

# TransferLab Seminar

## Second-Order Information and Applications

---

Kristof Schröder

11.04.2024

appliedAI Institute for Europe

Introduction

LLM Pre-Training

LLM Fine-Tuning

Physics Informed Neural Networks (PINNs)

Influence Functions

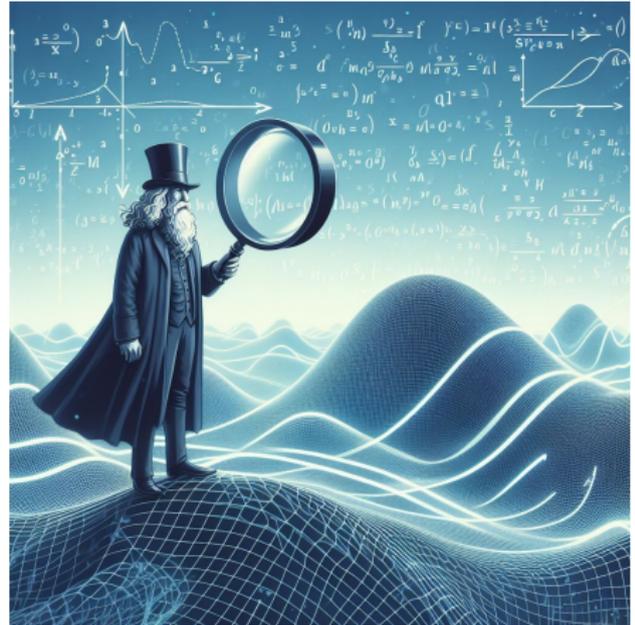
Software

# Introduction

---

# Why Second-Order?

- Encodes geometry/curvature information about the optimization objective
- Can reduce number of iterations
- Can improve the quality of the optimization result



Created with  
<https://www.bing.com/chat>

# Gradient descent

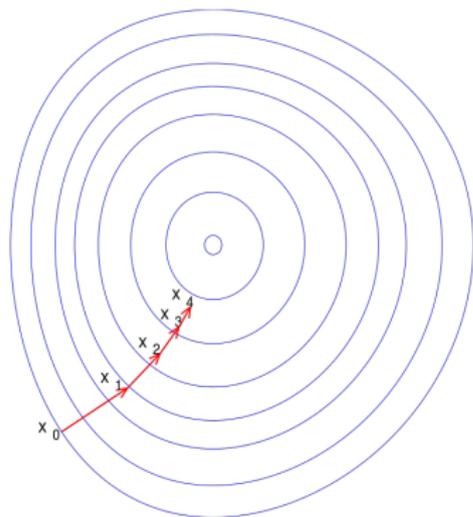
In order to solve (for  $\mathcal{L}$  smooth, strongly convex)

$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta),$$

iterate

$$\theta_{t+1} = \theta_t - \gamma_t \nabla_{\theta} \mathcal{L}(\theta_t).$$

- convergence rate is linear
- convergence depends on the condition number of the Hessian at the solution



Source: Wikipedia

# Newton's Method in Optimization

In order to solve

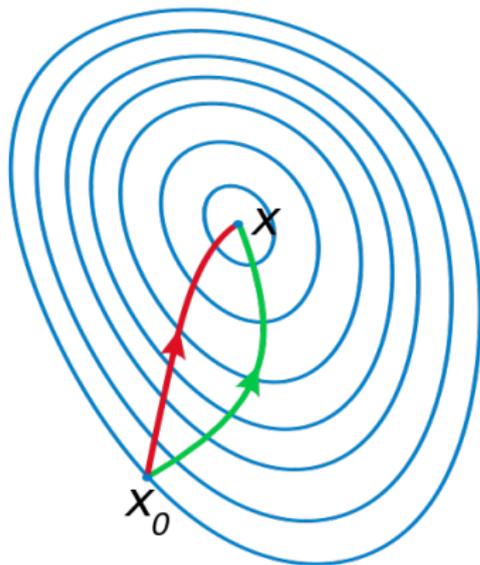
$$\min_{\theta \in \mathbb{R}^d} \mathcal{L}(\theta),$$

iterate

$$\theta_{t+1} = \theta_t - H_{\mathcal{L}}(\theta_t)^{-1} \nabla_{\theta} \mathcal{L}(\theta_t)$$

where  $H_{\mathcal{L}}(\theta_t) = \nabla^2_{\theta} \mathcal{L}(\theta_t)$ .

