

# Lecture 13: Optimization for Deep Learning

**Tao LIN**

SoE, Westlake University

December 16, 2025



1 Stochastic Gradient Descent (SGD) and Mini-batch SGD

2 Accelerated and Stabilized Optimization Methods

3 Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization
- Muon Optimizer

# Table of Contents

- ① Stochastic Gradient Descent (SGD) and Mini-batch SGD
- ② Accelerated and Stabilized Optimization Methods
- ③ Advanced Optimization Methods

# Background

- ① **Get data:**  $\xi_1, \dots, \xi_N$ , **where**  $\xi_i := (\mathbf{d}, y)_i$

# Background

- ① Get data:  $\xi_1, \dots, \xi_N$ , where  $\xi_i := (\mathbf{d}, y)_i$
- ② Choose a classifier

$$\begin{aligned} h_{\mathbf{x}}(\mathbf{d}) &\rightarrow y \\ h_{\mathbf{x}} \left( \begin{array}{c} \text{cat} \end{array} \right) &\rightarrow \boxed{\mathbf{cat}} \end{aligned} \tag{1}$$

# Background

- ① Get data:  $\xi_1, \dots, \xi_N$ , where  $\xi_i := (\mathbf{d}, y)_i$
- ② Choose a classifier

$$\begin{aligned} h_{\mathbf{x}}(\mathbf{d}) &\rightarrow y \\ h_{\mathbf{x}} \left( \begin{array}{c} \text{cat} \end{array} \right) &\rightarrow \text{cat} \end{aligned} \tag{1}$$

- ③ **Choose a loss function:**  $\ell(h_{\mathbf{x}}(\mathbf{d}, y)) \geq 0$

# Background

- 1 Get data:  $\xi_1, \dots, \xi_N$ , where  $\xi_i := (\mathbf{d}, y)_i$
- 2 Choose a classifier

$$\begin{aligned}
 h_{\mathbf{x}}(\mathbf{d}) &\rightarrow y \\
 h_{\mathbf{x}} \left( \begin{array}{c} \text{cat} \end{array} \right) &\rightarrow \text{cat}
 \end{aligned} \tag{1}$$

- 3 Choose a loss function:  $\ell(h_{\mathbf{x}}(\mathbf{d}, y)) \geq 0$
- 4 Solve the ***training problem***:

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N \ell(h_{\mathbf{x}}(\mathbf{d}_i), y_i) \tag{2}$$

# Background

Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

# Background

Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i) \leftarrow$

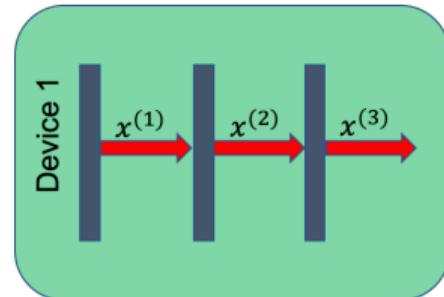
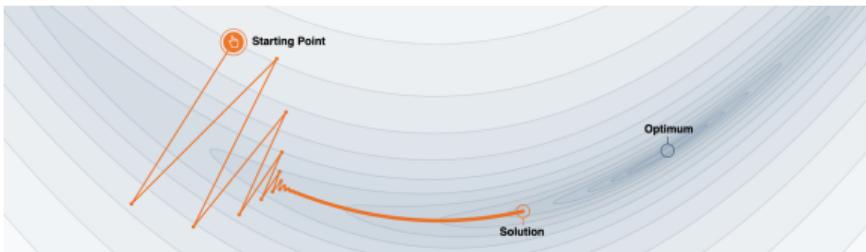
# Background

Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i) \leftarrow$
- Baseline method: Stochastic Gradient Descent (SGD)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla F(\mathbf{x}^{(t)}, \xi_i) \quad (4)$$



# Background

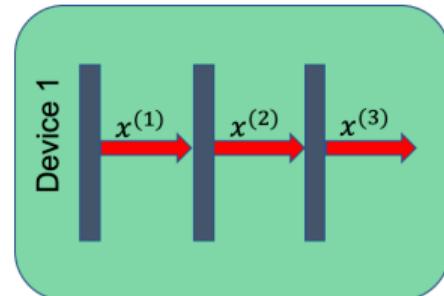
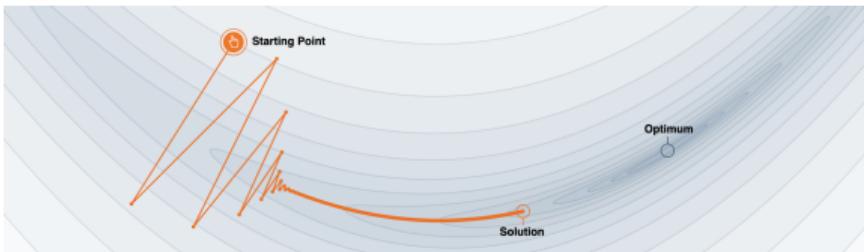
Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i)$
- Baseline method: Stochastic Gradient Descent (SGD)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla F(\mathbf{x}^{(t)}, \xi_i) \quad (4)$$

- $\eta$  is step-size/learning rate



# Background

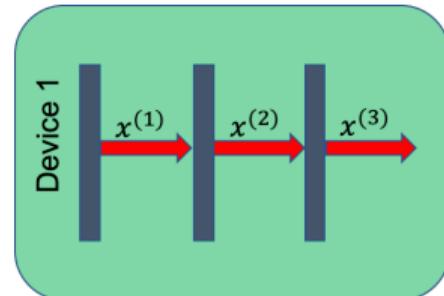
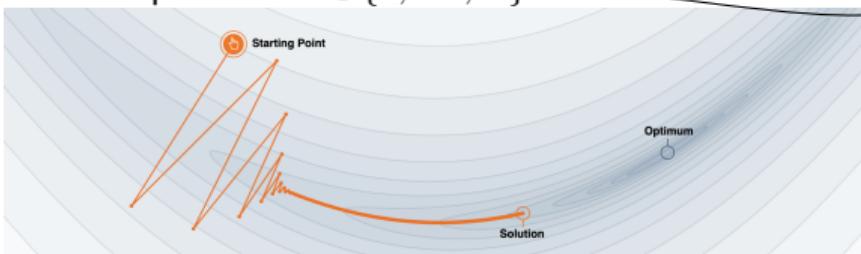
Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i) \leftarrow$
- Baseline method: Stochastic Gradient Descent (SGD)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta \nabla F(\mathbf{x}^{(t)}, \xi_i) \quad (4)$$

- $\eta$  is step-size/learning rate
- sampled i.i.d.  $i \in \{1, \dots, N\}$



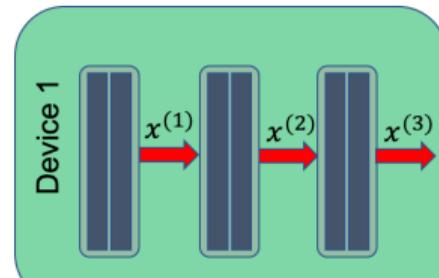
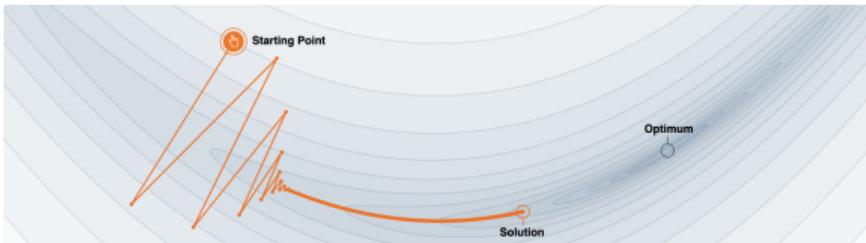
# Background

Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i) \leftarrow$
- Baseline method: Stochastic Gradient Descent (SGD)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{B} \sum_{i \in \mathcal{B}} \eta \nabla f_i(\mathbf{x}) \quad (\text{Using } \textit{mini-batching})$$



# Background

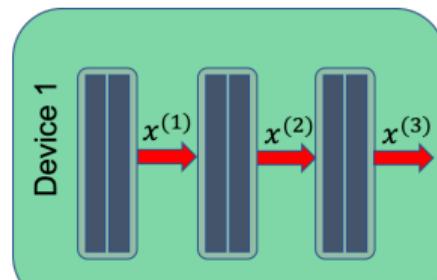
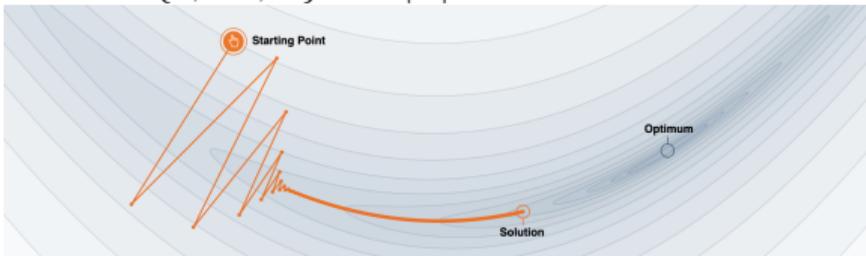
Finite-sum empirical risk minimization problem:

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N \left( f_i(\mathbf{x}) := F(\mathbf{x}, \xi_i) \right) \right\} \quad (3)$$

- The loss function of  $i$ -th data  $\xi_i := (\mathbf{d}_i, y_i) \leftarrow$
- Baseline method: Stochastic Gradient Descent (SGD)

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \frac{1}{B} \sum_{i \in \mathcal{B}} \eta \nabla f_i(\mathbf{x}) \quad (\text{Using } \textit{mini-batching})$$

- $\mathcal{B} \in \{1, \dots, N\}$  with  $|\mathcal{B}| = B$ .



## Background: Stochastic reformulation of finite-sum problems: SGD with arbitrary sampling

**Random sampling vector**  $\mathbf{v} = (v_1, \dots, v_N) \sim \mathcal{D}$  with  $\mathbb{E}[v_i] = 1 \quad \text{for } i = 1, \dots, N.$

## Background: Stochastic reformulation of finite-sum problems: SGD with arbitrary sampling

**Random sampling vector**  $\mathbf{v} = (v_1, \dots, v_N) \sim \mathcal{D}$  with  $\mathbb{E}[v_i] = 1 \quad \text{for } i = 1, \dots, N.$

$$f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[v_i] f_i(\mathbf{x}) = \mathbb{E} \left[ \underbrace{\frac{1}{N} \sum_{i=1}^N v_i f_i(\mathbf{x})}_{=: f_v(\mathbf{x})} \right] \quad (\text{Stochastic Reformulation})$$

## Background: Stochastic reformulation of finite-sum problems: SGD with arbitrary sampling

**Random sampling vector**  $\mathbf{v} = (v_1, \dots, v_N) \sim \mathcal{D}$  with  $\mathbb{E}[v_i] = 1$  for  $i = 1, \dots, N$ .

$$f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[v_i] f_i(\mathbf{x}) = \mathbb{E} \left[ \underbrace{\frac{1}{N} \sum_{i=1}^N v_i f_i(\mathbf{x})}_{=: f_v(\mathbf{x})} \right] \quad (\text{Stochastic Reformulation})$$

**Original Finite-sum problem**

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) \quad (4)$$

**Stochastic Reformulation**

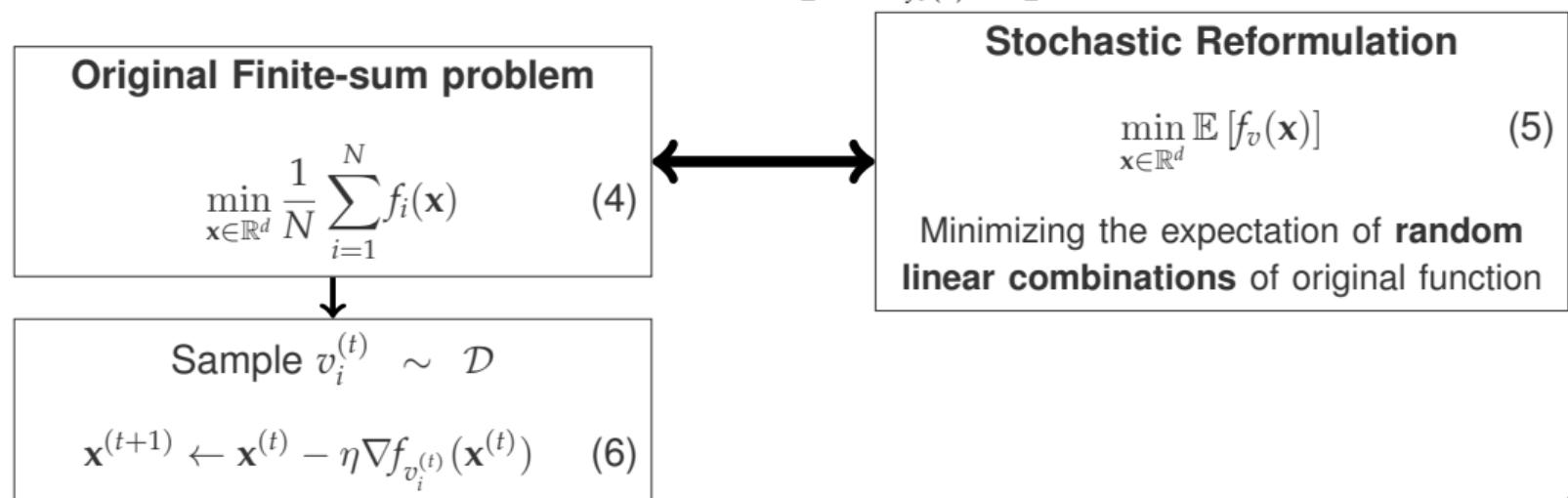
$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbb{E}[f_v(\mathbf{x})] \quad (5)$$

Minimizing the expectation of **random linear combinations** of original function

## Background: Stochastic reformulation of finite-sum problems: SGD with arbitrary sampling

**Random sampling vector**  $\mathbf{v} = (v_1, \dots, v_N) \sim \mathcal{D}$  with  $\mathbb{E}[v_i] = 1 \quad \text{for } i = 1, \dots, N.$

$$f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[v_i] f_i(\mathbf{x}) = \mathbb{E} \left[ \underbrace{\frac{1}{N} \sum_{i=1}^N v_i f_i(\mathbf{x})}_{=: f_v(\mathbf{x})} \right] \quad (\text{Stochastic Reformulation})$$



## Background: Stochastic reformulation of finite-sum problems: SGD with arbitrary sampling

**Random sampling vector**  $\mathbf{v} = (v_1, \dots, v_N) \sim \mathcal{D}$  with  $\mathbb{E}[v_i] = 1$  for  $i = 1, \dots, N$ .

$$f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbb{E}[v_i] f_i(\mathbf{x}) = \mathbb{E} \left[ \underbrace{\frac{1}{N} \sum_{i=1}^N v_i f_i(\mathbf{x})}_{=: f_v(\mathbf{x})} \right] \quad (\text{Stochastic Reformulation})$$

### Original Finite-sum problem

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) \quad (4)$$

### Stochastic Reformulation

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbb{E}[f_v(\mathbf{x})] \quad (5)$$

Minimizing the expectation of **random linear combinations** of original function

↓  
Sample  $v_i^{(t)} \sim \mathcal{D}$

$$\mathbf{x}^{(t+1)} \leftarrow \mathbf{x}^{(t)} - \eta \nabla f_{v_i^{(t)}}(\mathbf{x}^{(t)}) \quad (6)$$

The distribution  $\mathcal{D}$  encodes any form of mini-batching / non-uniform sampling.

# The convergence of mini-batch SGD

## Assumption 1

- ***The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$***

# The convergence of mini-batch SGD

## Assumption 1

- *The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$*
- *The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .*

# The convergence of mini-batch SGD

## Assumption 1

- The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$
- The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .

Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{x}^{(t)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{\frac{L}{T} (f(\mathbf{x}_0) - f^*)}{\sqrt{\frac{B}{T}}} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L (f(\mathbf{x}_0) - f^*)}{T}} \right)$$

# The convergence of mini-batch SGD

## Assumption 1

- The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$
- The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .

## Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- $L$ -smoothness

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{x}^{(t)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{L (f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L (f(\mathbf{x}_0) - f^*)}{T}} \right)$$


# The convergence of mini-batch SGD

## Assumption 1

- The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$
- The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .

## Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- $L$ -smoothness

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{x}^{(t)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{L (f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L (f(\mathbf{x}_0) - f^*)}{T}} \right)$$

- $T$ : number of iterations

# The convergence of mini-batch SGD

## Assumption 1

- The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$
- The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .

## Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- $L$ -smoothness

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{x}^{(t)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{L (f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L (f(\mathbf{x}_0) - f^*)}{T}} \right)$$

- $T$ : number of iterations

- $\sigma$ : stochastic gradient variance

# The convergence of mini-batch SGD

## Assumption 1

- The function  $f(\mathbf{x})$  we are minimizing is lower bounded from below by  $f^* := f(\mathbf{x}^*)$ , and each  $f_i$  is  $L$ -smooth satisfying  $\|\nabla f_i(\mathbf{y}) - \nabla f_i(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$
- The stochastic gradients satisfy  $\mathbb{E} [\nabla f_i(\mathbf{x})] = \nabla f(\mathbf{x})$  and  $\mathbb{E} \|\nabla f_i(\mathbf{x}) - \nabla f(\mathbf{x})\|^2 \leq \sigma^2$ .

## Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- $L$ -smoothness

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \left\| \nabla f(\mathbf{x}^{(t)}) \right\|^2 \right] \leq \mathcal{O} \left( \frac{L (f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L (f(\mathbf{x}_0) - f^*)}{T}} \right)$$

- $T$ : number of iterations

- $\sigma$ : stochastic gradient variance

- $B$ : mini-batch size of  $\mathcal{B}$

# The convergence of mini-batch SGD

Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- *L-smoothness*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \leq \mathcal{O} \left( \frac{L(f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L(f(\mathbf{x}_0) - f^*)}{T}} \right)$$

- *T: number of iterations*

- *$\sigma$ : stochastic gradient variance*

- *B: mini-batch size of  $\mathcal{B}$*

- When iterations  $T \rightarrow \infty$ , it holds that  $\mathbb{E} \left[ \|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \rightarrow 0$

# The convergence of mini-batch SGD

Theorem 1 (Convergence rate of mini-batch SGD for non-convex functions)

- *L-smoothness*

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[ \|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \leq \mathcal{O} \left( \frac{L(f(\mathbf{x}_0) - f^*)}{T} + \frac{\sigma}{\sqrt{B}} \sqrt{\frac{L(f(\mathbf{x}_0) - f^*)}{T}} \right)$$

- *T: number of iterations*

- *$\sigma$ : stochastic gradient variance*

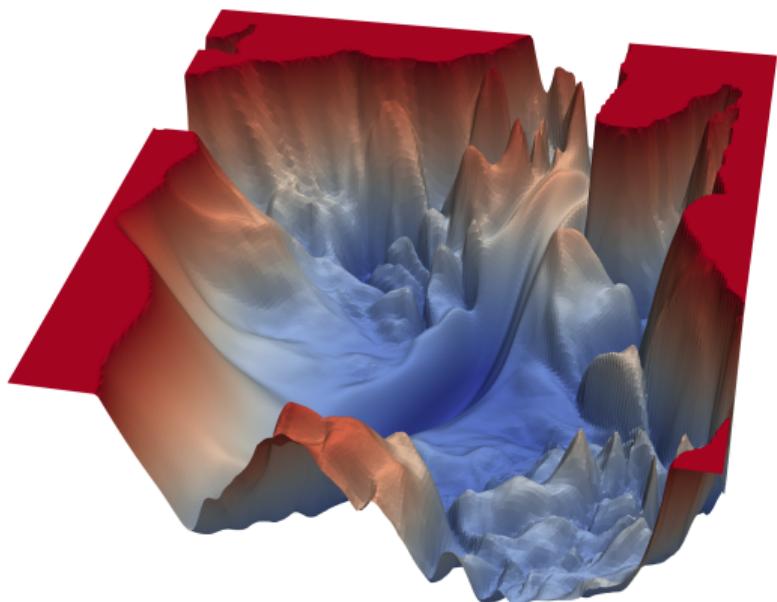
- *B: mini-batch size of  $\mathcal{B}$*

- When iterations  $T \rightarrow \infty$ , it holds that  $\mathbb{E} \left[ \|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \rightarrow 0$
- $\mathbb{E} \left[ \|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \rightarrow 0$  implies the sequence converges to a stationary solution

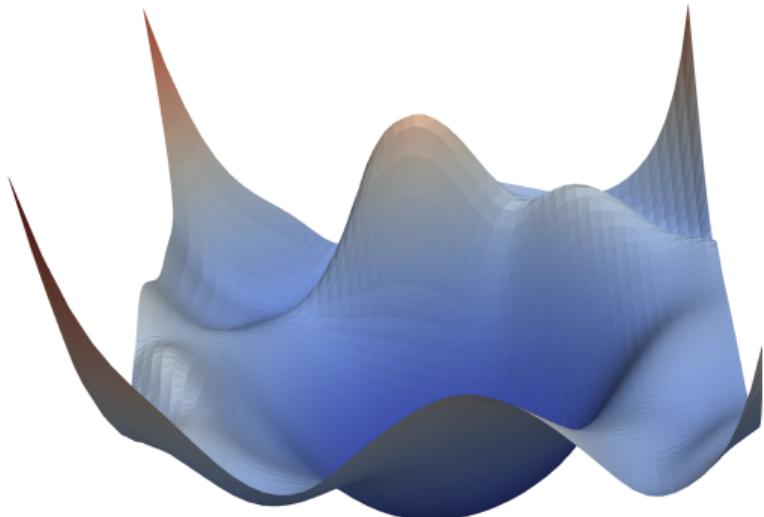
# Table of Contents

- ① Stochastic Gradient Descent (SGD) and Mini-batch SGD
- ② Accelerated and Stabilized Optimization Methods
- ③ Advanced Optimization Methods

# Issues of SGD—from the perspective of loss landscape



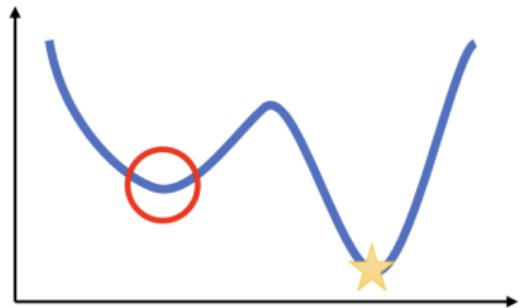
(a) ResNet w/o skip connections.



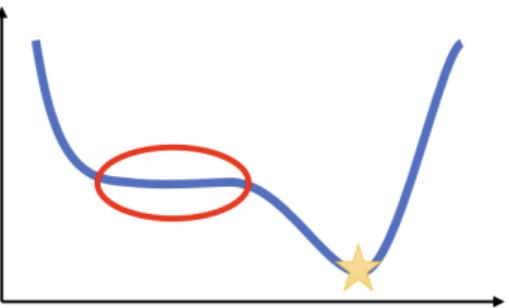
(b) ResNet w/ skip connections.

Figure: The surfaces of ResNet-56 w/ and w/o skip connections [4].

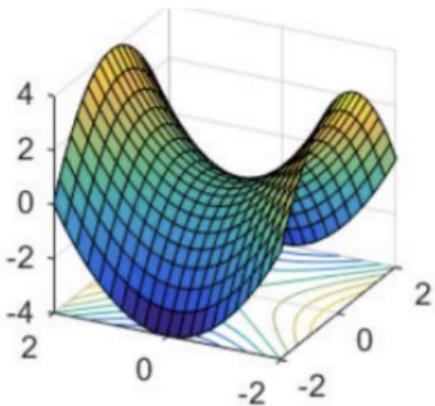
# Issues of SGD—from the perspective of loss landscape



the local optimum



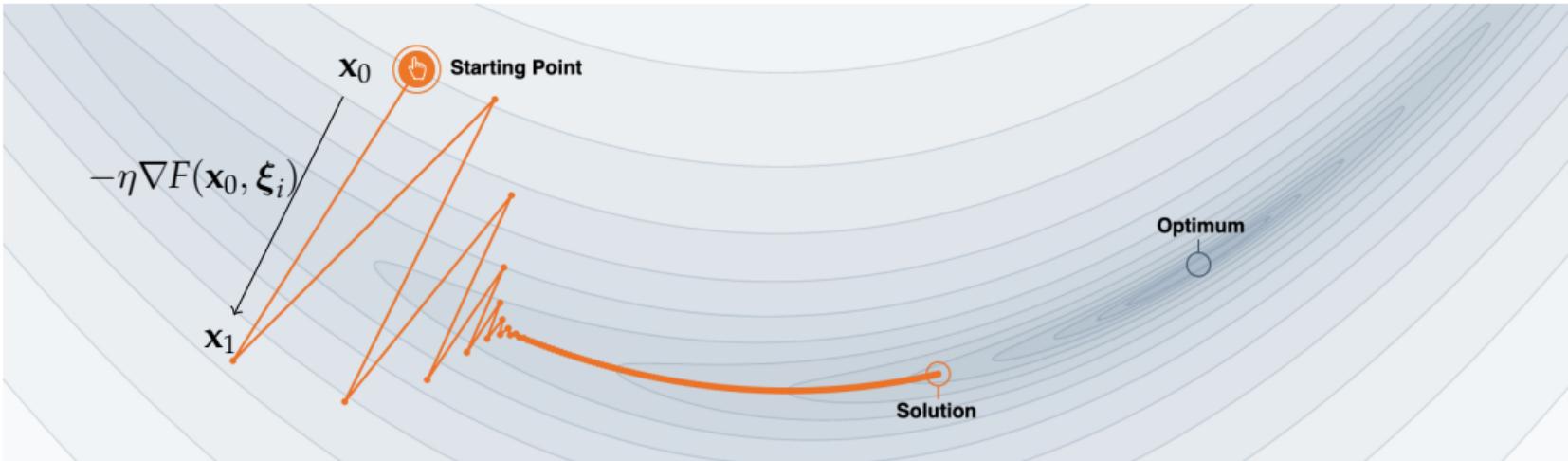
the plateau



the saddle point

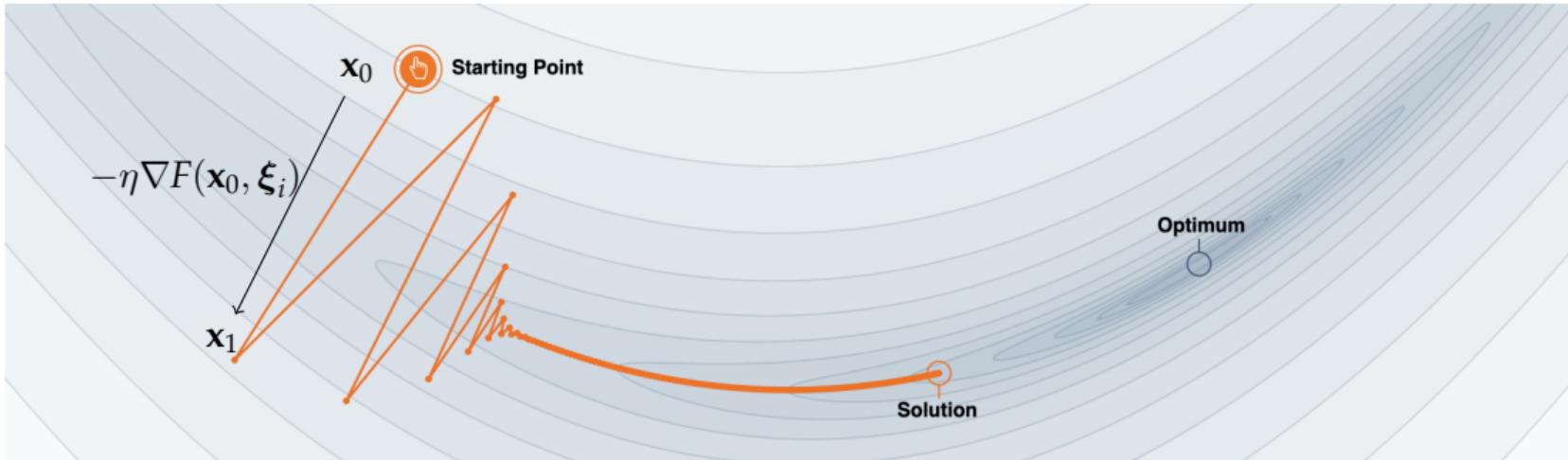
**Challenging optimization loss landscape!**

# Issues of SGD—from the perspective of loss landscape



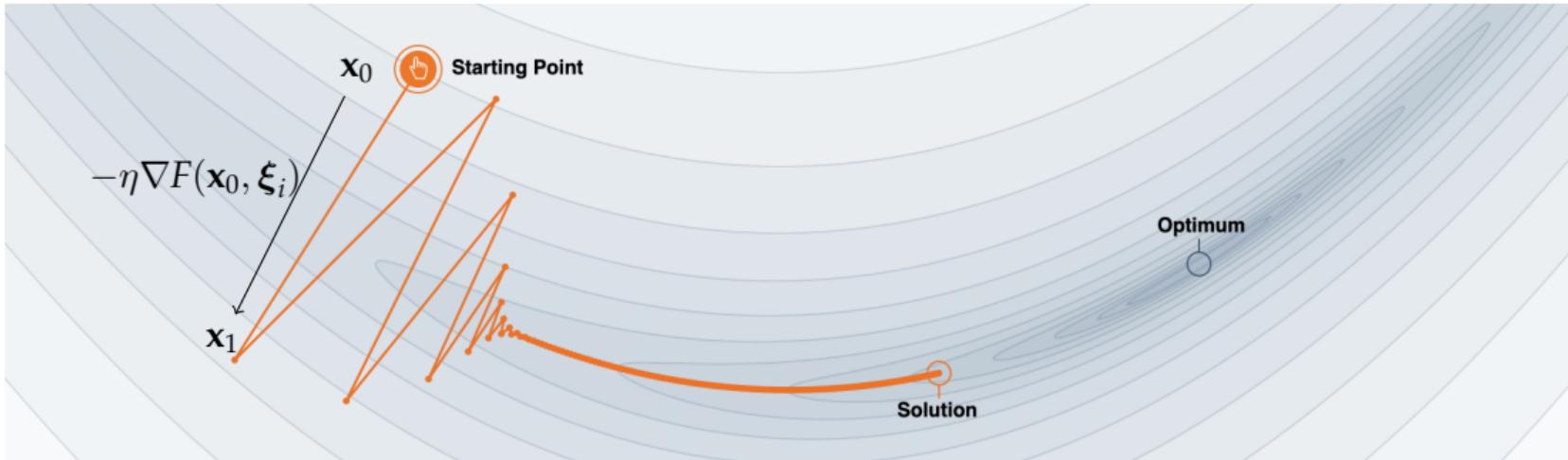
- **Challenges # 1:** loss function has high condition number.  
→ very slow progress along shallow dimension, jitter along steep direction.

# Issues of SGD—from the perspective of loss landscape



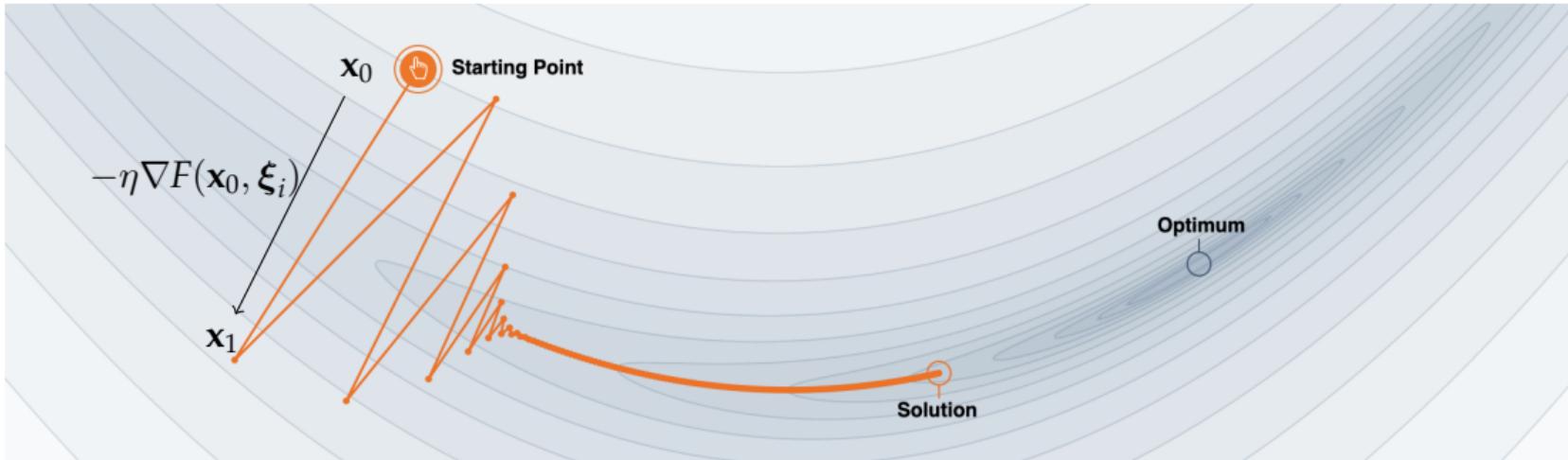
- **Challenges # 1:** loss function has high condition number.  
→ very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.

# Issues of SGD—from the perspective of loss landscape



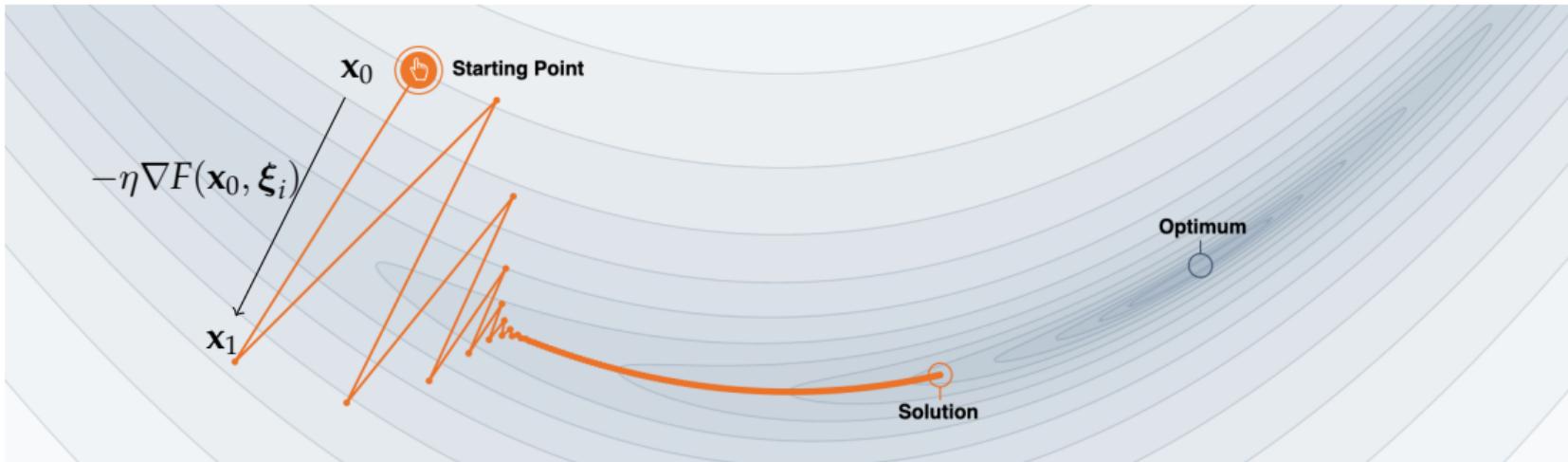
- **Challenges # 1:** loss function has high condition number.  
→ very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.  
→ cannot just choose tiny learning rates to prevent oscillation!

