

On the Iteration Complexity of Hypergradient Computation

Riccardo Grazi

Computational Statistics and Machine Learning, Istituto Italiano di Tecnologia.
Department of Computer Science, University College London.

riccardo.grazzi@iit.it

Joint work with Luca Franceschi, Massimiliano Pontil and Saverio Salzo.



ISTITUTO ITALIANO
DI TECNOLOGIA



Bilevel Optimization Problem

$\min_{\theta} f(\theta) \quad (\text{upper-level})$

$\min_{\theta} f(\theta) \quad (\text{lower-level})$

- Hyperparameter optimization, meta-learning.
- Graph and recurrent neural networks.

How can we solve this optimization problem?

- Black-box methods (random/grid search, Bayesian optimization, ...).

Bilevel Optimization Problem

$\min_{\theta} f(\theta) \quad \text{(upper-level)}$

$\min_{\phi} f(\phi) \quad \text{(lower-level)}$

- Hyperparameter optimization, meta-learning.
- Graph and recurrent neural networks.

How can we solve this optimization problem?

- Black-box methods (random/grid search, Bayesian optimization, ...).
- **Gradient-based methods** exploiting the *hypergradient* $\nabla_{\theta} f(\theta)$

Can be really expensive or even impossible to compute Exactly.

Two common approximation strategies are

1. Iterative Differentiation (ITD).
2. Approximate Implicit Differentiation (AID).

Can be really expensive or even impossible to compute Exactly.

Two common approximation strategies are

1. Iterative Differentiation (ITD).
2. Approximate Implicit Differentiation (AID).

Which one is the best?

- Previous works provide mostly qualitative and empirical results.

Can we have quantitative results on the approximation error?

- Yes! If the fixed point map γ_{f_i} is a **contraction**.

Our Contributions

Upper bounds on the approximation error for both ITD and AID

- Both methods achieve **non-asymptotic linear convergence rates**.
- We prove that ITD is generally worse than AID in terms of upper bounds.

Extensive experimental comparison among different AID strategies and ITD

- If f is a contraction, the results confirm the theory.
- If f is *NOT* a contraction, ITD can be still a reliable strategy.

Source: S.Ravi, H. Larochelle (2016).

- Hyperparameter optimization
(learn the kernel/regularization, ...).
- Meta-learning (MAML, L2LOpt, ...).

Source: S.Ravi, H. Larochelle (2016).

- Hyperparameter optimization (learn the kernel/regularization, ...).
- Meta-learning (MAML, L2LOpt, ...).

Source: snap.stanford.edu/proj/embeddings-www

- Graph Neural Networks.
- Some Recurrent Models.
- Deep Equilibrium Models.

Source: S.Ravi, H. Larochelle (2016).

Source: snap.stanford.edu/proj/embeddings-www

- Hyperparameter optimization (learn the kernel/regularization, ...).
- Meta-learning (MAML, L2LOpt, ...).
- Graph Neural Networks.
- Some Recurrent Models.
- Deep Equilibrium Models.

All can be cast into the same **bilevel framework** where at the lower-level we seek for the solution to a **parametric fixed point equation**.

Example: Optimizing the Regularization Hyperparameter in Ridge Regression

$$T_{\lambda} = \left(\frac{1}{n} X^T X + \lambda I \right)^{-1} X^T Y$$

λ is the **unique fixed point** of the *one step GD* map

$$\lambda \mapsto \frac{1}{n} \text{tr} \left(X^T X (\lambda I + X^T X)^{-1} \right)$$

If the step size is sufficiently small, $\lambda \mapsto \frac{1}{n} \text{tr} \left(X^T X (\lambda I + X^T X)^{-1} \right)$ is also a **contraction**.

$$\begin{aligned} & \min_{x \in X} \left[\min_{y \in Y(x)} f(x, y) \right] \quad (\text{upper-level}) \\ & \min_{y \in Y(x)} f(x, y) \quad (\text{lower-level}) \end{aligned}$$