- quadratic convergence rate
- think of  $H_{\mathcal{L}}(\theta_t)^{-1}$  as a preconditioner, i.e.  $\text{cond}(H_{\mathcal{L}}(\theta_t)^{-1} H_{\mathcal{L}}(\theta^*))$  is small



Source: Wikipedia

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

**Problems:**

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

**Problems:**

- $\mathcal{L}$  is non-convex, Newton could even converge to a maximum

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

## Problems:

- $\mathcal{L}$  is non-convex, Newton could even converge to a maximum
- loss-landscape is highly ill-conditioned, i.e. very heterogeneous curvature

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

## Problems:

- $\mathcal{L}$  is non-convex, Newton could even converge to a maximum
- loss-landscape is highly ill-conditioned, i.e. very heterogeneous curvature
- memory cost is  $\mathcal{O}(d^2)$  and naive inversion has complexity  $\mathcal{O}(d^3)$

In this case

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \ell(f(x_i; \theta), y_i)$$

where  $f(x, \theta)$  is a parametrized model and  $\ell$  is some loss function.

## Problems:

- $\mathcal{L}$  is non-convex, Newton could even converge to a maximum
- loss-landscape is highly ill-conditioned, i.e. very heterogeneous curvature
- memory cost is  $\mathcal{O}(d^2)$  and naive inversion has complexity  $\mathcal{O}(d^3)$



Approximate the preconditioned gradients

$$H_{\mathcal{L}}(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

What to do?

Approximate the preconditioned gradients

$$H_{\mathcal{L}}(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

What to do?

- Implicit inversion, block approximation, diagonal approximation,

Approximate the preconditioned gradients

$$H_{\mathcal{L}}(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

What to do?

- Implicit inversion, block approximation, diagonal approximation,
- Randomization: stochastic estimators + smoothing (EMA)

Approximate the preconditioned gradients

$$H_{\mathcal{L}}(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

What to do?

- Implicit inversion, block approximation, diagonal approximation,
- Randomization: stochastic estimators + smoothing (EMA)
- Low-rank approximations

Approximate the preconditioned gradients

$$H_{\mathcal{L}}(\theta)^{-1} \nabla_{\theta} \mathcal{L}(\theta)$$

What to do?

- Implicit inversion, block approximation, diagonal approximation,
- Randomization: stochastic estimators + smoothing (EMA)
- Low-rank approximations
- Decomposition based approximations: Gauss-Newton approximation, Fisher information matrix (Natural Gradient Method)

# LLM Pre-Training

---

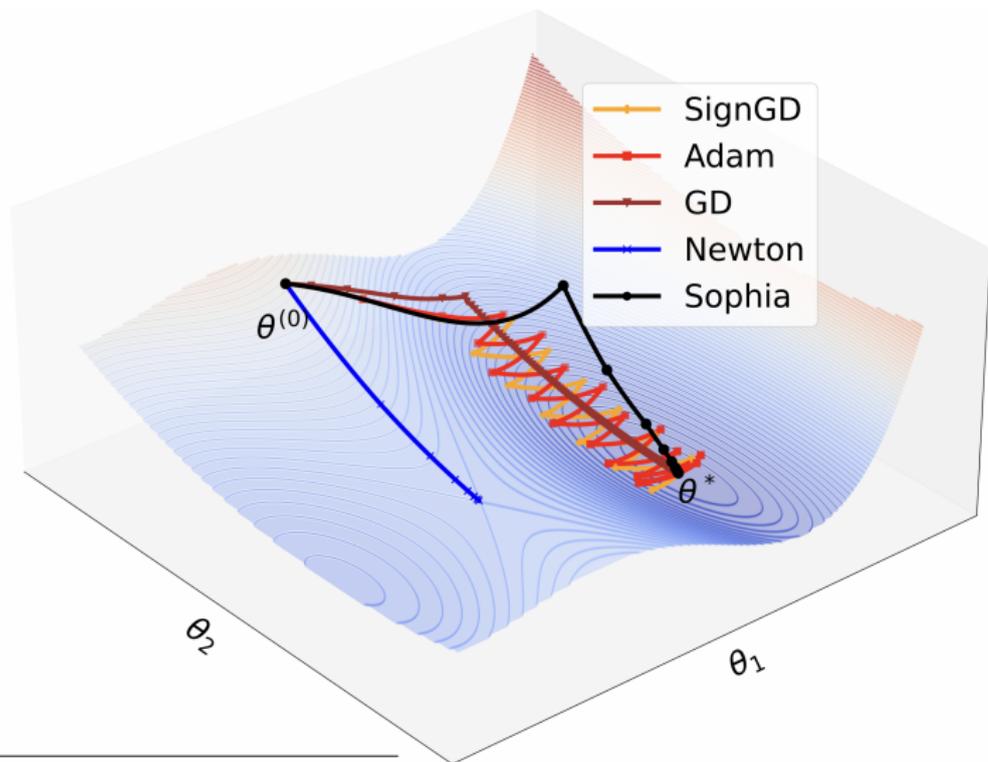
Pre-training of a large language model typically incur significant expenses. Often the chosen state of the state-of-the-art solver is **Adam**.

