

Optimization for Deep Learning

Tao LIN

Learning and INference Systems (LINs) Lab, Westlake University

March 13, 2023



Table of Contents

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

Background

- ➊ Get data: ξ_1, \dots, ξ_N , where $\xi_i := (\mathbf{d}, y)_i$

Background

- ① Get data: ξ_1, \dots, ξ_N , where $\xi_i := (\mathbf{d}, y)_i$
- ② Choose a classifier

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

Background

- ① Get data: ξ_1, \dots, ξ_N , where $\xi_i := (\mathbf{d}, y)_i$
- ② Choose a classifier

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

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

Background

- ① Get data: ξ_1, \dots, ξ_N , where $\xi_i := (\mathbf{d}, y)_i$
- ② Choose a classifier

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

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

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

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 $\xi_i := (\mathbf{d}_i, y_i)$

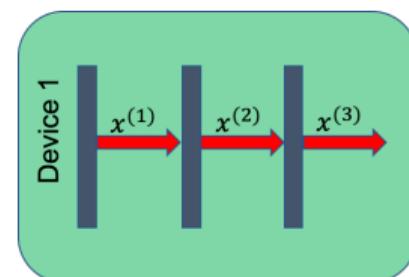
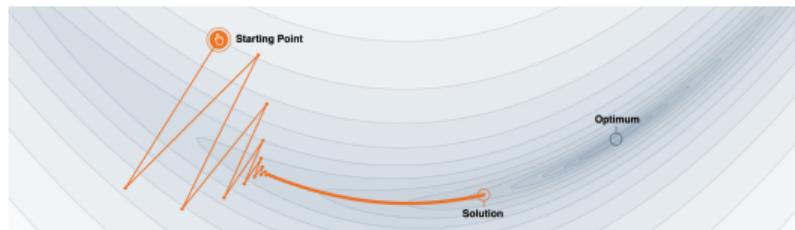
Background

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)

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



Background

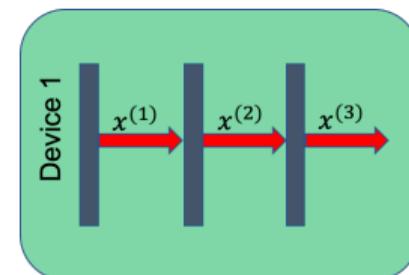
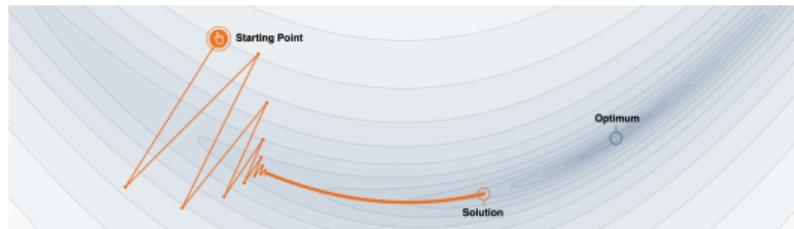
Finite-sum empirical risk minimization problem:

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

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

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

- η is step-size/learning rate



Background

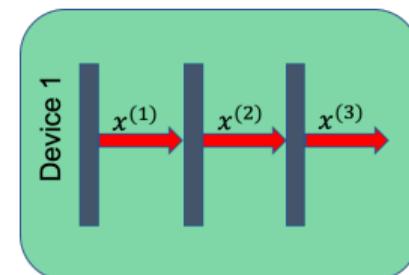
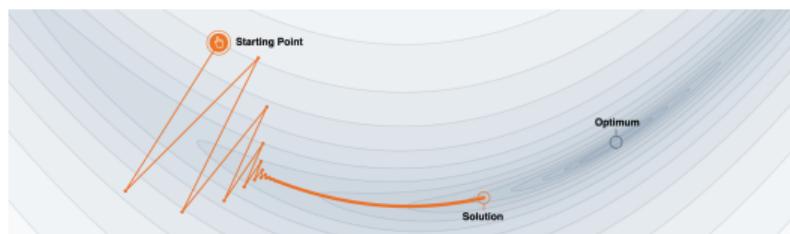
Finite-sum empirical risk minimization problem:

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

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

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

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



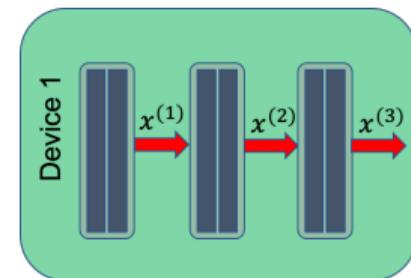
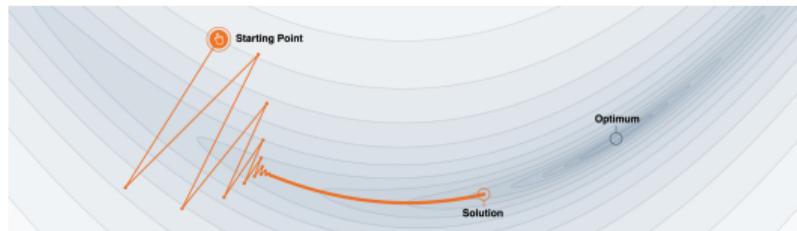
Background

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 $\xi_i := (\mathbf{d}_i, y_i)$
- 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

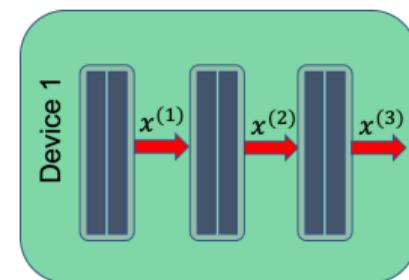
Finite-sum empirical risk minimization problem:

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

- The loss function of i -th data $\xi_i := (\mathbf{d}_i, y_i)$
- 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 $\mathbb{E}[v_i] = 1$ 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 $\mathbb{E}[v_i] = 1 \quad \text{for } i = 1, \dots, N.$

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

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 $\mathbb{E}[v_i] = 1$ 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

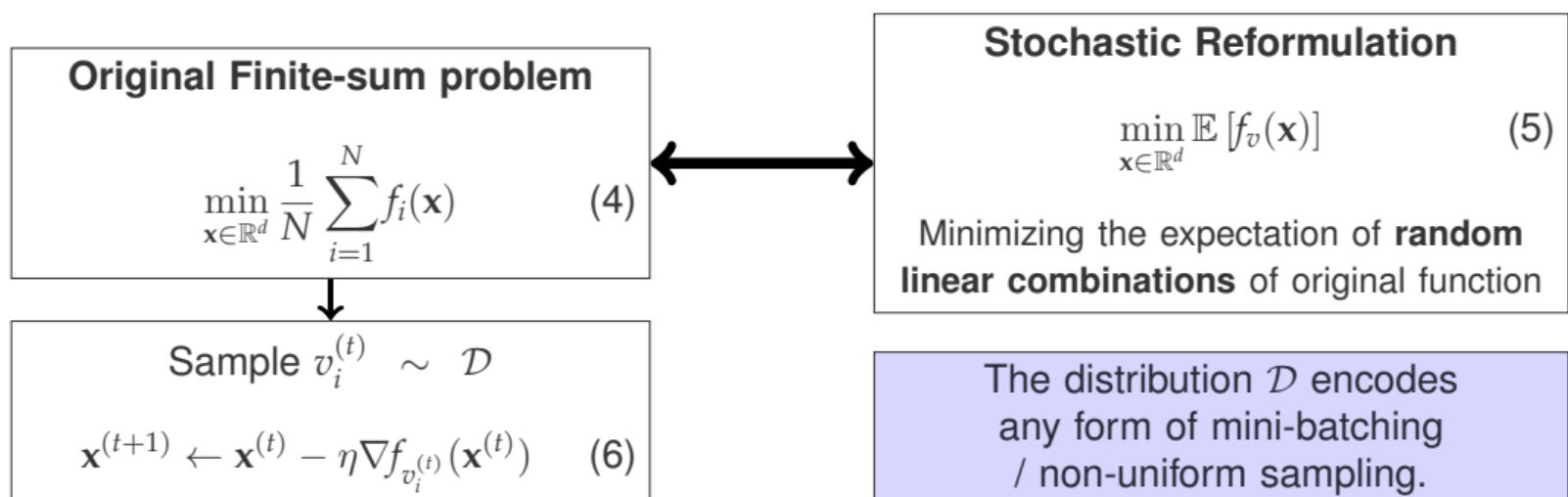
↓
Sample $v_i^{(t)} \sim \mathcal{D}$

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

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 $i = 1, \dots, N$.

