definition.** The Jacobian $\partial \mathbf{h}^* / \partial \boldsymbol{\theta}$ is approximated by a matrix \mathbf{A} :

$$\frac{\partial \widehat{\mathcal{L}}}{\partial \boldsymbol{\theta}} \triangleq \mathbf{A} \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}.$$

- **General Descent Condition.**

Theorem 1. Let σ_{\max} and σ_{\min} be the maximal and minimal singular value of $\partial \mathcal{F} / \partial \boldsymbol{\theta}$. If

$$\left\| \mathbf{A} \left(\mathbf{I} - \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \right) - \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \right\| \leq \frac{\sigma_{\min}^2}{\sigma_{\max}},$$

then the phantom gradient provides an ascent direction of the function \mathcal{F} , *i.e.*,

$$\left\langle \frac{\partial \widehat{\mathcal{L}}}{\partial \boldsymbol{\theta}}, \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\rangle \geq 0.$$

- **Instantiations.**

- a. **Unrolling-based Phantom Gradient (UPG).** Consider the damped fixed-point iteration:

$$\mathbf{h}_{t+1} = \lambda \mathcal{F}(\mathbf{h}_t, \mathbf{z}) + (1 - \lambda) \mathbf{h}_t, \quad t = 0, 1, \dots, T - 1.$$

Then, the matrix \mathbf{A} is given by

$$\mathbf{A}_{k,\lambda}^{\text{unr}} = \lambda \sum_{t=0}^{k-1} \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \bigg|_{\mathbf{h}_t} \prod_{s=t+1}^{k-1} \left(\lambda \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \bigg|_{\mathbf{h}_s} + (1 - \lambda) \mathbf{I} \right).$$

- b. **Neumann-series-based Phantom Gradient (NPG).** The matrix \mathbf{A} is given by

$$\mathbf{A}_{k,\lambda}^{\text{neu}} = \lambda \frac{\partial \mathcal{F}}{\partial \boldsymbol{\theta}} \bigg|_{\mathbf{h}^*} (\mathbf{I} + \mathbf{B} + \mathbf{B}^2 + \dots + \mathbf{B}^{k-1}), \quad \text{where } \mathbf{B} = \lambda \frac{\partial \mathcal{F}}{\partial \mathbf{h}} \bigg|_{\mathbf{h}^*} + (1 - \lambda) \mathbf{I}.$$

- **Convergence Theory.**

Theorem 3. Suppose the loss function \mathcal{R} is ℓ -smooth, lower-bounded, and has bounded gradient almost surely in the training process. Besides, assume the gradient $\partial \mathcal{L} / \partial \boldsymbol{\theta}$ is an unbiased estimator of $\nabla \mathcal{R}(\boldsymbol{\theta})$ with a bounded covariance. If the phantom gradient in is an ε -approximation to $\partial \mathcal{L} / \partial \boldsymbol{\theta}$, *i.e.*,

$$\left\| \frac{\partial \widehat{\mathcal{L}}}{\partial \boldsymbol{\theta}} - \frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} \right\| \leq \varepsilon, \quad \text{almost surely,}$$

then using the phantom gradient as a stochastic first-order oracle with a step size of $\eta_n = \mathcal{O}(1/\sqrt{n})$ to update $\boldsymbol{\theta}$ with gradient descent, it follows after N iterations that

$$\mathbb{E} \left[\frac{\sum_{n=1}^N \eta_n \|\nabla \mathcal{R}(\boldsymbol{\theta}_n)\|^2}{\sum_{n=1}^N \eta_n} \right] \leq \mathcal{O} \left(\varepsilon + \frac{\log N}{\sqrt{N}} \right).$$