# Issues of SGD—from the perspective of loss landscape



- **Challenges # 1:** loss function has high condition number.
  - very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.
  - cannot just choose tiny learning rates to prevent oscillation!
  - need learning rates to be large enough not to get stuck in a plateau.

# Issues of SGD—from the perspective of loss landscape



- **Challenges # 1:** loss function has high condition number.
  - very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.
  - cannot just choose tiny learning rates to prevent oscillation!
  - need learning rates to be large enough not to get stuck in a plateau.
  - saddle points have very small gradients: but much more common in high dimension.

# Improvement directions: leveraging the curvature information

Can we find a better descent direction in the loss landscape?

Yes! **By leveraging the curvature information** through Newton's method.

# Improvement directions: leveraging the curvature information

Can we find a better descent direction in the loss landscape?

Yes! By leveraging the curvature information through Newton's method.

Taylor expansion:

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \quad (7)$$



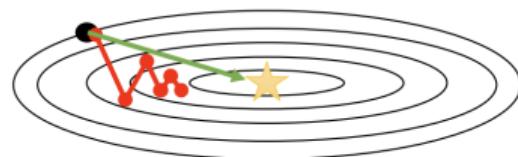
# Improvement directions: leveraging the curvature information

Can we find a better descent direction in the loss landscape?

Yes! By leveraging the curvature information through Newton's method.

**Taylor expansion:**

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \quad (7)$$



**Multivariate case:**

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_0)}_{\text{gradient}} (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \underbrace{\nabla_{\mathbf{x}}^2 f(\mathbf{x}_0)}_{\text{Hessian}} (\mathbf{x} - \mathbf{x}_0) \quad (8)$$

# Improvement directions: leveraging the curvature information

Can we find a better descent direction in the loss landscape?

Yes! By leveraging the curvature information through Newton's method.

**Taylor expansion:**

$$f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{1}{2}f''(x_0)(x - x_0)^2 \quad (7)$$



**Multivariate case:**

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \underbrace{\nabla_{\mathbf{x}} f(\mathbf{x}_0)}_{\text{gradient}} (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^\top \underbrace{\nabla_{\mathbf{x}}^2 f(\mathbf{x}_0)}_{\text{Hessian}} (\mathbf{x} - \mathbf{x}_0) \quad (8)$$

**Solution** (can optimize this analytically!):

$$\mathbf{x}^* \leftarrow \mathbf{x}_0 - (\nabla_{\mathbf{x}}^2 f(\mathbf{x}_0))^{-1} \nabla_{\mathbf{x}} f(\mathbf{x}_0) \quad (9)$$

# Improvement directions: trade-offs and approximations

**Q:** Why is Newton's method not a viable way to improve neural network optimization?

---

<sup>1</sup>if using naive approach, though fancy methods can be much faster if they avoid forming the Hessian explicitly.

# Improvement directions: trade-offs and approximations

**Q:** Why is Newton's method not a viable way to improve neural network optimization?

GD (w/o Hessian):  $\mathcal{O}(N)$

v.s.

GD (w/ Hessian)<sup>1</sup>:  $\mathcal{O}(N^3)$

---

<sup>1</sup>if using naive approach, though fancy methods can be much faster if they avoid forming the Hessian explicitly.

# Improvement directions: trade-offs and approximations

**Q:** Why is Newton's method not a viable way to improve neural network optimization?

GD (w/o Hessian):  $\mathcal{O}(N)$

v.s.

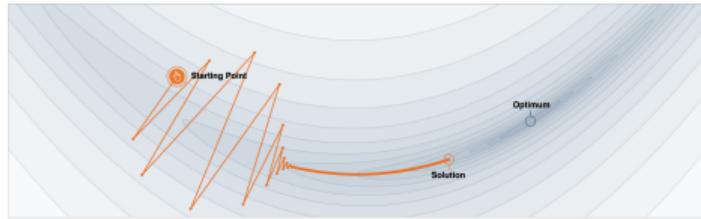
GD (w/ Hessian)<sup>1</sup>:  $\mathcal{O}(N^3)$

We would prefer methods that don't require second derivatives, but somehow "stabilize" / "accelerate" gradient descent instead.

---

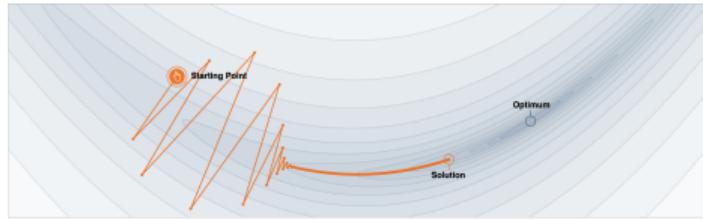
<sup>1</sup>if using naive approach, though fancy methods can be much faster if they avoid forming the Hessian explicitly.

# Momentum method

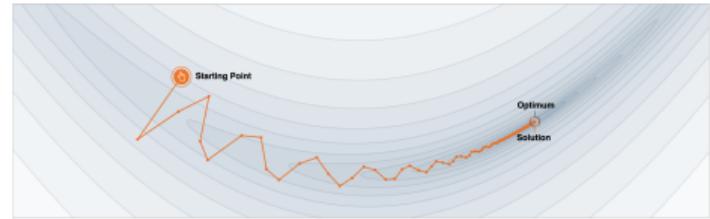


w/o momentum

# Momentum method

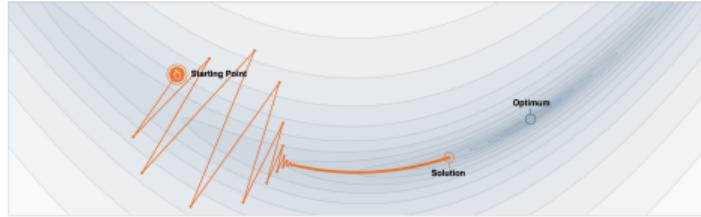


w/o momentum

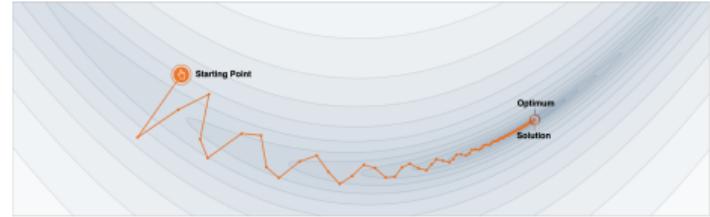


w/ momentum

# Momentum method



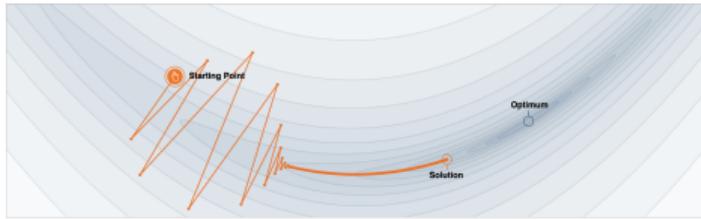
w/o momentum



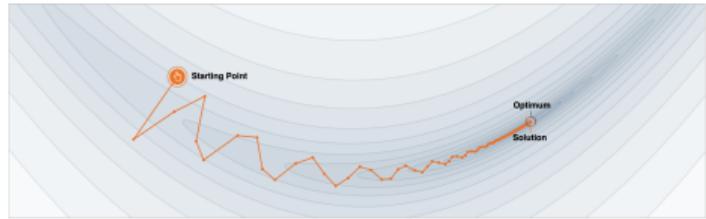
w/ momentum

**Intuition:** averaging together successive gradients yield a much better direction!

# Momentum method



w/o momentum

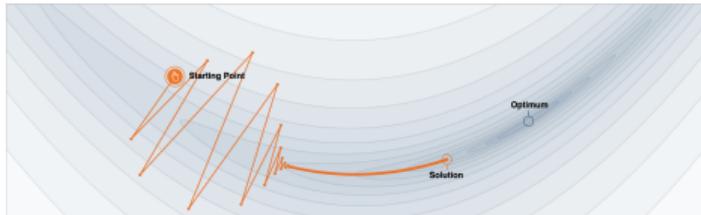


w/ momentum

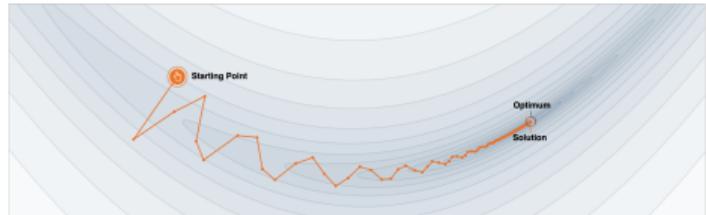
**Intuition:** averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions

# Momentum method



w/o momentum

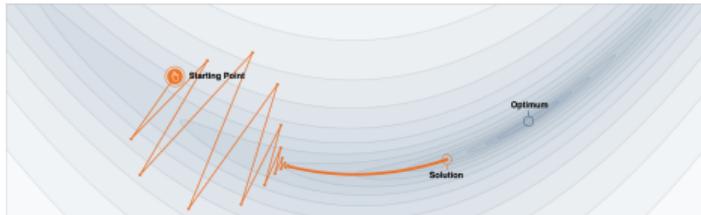


w/ momentum

**Intuition:** averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions  
→ we should cancel off the directions that disagree

# Momentum method



w/o momentum

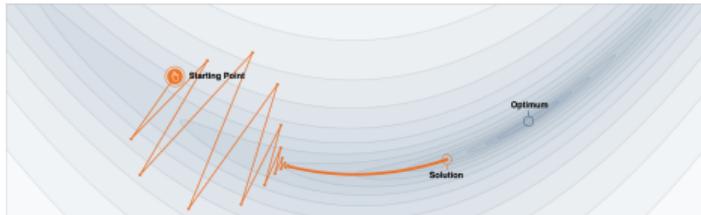


w/ momentum

**Intuition:** averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions  
→ we should cancel off the directions that disagree
- if successive gradient step point in similar directions

# Momentum method



w/o momentum

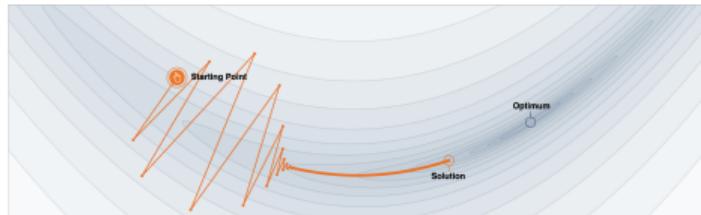


w/ momentum

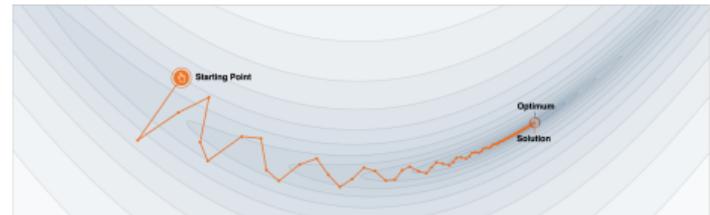
**Intuition:** averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions  
→ we should cancel off the directions that disagree
- if successive gradient step point in similar directions  
→ we should go faster in that direction

# Momentum method



w/o momentum



w/ momentum

**Intuition:** averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions  
→ we should cancel off the directions that disagree
- if successive gradient step point in similar directions  
→ we should go faster in that direction

$$\mathbf{m}_t = \beta \mathbf{m}_{t-1} + \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t), \quad \mathbf{x}_{t+1} = \mathbf{x}_t - \eta \mathbf{m}_t \quad (\text{SGD w/ momentum})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_0 - \eta \sum_{i=1}^t \nabla F(\mathbf{x}_i, \boldsymbol{\xi}_i) \quad (\text{Unroll SGD w/o momentum})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_0 - \eta \sum_{i=1}^t \frac{1 - \beta^{t+1-i}}{\beta - 1} \nabla F(\mathbf{x}_i, \boldsymbol{\xi}_i) \quad (\text{Unroll SGD w/ momentum})$$

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign:*

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign:*
- *magnitude:*

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign*: the sign of the gradient tells us which way to go along each dimension;
- *magnitude*:

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign*: the sign of the gradient tells us which way to go along each dimension;
- *magnitude*: the magnitude is not so great, and could be even worse:

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign*: the sign of the gradient tells us which way to go along each dimension;
- *magnitude*: the magnitude is not so great, and could be even worse:  
→ overall magnitude of the gradient can change drastically during the optimization,

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign*: the sign of the gradient tells us which way to go along each dimension;
- *magnitude*: the magnitude is not so great, and could be even worse:  
→ overall magnitude of the gradient can change drastically during the optimization, making learning rates hard to tune.

# Methods that manipulate gradient scale

**Intuition behind  $\nabla F(\mathbf{x}_i, \xi_i)$ :**

- *sign*: the sign of the gradient tells us which way to go along each dimension;
- *magnitude*: the magnitude is not so great, and could be even worse:  
→ overall magnitude of the gradient can change drastically during the optimization, making learning rates hard to tune.

**Idea:** *normalize* out the magnitude of the gradient along each dimension.

## Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \xi_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \xi_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

- AdaGrad has some appealing guarantees for convex problems.

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \xi_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \xi_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \xi_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \xi_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

- AdaGrad has some appealing guarantees for convex problems.  
→ AdaGrad originally proposed to benefit from sparse data.

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

- AdaGrad has some appealing guarantees for convex problems.
  - AdaGrad originally proposed to benefit from sparse data.
  - The accumulated sum of squares  $\mathbf{v}_t$  in the denominator grows monotonically throughout training.

## Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

- AdaGrad has some appealing guarantees for convex problems.
  - AdaGrad originally proposed to benefit from sparse data.
  - The accumulated sum of squares  $\mathbf{v}_t$  in the denominator grows monotonically throughout training.
  - Learning rate effectively “decreases” over time: good for convex (bad for non-convex).

# Methods that manipulate gradient scale: RMSProp, AdaGrad, and their differences

**AdaGrad [1]** (estimate per-dimension cumulative magnitude):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

**RMSProp** (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2$$

(roughly the squared length of each dimension)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}}$$

(each dimension is divided by its magnitude)

Remarks:

- AdaGrad has some appealing guarantees for convex problems.
  - AdaGrad originally proposed to benefit from sparse data.
  - The accumulated sum of squares  $\mathbf{v}_t$  in the denominator grows monotonically throughout training.
  - Learning rate effectively “decreases” over time: good for convex (bad for non-convex).
- RMSProp tends to be much better for deep learning (and most non-convex problems)

# Adam: combining momentum and RMSProp

Idea:

# Adam: combining momentum and RMSProp

## Idea:

- Maintain exponential moving averages of gradient and its square

# Adam: combining momentum and RMSProp

## Idea:

- Maintain exponential moving averages of gradient and its square
- Update proportional to  $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

# Adam: combining momentum and RMSProp

**Idea:**

- Maintain exponential moving averages of gradient and its square
- Update proportional to  $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) \quad (\text{first moment estimate})$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{second moment estimate})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} = \mathbf{x}_t - \underbrace{\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}}}_{\text{element-wise stepsize}} \mathbf{m}_t \quad (\text{update step})$$

# Adam: combining momentum and RMSProp

**Idea:**

- Maintain exponential moving averages of gradient and its square
- Update proportional to  $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) \quad (\text{first moment estimate})$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{second moment estimate})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} = \mathbf{x}_t - \underbrace{\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}}}_{\text{element-wise stepsize}} \mathbf{m}_t \quad (\text{update step})$$

where compared to RMSProp, Adam

# Adam: combining momentum and RMSProp

## Idea:

- Maintain exponential moving averages of gradient and its square
- Update proportional to  $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) \quad (\text{first moment estimate})$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{second moment estimate})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} = \mathbf{x}_t - \underbrace{\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}}}_{\text{element-wise stepsize}} \mathbf{m}_t \quad (\text{update step})$$

where compared to RMSProp, Adam

- replaces  $\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)$  by  $\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \mathbf{m}_t$ .

# Adam: combining momentum and RMSProp

## Idea:

- Maintain exponential moving averages of gradient and its square
- Update proportional to  $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) \quad (\text{first moment estimate})$$

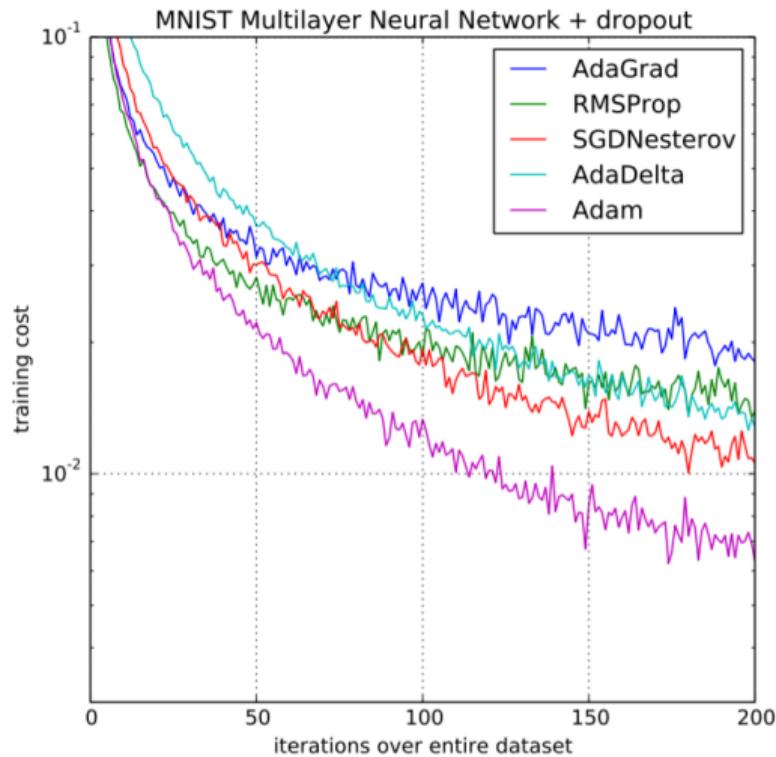
$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{second moment estimate})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\mathbf{m}_t}{\sqrt{\mathbf{v}_t + \epsilon}} = \mathbf{x}_t - \underbrace{\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}}}_{\text{element-wise stepsize}} \mathbf{m}_t \quad (\text{update step})$$

where compared to RMSProp, Adam

- replaces  $\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)$  by  $\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \mathbf{m}_t$ .
- adds bias correction (omitted in the expression above): it avoids large stepsizes in early stages of run (especially when  $\beta_2$  is close to 1).