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



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)

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

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

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

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

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

- T : number of iterations

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

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

- T : number of iterations

- σ : 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

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

- T : number of iterations

- σ : 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*

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

- *T: number of iterations*

- *σ : stochastic gradient variance*

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

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

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*

- *σ : stochastic gradient variance*

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

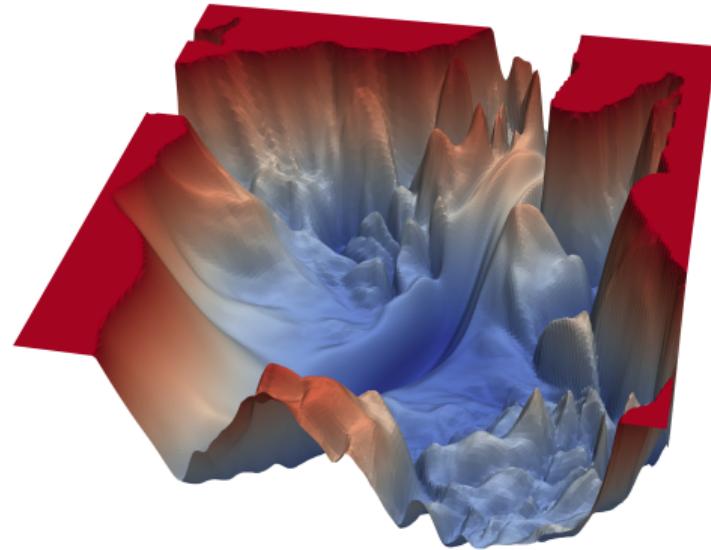
- When iterations $T \rightarrow \infty$, it holds that $\mathbb{E} \left[\|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \rightarrow 0$

- $\mathbb{E} \left[\|\nabla f(\mathbf{x}^{(t)})\|^2 \right] \rightarrow 0$ implies the sequence converges to a stationary solution

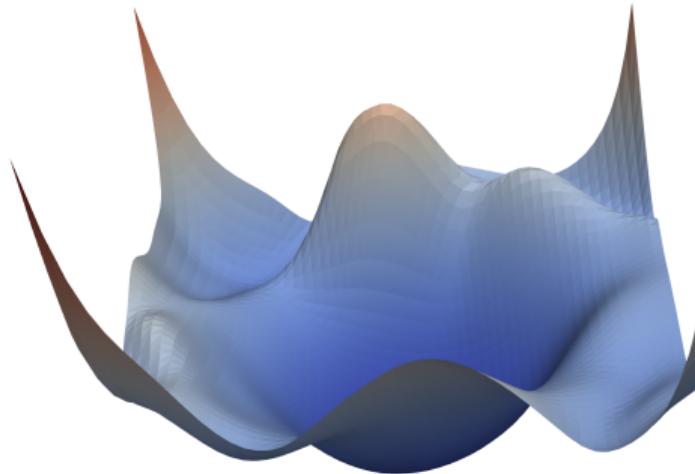
Table of Contents

- 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



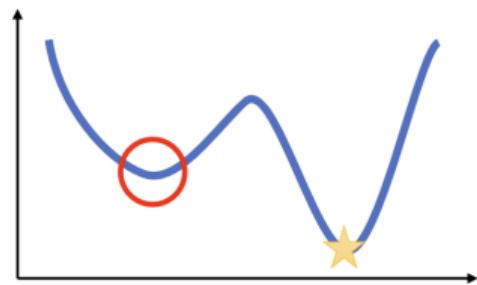
(a) ResNet w/o skip connections.



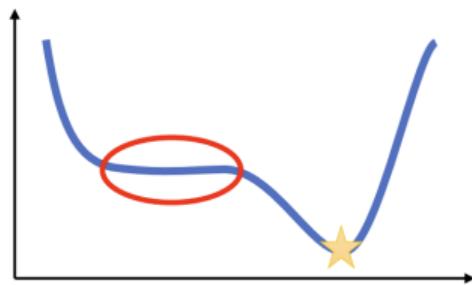
(b) ResNet w/ skip connections.

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

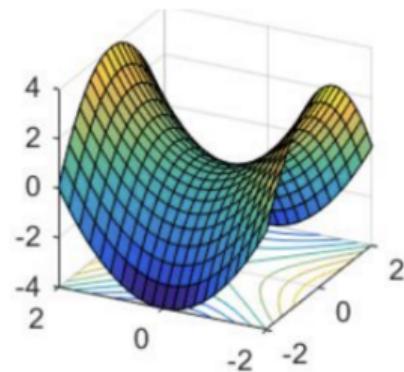
Issues of SGD—from the perspective of loss landscape



the local optimum



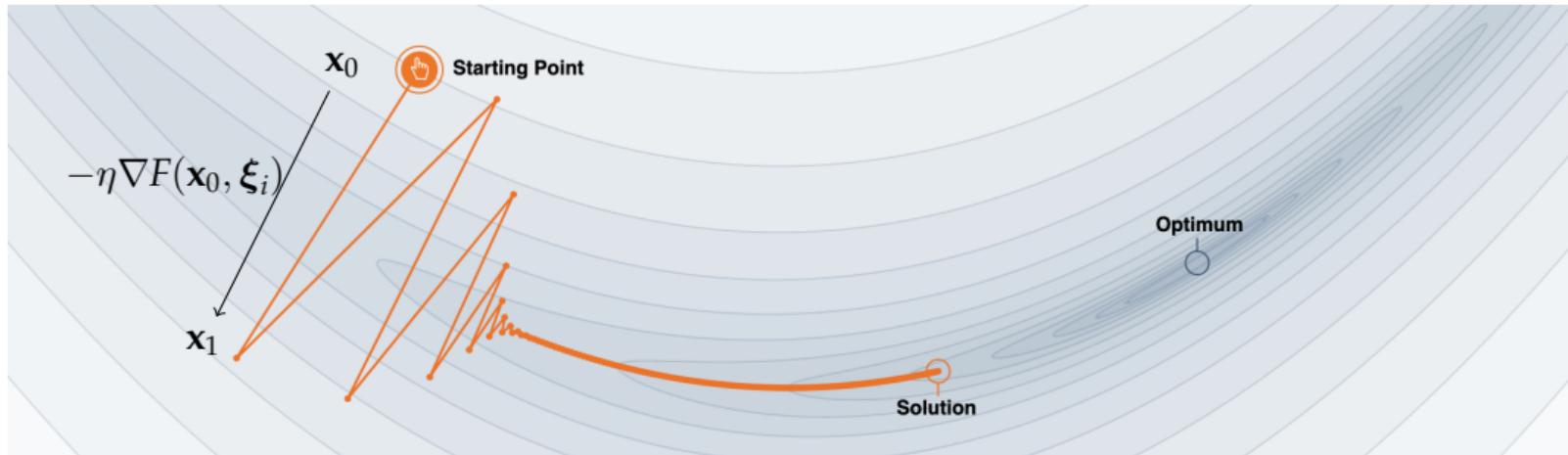
the plateau



the saddle point

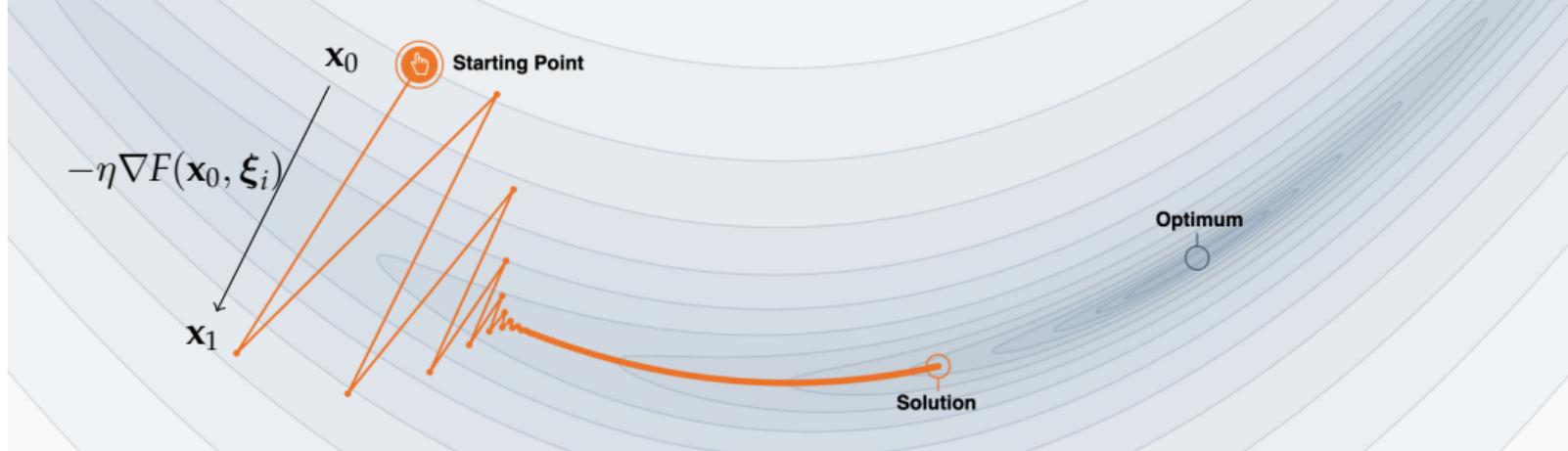
Challenging optimization loss landscape!