Pre-training of a large language model typically incur significant expenses. Often the chosen state of the state-of-the-art solver is **Adam**.

How to incorporate second-order information in order to

- reduce the number of iterations needed,
- only adding a small over-head per iteration,
- keep memory footprint comparable?

# Motivating Example



1

<sup>1</sup>Liu et al., Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, 2023 [LLH<sup>+</sup>23]

# Sophia: Second-order Clipped Stochastic Optimization<sup>2</sup>

- Construct stochastic estimators for the diagonal of the Hessian.
- Exponential smoothing of minibatch gradients at each iteration.
- Exponential smoothing of the second-order information every  $k = 10$  iteration.
- Per-coordinate clipping to handle negative curvature.

---

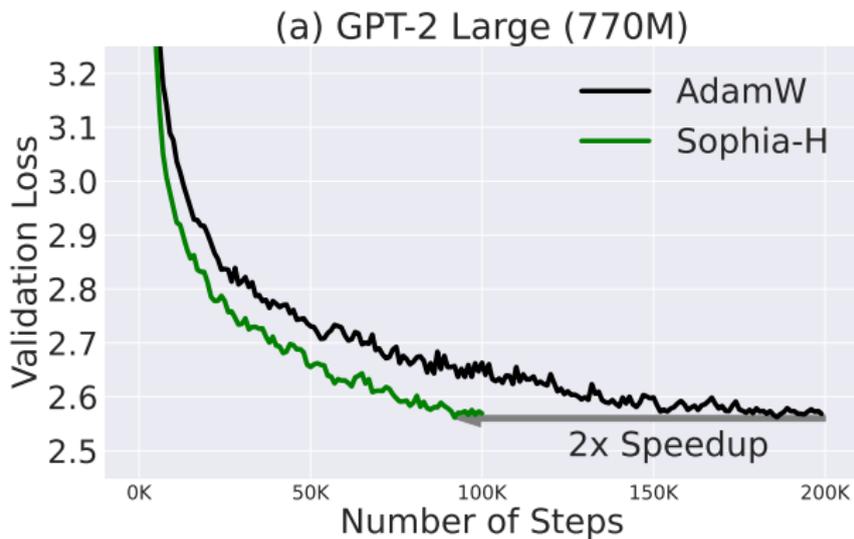
## Algorithm 3 Sophia

---

```
1: Input:  $\theta_1$ , learning rate  $\{\eta_t\}_{t=1}^T$ , hyperparameters  $\lambda, \beta_1, \beta_2, \epsilon$ , and estimator choice Estimator  $\in \{\text{Hutchinson}, \text{Gauss-Newton-Bartlett}\}$   
2: Set  $m_0 = 0, v_0 = 0, h_{1-k} = 0$   
3: for  $t = 1$  to  $T$  do  
4:   Compute minibatch loss  $L_t(\theta_t)$ .  
5:   Compute  $g_t = \nabla L_t(\theta_t)$ .  
6:    $m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$   
7:   if  $t \bmod k = 1$  then  
8:     Compute  $\hat{h}_t = \text{Estimator}(\theta_t)$ .  
9:      $h_t = \beta_2 h_{t-k} + (1 - \beta_2) \hat{h}_t$   
10:  else  
11:     $h_t = h_{t-1}$   
12:     $\theta_t = \theta_t - \eta_t \lambda \theta_t$  (weight decay)  
13:     $\theta_{t+1} = \theta_t - \eta_t \cdot \text{clip}(m_t / \max\{h_t, \epsilon\}, \rho)$ 
```