# Adam: combining momentum and RMSProp



# Weight decay in SGD and Adam: why AdamW matters

Many learning problems optimize the loss with  $L_2$  norm penalty:

$$\tilde{f}(\mathbf{x}) = f(\mathbf{x}) + \lambda \|\mathbf{x}\|_2^2, \quad (10)$$

# Weight decay in SGD and Adam: why AdamW matters

Many learning problems optimize the loss with  $L_2$  norm penalty:

$$\tilde{f}(\mathbf{x}) = f(\mathbf{x}) + \lambda \|\mathbf{x}\|_2^2, \quad (10)$$

where it is sometimes called “weight decay” in SGD, since its gradient decays weight:

$$\begin{array}{ccc} \mathbf{x} - \eta \nabla_{\mathbf{x}} \left( f(\mathbf{x}) + \lambda \|\mathbf{x}\|_2^2 \right) & \xleftrightarrow{\nabla_{\mathbf{x}} \|\mathbf{x}\|_2^2 = 2\mathbf{x}} & (1 - 2\eta\lambda)\mathbf{x} - \eta \nabla_{\mathbf{x}} f(\mathbf{x}) \\ \text{SGD on } L_2\text{-norm penalty} & & \text{weight decay} \end{array} \quad (11)$$

# Weight decay in SGD and Adam: why AdamW matters

On the discrepancy between  $L_2$  regularization and weight decay:

---

[5] Loshchilov et al. Decoupled weight decay regularization. ICLR 2018 (Google Scholar 4800+)

# Weight decay in SGD and Adam: why AdamW matters

**On the discrepancy between  $L_2$  regularization and weight decay:**

- $L_2$  regularization and weight decay are not identical (for momentum/adaptive SGD).

# Weight decay in SGD and Adam: why AdamW matters

**On the discrepancy between  $L_2$  regularization and weight decay:**

- $L_2$  regularization and weight decay are not identical (for momentum/adaptive SGD).
- $L_2$  regularization is not effective in Adam.

# Weight decay in SGD and Adam: why AdamW matters

**On the discrepancy between  $L_2$  regularization and weight decay:**

- $L_2$  regularization and weight decay are not identical (for momentum/adaptive SGD).
- $L_2$  regularization is not effective in Adam.
- Weight decay is equally effective in both SGD and Adam.

# Weight decay in SGD and Adam: why AdamW matters

**On the discrepancy between  $L_2$  regularization and weight decay:**

- $L_2$  regularization and weight decay are not identical (for momentum/adaptive SGD).
- $L_2$  regularization is not effective in Adam.
- Weight decay is equally effective in both SGD and Adam.

**Decoupled SGD with momentum:** (same trick applies to Adam)

$$\mathbf{m}_{t+1} = \beta \mathbf{m}_t + (1 - \beta) \begin{pmatrix} \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) + \lambda \mathbf{x}_t \\ \text{gradient of loss with } L_2 \text{ penalty} \end{pmatrix} \quad (10)$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{m}_t - \frac{2\eta\lambda\mathbf{x}_t}{\text{weight decay}} \quad (11)$$

# Weight decay in SGD and Adam: why AdamW matters

**On the discrepancy between  $L_2$  regularization and weight decay:**

- $L_2$  regularization and weight decay are not identical (for momentum/adaptive SGD).
- $L_2$  regularization is not effective in Adam.
- Weight decay is equally effective in both SGD and Adam.

**Decoupled SGD with momentum:** (same trick applies to Adam)

$$\mathbf{m}_{t+1} = \beta \mathbf{m}_t + (1 - \beta) \left( \nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t) + \lambda \mathbf{x}_t \right) \quad (10)$$

gradient of loss with  $L_2$  penalty

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{m}_t - \frac{2\eta\lambda\mathbf{x}_t}{\text{weight decay}} \quad (11)$$

AdamW is widely used in training STOA NNs from scratch or fine-tuning on downstream tasks.

# Second-Order Optimization Methods

- **Idea:** Utilize second-derivative information (the Hessian matrix  $\mathbf{H} = \nabla_x^2 F$ ) to get a better sense of the loss surface curvature.
- **Newton's Method Update:**  $\mathbf{x}_{t+1} = \mathbf{x}_t - \mathbf{H}_t^{-1} \nabla_x F(\mathbf{x}_t)$
- **Examples:** Newton's method, Quasi-Newton methods (e.g., L-BFGS which approximates  $\mathbf{H}_t^{-1}$ ).
- **Pros:** Can converge very rapidly (quadratically) near a good minimum.
- **Cons for Deep Learning:**
  - **Computationally Prohibitive:**
    - Computing the Hessian:  $\mathcal{O}(N_{\text{params}}^2)$  operations.
    - Storing the Hessian:  $\mathcal{O}(N_{\text{params}}^2)$  memory.
    - Inverting the Hessian:  $\mathcal{O}(N_{\text{params}}^3)$  operations.
  - Impractical for typical DNNs with millions/billions of parameters ( $N_{\text{params}}$ ).

# Optimizer Comparison and Practical Choices

- **No Single Best Optimizer:** The performance of an optimizer can vary significantly depending on the specific problem, model architecture, dataset, and hyper-parameter tuning.

# Optimizer Comparison and Practical Choices

- **No Single Best Optimizer:** The performance of an optimizer can vary significantly depending on the specific problem, model architecture, dataset, and hyper-parameter tuning.
- **Common Starting Points:**
  - **Adam / AdamW:** Often good default choices. They are generally robust, converge relatively fast, and work well across a wide range of problems.
  - Typical Adam hyper-parameters:  $\eta = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ .

# Optimizer Comparison and Practical Choices

- **No Single Best Optimizer:** The performance of an optimizer can vary significantly depending on the specific problem, model architecture, dataset, and hyper-parameter tuning.
- **Common Starting Points:**
  - **Adam / AdamW:** Often good default choices. They are generally robust, converge relatively fast, and work well across a wide range of problems.
  - Typical Adam hyper-parameters:  $\eta = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ .
- **SGD with Momentum:**
  - Can sometimes achieve better generalization performance than adaptive methods, especially with careful and extensive tuning of learning rate schedules and momentum.
  - May require more effort to tune and might converge slower initially.

# Optimizer Comparison and Practical Choices

- **No Single Best Optimizer:** The performance of an optimizer can vary significantly depending on the specific problem, model architecture, dataset, and hyper-parameter tuning.
- **Common Starting Points:**
  - **Adam / AdamW:** Often good default choices. They are generally robust, converge relatively fast, and work well across a wide range of problems.
  - Typical Adam hyper-parameters:  $\eta = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 10^{-8}$ .
- **SGD with Momentum:**
  - Can sometimes achieve better generalization performance than adaptive methods, especially with careful and extensive tuning of learning rate schedules and momentum.
  - May require more effort to tune and might converge slower initially.
- **Experimentation is Key:** It's often beneficial to try a few different optimizers and learning rate schedules.

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.
  - **Exponential Decay:**  $\eta_t = \eta_0 e^{-kt}$ .

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.
  - **Exponential Decay:**  $\eta_t = \eta_0 e^{-kt}$ .
  - **Cosine Annealing:**  $\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{t_{curr}}{T_{total}}\pi))$ . Smoothly varies  $\eta$  from  $\eta_{max}$  to  $\eta_{min}$ .

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.
  - **Exponential Decay:**  $\eta_t = \eta_0 e^{-kt}$ .
  - **Cosine Annealing:**  $\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{t_{curr}}{T_{total}}\pi))$ . Smoothly varies  $\eta$  from  $\eta_{max}$  to  $\eta_{min}$ .
  - **Linear Decay:** Decrease  $\eta$  linearly over epochs.

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.
  - **Exponential Decay:**  $\eta_t = \eta_0 e^{-kt}$ .
  - **Cosine Annealing:**  $\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{t_{curr}}{T_{total}}\pi))$ . Smoothly varies  $\eta$  from  $\eta_{max}$  to  $\eta_{min}$ .
  - **Linear Decay:** Decrease  $\eta$  linearly over epochs.
  - **Warmup:** Start with a very small  $\eta$  for a few initial epochs and gradually increase it to the target  $\eta$ . Helps stabilize training in the early phase, especially for large models or batches.

# Learning Rate Schedules: Guiding the Descent

- **Motivation:** The optimal learning rate often changes during training. Starting large can help escape poor local minima, while decreasing it later can help converge to a good solution.
- **Common Schedules:**
  - **Step Decay:** Reduce  $\eta$  by a factor (e.g., 0.1) at predefined epochs.
  - **Exponential Decay:**  $\eta_t = \eta_0 e^{-kt}$ .
  - **Cosine Annealing:**  $\eta_t = \eta_{min} + \frac{1}{2}(\eta_{max} - \eta_{min})(1 + \cos(\frac{t_{curr}}{T_{total}}\pi))$ . Smoothly varies  $\eta$  from  $\eta_{max}$  to  $\eta_{min}$ .
  - **Linear Decay:** Decrease  $\eta$  linearly over epochs.
  - **Warmup:** Start with a very small  $\eta$  for a few initial epochs and gradually increase it to the target  $\eta$ . Helps stabilize training in the early phase, especially for large models or batches.
- Learning rate schedules are commonly used in conjunction with all types of optimizers.

## Broader Context: Training Aids and Summary

- Optimizers are a critical component, but successful deep learning training relies on a combination of factors:

## Broader Context: Training Aids and Summary

- Optimizers are a critical component, but successful deep learning training relies on a combination of factors:
  - **Initialization:** Proper weight initialization (e.g., Xavier/Glorot, He) is crucial to prevent early vanishing/exploding gradients.

## Broader Context: Training Aids and Summary

- Optimizers are a critical component, but successful deep learning training relies on a combination of factors:
  - **Initialization:** Proper weight initialization (e.g., Xavier/Glorot, He) is crucial to prevent early vanishing/exploding gradients.
  - **Regularization:** Techniques to prevent overfitting and improve generalization:
    - L2 Regularization (Weight Decay)
    - Dropout
    - Early Stopping

# Broader Context: Training Aids and Summary

- Optimizers are a critical component, but successful deep learning training relies on a combination of factors:
  - **Initialization:** Proper weight initialization (e.g., Xavier/Glorot, He) is crucial to prevent early vanishing/exploding gradients.
  - **Regularization:** Techniques to prevent overfitting and improve generalization:
    - L2 Regularization (Weight Decay)
    - Dropout
    - Early Stopping
  - **Batch Normalization:**
    - Standardizes inputs to layers, stabilizing training.
    - Smooths the loss landscape.
    - Allows for higher learning rates.
    - Can act as a regularizer.

# Broader Context: Training Aids and Summary

- Optimizers are a critical component, but successful deep learning training relies on a combination of factors:
  - **Initialization:** Proper weight initialization (e.g., Xavier/Glorot, He) is crucial to prevent early vanishing/exploding gradients.
  - **Regularization:** Techniques to prevent overfitting and improve generalization:
    - L2 Regularization (Weight Decay)
    - Dropout
    - Early Stopping
  - **Batch Normalization:**
    - Standardizes inputs to layers, stabilizing training.
    - Smooths the loss landscape.
    - Allows for higher learning rates.
    - Can act as a regularizer.

**Summary:** Choosing and tuning an optimizer is an art and science. It involves understanding the trade-offs and often requires empirical validation. The interplay between the optimizer, learning rate schedule, regularization, and model architecture determines the final performance.

# Table of Contents

① Stochastic Gradient Descent (SGD) and Mini-batch SGD

② Accelerated and Stabilized Optimization Methods

③ Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization
- Muon Optimizer

# Table of Contents

① Stochastic Gradient Descent (SGD) and Mini-batch SGD

② Accelerated and Stabilized Optimization Methods

③ Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization
- Muon Optimizer

# Lookahead Optimizer: $k$ steps forward, 1 step back

Different from the ideas e.g.,

# Lookahead Optimizer: $k$ steps forward, 1 step back

Different from the ideas e.g.,

- **Adaptive element-wise learning rate**, e.g., AdaGrad and Adam

# Lookahead Optimizer: $k$ steps forward, 1 step back

Different from the ideas e.g.,

- **Adaptive element-wise learning rate**, e.g., AdaGrad and Adam
- **Accelerated optimization**, e.g., Heavy-ball momentum and Nesterov momentum.

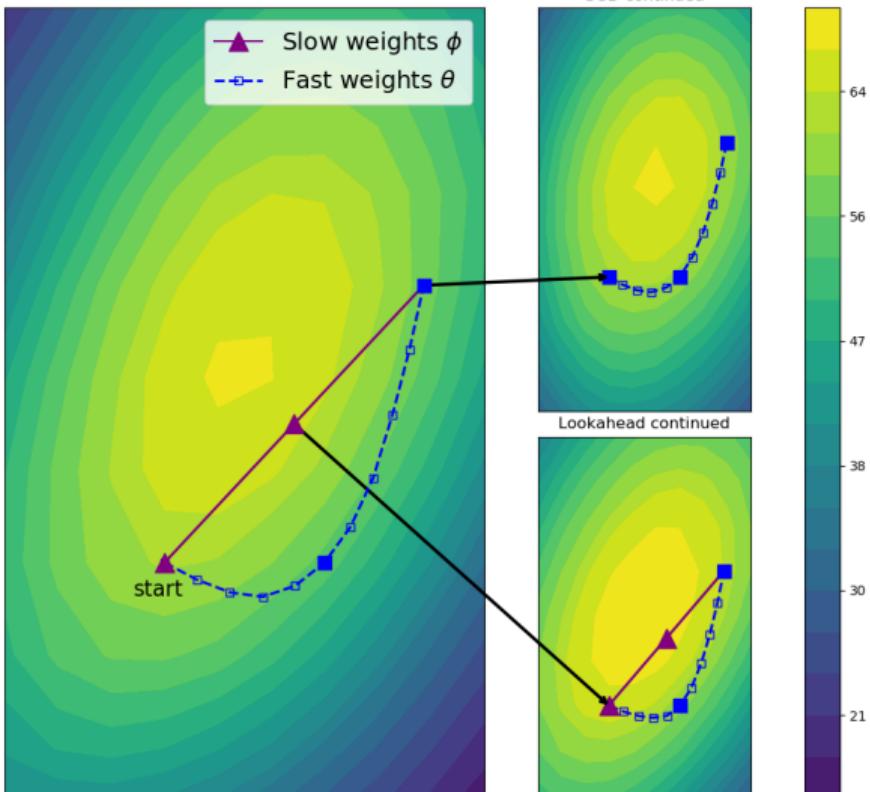
# Lookahead Optimizer: $k$ steps forward, 1 step back

Different from the ideas e.g.,

- **Adaptive element-wise learning rate**, e.g., AdaGrad and Adam
- **Accelerated optimization**, e.g., Heavy-ball momentum and Nesterov momentum.

# Lookahead Optimizer: $k$ steps forward, 1 step back

CIFAR-100 accuracy surface with Lookahead interpolation



**Algorithm 1:** SGD

**Input** : Objective  $F_S(\theta)$ , dataset  $S$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_\tau^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**Inner-loop optimization**

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, S) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = \mathbf{v}_k^{(t)} = (1 - 1)\theta_{t-1} + 1 * \mathbf{v}_k^{(t)} (\alpha = 1)$$

**end**

**outer-loop optimization**

**Output** :  $\theta_{\mathcal{A}, S} = \theta_T$

**Algorithm 2:** Lookahead

**Input** : Objective  $F_S(\theta)$ , dataset  $S$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_\tau^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**Inner-loop optimization**

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\mathbf{v}_\tau^{(t)} = \mathcal{A}(F_S(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, S) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}.$$

**end**

**outer-loop optimization**

**Output** :  $\theta_{\mathcal{A}, S} = \theta_T$

**Algorithm 1:** SGD

**Input** : Objective  $F_{\mathcal{S}}(\theta)$ , dataset  $\mathcal{S}$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_{\tau}^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**Inner-loop optimization**

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\mathbf{v}_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = \mathbf{v}_k^{(t)} = (1 - \alpha)\theta_{t-1} + 1 * \mathbf{v}_k^{(t)} \ (\alpha = 1)$$

**end**

**Output** :  $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

**Algorithm 2:** Lookahead

**Input** : Objective  $F_{\mathcal{S}}(\theta)$ , dataset  $\mathcal{S}$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_{\tau}^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**Inner-loop optimization**

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\mathbf{v}_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}.$$

**end**

**Output** :  $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

- inner-loop optimization:  $k$  steps forward in SGD & LA
- outer-loop optimization: 1 step back in LA, while no step back in SGD

**Algorithm 1:** SGD

**Input** : Objective  $F_{\mathcal{S}}(\theta)$ , dataset  $\mathcal{S}$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_{\tau}^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\quad \mathbf{v}_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = \mathbf{v}_k^{(t)} = (1 - \alpha)\theta_{t-1} + \alpha * \mathbf{v}_k^{(t)} \quad (\alpha = 1)$$

**end**

**Output:**  $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

outer-loop optimization

**Algorithm 2:** Lookahead

**Input** : Objective  $F_{\mathcal{S}}(\theta)$ , dataset  $\mathcal{S}$ , inner-loop optimizer  $\mathcal{A}$ , inner-loop step number  $k$  and learning rate  $\{\eta_{\tau}^{(t)}\}$ , outer-loop learning rate  $\alpha \in (0, 1)$ .