Issues of SGD—from the perspective of loss landscape



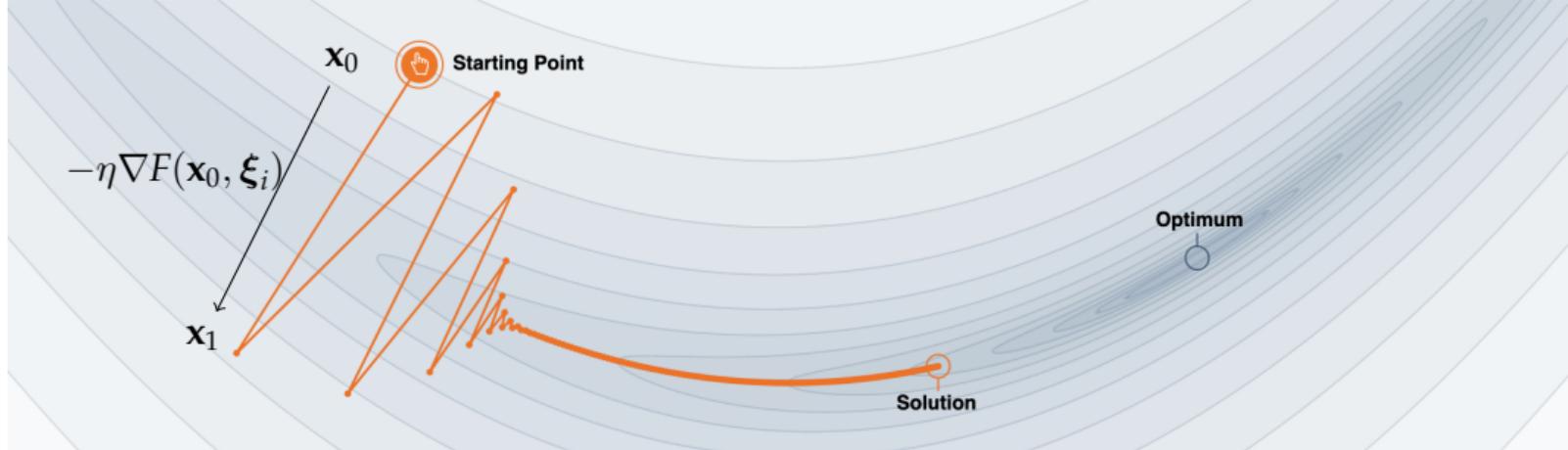
- **Challenges # 1:** loss function has high condition number.
→ very slow progress along shallow dimension, jitter along steep direction.

Issues of SGD—from the perspective of loss landscape



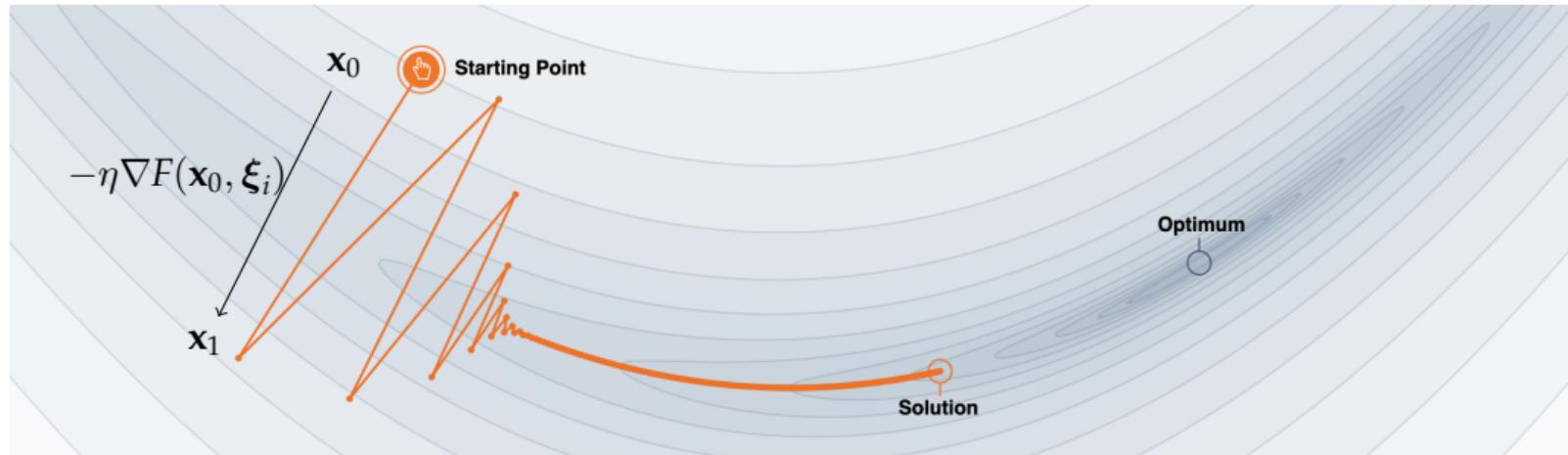
- **Challenges # 1:** loss function has high condition number.
→ very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.

Issues of SGD—from the perspective of loss landscape



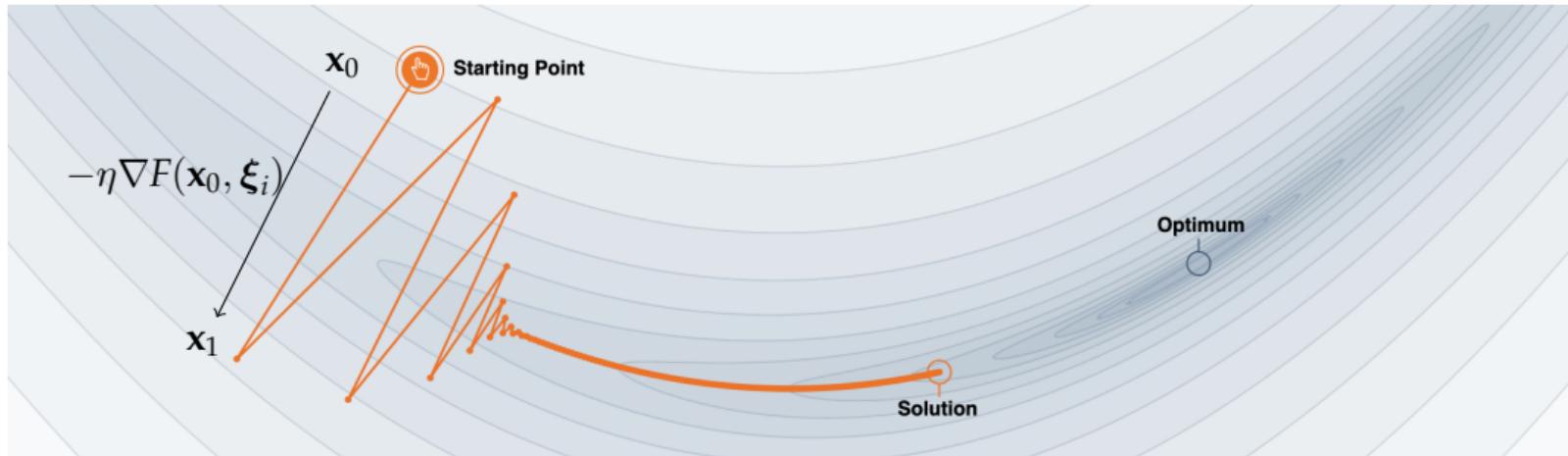
- **Challenges # 1:** loss function has high condition number.
→ very slow progress along shallow dimension, jitter along steep direction.
- **Challenges # 2 & more:** plateaus & saddle points.
→ cannot just choose tiny learning rates to prevent oscillation!

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.

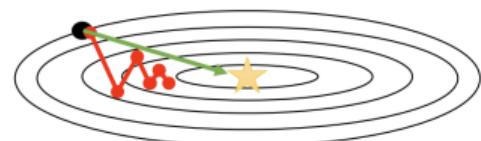
Improvement directions: leveraging the curvature information

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?

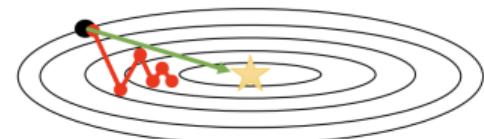
Yes! By leveraging the curvature information through Newton's method.

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:

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



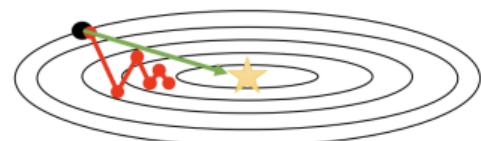
Improvement directions: leveraging the curvature information

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?

¹if 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)¹: $\mathcal{O}(N^3)$

¹if 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): $\mathcal{O}(N)$

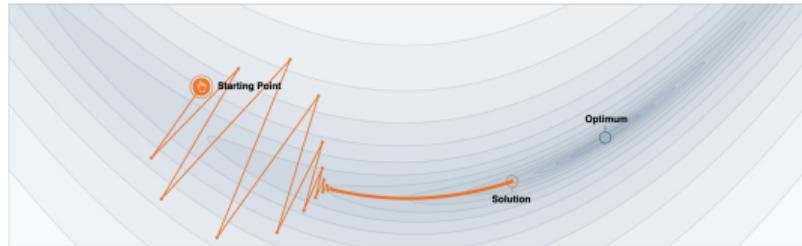
v.s.

GD (w/ Hessian)¹: $\mathcal{O}(N^3)$

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

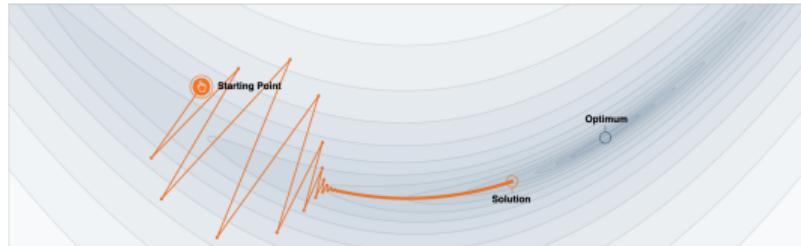
¹if using naive approach, though fancy methods can be much faster if they avoid forming the Hessian explicitly.