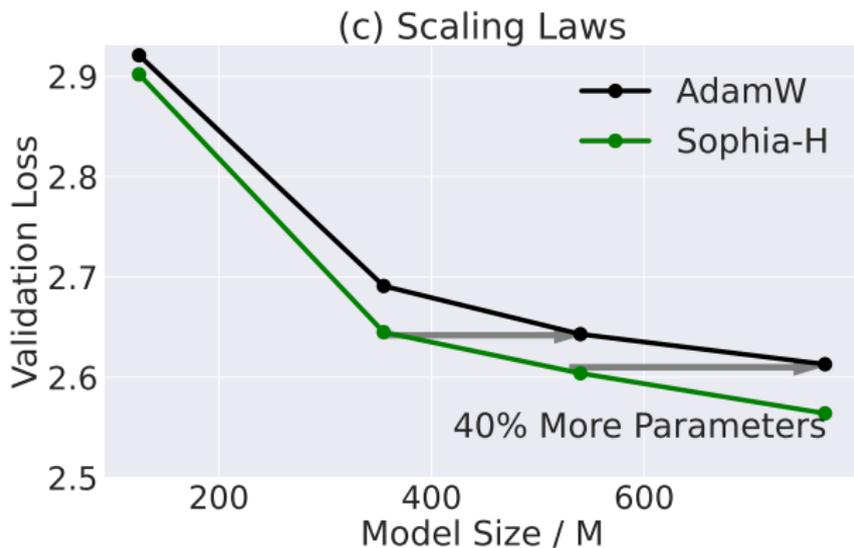
---

<sup>2</sup>Liu et al., Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, 2023 [LLH<sup>+</sup>23]



**Figure 4:** Loss evolution for training GPT-2 on [OpenWebText](#)<sup>3</sup>

<sup>3</sup>Liu et al., Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, 2023 [LLH<sup>+</sup>23]



**Figure 5:** Validation loss vs. number of parameters<sup>4</sup>

<sup>4</sup>Liu et al., Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, 2023 [LLH<sup>+</sup>23]

- same memory cost as AdamW,
- overall wall-clock time overhead less than 5%
- halving the number of iteration results in almost halving the total wall-clock time

Table 1: Wall-clock time and compute.

| Algorithm | Model Size | T(step) | T(Hessian) | Compute |
|-----------|------------|---------|------------|---------|
| AdamW     | 770M       | 3.25s   | –          | 2550    |
| Sophia-H  | 770M       | 3.40s   | 0.12s      | 2708    |
| Sophia-G  | 770M       | 3.42s   | 0.17s      | 2678    |
| AdamW     | 355M       | 1.77s   | –          | 1195    |
| Sophia-H  | 355M       | 1.88s   | 0.09s      | 1249    |
| Sophia-G  | 355M       | 1.86s   | 0.09s      | 1255    |

**Figure 6:** Computation Time<sup>a</sup>

<sup>a</sup>Liu et al., Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, 2023 [LLH<sup>+</sup>23]

# LLM Fine-Tuning

---

- Fine-tuning language models has been effective for various tasks, but the memory demands of backpropagation in large models are still significant.
- The limitations of available compute environments typically constrain the fine-tuning.
- Recently, zero-order methods have been adapted to fine-tuning LLMs<sup>5</sup>(with up to 12x reduced memory consumption), but are known to converge slowly.

---

<sup>5</sup>Malladi et al., Fine-Tuning Language Models with Just Forward Passes, 2023  
[MGN<sup>+</sup>23]

### Definition (Simultaneous Perturbation Stochastic Approximation or SPSA)

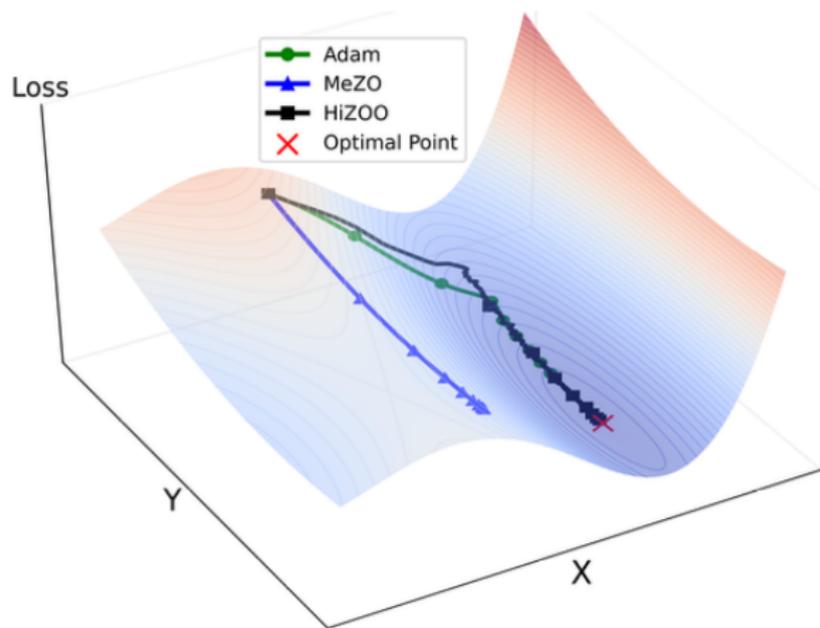
For  $\theta \in \mathbb{R}^d$  and loss function  $L$  estimate the gradient on a mini-batch  $B$  as