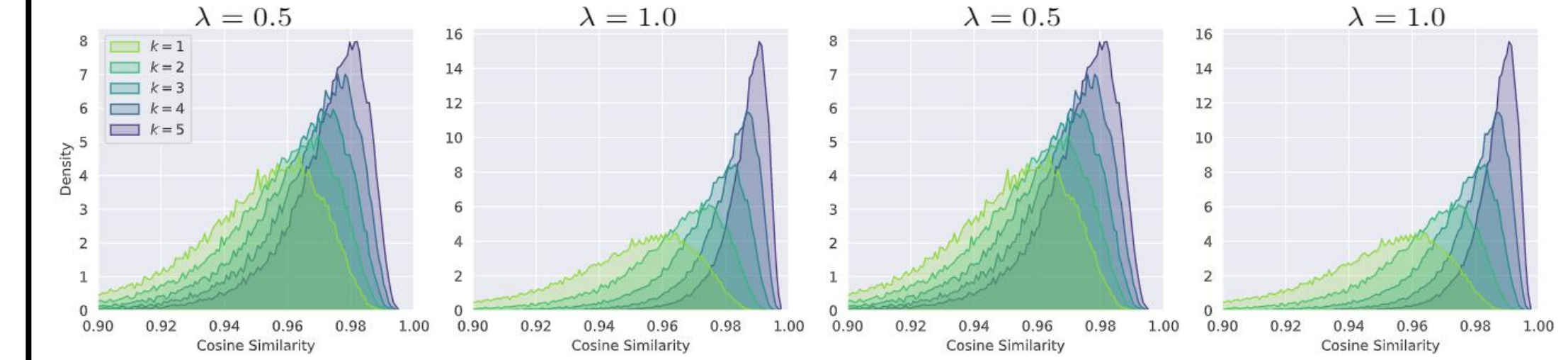
- **Complexity.**

Let **Mem** denote the memory cost, and K and k be the solver's steps and the unrolling/Neumann steps, respectively. Here, $K \gg k \approx 1$.

Method	Time	Mem	Peak Mem
Implicit	$\mathcal{O}(K)$	$\mathcal{O}(1)$	$\mathcal{O}(k)$
UPG	$\mathcal{O}(k)$	$\mathcal{O}(k)$	$\mathcal{O}(k)$
NPG	$\mathcal{O}(k)$	$\mathcal{O}(1)$	$\mathcal{O}(1)$