Momentum method

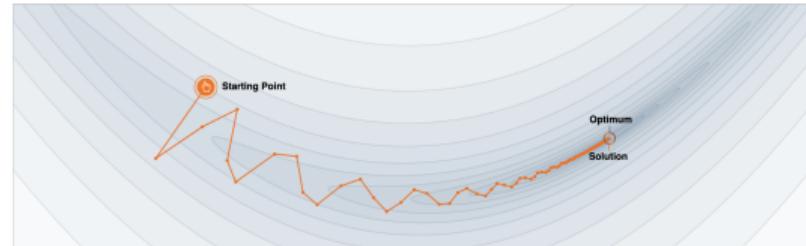


w/o momentum

Momentum method

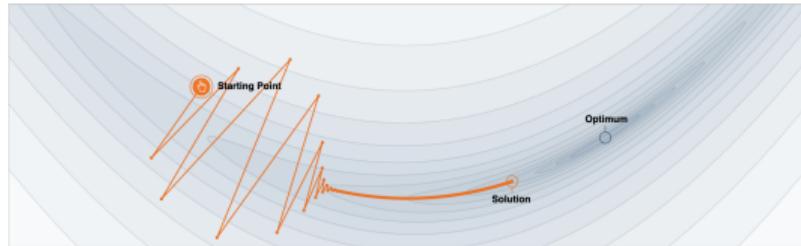


w/o momentum

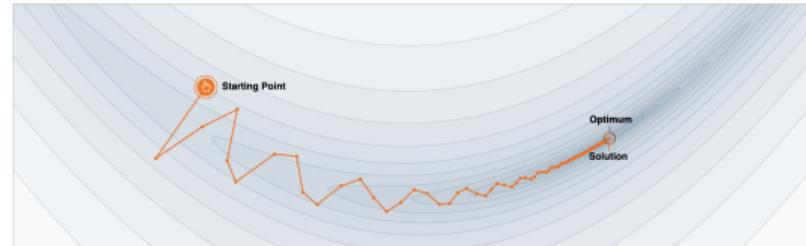


w/ momentum

Momentum method



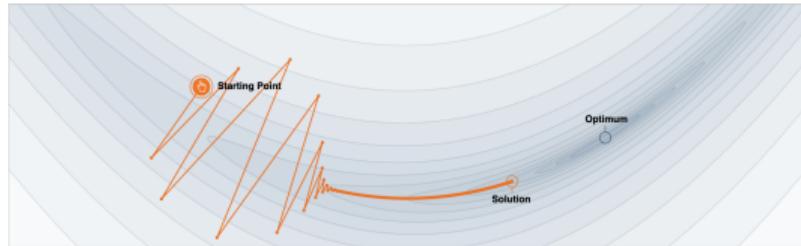
w/o momentum



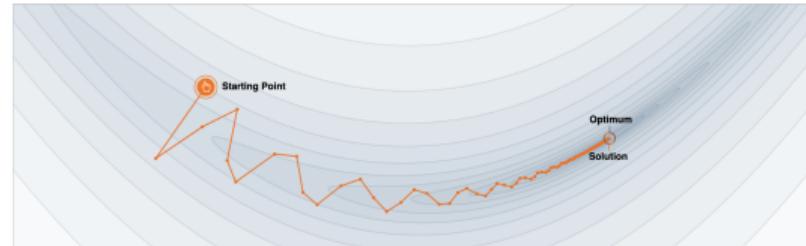
w/ momentum

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

Momentum method



w/o momentum

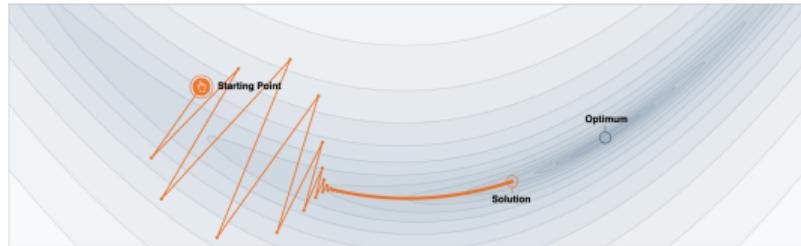


w/ momentum

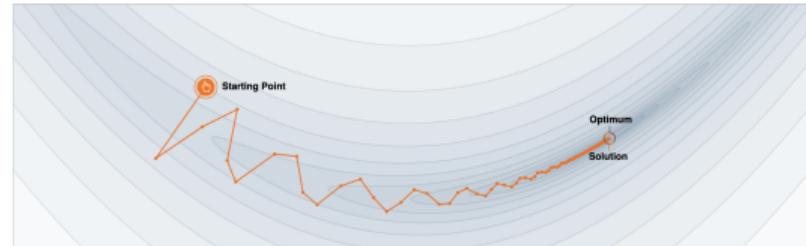
Intuition: averaging together successive gradients yield a much better direction!

- if successive gradient step point in different directions

Momentum method



w/o momentum

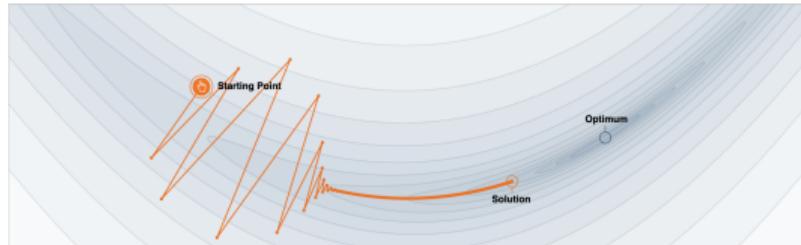


w/ momentum

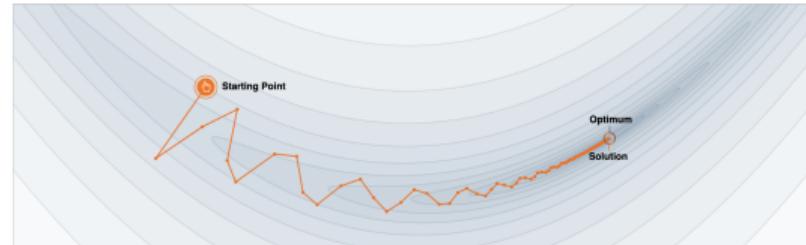
Intuition: averaging together successive gradients yield a much better direction!

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

Momentum method



w/o momentum

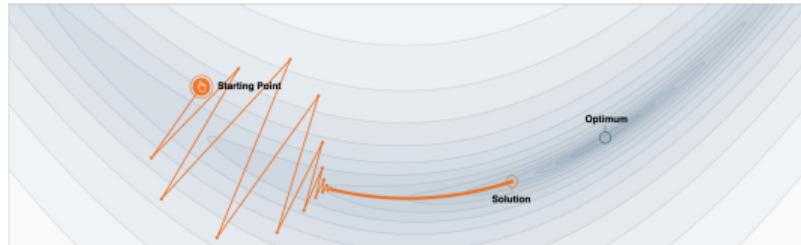


w/ momentum

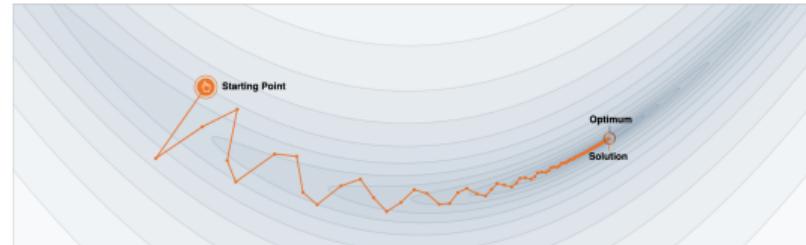
Intuition: averaging together successive gradients yield a much better direction!

- 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

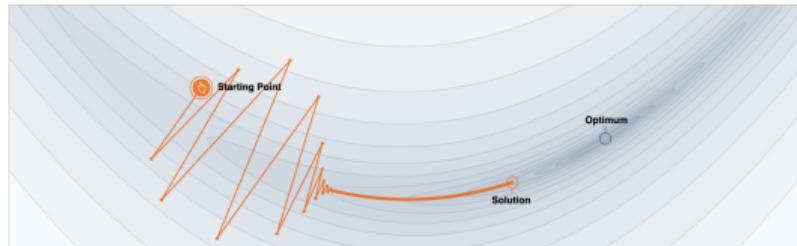


w/ momentum

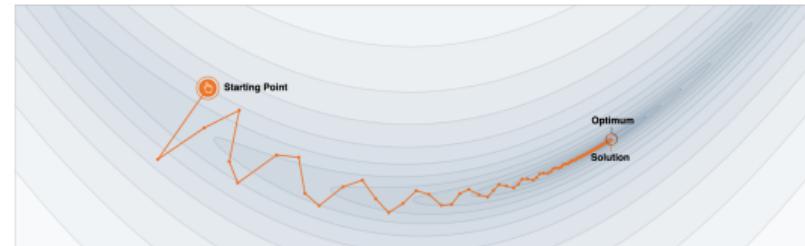
Intuition: averaging together successive gradients yield a much better direction!