$$g_\varepsilon(\theta) = \frac{L(\theta + \varepsilon z, B) - L(\theta - \varepsilon z, B)}{2\varepsilon} \cdot z \approx zz^T \nabla_\theta L(\theta, B),$$

for  $z \sim \mathcal{N}(0, \text{Id}_d)$ .

- Requires two forward passes per mini-batch, no back-propagation.
- Can be averaged over several samplings from  $\mathcal{N}(0, \text{Id}_d)$  ( $n$ -SPSA).

# Motivation



**Figure 7:** Toy example in figure 1<sup>6</sup>

<sup>6</sup>Zhao et al., Second-Order Fine-Tuning without Pain for LLMs: A Hessian Informed Zeroth-Order Optimizer, 2024 [ZDY<sup>+</sup>24]

# Hessian Informed Zero-Order Optimization (HiZOO)

For  $\theta \in \mathbb{R}^d$  and loss function  $L$  estimate the gradient on a mini-batch  $B$  as

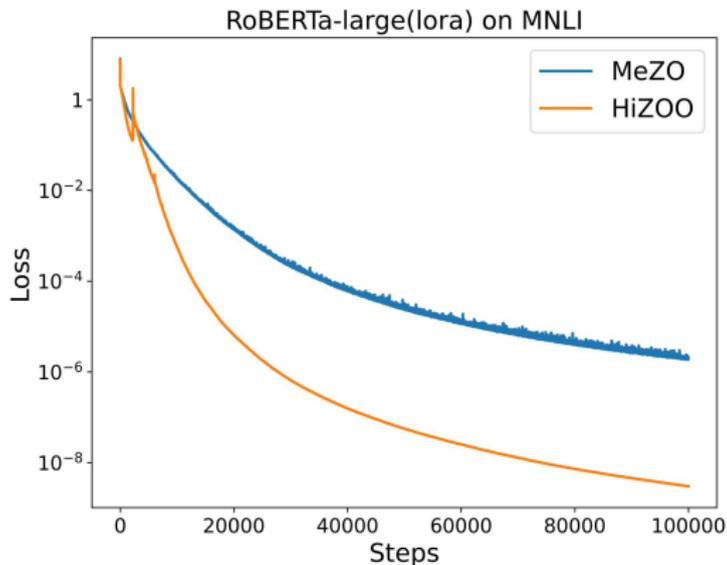
$$g_\epsilon(\theta) = \frac{L(\theta + \epsilon \Sigma_t^{1/2} z, B) - L(\theta - \epsilon \Sigma_t^{1/2} z, B)}{2\epsilon} \cdot \Sigma_t^{1/2} z,$$

for  $z \sim \mathcal{N}(0, \text{Id}_d)$  and  $\Sigma_t$  is an approximation to the diagonal of the inverse Hessian matrix, which gets updated simultaneously during iteration.

Overall computation requires three forward passes of the model<sup>7</sup>.

---

<sup>7</sup>Zhao et al., Second-Order Fine-Tuning without Pain for LLMs: A Hessian Informed Zeroth-Order Optimizer, 2024 [ZDY<sup>+</sup>24]



**Figure 8:** Loss evolution for LoRA training RoBERTa on the MultiNLI dataset<sup>8</sup>, see figure 3<sup>9</sup>

<sup>8</sup><https://paperswithcode.com/dataset/multinli>

<sup>9</sup>Zhao et al., Second-Order Fine-Tuning without Pain for LLMs: A Hessian Informed Zeroth-Order Optimizer, 2024 [ZDY<sup>+</sup>24]

# Physics Informed Neural Networks (PINNs)

---

Consider a partial differential equation of the form

$$\mathcal{D}[u(x), x] = 0, \quad x \in \Omega$$

$$\mathcal{B}[u(x), x] = 0, \quad x \in \partial\Omega$$

where  $\Omega \subseteq \mathbb{R}^d$ ,  $\mathcal{D}$  is a differential operator and  $\mathcal{B}$  is the boundary value operator.

