

SHINE: Sharing the Inverse Estimate for bi-level optimization.

Thomas Moreau INRIA Saclay

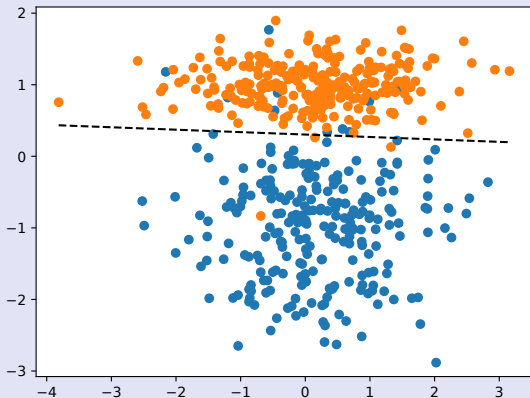
Joint work with Z. Ramzi, S. Bai, F. Mannel, J.-L. Starck & P. Ciuciu



Learning a linear ML model

Setup:

- ▶ Binary classification task $(X_i, y_i)_{i=1}^N \in \mathbb{R}^p \times \{-1, 1\}$
- ▶ Linear model: predict y from X with $\text{sign}(\langle \theta, X \rangle)$.



Regularized Logistic loss:

$$G(\lambda, \theta) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-y_i \langle \theta, X_i \rangle}) + \lambda \|\theta\|_2^2$$

Training the model:

$$\theta^*(\lambda) = \underset{\theta}{\operatorname{argmin}} G(\lambda, \theta)$$

Regularized Logistic loss:

$$G(\lambda, \theta) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-y_i \langle \theta, X_i \rangle}) + \lambda \|\theta\|_2^2$$

Training the model:

$$\theta^*(\lambda) = \underset{\theta}{\operatorname{argmin}} G(\lambda, \theta)$$

⇒ Choose λ using validation data.

Evaluating the generalization

We want to find λ that ensure the best *generalization* of $\theta^*(\lambda)$.

Validation loss: use held out data $(X_i^{val}, y_i^{val})_{i=1}^M$

$$F(\theta) = \frac{1}{M} \sum_{i=1}^M \log(1 + e^{-y_i^{val} \langle \theta, X_i^{val} \rangle})$$

Independent estimate of the risk of the model.

Evaluating the generalization

We want to find λ that ensure the best *generalization* of $\theta^*(\lambda)$.

Validation loss: use held out data $(X_i^{val}, y_i^{val})_{i=1}^M$

$$F(\theta) = \frac{1}{M} \sum_{i=1}^M \log(1 + e^{-y_i^{val} \langle \theta, X_i^{val} \rangle})$$

Independent estimate of the risk of the model.

\Rightarrow Find λ that gives a model $\theta^*(\lambda)$ with a good validation loss.

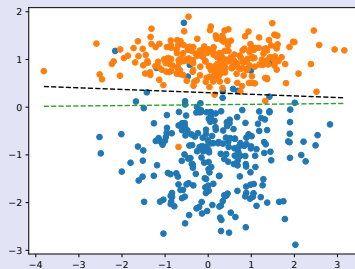
Mathematical rewriting:
$$\begin{cases} \min_{\lambda} F(\theta^*(\lambda)) \\ \text{s.t. } \theta^*(\lambda) = \operatorname{argmin}_{\theta} G(\lambda, \theta) \end{cases}$$

Bi-level optimization problems: Model selection

Selecting the best model:

- ▶ G is the training loss and θ are the parameters of the model.
- ▶ Select the hyper-parameter λ to get the best validation loss F .

Hyperparameter optimization: λ is a regularization parameter:

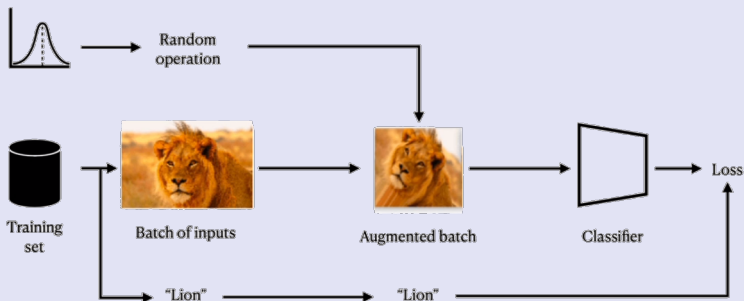


Bi-level optimization problems: Model selection

Selecting the best model:

- ▶ G is the training loss and θ are the parameters of the model.
- ▶ Select the hyper-parameter λ to get the best validation loss F .