**for**  $t = 1, 2, \dots, T$  **do**

$$\mathbf{v}_0^{(t)} = \theta_{t-1};$$

**for**  $\tau = 1, 2, \dots, k$  **do**

$$\quad \mathbf{v}_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), \mathbf{v}_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = \mathbf{v}_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} \mathbf{g}_{\tau-1}^{(t)}$$

**end**

$$\theta_t = (1 - \alpha)\theta_{t-1} + \alpha \mathbf{v}_k^{(t)}.$$

**end**

**Output:**  $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$

outer-loop optimization

- inner-loop optimization:  $k$  steps forward in SGD & LA
- outer-loop optimization: 1 step back in LA, while no step back in SGD

# Table of Contents

① Stochastic Gradient Descent (SGD) and Mini-batch SGD

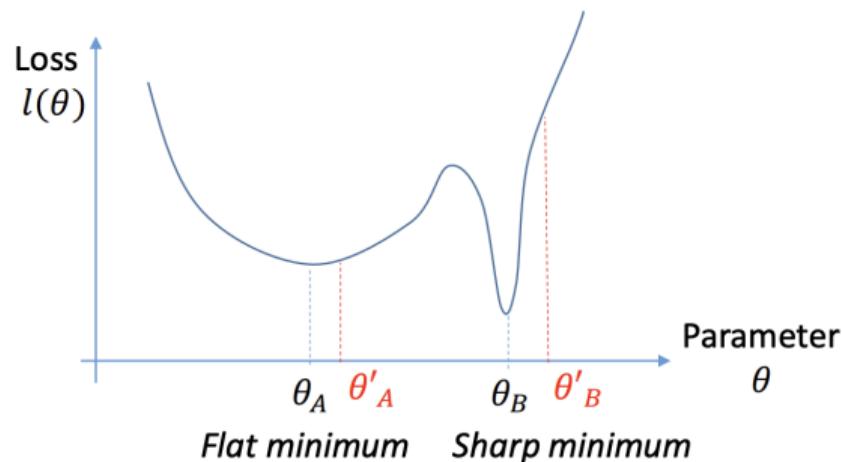
② Accelerated and Stabilized Optimization Methods

③ Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization
- Muon Optimizer

# Flat Minima in Deep Learning

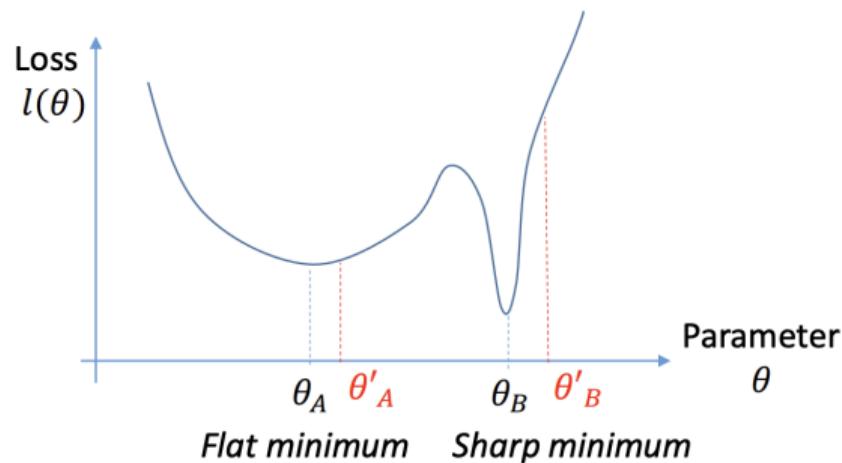
In many cases,



# Flat Minima in Deep Learning

In many cases,

DL → minimizing a loss function  $\ell(\theta)$

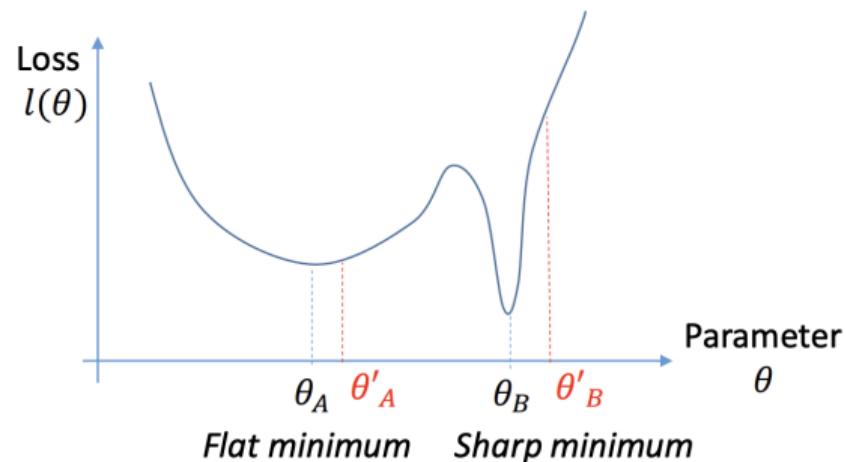


# Flat Minima in Deep Learning

In many cases,

DL → minimizing a loss function  $\ell(\theta)$

*Highly non-convex (many local minima)!*

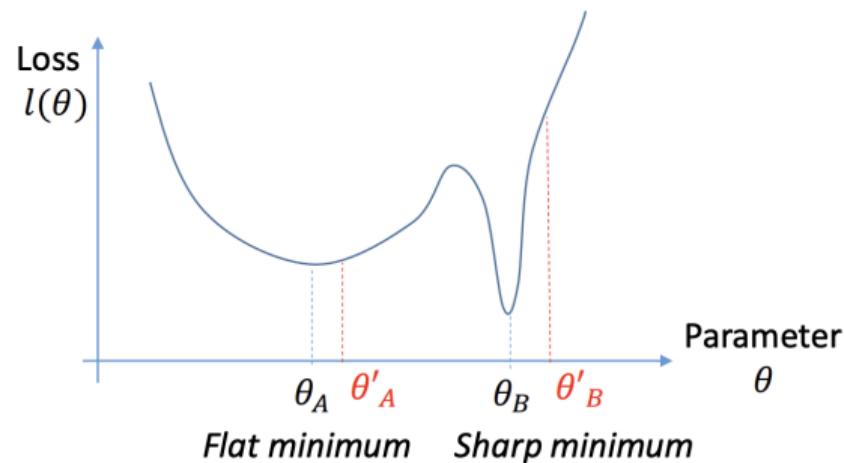


# Flat Minima in Deep Learning

In many cases,

DL → minimizing a loss function  $\ell(\theta)$

*Highly non-convex (many local minima)!*



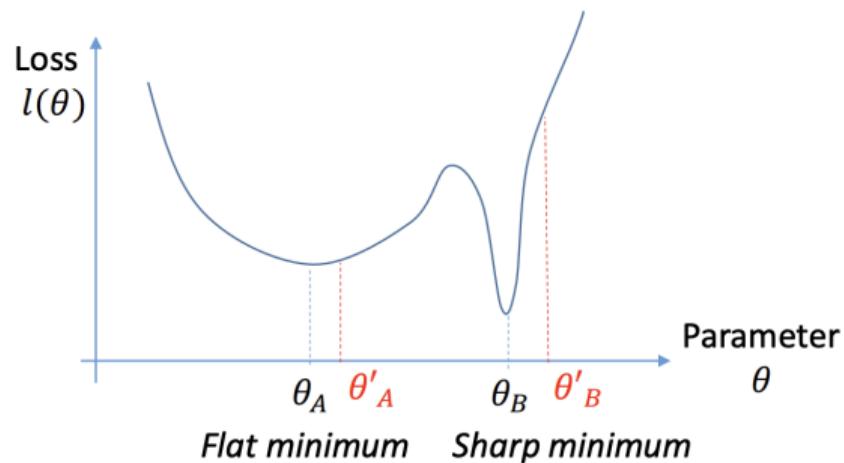
- Q) Which is better,  $\theta_A$  or  $\theta_B$ ?

# Flat Minima in Deep Learning

In many cases,

DL → minimizing a loss function  $\ell(\theta)$

*Highly non-convex (many local minima)!*



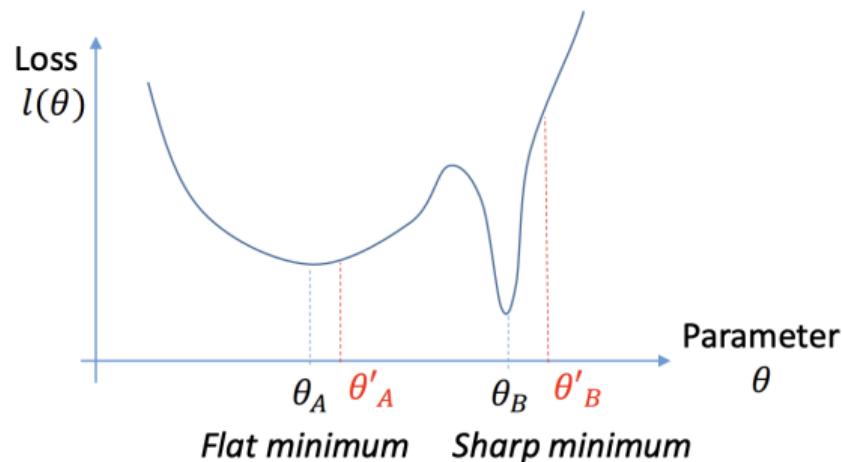
- Q) Which is better,  $\theta_A$  or  $\theta_B$ ?
- A) We prefer  $\theta_A$  to  $\theta_B$  even though  $\ell(\theta_A) > \ell(\theta_B)$

# Flat Minima in Deep Learning

In many cases,

DL → minimizing a loss function  $\ell(\theta)$

*Highly non-convex (many local minima)!*



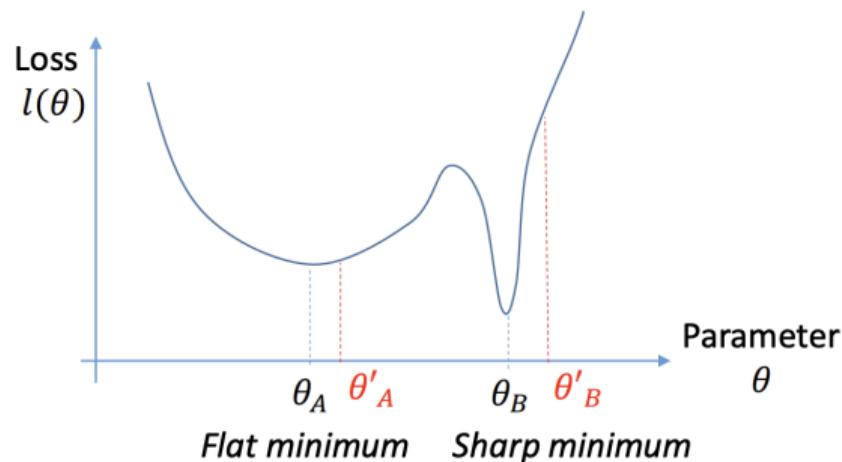
- Q) Which is better,  $\theta_A$  or  $\theta_B$ ?
- A) We prefer  $\theta_A$  to  $\theta_B$  even though  $\ell(\theta_A) > \ell(\theta_B)$ 
  - Why? Because  $\theta_A$  is more robust.

# Flat Minima in Deep Learning

In many cases,

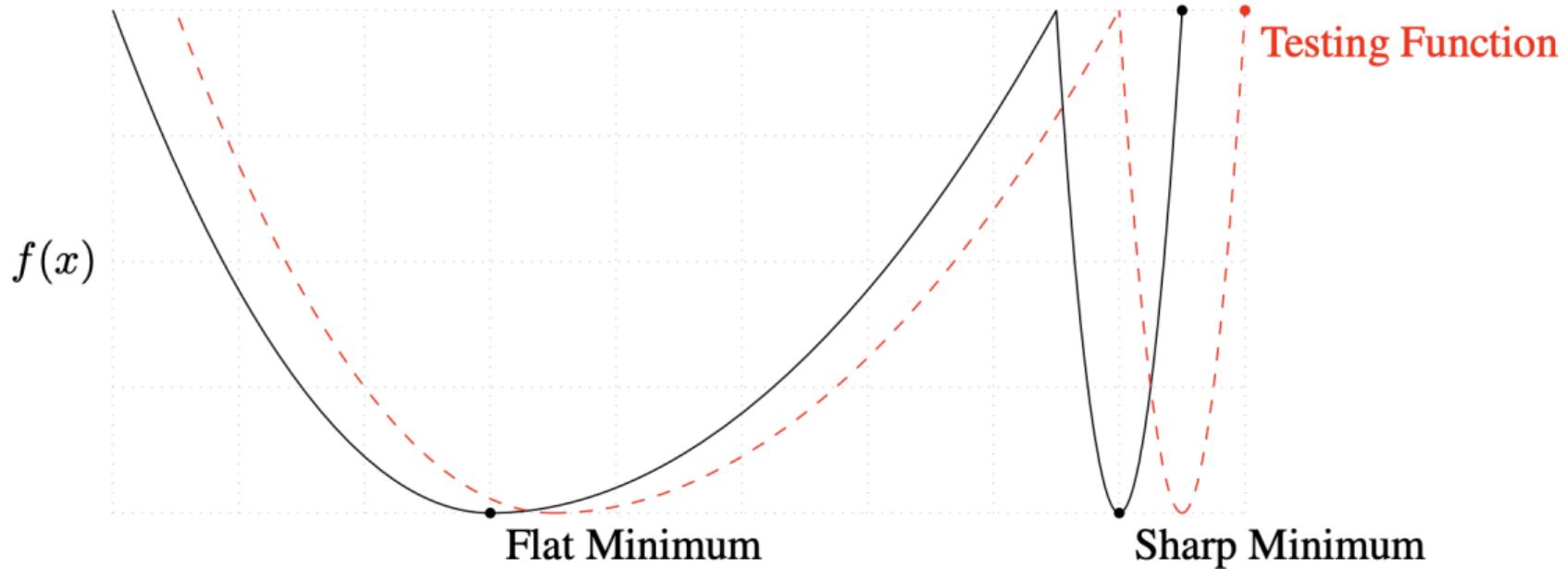
DL → minimizing a loss function  $\ell(\theta)$

*Highly non-convex (many local minima)!*



- Q) Which is better,  $\theta_A$  or  $\theta_B$ ?
- A) We prefer  $\theta_A$  to  $\theta_B$  even though  $\ell(\theta_A) > \ell(\theta_B)$ 
  - Why? Because  $\theta_A$  is more robust.
  - Imagine some perturbation:  $\theta_A \rightarrow \theta'_A, \theta_B \rightarrow \theta'_B \Rightarrow \ell(\theta'_A) \ll \ell(\theta'_B)$ .

## Training Function



# Let's seek for a Flat Minimum

Flat Minima = Robust models (12)

= Resilient to data noise or model corruption (13)

= (often encountered in AI applications) (14)

# Let's seek for a Flat Minimum

Flat Minima = Robust models (12)

= Resilient to data noise or model corruption (13)

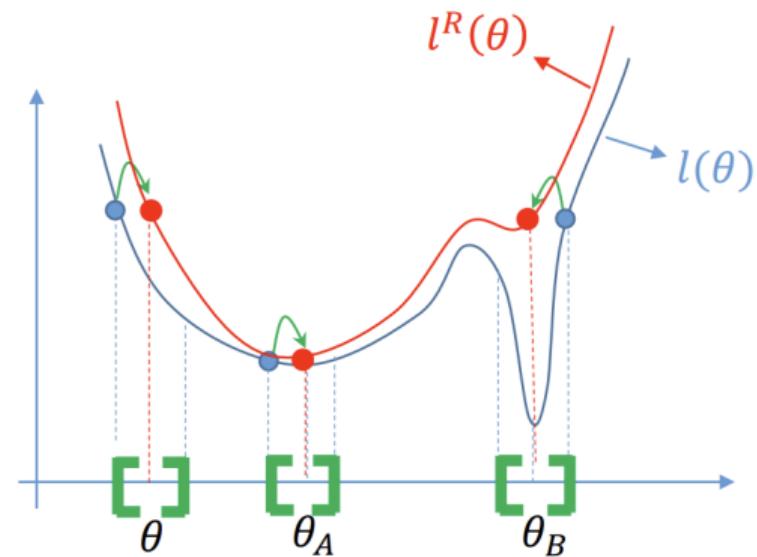
= (often encountered in AI applications) (14)

But, how?

# Sharpness-Aware Minimization (SAM)

## Idea of SAM:

Define a robust loss  $\ell^R(\theta)$  as worst-case loss within a neighborhood of  $\theta$ .