- $\min_{y \in Y(x)} f(x, y)$ is often not available in closed form.
- $f(x, y)$ is usually non convex and **expensive** or **impossible** to evaluate exactly.

$$\begin{aligned} \min_{x \in X} & \quad f(x) \quad (\text{upper-level}) \\ \text{s.t.} & \quad x \in \arg \min_{y \in Y} f(x, y) \quad (\text{lower-level}) \end{aligned}$$

- $\arg \min_{y \in Y} f(x, y)$ is often not available in closed form.
- $f(x, y)$ is usually non convex and **expensive** or **impossible** to evaluate exactly.
- **is even harder to evaluate.**

Iterative Differentiation (ITD)

1. Set \mathbf{f}_i and compute,

$$\text{for } \mathbf{S} \\ \mathbf{f}_i / \mathbf{f}_i / \mathbf{f}_i$$

2. Compute $\mathbf{f}_i / \mathbf{f}_i / \mathbf{f}_i$

3. Compute \mathbf{f}_i *efficiently* using reverse (RMAD) or forward (FMAD) mode automatic differentiation.

Iterative Differentiation (ITD)

1. Set $\mathbf{x} / \mathbf{f}^*$ and compute,

$$\text{for } \mathbf{x} \in \mathcal{S} \\ \nabla \mathbf{f} / \nabla \mathbf{f} \mathbf{f}^*$$

2. Compute $\nabla \mathbf{f} / \nabla \mathbf{f} \mathbf{f}^*$

3. Compute $\nabla \mathbf{f} / \mathbf{f}^*$ efficiently using reverse (RMAD) or forward (FMAD) mode automatic differentiation.

Approximate Implicit Differentiation (AID)

1. Get $\mathbf{x} / \mathbf{f}^*$ with **steps** of a lower-level solver.

2. Compute $\nabla \mathbf{f} / \mathbf{f}^*$ with **steps** of a solver for the *linear system*

$$\nabla \mathbf{f} / \nabla \mathbf{f} \mathbf{f}^* / \nabla \mathbf{f} \mathbf{f}^*$$

3. Compute the approximate gradient as

$$\nabla \mathbf{f} / \nabla \mathbf{f} \mathbf{f}^* / \nabla \mathbf{f} \mathbf{f}^* \\ \mathbb{L} / \nabla \mathbf{f} \mathbf{f}^* / \nabla \mathbf{f} \mathbf{f}^*$$

Iterative Differentiation (ITD)

1. Set \mathbf{x} / f_i and compute,

$$\text{for } \mathbf{x} / \mathbf{S} \\ \mathbf{x} / f_i / \mathbf{x} / f_i f_i$$

2. Compute $\mathbf{x} / f_i / \mathbf{x} / f_i f_i$

3. Compute \mathbf{x} / f_i efficiently using reverse (RMAD) or forward (FMAD) mode automatic differentiation.

Approximate Implicit Differentiation (AID)

1. Get \mathbf{x} / f_i with **steps** of a lower-level solver.

2. Compute \mathbf{x} / f_i with **steps** of a solver for the *linear system*

$$\mathbf{x} / \mathbf{x} / f_i f_i f_i / \mathbf{x} / f_i f_i$$

3. Compute the approximate gradient as

$$\mathbf{x} / f_i / \mathbf{x} / f_i f_i \\ \mathbf{L} / f_i f_i / f_i$$

Which one is the best?

A First Comparison

ITD

- Ignores the bilevel structure.
- Cost in time (RMAD): $\mathcal{O}(n^2)$ / $\mathcal{O}(n^3)$
- Cost in memory (RMAD): $\mathcal{O}(n)$ / $\mathcal{O}(n)$
- **Can we control** $\mathcal{O}(n)$ / $\mathcal{O}(n)$?

AID

- Can use any lower-level solver.
- Cost in time (/): $\mathcal{O}(n^2)$ / $\mathcal{O}(n^3)$
- Cost in memory: $\mathcal{O}(n)$ / $\mathcal{O}(n)$
- **Can we control** $\mathcal{O}(n)$ / $\mathcal{O}(n)$?