Data augmentation: λ parametrizes the transformations distribution.

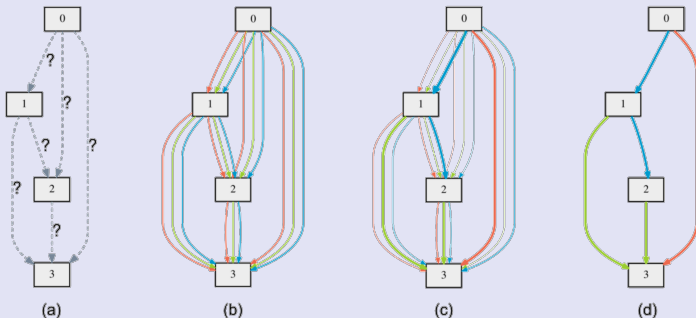


Bi-level optimization problems: Model selection

Selecting the best model:

- ▶ G is the training loss and θ are the parameters of the model.
- ▶ Select the hyper-parameter λ to get the best validation loss F .

Neural Architecture Search: λ parametrizes the architecture.



Black box methods: Take $\{\lambda_k\}_k$ and compute $\min_k h(\lambda_k)$

- ▶ Grid-Search
- ▶ Random-Search
- ▶ Bayesian-Optimization

⇒ Do not scale well with the dimension

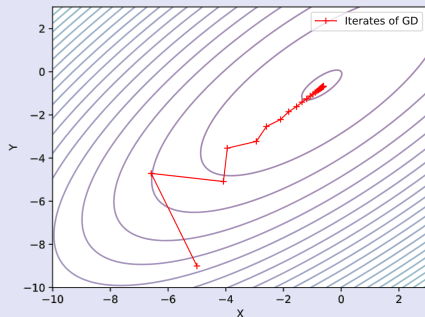
Solving bi-level optimization

First order methods: Gradient descent on h

Iterate in the steepest direction:

$$\lambda^{t+1} = \lambda^t - \rho^t \nabla h(\lambda)$$

- ▶ Gradient $\nabla h(\lambda) = \frac{d F(\lambda, \theta^*(\lambda))}{d \lambda}$
- ▶ Step size ρ^t .



Value function definition:

$$h(\lambda) = F(\lambda, \theta^*(\lambda))$$

Chain rule:

$$\nabla_{\lambda} h(\lambda) = \nabla_1 F(\lambda, \theta^*(\lambda)) + (d\theta^*(\lambda))^T \nabla_2 F(\lambda, \theta^*(\lambda))$$

Optimality condition for θ^*

$$\nabla_2 G(\lambda, \theta^*(\lambda)) = 0$$

Optimality condition for θ^*

$$\nabla_2 G(\lambda, \theta^*(\lambda)) = 0$$

Derivating this equation relative to λ gives:

$$\nabla_{22}^2 G(\lambda, \theta^*(\lambda)) d\theta^*(\lambda) + \nabla_{21}^2 G(\lambda, \theta^*(\lambda)) = 0,$$

Optimality condition for θ^*

$$\nabla_2 G(\lambda, \theta^*(\lambda)) = 0$$

Derivating this equation relative to λ gives:

$$\nabla_{22}^2 G(\lambda, \theta^*(\lambda)) d\theta^*(\lambda) + \nabla_{21}^2 G(\lambda, \theta^*(\lambda)) = 0,$$

Implicit function theorem

$$d\theta^*(\lambda) = -[\nabla_{22}^2 G(\lambda, \theta^*(\lambda))]^{-1} \nabla_{21}^2 G(\lambda, \theta^*(\lambda)),$$

Value function gradient:

$$\nabla h(\lambda) = \nabla_1 F(\lambda, \theta^*) - \nabla_{21}^2 G(\lambda, \theta^*) [\nabla_{22}^2 G(\lambda, \theta^*)]^{-1} \nabla_2 F(\lambda, \theta^*)$$

Value function gradient:

$$\nabla h(\lambda) = \nabla_1 F(\lambda, \theta^*) - \nabla_{21}^2 G(\lambda, \theta^*) [\nabla_{22}^2 G(\lambda, \theta^*)]^{-1} \nabla_2 F(\lambda, \theta^*)$$

- ▶ Need to compute the solution of the inner

Value function gradient:

$$\nabla h(\lambda) = \nabla_1 F(\lambda, \theta^*) - \nabla_{21}^2 G(\lambda, \theta^*) [\nabla_{22}^2 G(\lambda, \theta^*)]^{-1} \nabla_2 F(\lambda, \theta^*)$$