- 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

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

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

Methods that manipulate gradient scale

Intuition behind $\nabla F(\mathbf{x}_i, \boldsymbol{\xi}_i)$:

- *sign:*

Methods that manipulate gradient scale

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

- *sign:*
- *magnitude:*

Methods that manipulate gradient scale

Intuition behind $\nabla F(\mathbf{x}_i, \boldsymbol{\xi}_i)$:

- *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):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad \text{(roughly the squared length of each dimension)}$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad \text{(each dimension is divided by its magnitude)}$$

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

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

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

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

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

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

Remarks:

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

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

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

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):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

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):

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

Remarks:

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

$$\mathbf{v}_t = \mathbf{v}_{t-1} + (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

RMSProp (estimate per-dimension magnitude):

$$\mathbf{v}_t = \beta \mathbf{v}_{t-1} + (1 - \beta) (\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t))^2 \quad (\text{roughly the squared length of each dimension})$$

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta \frac{\nabla F(\mathbf{x}_t, \boldsymbol{\xi}_t)}{\sqrt{\mathbf{v}_t}} \quad (\text{each dimension is divided by its magnitude})$$

Remarks:

- AdaGrad has some appealing guarantees for convex problems.
 - AdaGrad originally proposed to benefit from sparse data.
 - 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 $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

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

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

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

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 $\frac{\text{average gradient}}{\sqrt{\text{average squared gradient}}}$

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

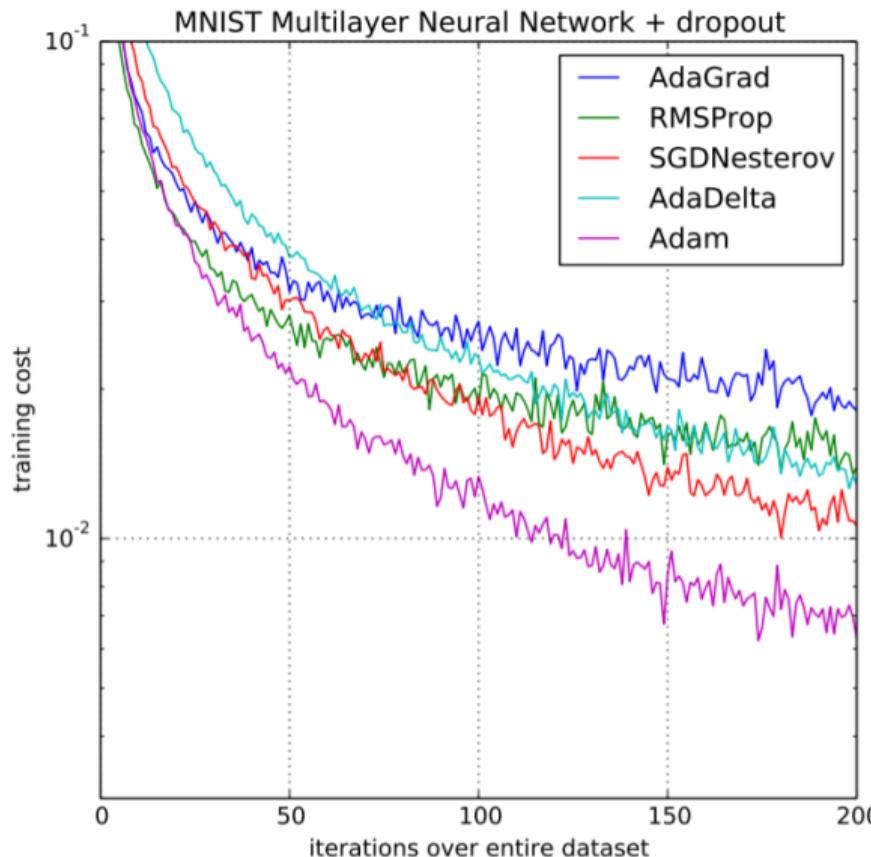
$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla F(\mathbf{x}_t, \xi_t))^2 \quad (\text{second moment estimate})$$

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

where compared to RMSProp, Adam

- replaces $\frac{\eta}{\sqrt{\mathbf{v}_t + \epsilon}} \nabla F(\mathbf{x}_t, \xi_t)$ 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 β_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:

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)

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

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

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)

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

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

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

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

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

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

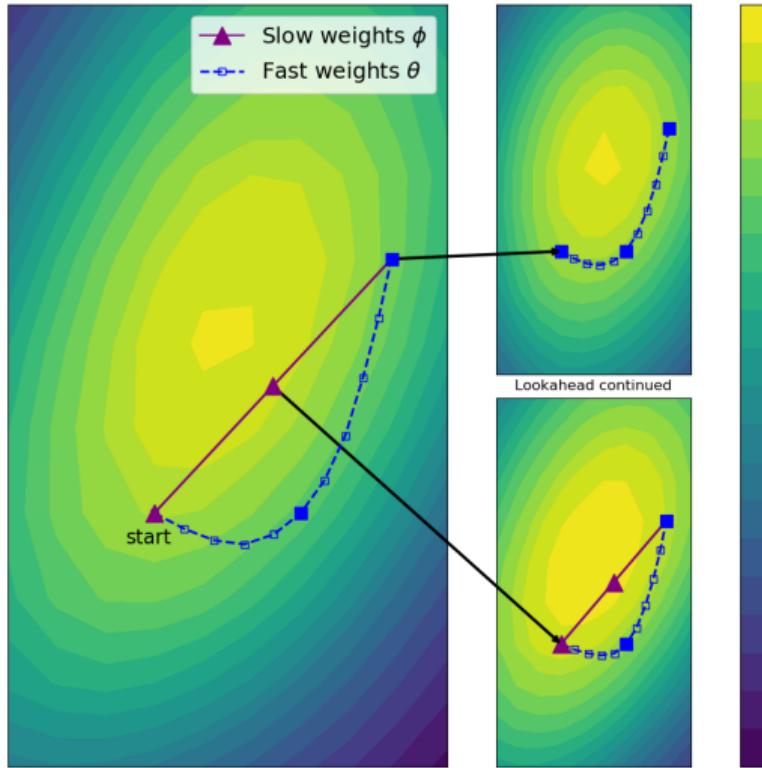
Different from the ideas e.g.,

- **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

SGD continued



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**

$v_0^{(t)} = \theta_{t-1}$; Inner-loop optimization

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

$v_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

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

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**

$v_0^{(t)} = \theta_{t-1}$; Inner-loop optimization

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

$v_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

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

end

Output : $\theta_{\mathcal{A}, \mathcal{S}} = \theta_T$ outer-loop optimization

Algorithm 1: SGD

Input : Objective $F_S(\theta)$, 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**

$v_0^{(t)} = \theta_{t-1}$; Inner-loop optimization

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

$v_\tau^{(t)} = \mathcal{A}(F_S(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

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

end

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

Algorithm 2: Lookahead

Input : Objective $F_S(\theta)$, 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**

$v_0^{(t)} = \theta_{t-1}$; Inner-loop optimization

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

$v_\tau^{(t)} = \mathcal{A}(F_S(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

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

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**

$v_0^{(t)} = \theta_{t-1}$;

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

$v_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

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

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**

$v_0^{(t)} = \theta_{t-1}$;

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

$v_{\tau}^{(t)} = \mathcal{A}(F_{\mathcal{S}}(\theta), v_{\tau-1}^{(t)}, \eta_{\tau-1}^{(t)}, \mathcal{S}) = v_{\tau-1}^{(t)} - \eta_{\tau-1}^{(t)} g_{\tau-1}^{(t)}$

end

$\theta_t = (1 - \alpha)\theta_{t-1} + \alpha v_k^{(t)}$

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

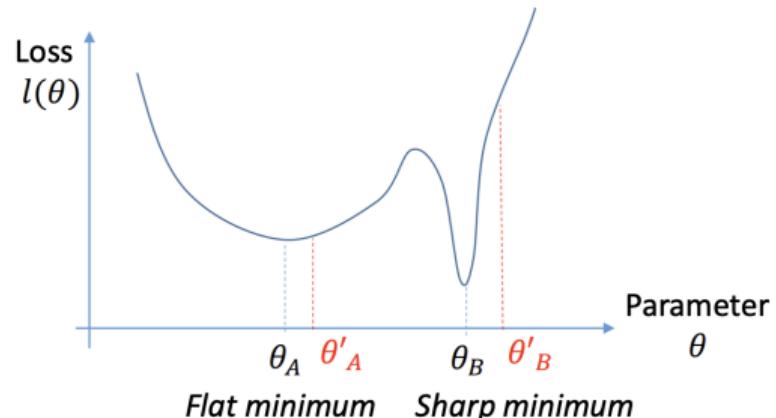
2 Accelerated and Stabilized Optimization Methods