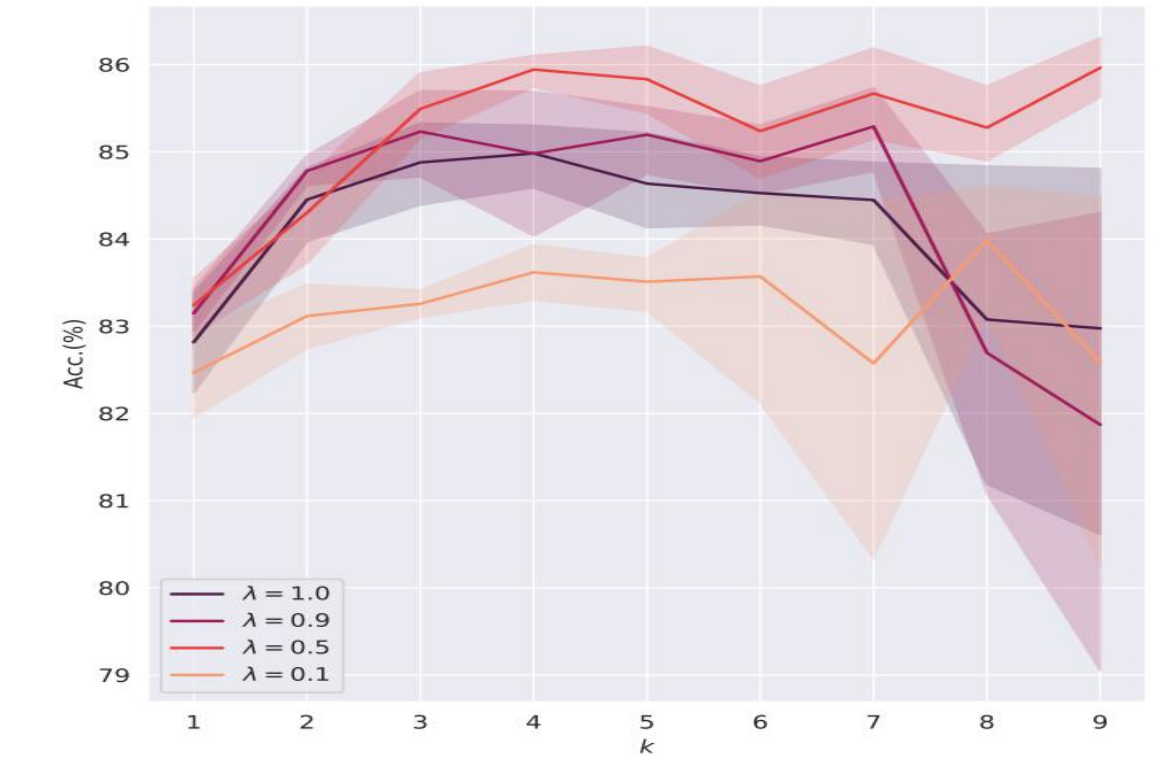
Experiments

- Cosine similarity between the phantom gradient and the exact gradient in the synthetic setting (see paper for details);



(a) Phantom gradient in the Neumann form (b) Phantom gradient in the unrolling form

- Impact of hyperparameters λ and k on the CIFAR-10 classification accuracy;



- Large-scale experiments;

Datasets	Model	Method	Params	Metrics	Speed
CIFAR-10	MDEQ	Implicit	10M	93.8 ± 0.17	1.0×
CIFAR-10	MDEQ	UPG $\mathbf{A}_{5,0.5}$	10M	95.0 ± 0.16	1.4×
ImageNet	MDEQ	Implicit	18M	75.3	1.0×
ImageNet	MDEQ	UPG $\mathbf{A}_{5,0.6}$	18M	75.7	1.7×
Wikitext-103	DEQ (PostLN)	Implicit	98M	24.0	1.0×
Wikitext-103	DEQ (PostLN)	UPG $\mathbf{A}_{5,0.8}$	98M	25.7	1.7×
Wikitext-103	DEQ (PreLN)	JR + Implicit	98M	24.5	1.7×
Wikitext-103	DEQ (PreLN)	JR + UPG $\mathbf{A}_{5,0.8}$	98M	24.4	2.2×
Wikitext-103	DEQ (PreLN)	JR + UPG $\mathbf{A}_{5,0.8}$	98M	24.0 [†]	1.7×

- Implicit GNN [2] model on graph tasks.

Datasets	Model	Method	Params	Metrics (%)
COX2	IGNN	Implicit	38K	84.1 ± 2.9
COX2	IGNN	UPG $\mathbf{A}_{5,0.5}$	38K	83.9 ± 3.0
COX2	IGNN	UPG $\mathbf{A}_{5,0.8}$	38K	83.9 ± 2.7
COX2	IGNN	UPG $\mathbf{A}_{5,1.0}$	38K	83.0 ± 2.9
PROTEINS	IGNN	Implicit	34K	78.6 ± 4.1
PROTEINS	IGNN	UPG $\mathbf{A}_{5,0.5}$	34K	78.4 ± 4.2
PROTEINS	IGNN	UPG $\mathbf{A}_{5,0.8}$	34K	78.6 ± 4.2
PROTEINS	IGNN	UPG $\mathbf{A}_{5,1.0}$	34K	78.8 ± 4.2
PPI	IGNN	Implicit	4.7M	97.6
PPI	IGNN	UPG $\mathbf{A}_{5,0.5}$	4.7M	98.2
PPI	IGNN	UPG $\mathbf{A}_{5,0.8}$	4.7M	97.4
PPI	IGNN	UPG $\mathbf{A}_{5,1.0}$	4.7M	96.2

References

- [1] Shaojie Bai, J. Zico Kolter, Vladlen Koltun. Deep Equilibrium Models.
- [2] Fangda Gu, Heng Chang, Wenwu Zhu, Somayeh Sojoudi, Laurent El Ghaoui. Implicit Graph Neural Netowrks.