- ▶ Need to compute the solution of the inner
- ▶ Need to solve a $p \times p$ linear system

$$v^*(\lambda) = [\nabla_{22}^2 G(\lambda, \theta^*)]^{-1} \nabla_2 F(\lambda, \theta^*)$$

Approximate bi-level optimization

References

- ▶ Pedregosa, F. (2016). [Hyperparameter optimization with approximate gradient](#). In *International Conference on Machine Learning (ICML)*, pages 737–746, New-York, NY, USA

Hyperparameter optimization with Approximate Gradient HOAG

[Pedregosa 2016]

Do we need to compute θ^ and v^* precisely?*

Idea: Approximate $\theta^*(\lambda^t)$ and $v^*(\lambda^t) = [\nabla_{22}^2 G(\lambda^t, \theta^*)]^{-1} \nabla_2 F(\lambda^t, \theta^*)$

Hyperparameter optimization with Approximate Gradient HOAG

[Pedregosa 2016]

Do we need to compute θ^ and v^* precisely?*

Idea: Approximate $\theta^*(\lambda^t)$ and $v^*(\lambda^t) = [\nabla_{\theta\theta}^2 G(\lambda^t, \theta^*)]^{-1} \nabla_{\theta} F(\lambda^t, \theta^*)$

- ▶ Compute θ^t such that $\|\theta^t - \theta^*(\lambda^t)\|_2 \leq \epsilon_t$,
iterative solver e.g. L-BFGS
- ▶ Compute v^t such that $\|\frac{\partial^2 G}{\partial \theta^2}(\lambda^t, \theta^t)v^t + \frac{\partial F}{\partial \theta}(\lambda^t, \theta^t)\|_2 \leq \epsilon_t$,
L-BFGS or CG
- ▶ Compute the approximate gradient $g_t = \frac{\partial F}{\partial \lambda}(\lambda^t, \theta^t) + \frac{\partial^2 G}{\partial \theta \partial \lambda}(\lambda^t, \theta^t)v^t$
- ▶ Update the outer variable $\lambda^{t+1} = \lambda^t - \rho^t g^t$

Theorem: If $\sum_t \epsilon_t < \infty$ and the step-sizes are chosen appropriately, then the algorithm converges to a stationary point *i.e.*

$$\|\nabla h(\lambda^t)\|_2 \rightarrow 0 .$$

Further linear system approximation v^*

Linear system solution $v^*(\lambda^t)$ is a by product.

⇒ Avoid computing it as much as possible.

Proposed Methods:

▶ L-BFGS

▶ Conjugate Gradient

▶ Jacobian-Free method

▶ Neumann iterations

$$\nabla_{22}^2 G(\lambda^t, \theta^t) \approx Id$$

$$\nabla_{22}^2 G(\lambda^t, \theta^t)^{-1} \approx \sum_k (Id - \nabla_{22}^2 G(\lambda^t, \theta^t))^k$$

▶ Algorithm unrolling

▶ Use Quasi-newton hessian approximation

[Pedregosa 2016, Lorraine et al. 2020, Luketina et al. 2016, Ramzi et al. 2022]

SHINE - Sharing the INverse Estimate

References

- ▶ Ramzi, Z., Mannel, F., Bai, S., Starck, J.-L., Ciuciu, P., and Moreau, T. (2022). [SHINE: SHaring the INverse Estimate from the forward pass for bi-level optimization and implicit models](#). In *International Conference on Learning Representations (ICLR)*, online

Quasi Newton 101:

Solving $\theta^* = \operatorname{argmin}_{\theta} G(\theta)$

Newton Method

$$\theta^{t+1} = \theta^t - [\nabla^2 G(\theta^t)]^{-1} \nabla G(\theta^t)$$

Quasi-Newton Method

$$\theta^{t+1} = \theta^t - B_t^{-1} \nabla G(\theta^t)$$

B_t : low-rank approx. of $\nabla^2 G(\theta^t)$.

Inverse with Sherman-Morrison

Quasi Newton 101:

Solving $\theta^* = \operatorname{argmin}_{\theta} G(\theta)$

Newton Method

$$\theta^{t+1} = \theta^t - [\nabla^2 G(\theta^t)]^{-1} \nabla G(\theta^t)$$

Quasi-Newton Method

$$\theta^{t+1} = \theta^t - B_t^{-1} \nabla G(\theta^t)$$

B_t : low-rank approx. of $\nabla^2 G(\theta^t)$.