3 Advanced Optimization Methods

- Lookahead
- Sharpness-aware Minimization

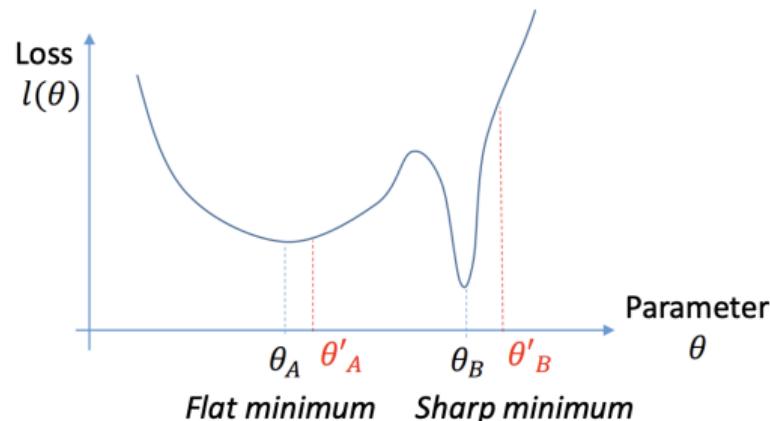
Flat Minima in Deep Learning

In many cases,



Flat Minima in Deep Learning

In many cases,
DL → minimizing a loss function $\ell(\theta)$

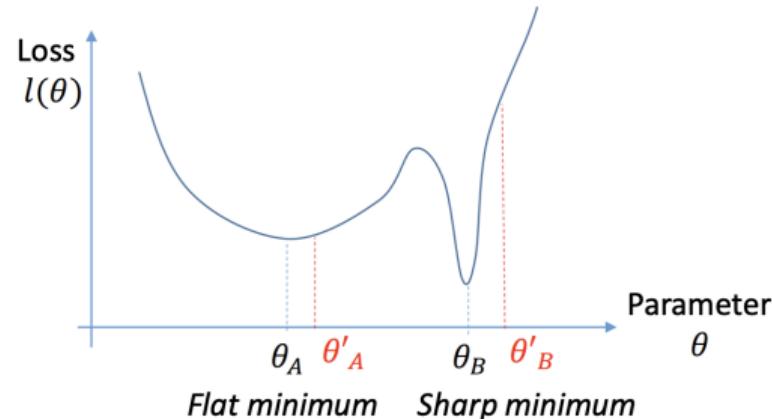


Flat Minima in Deep Learning

In many cases,

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

Highly non-convex (many local minima)!

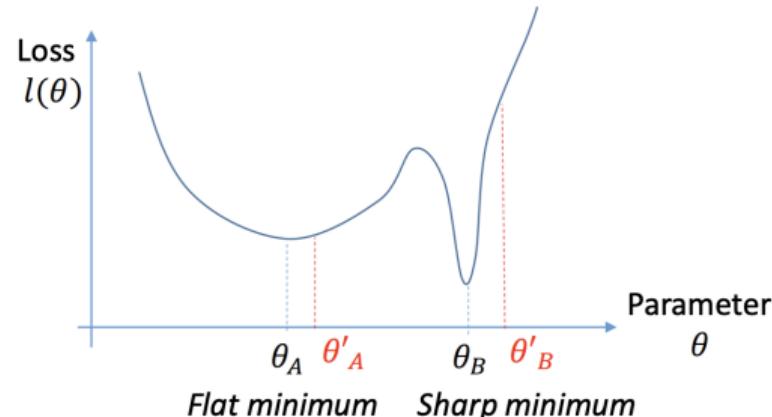


Flat Minima in Deep Learning

In many cases,

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

Highly non-convex (many local minima)!



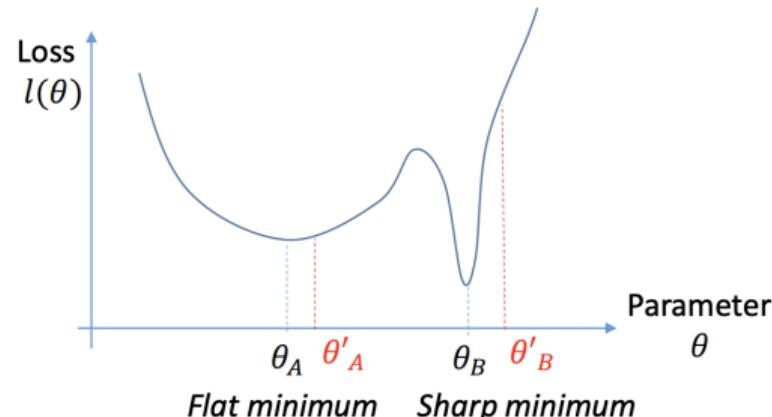
- Q) Which is better, θ_A or θ_B ?

Flat Minima in Deep Learning

In many cases,

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

Highly non-convex (many local minima)!



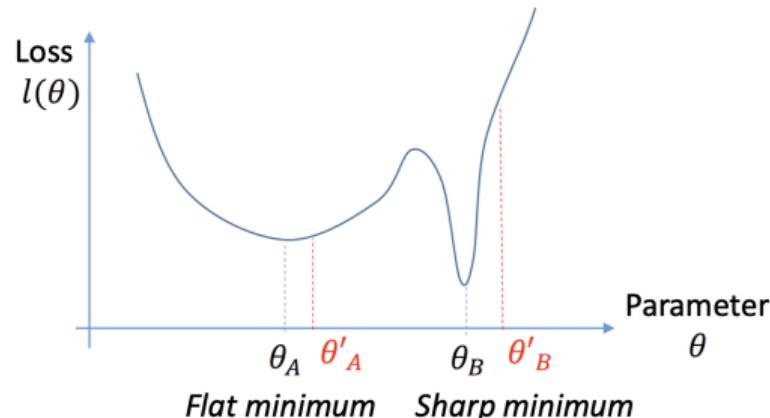
- Q) Which is better, θ_A or θ_B ?
- A) We prefer θ_A to θ_B even though $\ell(\theta_A) > \ell(\theta_B)$

Flat Minima in Deep Learning

In many cases,

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

Highly non-convex (many local minima)!



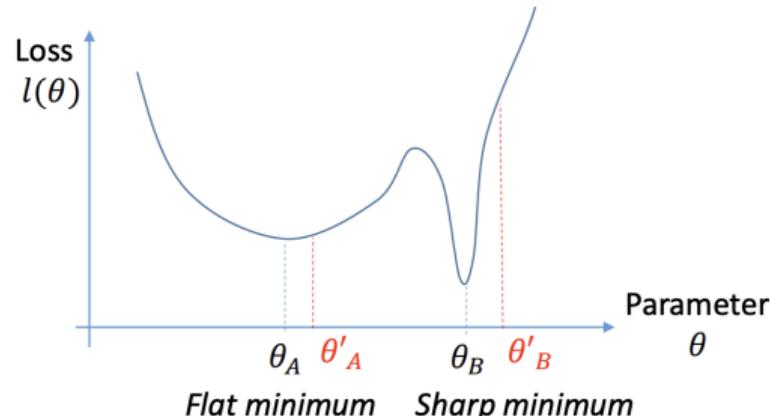
- Q) Which is better, θ_A or θ_B ?
- A) We prefer θ_A to θ_B even though $\ell(\theta_A) > \ell(\theta_B)$
 - Why? Because θ_A is more robust.

Flat Minima in Deep Learning

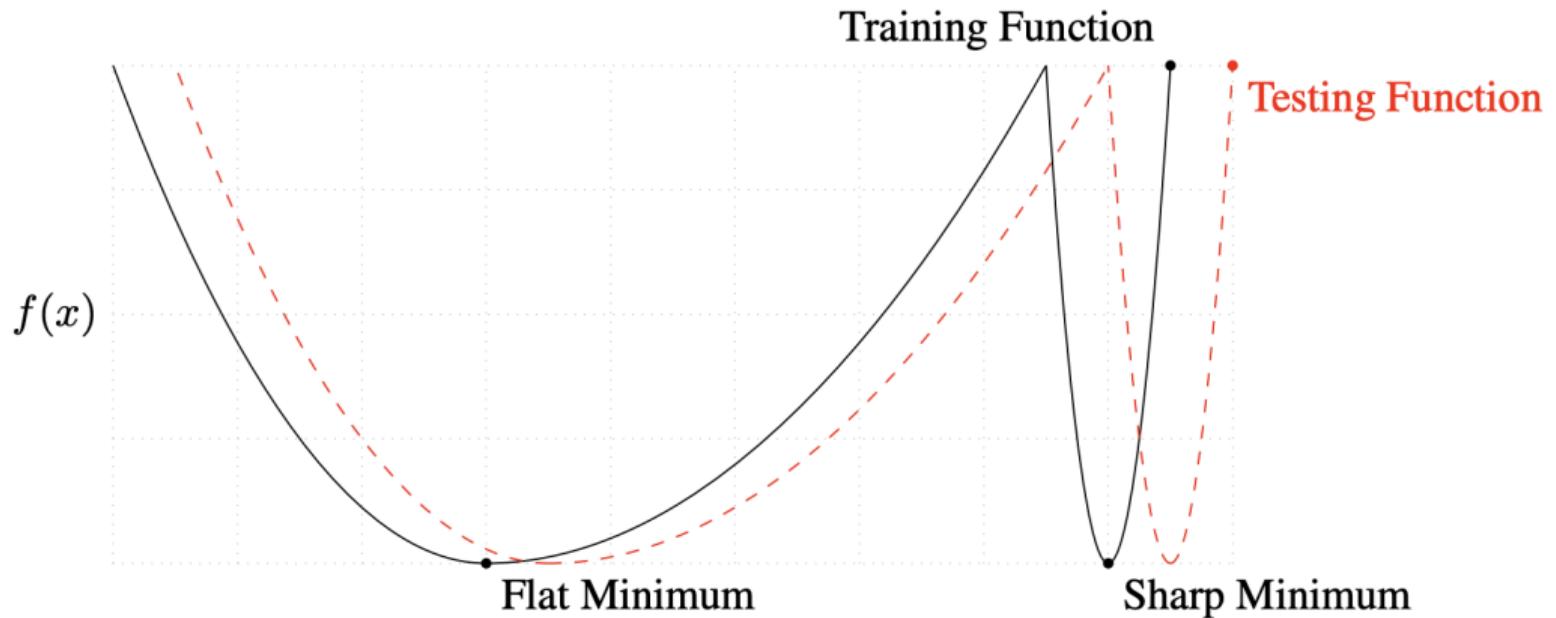
In many cases,

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

Highly non-convex (many local minima)!