## Physics Informed Neural Networks (PINNs)

Consider a partial differential equation of the form

$$\begin{aligned}\mathcal{D}[u(x), x] &= 0, & x \in \Omega \\ \mathcal{B}[u(x), x] &= 0, & x \in \partial\Omega\end{aligned}$$

where  $\Omega \subseteq \mathbb{R}^d$ ,  $\mathcal{D}$  is a differential operator and  $\mathcal{B}$  is the boundary value operator. Let  $u(x; \theta)$  be a neural network and minimize

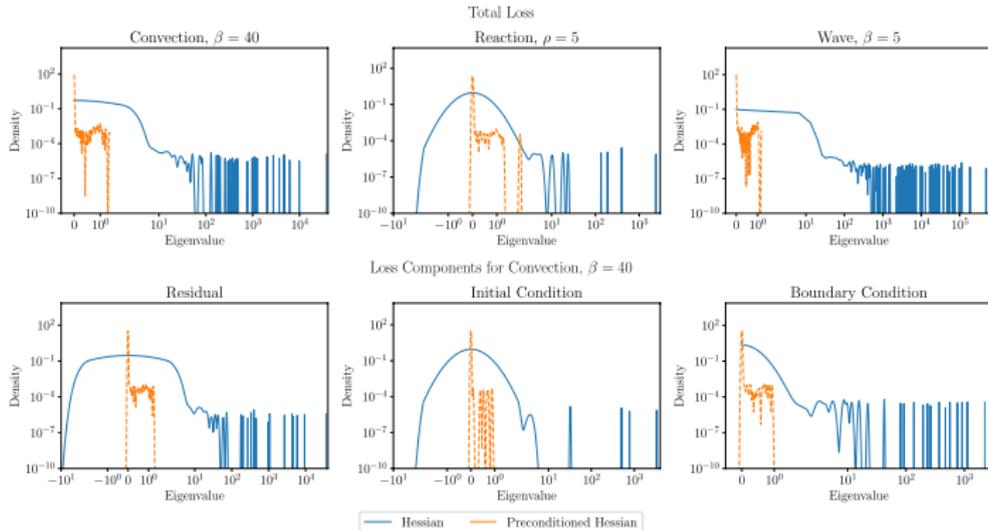
$$\begin{aligned}\mathcal{L}(\theta) &= \frac{1}{2n_{res}} \sum_{i=1}^{n_{res}} (\mathcal{D}[u(x_i^{res}; \theta), x_i^{res}])^2 \\ &+ \frac{1}{2n_{bc}} \sum_{i=1}^{n_{bc}} (\mathcal{B}[u(x_i^{bc}; \theta), x_i^{bc}])^2\end{aligned}$$

- PINNs must be trained to near-zero loss to obtain an adequate solution (in  $\ell_2$  sense) to the PDE.
- The loss-landscape is ill-conditioned, i.e. the Hessian  $H_{\mathcal{L}}(\theta)$  has a large condition number. In other words, the loss is very steep in some directions and very flat in others.
- Pre-conditioning with second-order information improves the conditioning significantly.

---

<sup>10</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

# Spectral Density



**Figure 9:** Optimizing with Adam + L-BFGS, Hessian spectral density after 41k iterations with and without preconditioning with quasi-Newton matrix<sup>11</sup>

<sup>11</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

**Motivation:** L-BFGS may stop early and leaves the loss under-optimized (see Section 7.1<sup>12</sup>).

---

<sup>12</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

<sup>13</sup>Frangella, et al., Randomized Nyström Preconditioning, 2023 [FTU23]

**Motivation:** L-BFGS may stop early and leaves the loss under-optimized (see Section 7.1<sup>12</sup>).

**Idea:** Use a different approximate Newton update step, which allow for further progress.

---

<sup>12</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

<sup>13</sup>Frangella, et al., Randomized Nyström Preconditioning, 2023 [FTU23]