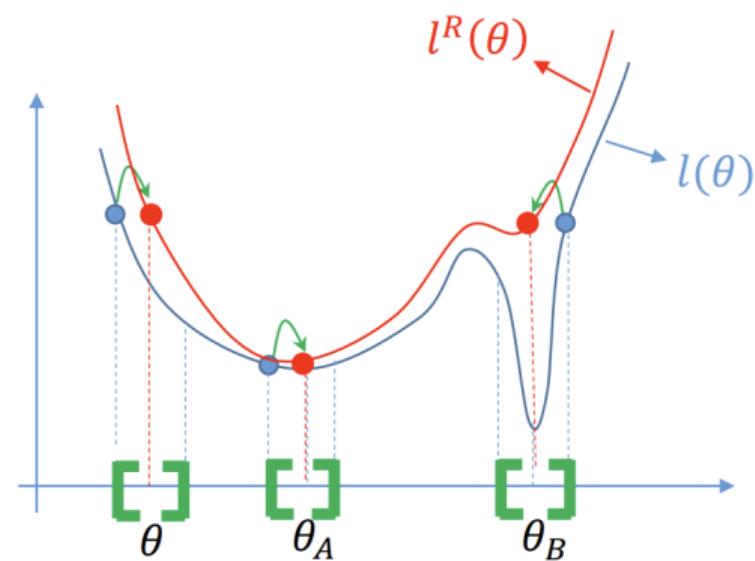


# Sharpness-Aware Minimization (SAM)

## Idea of SAM:

Define a robust loss  $\ell^R(\theta)$  as worst-case loss within a neighborhood of  $\theta$ .

$$\ell^R(\theta) = \max_{\epsilon \in N_\theta} \ell(\theta + \epsilon) \quad (15)$$



# Intuition behind SAM

- **Goal:** Find the local minima  $\theta$  that are generalizable to test samples

# Intuition behind SAM

- **Goal:** Find the local minima  $\theta$  that are generalizable to test samples
- **Theorem (Flatness-based generalization bounds):**

# Intuition behind SAM

- **Goal:** Find the local minima  $\theta$  that are generalizable to test samples
- **Theorem (Flatness-based generalization bounds):**

Theorem 2

*With high probability over  $\mathcal{S}$ , the flatness-based bound says:*

$$\mathcal{L}_{\mathcal{D}}(\theta) \leq \mathcal{L}_{\mathcal{S}}(\theta) + \underbrace{\left[ \max_{\|\epsilon\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta + \epsilon) - \mathcal{L}_{\mathcal{S}}(\theta) \right]}_{\text{flatness measure}} + h(\|\theta\|_2^2 / \rho^2), \quad (16)$$

# Intuition behind SAM

- **Goal:** Find the local minima  $\theta$  that are generalizable to test samples
- **Theorem (Flatness-based generalization bounds):**

Theorem 2

*With high probability over  $\mathcal{S}$ , the flatness-based bound says:*

$$\mathcal{L}_{\mathcal{D}}(\theta) \leq \mathcal{L}_{\mathcal{S}}(\theta) + \underbrace{\left[ \max_{\|\epsilon\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\theta + \epsilon) - \mathcal{L}_{\mathcal{S}}(\theta) \right]}_{\text{flatness measure}} + h(\|\theta\|_2^2 / \rho^2), \quad (16)$$

*where  $h(\|\theta\|_2^2 / \rho^2)$  is a strictly increasing function of  $\theta$ . It decreases as the number of samples  $n = |\mathcal{S}|$  increases.*

# Sharpness-Aware Minimization (SAM)

$$\mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) \leq \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) + \underbrace{\left[ \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}) - \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) \right]}_{\text{flatness measure}} + h(\|\boldsymbol{\theta}\|_2^2 / \rho^2) \quad (17)$$

$$= \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}) + h(\|\boldsymbol{\theta}\|_2^2 / \rho^2) \quad (18)$$

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

- the optimal  $\hat{\boldsymbol{\epsilon}}$  is given by (linear approximation through first-order Taylor expansion)

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

- the optimal  $\hat{\boldsymbol{\epsilon}}$  is given by (linear approximation through first-order Taylor expansion)

$$\hat{\boldsymbol{\epsilon}} = \arg \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} \boldsymbol{\epsilon}^\top \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \rho \cdot \text{sign}(\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})) \frac{|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})|^{q-1}}{\left( \|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})\|_q^q \right)^{1/p}}, \quad (18)$$

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

- the optimal  $\hat{\boldsymbol{\epsilon}}$  is given by (linear approximation through first-order Taylor expansion)

$$\hat{\boldsymbol{\epsilon}} = \arg \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} \boldsymbol{\epsilon}^\top \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \rho \cdot \text{sign}(\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})) \frac{|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})|^{q-1}}{\left( \|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})\|_q^q \right)^{1/p}}, \quad (18)$$

where  $1/p + 1/q = 1$

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

- the optimal  $\hat{\boldsymbol{\epsilon}}$  is given by (linear approximation through first-order Taylor expansion)

$$\hat{\boldsymbol{\epsilon}} = \arg \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} \boldsymbol{\epsilon}^\top \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \rho \cdot \text{sign}(\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})) \frac{|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})|^{q-1}}{\left( \|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})\|_q^q \right)^{1/p}}, \quad (18)$$

where  $1/p + 1/q = 1$  and the solution to a classical dual norm problem can solve this approximation.

# Sharpness-Aware Minimization (SAM)

The objective of SAM becomes:

$$\boldsymbol{\theta}^* := \arg \min_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{D}}(\boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \max_{\|\boldsymbol{\epsilon}\|_2 \leq \rho} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \boldsymbol{\epsilon}). \quad (17)$$

- the optimal  $\hat{\boldsymbol{\epsilon}}$  is given by (linear approximation through first-order Taylor expansion)

$$\hat{\boldsymbol{\epsilon}} = \arg \max_{\|\boldsymbol{\epsilon}\|_p \leq \rho} \boldsymbol{\epsilon}^\top \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta}) = \rho \cdot \text{sign}(\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})) \frac{|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})|^{q-1}}{\left(\|\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})\|_q^q\right)^{1/p}}, \quad (18)$$

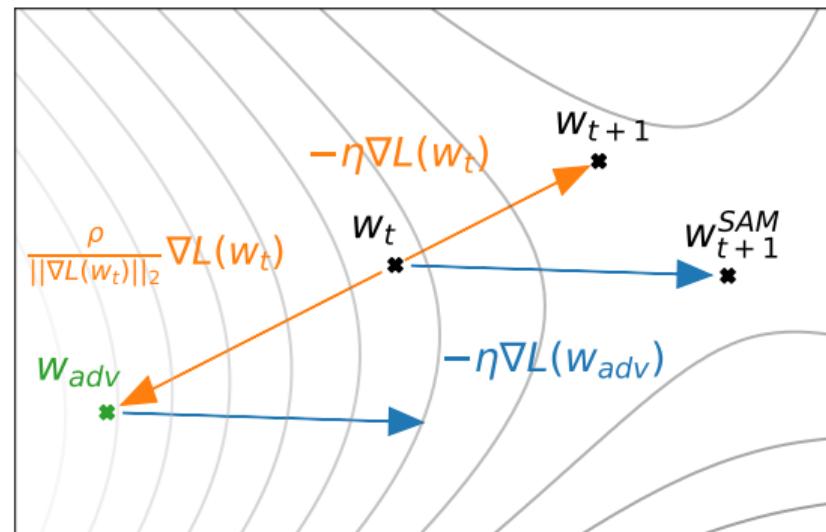
where  $1/p + 1/q = 1$  and the solution to a classical dual norm problem can solve this approximation.

- substituting  $\hat{\boldsymbol{\epsilon}}$  gives a gradient estimator

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}^{\text{SAM}}(\boldsymbol{\theta}) := \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta} + \hat{\boldsymbol{\epsilon}}) \approx \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{S}}(\boldsymbol{\theta})|_{\boldsymbol{\theta} + \hat{\boldsymbol{\epsilon}}} \quad (19)$$

- The gradient estimator of SAM is given by:

$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

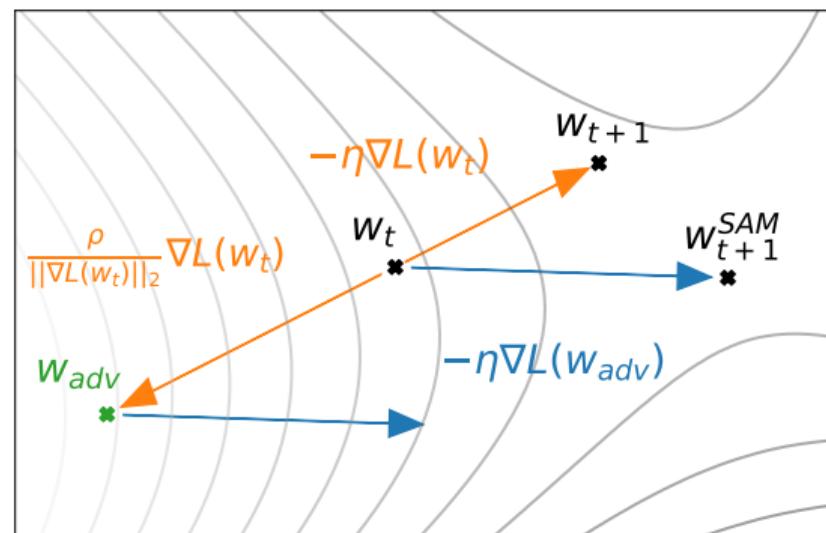


From original SAM paper.

- The gradient estimator of SAM is given by:

$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

- Recall that SAM requires 2-step gradient descent (thus, twice slow)

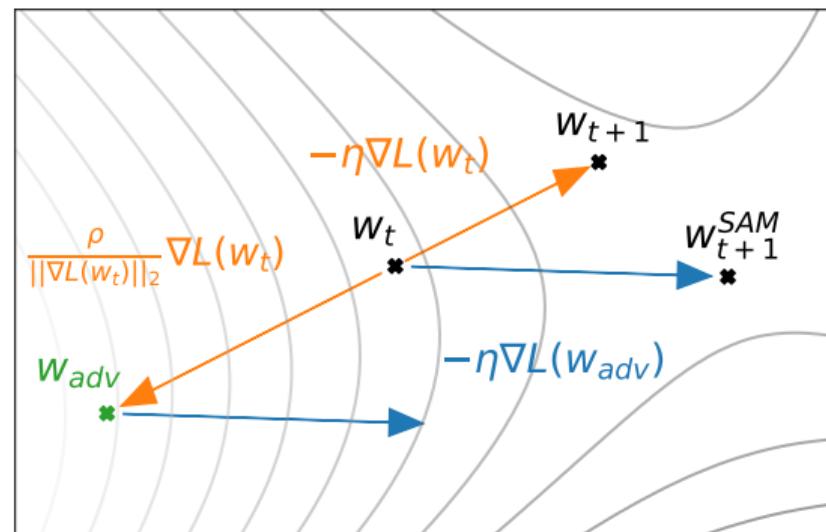


From original SAM paper.

- The gradient estimator of SAM is given by:

$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

- Recall that SAM requires 2-step gradient descent (thus, twice slow)
  - 1st for computing  $\hat{\epsilon}$  using  $\nabla_{\theta} \mathcal{L}_S(\theta)$

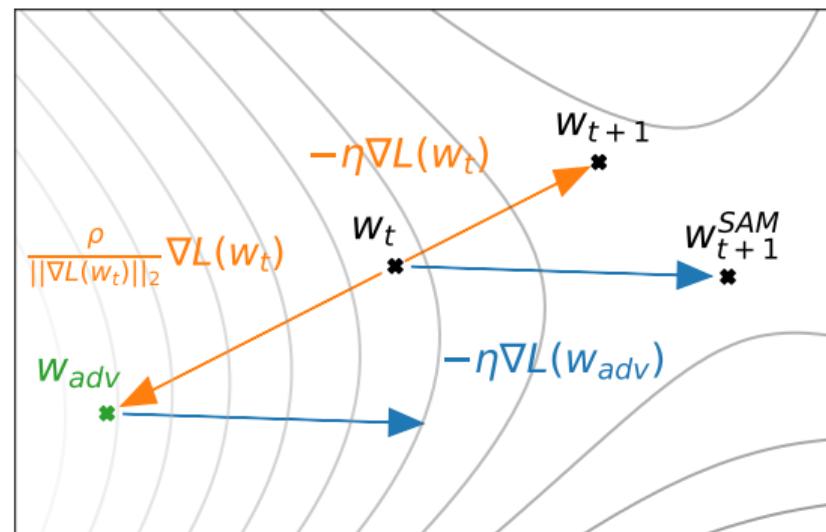


From original SAM paper.

- The gradient estimator of SAM is given by:

$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

- Recall that SAM requires 2-step gradient descent (thus, twice slow)
  - 1st for computing  $\hat{\epsilon}$  using  $\nabla_{\theta} \mathcal{L}_S(\theta)$
  - 2nd for computing  $\nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}}$

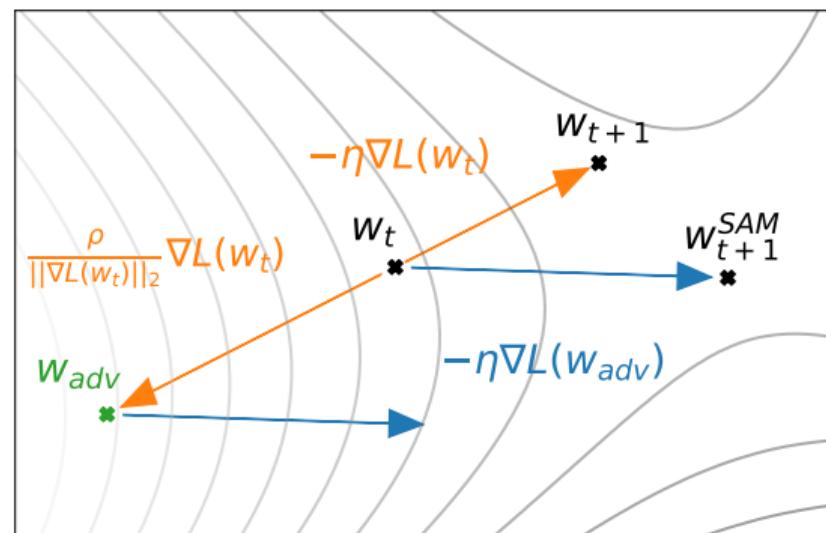


From original SAM paper.

- The gradient estimator of SAM is given by:

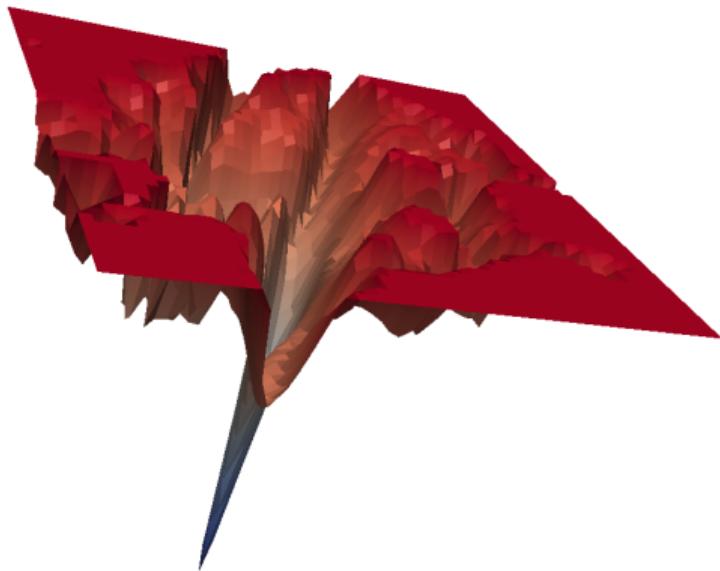
$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

- Recall that SAM requires *2-step* gradient descent (thus, twice slow)
  - 1st for computing  $\hat{\epsilon}$  using  $\nabla_{\theta} \mathcal{L}_S(\theta)$
  - 2nd for computing  $\nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}}$
- Set  $p = 2$ -norm and neighborhood-size  $\rho = 0.05$  as a default setup.

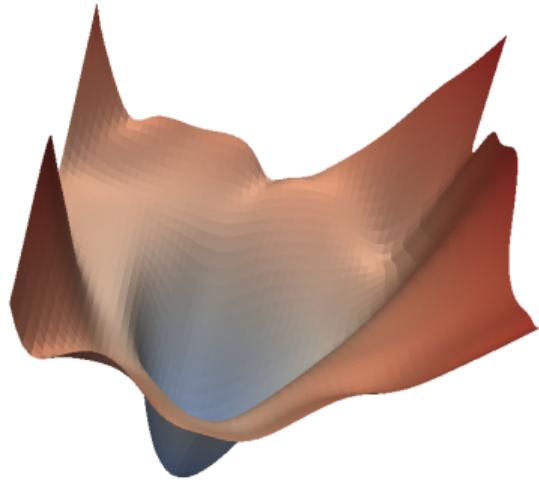


From original SAM paper.

# Verification of the flatness



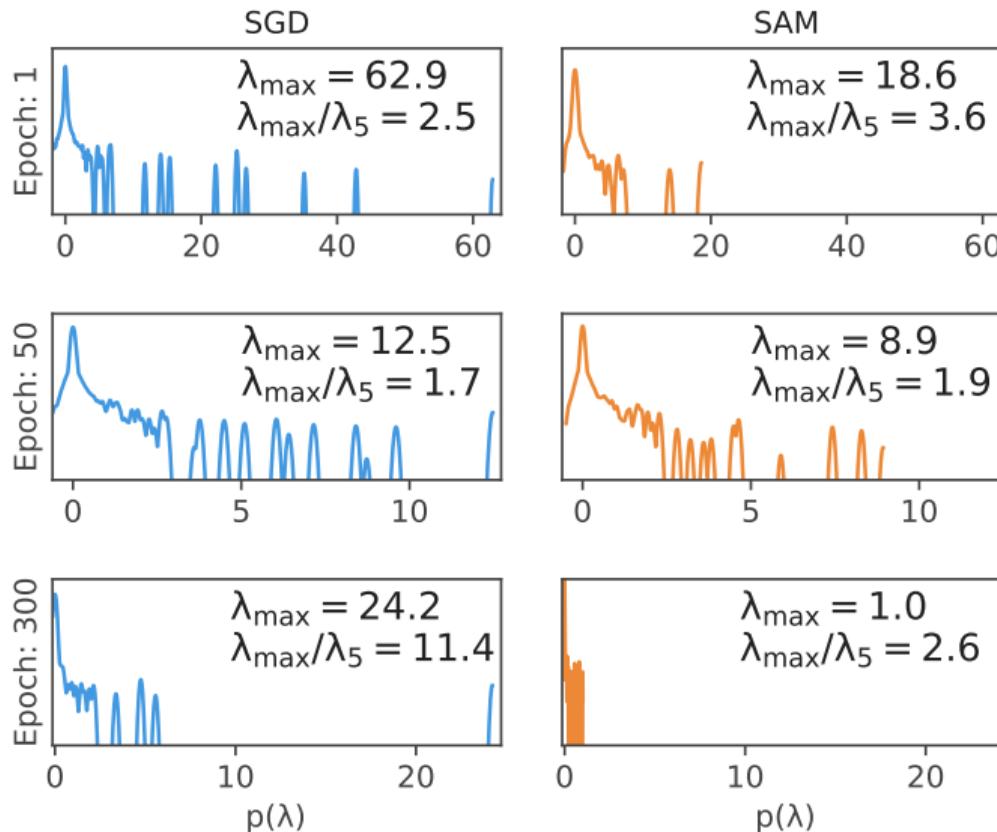
(a) ERM.