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



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)

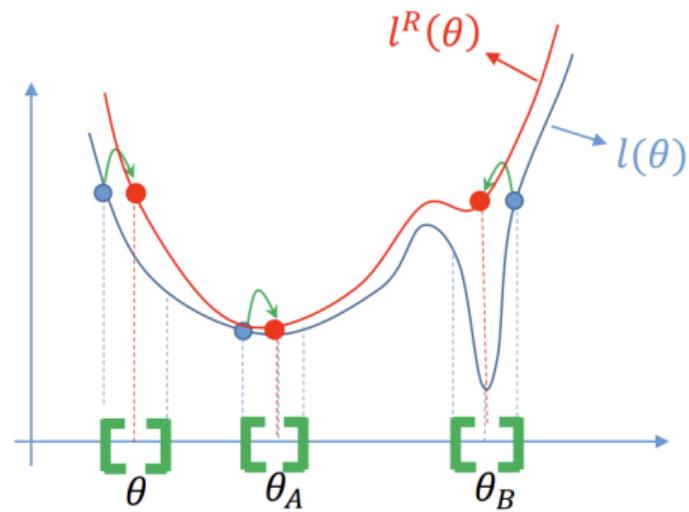
= (often encountered in AI applications) (14)

But, how?

Sharpness-Aware Minimization (SAM)

Idea of SAM:

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

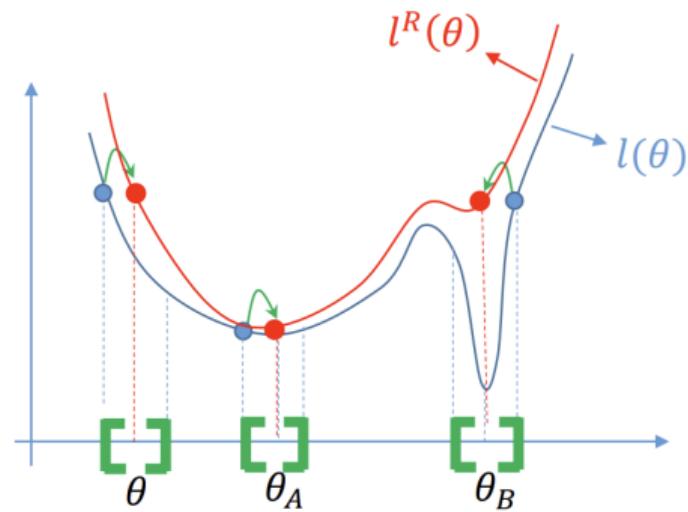


Sharpness-Aware Minimization (SAM)

Idea of SAM:

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

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



Intuition behind SAM

- **Goal:** Find the local minima θ that are generalizable to test samples

Intuition behind SAM

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

Intuition behind SAM

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

Theorem 2

With high probability over S , the flatness-based bound says:

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

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

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

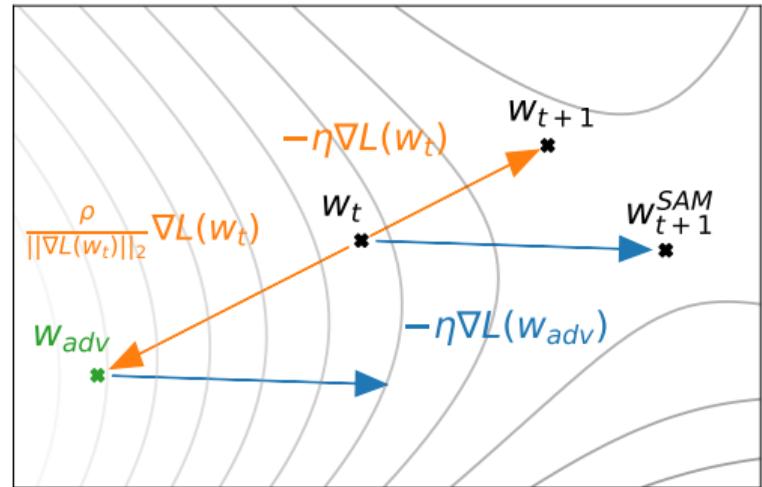
where $1/p + 1/q = 1$ and the solution to a classical dual norm problem can solve this approximation.

- 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:

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

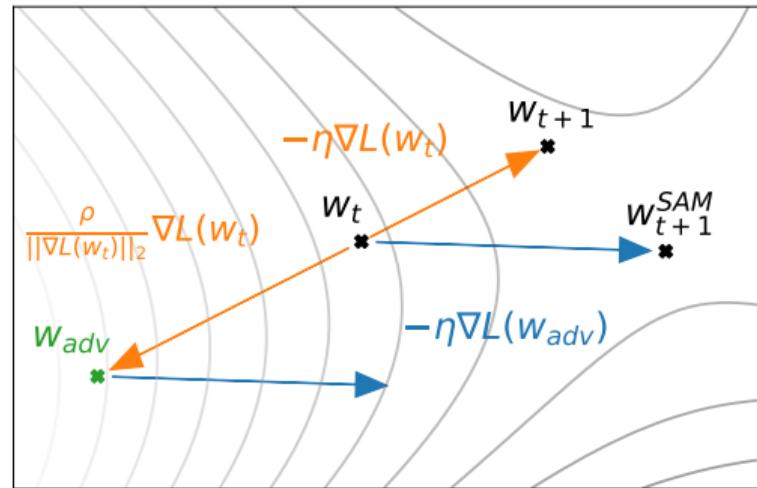


From original SAM paper.

- 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)

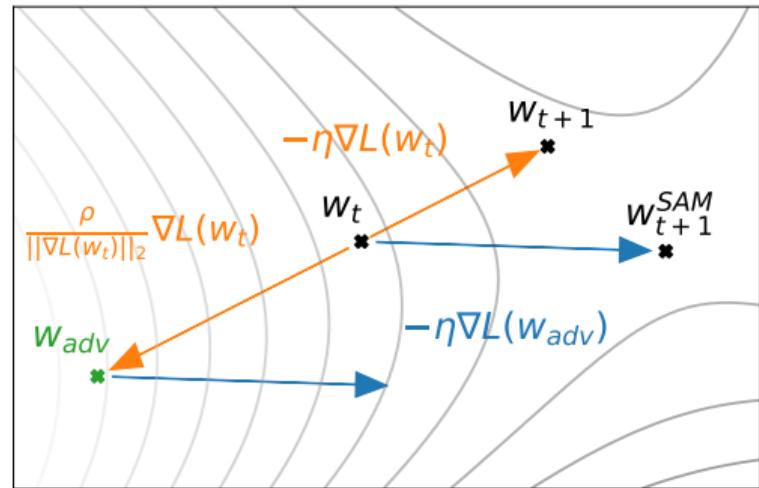


From original SAM paper.

- 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 $\nabla_{\theta} \mathcal{L}_S(\theta)$

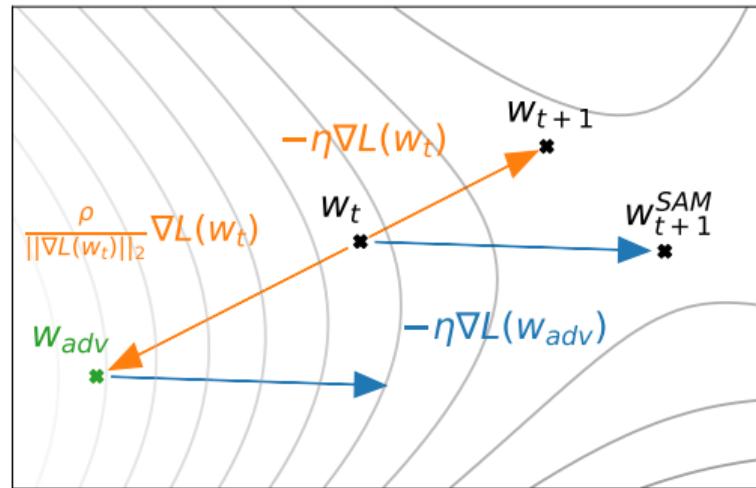


From original SAM paper.