**Motivation:** L-BFGS may stop early and leaves the loss under-optimized (see Section 7.1<sup>12</sup>).

**Idea:** Use a different approximate Newton update step, which allow for further progress.

Introduce a positive-definite rank- $r$  approximation (Nyström approximation)

$$H_{nys} = (H_{\mathcal{L}}S)(S^T H_{\mathcal{L}}S)^{\dagger}(H_{\mathcal{L}}S)^T, \quad S \in \mathbb{R}^{d \times r} \text{ standard normal,}$$

of the Hessian as preconditioner and use Conjugate Gradient. <sup>13</sup>

---

<sup>12</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

<sup>13</sup>Frangella, et al., Randomized Nyström Preconditioning, 2023 [FTU23]

---

**Algorithm 4** NysNewton-CG (NnCG)

---

**input** Initialization  $w_0$ , max. learning rate  $\eta$ , number of iterations  $K$ , preconditioner sketch size  $s$ , preconditioner update frequency  $F$ , damping parameter  $\mu$ , CG tolerance  $\epsilon$ , CG max. iterations  $M$ , backtracking parameters  $\alpha, \beta$

$d_{-1} = 0$

**for**  $k = 0, \dots, K - 1$  **do**

**if**  $k$  is a multiple of  $F$  **then**

$[U, \hat{\Lambda}] = \text{RandomizedNyströmApproximation}(H_L(w_k), s)$    ▷ Update Nyström preconditioner every  $F$  iterations

**end if**

$d_k = \text{NyströmPCG}(H_L(w_k), \nabla L(w_k), d_{k-1}, U, \hat{\Lambda}, s, \mu, \epsilon, M)$    ▷ Damped Newton step  $(H_L(w_k) + \mu I)^{-1} \nabla L(w_k)$

$\eta_k = \text{Armijo}(L, w_k, \nabla L(w_k), -d_k, \eta)$    ▷ Compute step size via line search

$w_{k+1} = w_k - \eta_k d_k$    ▷ Update parameters

**end for**

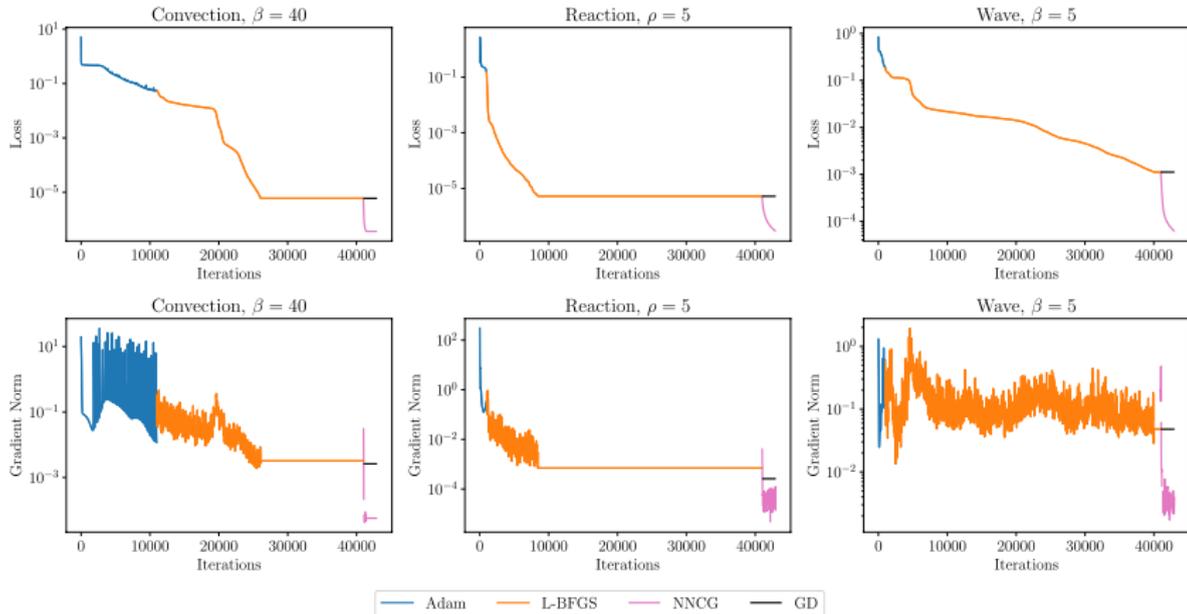
---

**Figure 10:** NysNewton-CG optimizer<sup>14</sup>

---

<sup>14</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

# NysNewton-CG (NNCG)



**Figure 11:** Loss evolution Adam + L-BFGS + NNCG<sup>15</sup>

<sup>15</sup>Rathore et al., Challenges in Training PINNs: A Loss Landscape Perspective, 2024 [RLF<sup>+</sup>24]

# Influence Functions

---

How to estimate the influence of single training points on model parameters or model output on test points?