Inverse with Sherman-Morrison

\Rightarrow The Hessian for v^* is the same as the one from the inner problem.

Idea: reuse the approximation of the Hessian B_t computed by L-BFGS for the inner problem.

$$\begin{cases} \tilde{v}_t = B_t^{-1} \nabla_2 F(\lambda, \theta^t) \\ \tilde{\nabla} h(\lambda) = \nabla_1 F(\lambda, \theta^t) + \nabla_{12}^2 G(\lambda, \theta^t) \tilde{v}_t \end{cases}$$

Properties of B :

Idea: reuse the approximation of the Hessian B_t computed by L-BFGS for the inner problem.

$$\begin{cases} \tilde{v}_t = B_t^{-1} \nabla_2 F(\lambda, \theta^t) \\ \tilde{\nabla} h(\lambda) = \nabla_1 F(\lambda, \theta^t) + \nabla_{12}^2 G(\lambda, \theta^t) \tilde{v}_t \end{cases}$$

Properties of B :

- ▶ It is computed when solving $\theta^* = \operatorname{argmin}_{\theta} G(\theta)$ using a quasi-Newton method.

Idea: reuse the approximation of the Hessian B_t computed by L-BFGS for the inner problem.

$$\begin{cases} \tilde{v}_t = B_t^{-1} \nabla_2 F(\lambda, \theta^t) \\ \tilde{\nabla} h(\lambda) = \nabla_1 F(\lambda, \theta^t) + \nabla_{12}^2 G(\lambda, \theta^t) \tilde{v}_t \end{cases}$$

Properties of B :

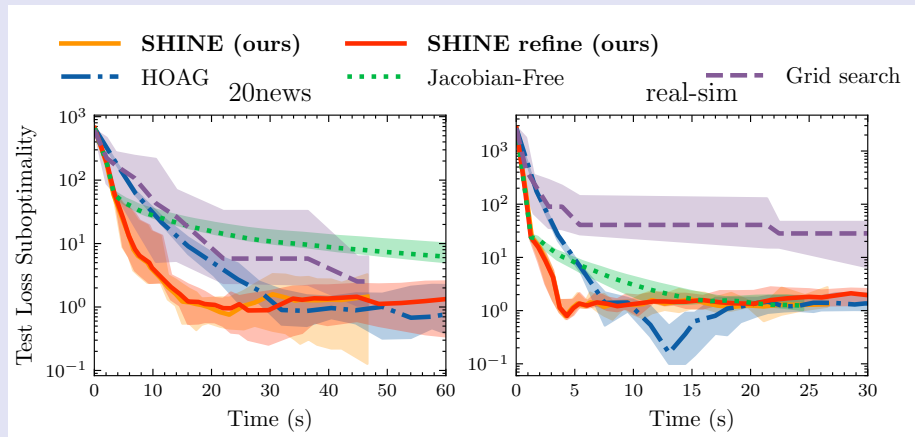
- ▶ It is computed when solving $\theta^* = \operatorname{argmin}_{\theta} G(\theta)$ using a quasi-Newton method.
- ▶ It is easily invertible using the Sherman-Morrison formula, because low-rank.

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 λ , (θ^t) converges q -superlinearly to θ^ and*

$$\lim_{t \rightarrow \infty} \nabla_1 F(\lambda, \theta^t) + \nabla_{12}^2 G(\lambda, \theta^t) \tilde{v}_t = \nabla h(\lambda).$$

Logistic Regression with ℓ_2 -regularisation on 2 datasets:



Application to Deep Equilibrium Networks (DEQs)

A recurrent expression of classical, explicit networks:

$$\theta_n = f_{\lambda_n}(\theta_{n-1}), \quad \forall n < N$$

Infinite depth neural networks

A recurrent expression of classical, explicit networks:

$$\theta_n = f_{\lambda_n}(\theta_{n-1}), \quad \forall n < N$$

What if $N \rightarrow \infty$?