- 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 $\nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}}$

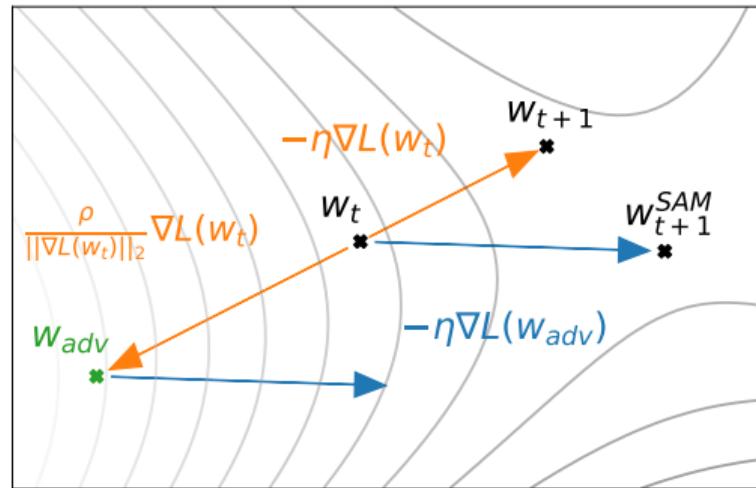


From original SAM paper.

- The gradient estimator of SAM is given by:

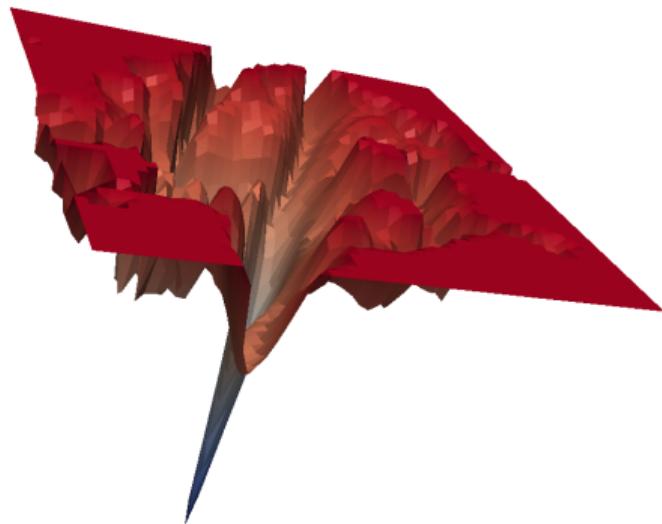
$$\nabla_{\theta} \mathcal{L}_S^{\text{SAM}}(\theta) \approx \nabla_{\theta} \mathcal{L}_S(\theta)|_{\theta+\hat{\epsilon}} \quad (20)$$

- 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.

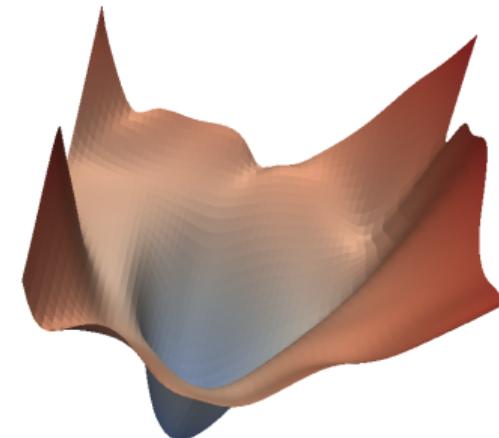


From original SAM paper.

Verification of the flatness



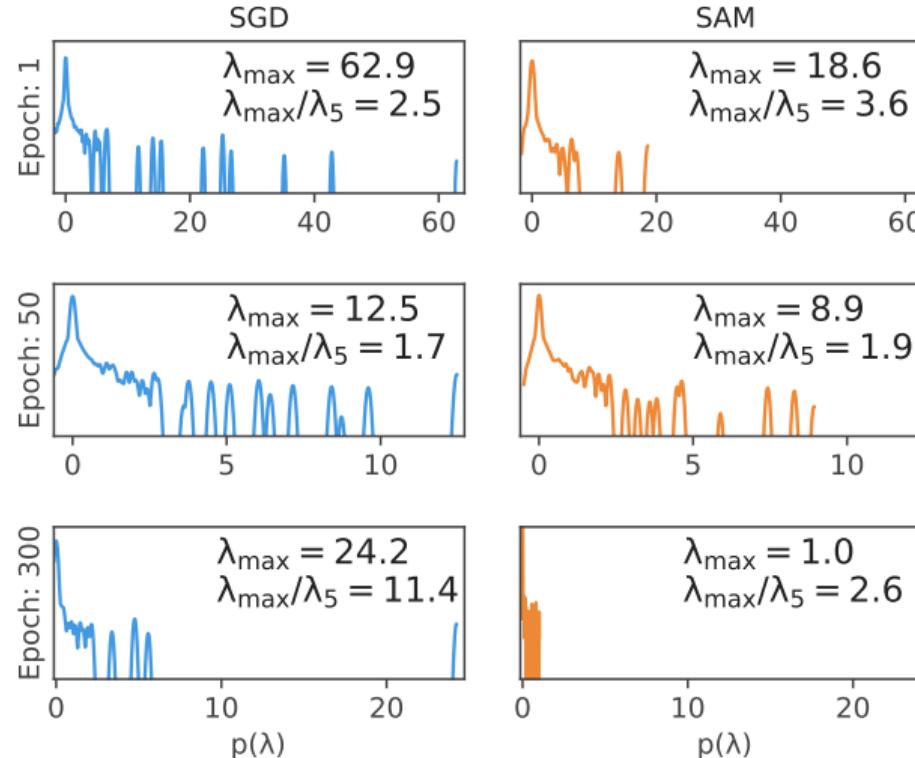
(a) ERM.



(b) SAM.

Loss surface visualization.

Verification of the flatness



Hessian spectra.

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	22.5 _{±0.1}	6.28 _{±0.08}	22.9 _{±0.1}	6.62 _{±0.11}
	200	21.4 _{±0.1}	5.82 _{±0.03}	22.3 _{±0.1}	6.37 _{±0.04}
	400	20.9 _{±0.1}	5.51 _{±0.03}	22.3 _{±0.1}	6.40 _{±0.06}
ResNet-101	100	20.2 _{±0.1}	5.12 _{±0.03}	21.2 _{±0.1}	5.66 _{±0.05}
	200	19.4 _{±0.1}	4.76 _{±0.03}	20.9 _{±0.1}	5.66 _{±0.04}
	400	19.0 _{±<0.01}	4.65 _{±0.05}	22.3 _{±0.1}	6.41 _{±0.06}
ResNet-152	100	19.2 _{±<0.01}	4.69 _{±0.04}	20.4 _{±<0.0}	5.39 _{±0.06}
	200	18.5 _{±0.1}	4.37 _{±0.03}	20.3 _{±0.2}	5.39 _{±0.07}
	400	18.4 _{±<0.01}	4.35 _{±0.04}	20.9 _{±<0.0}	5.84 _{±0.07}

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

Results on ViT (and MLP-Mixer)

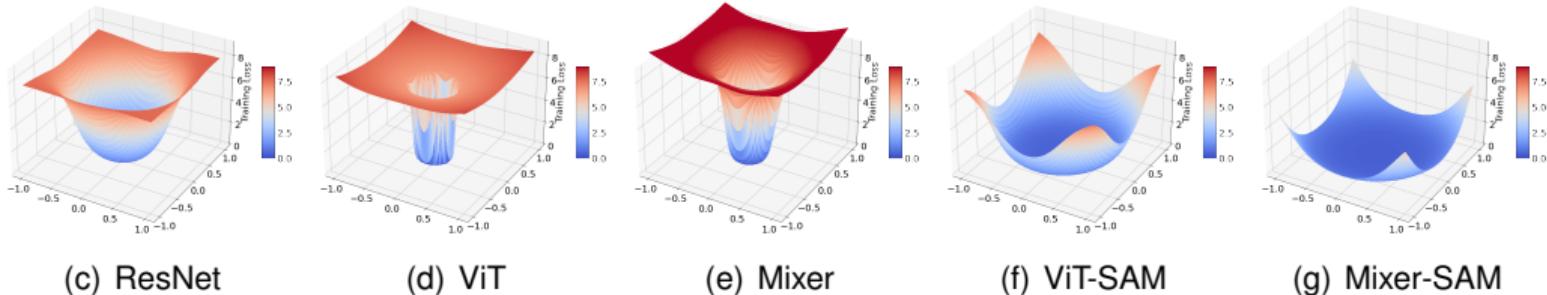


Figure: Cross-entropy loss landscapes of ResNet-152, ViT-B/16, and Mixer-B/16. ViT and MLP-Mixer converge to sharper regions than ResNet when trained on ImageNet with the basic Inception-style preprocessing. SAM significantly smooths the landscapes.

Results on ViT (and MLP-Mixer)

Table: Number of parameters, Hessian dominate eigenvalue λ_{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 κ 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
#Params	60M		87M		59M	
Hessian λ_{max}	179.8	42.0	738.8	20.9	1644.4	22.5
L_{train}	0.86	0.90	0.65	0.82	0.45	0.97
L_{train}^N *	2.39	2.16	6.66	0.96	7.78	1.01
ImageNet (%)	78.5	79.3	74.6	79.9	66.4	77.4
ImageNet-C (%)	50.0	52.2	46.6	56.5	33.8	48.8

Thanks & Question Time!

- [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.