$\mathcal{O}(n)$ / $\mathcal{O}(n)$ / $\mathcal{O}(n)$ / $\mathcal{O}(n)$

Previous Work on the Approximation Error

ITD

- $SdY_{\text{ITD}} \approx SdY_{\text{AID}}$
(Franceschi et al., 2018).

AID

- $\frac{1}{n} \sum_{i=1}^n \ell(f_i) \approx \frac{1}{n} \sum_{i=1}^n \ell(f_i^*)$
(Pedregosa, 2016).
- $\frac{1}{n} \sum_{i=1}^n \ell(f_i) \approx \frac{1}{n} \sum_{i=1}^n \ell(f_i^*)$ at a
linear rate in n and λ for
meta-learning with biased
regularization
(Rajeswaran et al., 2019).

$$\frac{1}{n} \sum_{i=1}^n \ell(f_i) \approx \frac{1}{n} \sum_{i=1}^n \ell(f_i^*)$$

Previous Work on the Approximation Error

ITD

- $SdY_{\cdot}[\cdot]$ $SdY_{\cdot}[\cdot]$
(Franceschi et al., 2018).
- **We provide non-asymptotic upper bounds on $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ and $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$.**

AID

- $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ "
(Pedregosa, 2016).
- $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ " at a
linear rate in n and $\frac{1}{n}$ for
meta-learning with biased
regularization
(Rajeswaran et al., 2019).
- **We provide non-asymptotic upper bounds on $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ and $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$.**

$\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$ $\frac{1}{n} \sum_{i=1}^n \ell(f_i)$

Assumptions

- f is a **contraction** with constant $\beta < 1$.
- f , f' , f'' and f''' are Lipschitz continuous

differentiable and

$$\begin{aligned} & \|f(x) - f(y)\| \leq \beta \|x - y\| \\ & \|f'(x) - f'(y)\| \leq L \|x - y\| \\ & \|f''(x) - f''(y)\| \leq L \|x - y\| \\ & \|f'''(x) - f'''(y)\| \leq L \|x - y\| \end{aligned}$$

Main Contribution

Theorem 1

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{1}{\sqrt{f_i}} \left(\frac{1}{\sqrt{f_i}} - \frac{1}{\sqrt{f_i}} \right) \xrightarrow{D} N\left(0, \frac{1}{\sqrt{f_i}}\right)$$

Theorem 2

Let $\frac{1}{\sqrt{f_i}}$ and $\frac{1}{\sqrt{f_i}}$ be independent and assume that

- $\frac{1}{\sqrt{f_i}}$ and $\frac{1}{\sqrt{f_i}}$ are independent,
- $\frac{1}{\sqrt{f_i}}$ and $\frac{1}{\sqrt{f_i}}$ are independent.

Then,

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{1}{\sqrt{f_i}} \left(\frac{1}{\sqrt{f_i}} - \frac{1}{\sqrt{f_i}} \right) \xrightarrow{D} N\left(0, \frac{1}{\sqrt{f_i}}\right)$$

Main Contribution (Part 2)

Efficient solvers for the **linear system** in AID:

- fixed point method (FP)
- conjugate gradient (CG)

Theorem ; r ã

Assume that the lower-level problem is solved as in ITD. Then

$$(FP) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) = \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \mathbb{L} \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*)$$

Moreover, when \mathbb{L} / \mathbb{L} f is symmetric,

$$(CG) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \mathbb{L} \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) / \mathbb{L} \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*) \quad \frac{\partial f}{\partial \mathbf{x}}(\mathbf{x}^*)$$

where . . .

So... Which method has the best approximation error?

From our analysis:

- ITD, CG and FP **converge linearly** (in n and m) to $\|f - f_n\|$
- FP bound \leq ITD bound for every n, m , when $\|f\| \leq 1$.
- CG bound \leq FP bound for n big enough when $\|f\| \leq 1$ / $\|f\|$ is symmetric.

So... Which method has the best approximation error?