How to estimate the influence of single training points on model parameters or model output on test points?

**Definition (Influence Function)**

$$\mathcal{I}(z_t, z) = \nabla_{\theta} \ell(z_t, \theta)^T (H_{\theta} + \lambda I_d)^{-1} \nabla_{\theta} \ell(z, \theta)$$

How to estimate the influence of single training points on model parameters or model output on test points?

## Definition (Influence Function)

$$\mathcal{I}(z_t, z) = \nabla_{\theta} \ell(z_t, \theta)^T (H_{\theta} + \lambda I_d)^{-1} \nabla_{\theta} \ell(z, \theta)$$

Again, the bottleneck is the computation of the preconditioned gradients

$$(H_{\theta} + \lambda I_d)^{-1} \nabla_{\theta} \ell(z, \theta)$$

## Recent Implementations

- EKFAC (Eigenvalue Corrected Kronecker Factorization)[GBA<sup>+</sup>23]:
  - Block inversion (per layer), Gauss-Newton approximation, Kronecker factorization
  - Implementation: [pyDVL](#)
  - [Seminar talk](#), [paper pill](#)
- Nyström Preconditioned CG[FTU23]
  - Randomized low-rank approximation as preconditioner, [Woodbury matrix identity](#)
  - Implementation: [pyDVL](#), implemented for the next minor release
- DataInf[KWWZ23]:
  - Block inversion (per layer), Gauss-Newton approximation, harmonic mean estimator, [Woodbury matrix identity](#)
  - Implementation: planned for next major release, [github issue](#)
  - [paper pill](#)

**Software**

---

- Sophia: Second-order Clipped Stochastic Optimization[LLH<sup>+</sup>23]
  - <https://github.com/Liuhong99/Sophia>
  - Torch optimizer implementation, no package
  - [paper](#) [pill](#)
- MeZO: Memory-efficient Zeroth-order[MGN<sup>+</sup>23]
  - <https://github.com/princeton-nlp/MeZO>
  - HuggingFace trainer implementation, no package
- HiZOO: Hessian informed zeroth-order optimization[ZDY<sup>+</sup>24]
  - <https://anonymous.4open.science/r/HiZOO-27F8>
  - HuggingFace trainer implementation, no package
- NysNewton-CG[RLF<sup>+</sup>24]
  - [https://anonymous.4open.science/r/opt\\_for\\_pinns-9246](https://anonymous.4open.science/r/opt_for_pinns-9246)
  - Torch optimizer implementation, no package

**Thank you!**

-  Zachary Frangella, Joel A. Tropp, and Madeleine Udell.  
**Randomized Nyström Preconditioning.**  
*SIAM Journal on Matrix Analysis and Applications*,  
44(2):718–752, June 2023.
-  Roger Grosse, Juhan Bae, Cem Anil, Nelson Elhage, Alex Tamkin, Amirhossein Tajdini, Benoit Steiner, Dustin Li, Esin Durmus, Ethan Perez, Evan Hubinger, Kamilė Lukošiuūtė, Karina Nguyen, Nicholas Joseph, Sam McCandlish, Jared Kaplan, and Samuel R. Bowman.  
**Studying Large Language Model Generalization with Influence Functions, August 2023.**

-  Yongchan Kwon, Eric Wu, Kevin Wu, and James Zou.  
**DataInf: Efficiently Estimating Data Influence in LoRA-tuned LLMs and Diffusion Models, October 2023.**
-  Hong Liu, Zhiyuan Li, David Hall, Percy Liang, and Tengyu Ma.  
**Sophia: A Scalable Stochastic Second-order Optimizer for Language Model Pre-training, May 2023.**
-  Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D. Lee, Danqi Chen, and Sanjeev Arora.  
**Fine-Tuning Language Models with Just Forward Passes, May 2023.**

 Pratik Rathore, Weimu Lei, Zachary Frangella, Lu Lu, and Madeleine Udell.

**Challenges in Training PINNs: A Loss Landscape Perspective, February 2024.**

 Yanjun Zhao, Sizhe Dang, Haishan Ye, Guang Dai, Yi Qian, and Ivor W. Tsang.

**Second-Order Fine-Tuning without Pain for LLMs: A Hessian Informed Zeroth-Order Optimizer, February 2024.**