(b) SAM.

Loss surface visualization.

# Verification of the flatness



# Results (i.e., SAM > ERM)

- SAM consistently improves classification tasks, particularly with label noises

Model	Epoch	SAM		Standard Training (No SAM)	
		Top-1	Top-5	Top-1	Top-5
ResNet-50	100	<b>22.5</b> <sub>±0.1</sub>	6.28 <sub>±0.08</sub>	22.9 <sub>±0.1</sub>	6.62 <sub>±0.11</sub>
	200	<b>21.4</b> <sub>±0.1</sub>	5.82 <sub>±0.03</sub>	22.3 <sub>±0.1</sub>	6.37 <sub>±0.04</sub>
	400	<b>20.9</b> <sub>±0.1</sub>	5.51 <sub>±0.03</sub>	22.3 <sub>±0.1</sub>	6.40 <sub>±0.06</sub>
ResNet-101	100	<b>20.2</b> <sub>±0.1</sub>	5.12 <sub>±0.03</sub>	21.2 <sub>±0.1</sub>	5.66 <sub>±0.05</sub>
	200	<b>19.4</b> <sub>±0.1</sub>	4.76 <sub>±0.03</sub>	20.9 <sub>±0.1</sub>	5.66 <sub>±0.04</sub>
	400	<b>19.0</b> <sub>±&lt;0.01</sub>	4.65 <sub>±0.05</sub>	22.3 <sub>±0.1</sub>	6.41 <sub>±0.06</sub>
ResNet-152	100	<b>19.2</b> <sub>±&lt;0.01</sub>	4.69 <sub>±0.04</sub>	20.4 <sub>±&lt;0.0</sub>	5.39 <sub>±0.06</sub>
	200	<b>18.5</b> <sub>±0.1</sub>	4.37 <sub>±0.03</sub>	20.3 <sub>±0.2</sub>	5.39 <sub>±0.07</sub>
	400	<b>18.4</b> <sub>±&lt;0.01</sub>	4.35 <sub>±0.04</sub>	20.9 <sub>±&lt;0.0</sub>	5.84 <sub>±0.07</sub>

Table: Test error rates for ResNets trained on ImageNet, with and without SAM.

# Results on ViT (and MLP-Mixer)

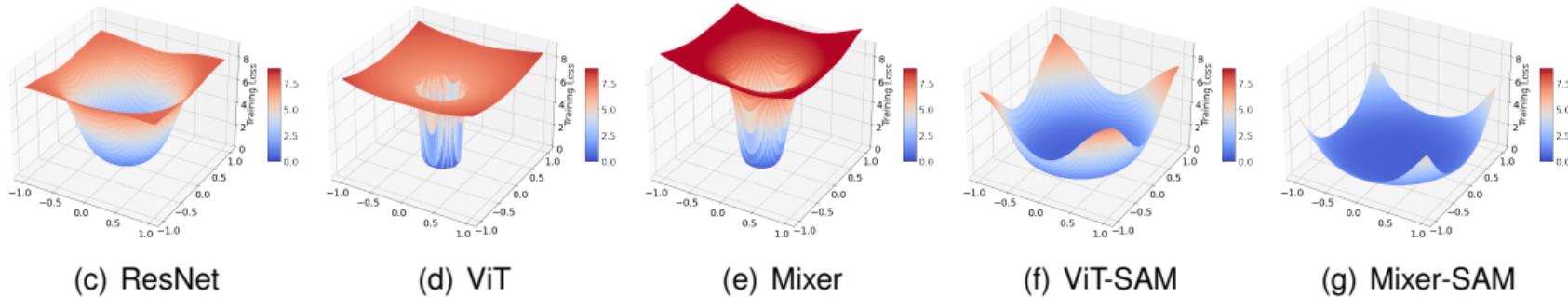


Figure: **Cross-entropy loss landscapes of ResNet-152, ViT-B/16, and Mixer-B/16.** ViT and MLP-Mixer converge to sharper regions than ResNet when trained on ImageNet with the basic Inception-style preprocessing. SAM significantly smooths the landscapes.

# Results on ViT (and MLP-Mixer)

Table: Number of parameters, Hessian dominate eigenvalue  $\lambda_{max}$ , training error at convergence  $L_{train}$ , average flatness  $L_{train}^N$ , accuracy on ImageNet, and accuracy/robustness on ImageNet-C. ViT and MLP-Mixer suffer divergent  $\kappa$  and converge at sharp regions; SAM rescues that and leads to better generalization.

	ResNet-152	ResNet-152-SAM	ViT-B/16	ViT-B/16-SAM	Mixer-B/16	Mixer-B/16-SAM
<b>#Params</b>	60M			87M		
<b>Hessian <math>\lambda_{max}</math></b>	179.8	<b>42.0</b>	738.8	<b>20.9</b>	1644.4	<b>22.5</b>
$L_{train}$	<b>0.86</b>	0.90	<b>0.65</b>	0.82	<b>0.45</b>	0.97
$L_{train}^N$ *	2.39	<b>2.16</b>	6.66	<b>0.96</b>	7.78	<b>1.01</b>
<b>ImageNet (%)</b>	78.5	<b>79.3</b>	74.6	<b>79.9</b>	66.4	<b>77.4</b>
<b>ImageNet-C (%)</b>	50.0	<b>52.2</b>	46.6	<b>56.5</b>	33.8	<b>48.8</b>

# Table of Contents

① Stochastic Gradient Descent (SGD) and Mini-batch SGD

② Accelerated and Stabilized Optimization Methods

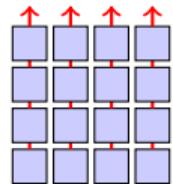
③ Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization
- Muon Optimizer

# Muon: From Vector to Matrix Optimization

## Element-wise Updates

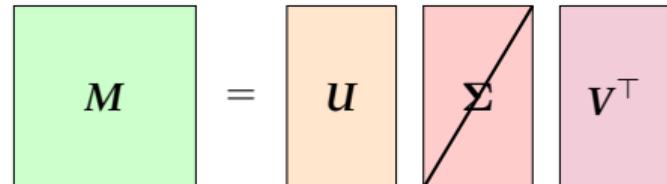
(SGD, Adam, AdaGrad)



Independent updates

## Matrix-aware Updates

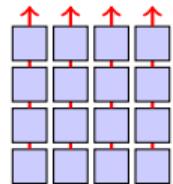
(Muon)

$$M = U \Sigma V^\top$$


# Muon: From Vector to Matrix Optimization

## Element-wise Updates

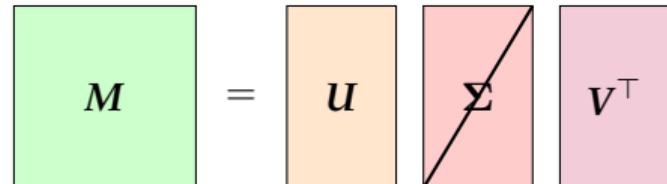
(SGD, Adam, AdaGrad)



Independent updates

## Matrix-aware Updates

(Muon)

$$M = U \Sigma V^\top$$


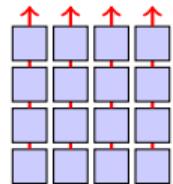
## Limitation

Ignores matrix structure;  
Treats all parameters as a flat vector

# Muon: From Vector to Matrix Optimization

## Element-wise Updates

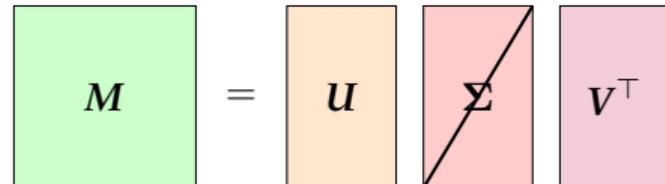
(SGD, Adam, AdaGrad)



Independent updates

## Matrix-aware Updates

(Muon)

$$M = U \Sigma V^T$$


### Limitation

Ignores matrix structure;  
Treats all parameters as a flat vector

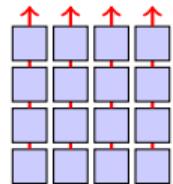
### Advantage

Leverages matrix structure;  
Captures correlations

# Muon: From Vector to Matrix Optimization

## Element-wise Updates

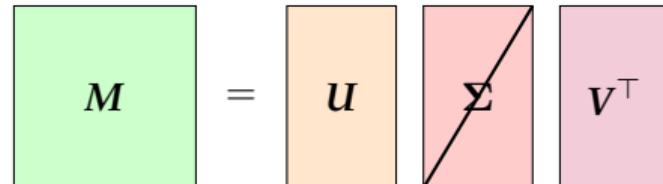
(SGD, Adam, AdaGrad)



Independent updates

## Matrix-aware Updates

(Muon)

$$M = U \Sigma V^T$$


### Limitation

Ignores matrix structure;  
Treats all parameters as a flat vector

### Advantage

Leverages matrix structure;  
Captures correlations

**Key Insight:** Neural networks are built on matrix multiplications.  
Matrix parameters have inherent structure that element-wise optimizers ignore!

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

Least Action

**Stable:** Small parameter change

**Fast:** Large loss decrease

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

Least Action

**Stable:** Small parameter change

**Fast:** Large loss decrease

Mathematically:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \text{ s.t. } \rho(\Delta W) \leq \eta$$

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

Least Action

**Stable:** Small parameter change

**Fast:** Large loss decrease

Mathematically:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \text{ s.t. } \rho(\Delta W) \leq \eta$$

**Question:** What is the right  $\rho(\cdot)$ ?

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

Least Action

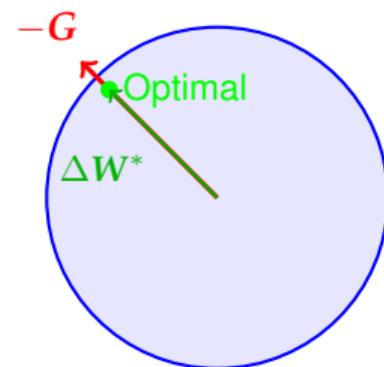
**Stable:** Small parameter change

**Fast:** Large loss decrease

Mathematically:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \text{ s.t. } \rho(\Delta W) \leq \eta$$

**Question:** What is the right  $\rho(\cdot)$ ?



$$\text{Constraint: } \|\Delta W\|_2 \leq \eta$$

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

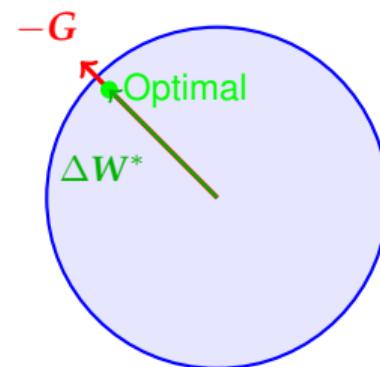
Least Action

**Stable:** Small parameter change

**Fast:** Large loss decrease

Mathematically:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \text{ s.t. } \rho(\Delta W) \leq \eta$$



$$\text{Constraint: } \|\Delta W\|_2 \leq \eta$$

**Question:** What is the right  $\rho(\cdot)$ ? **Answer:** For matrix parameters, use **spectral norm**  $\|\Delta W\|_2$

$$\|x \Delta W\| \leq \|\Delta W\|_2 \|x\| \quad (\text{tightest bound!})$$

# Why Spectral Norm? The Least Action Principle

## Optimization Principle:

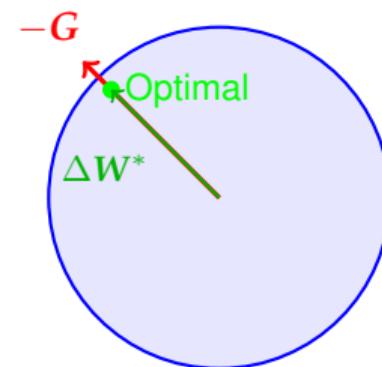
Least Action

**Stable:** Small parameter change

**Fast:** Large loss decrease

Mathematically:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \text{ s.t. } \rho(\Delta W) \leq \eta$$



Constraint:  $\|\Delta W\|_2 \leq \eta$

**Question:** What is the right  $\rho(\cdot)$ ? **Answer:** For matrix parameters, use **spectral norm**  $\|\Delta W\|_2$

$$\|x\Delta W\| \leq \|\Delta W\|_2 \|x\| \quad (\text{tightest bound!})$$

## Key Insight

Spectral norm precisely measures how much  $\Delta W$  changes the output  $y = xW$ , while Frobenius norm (used by SGD) is looser!

# From Spectral Norm to Matrix Sign Function

Recall the optimization problem:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \quad \text{s.t.} \quad \|\Delta W\|_2 \leq \eta$$

# From Spectral Norm to Matrix Sign Function

Recall the optimization problem:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \quad \text{s.t.} \quad \|\Delta W\|_2 \leq \eta$$

Solution via Lagrange multipliers:

$$\Delta W^* = -\eta \cdot \frac{G}{\|G\|_2} = -\eta \cdot \text{msign}(G)$$

where the optimal direction is the normalized gradient under spectral norm.

# From Spectral Norm to Matrix Sign Function

Recall the optimization problem:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \quad \text{s.t.} \quad \|\Delta W\|_2 \leq \eta$$

Solution via Lagrange multipliers:

$$\Delta W^* = -\eta \cdot \frac{G}{\|G\|_2} = -\eta \cdot \text{msign}(G)$$

where the optimal direction is the normalized gradient under spectral norm.

**With momentum:** Replace  $G$  with momentum  $M_t = \beta M_{t-1} + G_t$

$$\Delta W^* = -\eta \cdot \text{msign}(M_t)$$

# From Spectral Norm to Matrix Sign Function

Recall the optimization problem:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \quad \text{s.t.} \quad \|\Delta W\|_2 \leq \eta$$

Solution via Lagrange multipliers:

$$\Delta W^* = -\eta \cdot \frac{G}{\|G\|_2} = -\eta \cdot \text{msign}(G)$$

where the optimal direction is the normalized gradient under spectral norm.

**With momentum:** Replace  $G$  with momentum  $M_t = \beta M_{t-1} + G_t$

$$\Delta W^* = -\eta \cdot \text{msign}(M_t)$$

Add weight decay:

$$W_{t+1} = W_t + \Delta W^* - \eta \lambda W_t$$

# From Spectral Norm to Matrix Sign Function

Recall the optimization problem:

$$\min_{\Delta W} \text{Tr}(G^\top \Delta W) \quad \text{s.t.} \quad \|\Delta W\|_2 \leq \eta$$

Solution via Lagrange multipliers:

$$\Delta W^* = -\eta \cdot \frac{G}{\|G\|_2} = -\eta \cdot \text{msign}(G)$$

where the optimal direction is the normalized gradient under spectral norm.

**With momentum:** Replace  $G$  with momentum  $M_t = \beta M_{t-1} + G_t$

$$\Delta W^* = -\eta \cdot \text{msign}(M_t)$$

Add weight decay:

$$W_{t+1} = W_t + \Delta W^* - \eta \lambda W_t$$

# Understanding the Matrix Sign Function

**Definition via SVD:**

$$M = U\Sigma V^\top$$

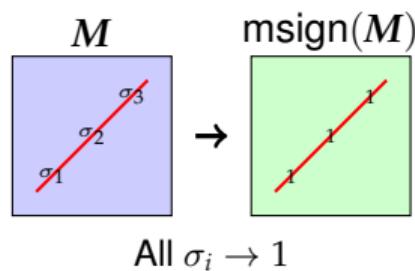
$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$

# Understanding the Matrix Sign Function

**Definition via SVD:**

$$M = U\Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$



# Understanding the Matrix Sign Function

Definition via SVD:

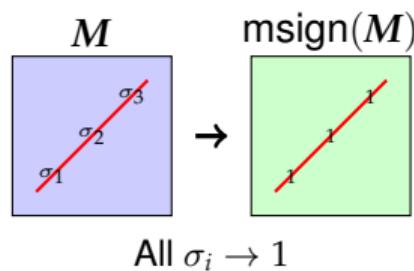
$$M = U \Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$

Special Cases:

- Scalar:

$$\text{sign}(x) = x/|x|$$



# Understanding the Matrix Sign Function

Definition via SVD:

$$M = U\Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$

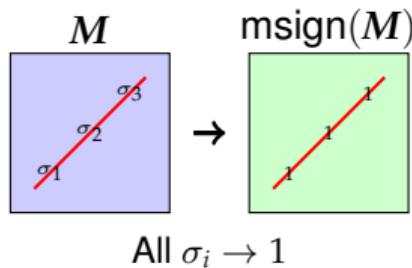
Special Cases:

- **Scalar:**

$$\text{sign}(x) = x/|x|$$

- **Vector ( $n \times 1$ ):**

$$\text{msign}(\mathbf{m}) = \frac{\mathbf{m}}{\|\mathbf{m}\|_2}$$

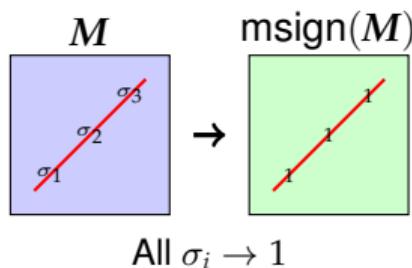


# Understanding the Matrix Sign Function

Definition via SVD:

$$M = U\Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$



Special Cases:

- **Scalar:**

$$\text{sign}(x) = x/|x|$$

- **Vector ( $n \times 1$ ):**

$$\text{msign}(\mathbf{m}) = \frac{\mathbf{m}}{\|\mathbf{m}\|_2}$$

- **Diagonal Matrix:**

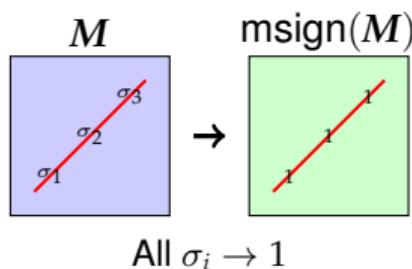
$$\text{msign}(\text{diag}(\mathbf{m})) = \text{diag}(\text{sign}(\mathbf{m}))$$