From our analysis:

- ITD, CG and FP **converge linearly** (in $\| \cdot \|$ and $\| \cdot \|_F$) to $\| \cdot \| / \mathbf{f}_1$
- FP bound \leq ITD bound for every $\| \cdot \|$, when $\| \cdot \| / \mathbf{f}_1$.
- CG bound \leq FP bound for $\| \cdot \|$ big enough when $\| \cdot \| / \mathbf{f}_1 / \mathbf{f}_1$ \mathbf{f}_1 is symmetric.

Is this true also for the actual error in practice?

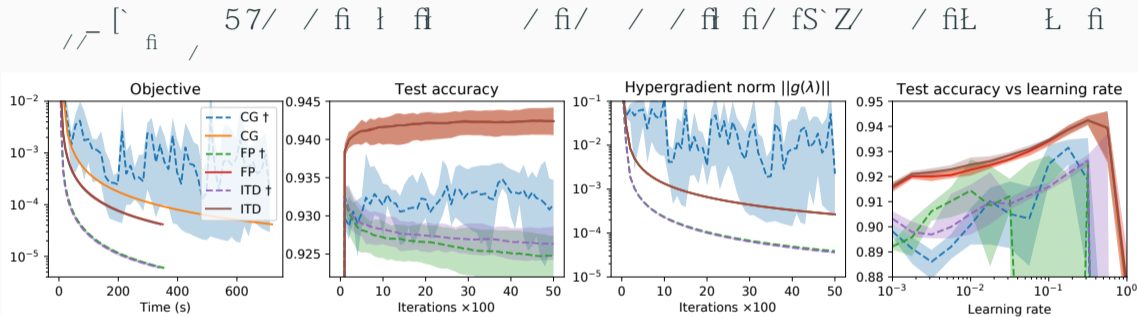
What happens when $\| \cdot \| / \mathbf{f}_1$ \mathbf{f}_1 is not a contraction?

Hypergradient approximation errors (mean/std on randomly drawn values of θ).

$\| \nabla_{\theta} f \|^2$ is equal to $\| \nabla_{\theta} f \|^2$ for ITD and to $\| \nabla_{\theta} f \|^2$ for CG and FP. In all settings $\| \nabla_{\theta} f \|^2$ is a contraction and $\| \nabla_{\theta} f \|^2$ is symmetric.

- After a while the error decreases linearly for all methods.
- Methods with lower error bounds have lower error on average.

Equilibrium Models¹ on MNIST (Proof of Concept)



ITD is NOT a contraction for some methods.

- When ITD is a contraction all the methods perform similarly.
- ITD is the most stable when the contraction assumption is not satisfied.

¹Shaojie Bai, J Zico Kolter, and Vladlen Koltun. "Deep equilibrium models". In Advances in Neural Information Processing Systems. 2019, pp. 688–699.

Conclusions

We studied the **iteration complexity** of two strategies used to approximate the *hypergradient* in bilevel problems: **iterative differentiation (ITD)** and **approximate implicit differentiation (AID)**.

We proved non-asymptotic upper bounds on the approximation error

- CG, FP and ITD converge linearly to the exact *hypergradient*.
- ITD is generally worse than AID in terms of upper bounds.

We conducted experiments comparing ITD and AID

- If $\|f\|$ is a contraction, the results confirm the theory.
- If $\|f\|$ is *NOT* a contraction, ITD can be still a reliable strategy.

Thank you for the attention

CODE (PyTorch): <https://github.com/prolearner/hypertorch>

References

Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., and Pontil, M. (2018). Bilevel programming for hyperparameter optimization and meta-learning. In *International Conference on Machine Learning*, pages 1563–1572.

Pedregosa, F. (2016). Hyperparameter optimization with approximate gradient. In *International Conference on Machine Learning*, pages 737–746.

Rajeswaran, A., Finn, C., Kakade, S. M., and Levine, S. (2019). Meta-learning with implicit gradients. In *Advances in Neural Information Processing Systems*, pages 113–124.