If we suppose $\lambda_n = \lambda, \forall n$:

$$\theta^* = f_{\lambda}(\theta^*)$$

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

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

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

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

This approximates an infinite depth network:

$$\theta_n = f_\lambda(\theta_{n-1}), \quad \forall n \rightarrow \infty$$

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

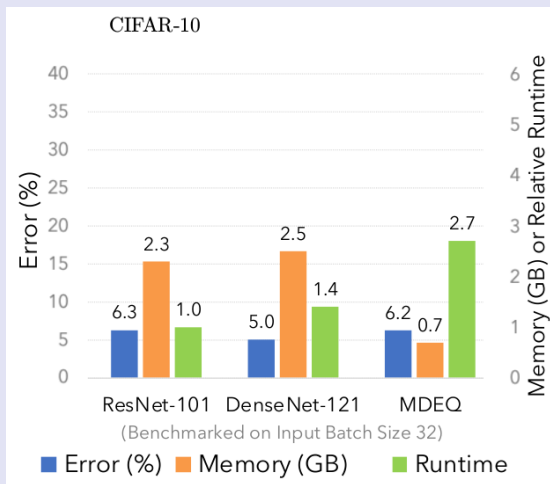
Deep Equilibrium Network:

$$\begin{cases} \min_{\lambda} h(\lambda) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, \theta^*(X_i, \lambda)) \\ \text{s.t. } \theta^*(X_i, \lambda) = g_{\lambda}(\theta^*(X_i, \lambda)) \end{cases}$$

Output of the network is the root of $G(\theta, \lambda) = \theta - g_{\lambda}(\theta) = 0$.

In practice, use quasi-Newton algorithm such as the Broyden method for both finding θ^* and its derivative v^* .

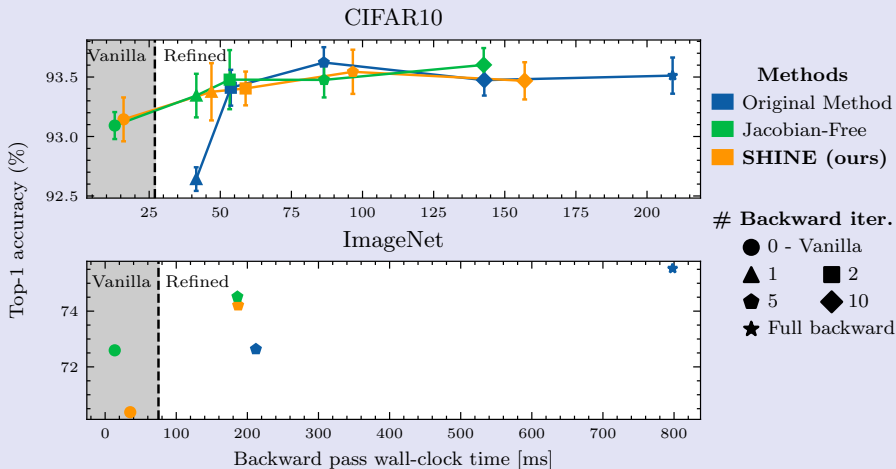
Deep Equilibrium Network



► Efficient memory

► Slow runtime

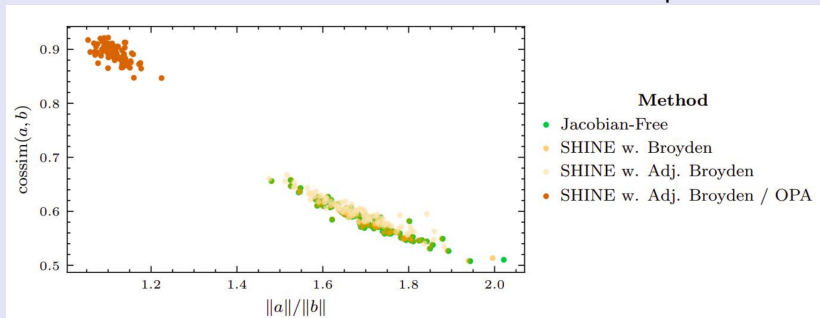
Multiscale DEQ on CIFAR10:



OPA - Outer Problem Awareness

B^{-1} is not a uniformly good approximation.

OPA: add additional secant conditions for B update.

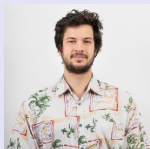


⇒ Better gradient approximation (theoretical and empirical)

However, this does not improve results for test error.

Conclusion

- ▶ Bi-level optimization is intrinsic in many ML problems.
- ▶ We propose to re-use by-product from θ^* computation to make it easy to get v^* .
- ▶ Good results for HO but still open questions for DEQs.



Z. Ramzi



S. Bai



F. Mannel





J.L. Starck



P. Ciuciu

Slides will be on my web page:

 [tommoral.github.io](https://github.com/tommoral)

 [@tomamoral](https://twitter.com/tomamoral)