# Understanding the Matrix Sign Function

**Definition via SVD:**

$$M = U\Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$



**Special Cases:**

- **Scalar:**

$$\text{sign}(x) = x/|x|$$

- **Vector ( $n \times 1$ ):**

$$\text{msign}(\mathbf{m}) = \frac{\mathbf{m}}{\|\mathbf{m}\|_2}$$

- **Diagonal Matrix:**

$$\text{msign}(\text{diag}(\mathbf{m})) = \text{diag}(\text{sign}(\mathbf{m}))$$

- **General Matrix:**

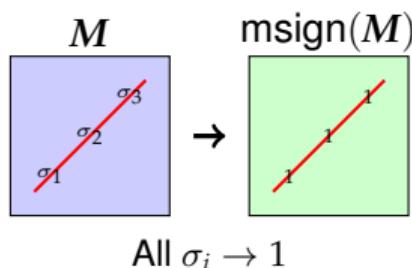
Optimal orthogonal approximation

# Understanding the Matrix Sign Function

**Definition via SVD:**

$$M = U \Sigma V^\top$$

$$\text{msign}(M) = U_{[:, :r]} V_{[:, :r]}^\top$$



**Special Cases:**

- **Scalar:**

$$\text{sign}(x) = x/|x|$$

- **Vector ( $n \times 1$ ):**

$$\text{msign}(\mathbf{m}) = \frac{\mathbf{m}}{\|\mathbf{m}\|_2}$$

- **Diagonal Matrix:**

$$\text{msign}(\text{diag}(\mathbf{m})) = \text{diag}(\text{sign}(\mathbf{m}))$$

- **General Matrix:**

Optimal orthogonal approximation

**Alternative Form:**  $\text{msign}(M) = M(M^\top M)^{-1/2}$

**Interpretation:** "Optimal orthogonal approximation" - closest orthogonal matrix to  $M$

# Muon Update Rule

**Muon** (MomentUm Orthogonalized by Newton-schulz):

$$\mathbf{M}_t = \beta \mathbf{M}_{t-1} + \mathbf{G}_t \quad (\text{Momentum})$$

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta_t [\text{msign}(\mathbf{M}_t) + \lambda \mathbf{W}_{t-1}] \quad (\text{Update})$$

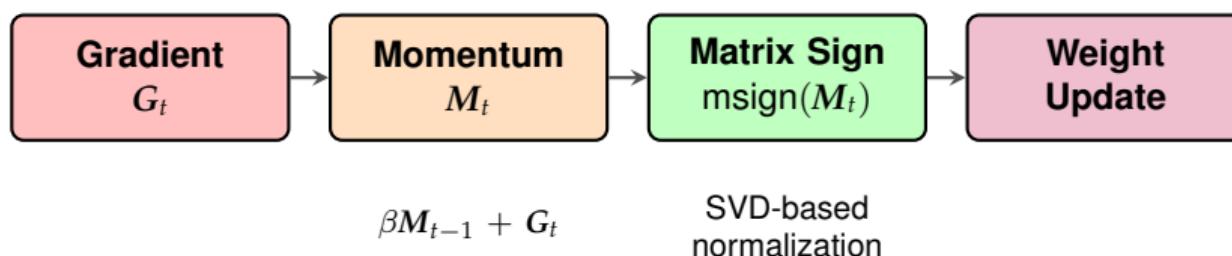
# Muon Update Rule

**Muon** (MomentUm Orthogonalized by Newton-schulz):

$$\mathbf{M}_t = \beta \mathbf{M}_{t-1} + \mathbf{G}_t \quad (\text{Momentum})$$

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta_t [\text{msign}(\mathbf{M}_t) + \lambda \mathbf{W}_{t-1}] \quad (\text{Update})$$

Update Pipeline:



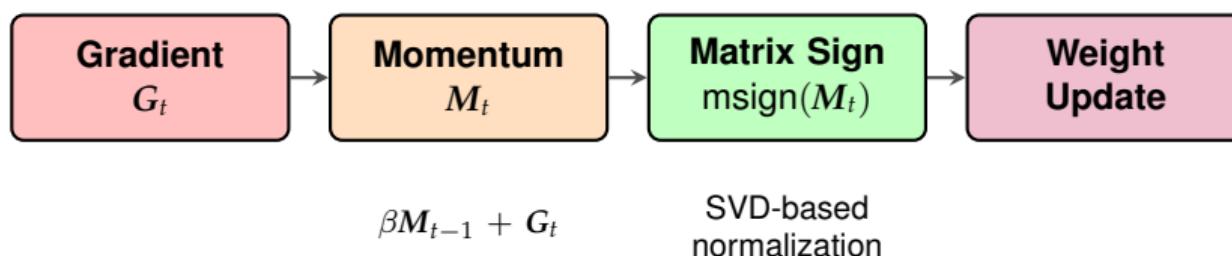
# Muon Update Rule

**Muon** (MomentUm Orthogonalized by Newton-schulz):

$$\mathbf{M}_t = \beta \mathbf{M}_{t-1} + \mathbf{G}_t \quad (\text{Momentum})$$

$$\mathbf{W}_t = \mathbf{W}_{t-1} - \eta_t [\text{msign}(\mathbf{M}_t) + \lambda \mathbf{W}_{t-1}] \quad (\text{Update})$$

Update Pipeline:



- $\text{msign}(\cdot)$ : **Matrix sign function** (not element-wise sign!)
- $\lambda$ : Weight decay coefficient
- $\beta$ : Momentum parameter (typically 0.95)

# Muon vs Other Optimizers

Optimizer	Update Rule	Constraint	Scale Inv.	Memory
SGD	$-\eta G$	$F$ -norm	✓	Low
Adam	$-\eta \frac{M}{\sqrt{V}}$	Element-wise	✓	High
SignSGD	$-\eta \text{sign}(M)$	Element-wise	✓	Med
<b>Muon</b>	$-\eta m \text{sign}(M)$	Spectral	✓	Med

# Muon vs Other Optimizers

Optimizer	Update Rule	Constraint	Scale Inv.	Memory
SGD	$-\eta G$	$F$ -norm	✓	Low
Adam	$-\eta \frac{M}{\sqrt{V}}$	Element-wise	✓	High
SignSGD	$-\eta \text{sign}(M)$	Element-wise	✓	Med
<b>Muon</b>	$-\eta m \text{sign}(M)$	Spectral	✓	Med

**Element-wise:**



Scattered

# Muon vs Other Optimizers

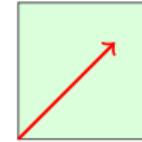
Optimizer	Update Rule	Constraint	Scale Inv.	Memory
SGD	$-\eta G$	$F$ -norm	✓	Low
Adam	$-\eta \frac{M}{\sqrt{V}}$	Element-wise	✓	High
SignSGD	$-\eta \text{sign}(M)$	Element-wise	✓	Med
<b>Muon</b>	$-\eta m \text{sign}(M)$	Spectral	✓	Med

Element-wise:



Scattered

Matrix-aware (Muon):



Structured

# Muon vs Other Optimizers

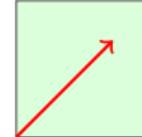
Optimizer	Update Rule	Constraint	Scale Inv.	Memory
SGD	$-\eta G$	$F$ -norm	✓	Low
Adam	$-\eta \frac{M}{\sqrt{V}}$	Element-wise	✓	High
SignSGD	$-\eta \text{sign}(M)$	Element-wise	✓	Med
<b>Muon</b>	$-\eta m \text{sign}(M)$	Spectral	✓	Med

**Element-wise:**



Scattered

**Matrix-aware (Muon):**



Structured

## Muon Advantages

- Spectral norm = tightest constraint for matrix stability
- Captures parameter correlations via SVD
- Lower memory than Adam (no second moment)

# Key Properties of Muon

1. Scale Invariance

2. Isotropy

# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path

## 2. Isotropy

# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$

## 2. Isotropy

# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
- $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$

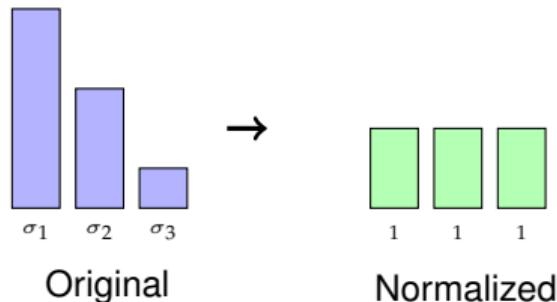
## 2. Isotropy

# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
- $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$

## 2. Isotropy



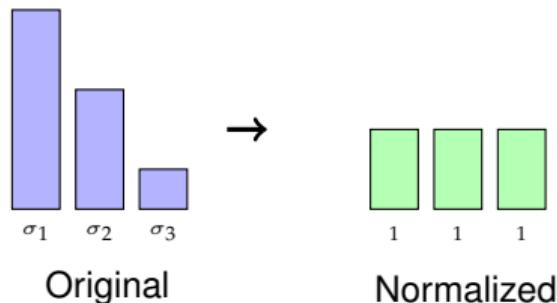
# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
- $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$

## 2. Isotropy

- Different singular values  $\rightarrow$  anisotropic



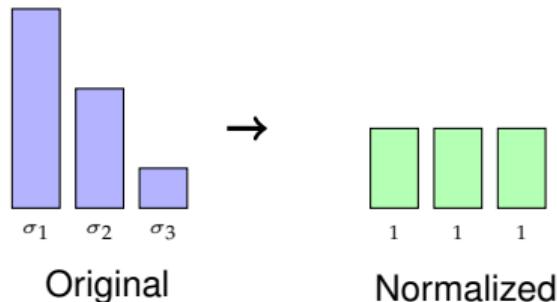
# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
- $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$

## 2. Isotropy

- Different singular values  $\rightarrow$  anisotropic
- All become 1  $\rightarrow$  isotropic updates



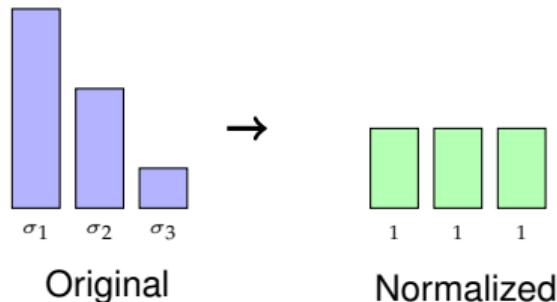
# Key Properties of Muon

## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
- SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
- $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$

## 2. Isotropy

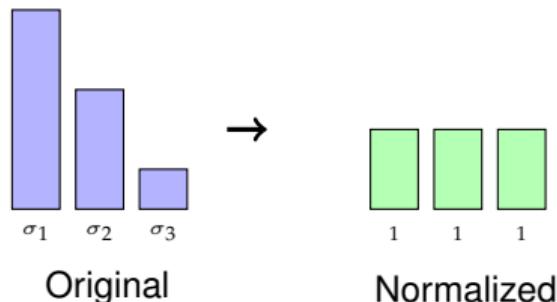
- Different singular values  $\rightarrow$  anisotropic
- All become 1  $\rightarrow$  isotropic updates
- Synchronizes magnitudes across directions



## Key Properties of Muon

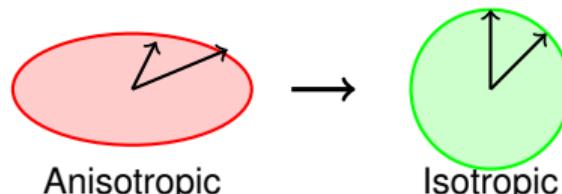
## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
  - SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
  - $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$



## 2. Isotropy

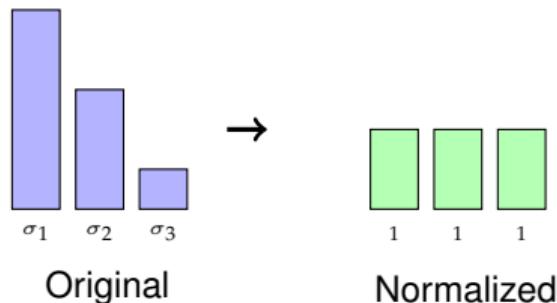
- Different singular values  $\rightarrow$  anisotropic
  - All become 1  $\rightarrow$  isotropic updates
  - Synchronizes magnitudes across directions



## Key Properties of Muon

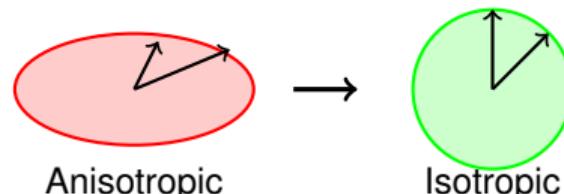
## 1. Scale Invariance

- Loss  $\times \lambda$  doesn't change optimization path
  - SVD:  $M \rightarrow \lambda M = U(\lambda \Sigma)V^\top$
  - $\text{msign}(\lambda M) = UV^\top = \text{msign}(M)$



## 2. Isotropy

- Different singular values  $\rightarrow$  anisotropic
  - All become 1  $\rightarrow$  isotropic updates
  - Synchronizes magnitudes across directions



**Connection:** Similar to BERT-whitening, Muon "whitens" the gradient distribution!

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

## Solution: Newton-Schulz Iteration

$$\mathbf{X}_0 = \mathbf{M}$$

$$\mathbf{X}_{t+1} = \frac{15}{8}\mathbf{X}_t - \frac{5}{4}\mathbf{X}_t(\mathbf{X}_t^\top \mathbf{X}_t) + \frac{3}{8}\mathbf{X}_t(\mathbf{X}_t^\top \mathbf{X}_t)^2$$

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

## Solution: Newton-Schulz Iteration

$$\mathbf{X}_0 = \mathbf{M}$$

$$\mathbf{X}_{t+1} = \frac{15}{8}\mathbf{X}_t - \frac{5}{4}\mathbf{X}_t(\mathbf{X}_t^\top \mathbf{X}_t) + \frac{3}{8}\mathbf{X}_t(\mathbf{X}_t^\top \mathbf{X}_t)^2$$

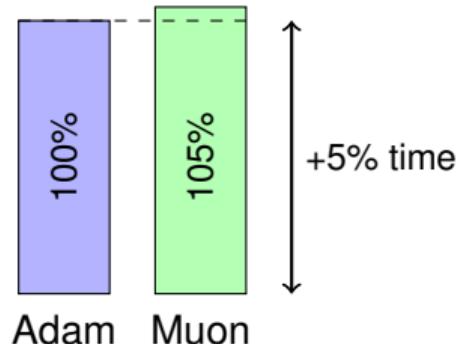
Converges to  $\text{msign}(\mathbf{M})$  in 5 iterations

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

## Computational Cost:



## Solution: Newton-Schulz Iteration

$$X_0 = M$$

$$X_{t+1} = \frac{15}{8}X_t - \frac{5}{4}X_t(X_t^\top X_t) + \frac{3}{8}X_t(X_t^\top X_t)^2$$

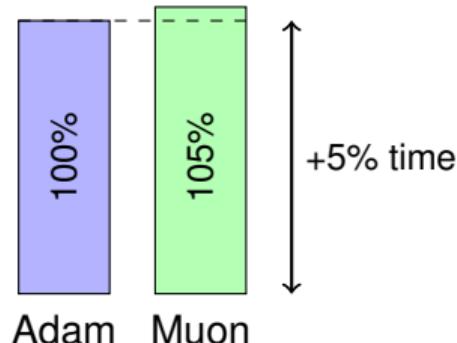
Converges to  $\text{msign}(M)$  in 5 iterations

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

## Computational Cost:



## Solution: Newton-Schulz Iteration

$$X_0 = M$$

$$X_{t+1} = \frac{15}{8}X_t - \frac{5}{4}X_t(X_t^\top X_t) + \frac{3}{8}X_t(X_t^\top X_t)^2$$

Converges to  $\text{msign}(M)$  in 5 iterations

## Memory Comparison:

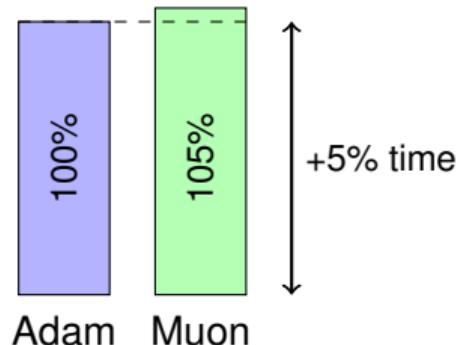
- Adam: 2 states ( $M, V$ )
- Muon: 1 state ( $M$ )

# Efficient Implementation via Newton-Schulz

## Challenge:

- Direct SVD is expensive:  $\mathcal{O}(n^3)$
- Need faster approximation

## Computational Cost:



## Solution: Newton-Schulz Iteration

$$\begin{aligned} X_0 &= M \\ X_{t+1} &= \frac{15}{8}X_t - \frac{5}{4}X_t(X_t^\top X_t) + \frac{3}{8}X_t(X_t^\top X_t)^2 \end{aligned}$$

Converges to  $\text{msign}(M)$  in 5 iterations

## Memory Comparison:

- Adam: 2 states ( $M, V$ )
- Muon: 1 state ( $M$ )

## Practical Consideration

Matrix multiplications happen during idle time (between gradient computations), so actual wall-clock overhead is minimal ( $\approx 2\text{-}5\%$ )!

- [1] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of machine learning research*, 12(7), 2011.
- [2] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur. Sharpness-aware minimization for efficiently improving generalization. In *International Conference on Learning Representations*, 2021.
- [3] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [4] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems*, 31, 2018.
- [5] I. Loshchilov and F. Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [6] M. Zhang, J. Lucas, J. Ba, and G. E. Hinton. Lookahead optimizer: k steps forward, 1 step back. *Advances in neural information processing systems*, 32, 2019.