

Optimization in Deep Learning



Bingjie YAN

bj.yan.pa@qq.com

ICT, CAS

Powered by @slidev

TOC

1. Optimization Problem & Deep Learning
2. Basic Definition (Convexity, Lipschitz Continuity, Smoothness, etc.)
3. Gradient Decent & Mini-batch Stochastic Gradient Decent
4. Convergence Analysis
5. (Optional) Convergence Analysis for Distributed Machine Learning
6. (Optional) Convergence Analysis for Federated Learning

Note:

- [illegible]

Optimization Problem & Deep Learning

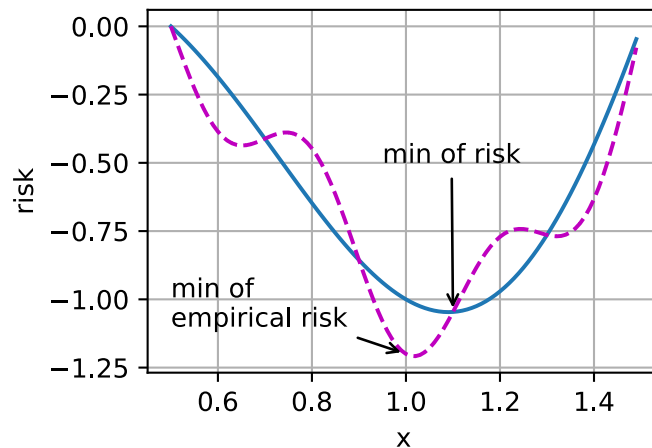
□ □

Optimization Problem Definition

□□□□□□□□□□

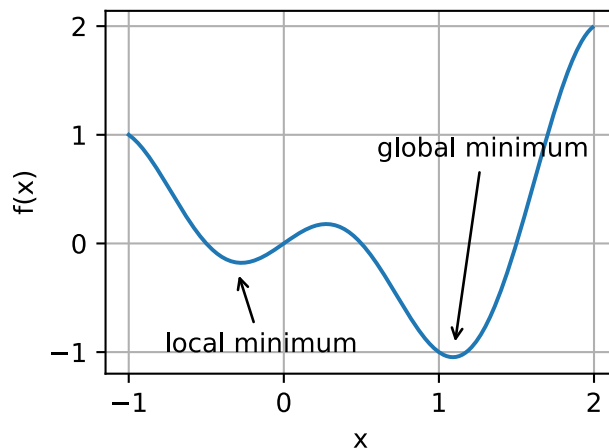
$$\begin{aligned} \min_{x \in \mathbb{R}^n} \quad & f(x) \\ \text{s.t.} \quad & g_i(x) \leq 0, \quad i = 1, \dots, m \\ & h_i(x) = 0, \quad i = 1, \dots, p \end{aligned}$$

×× $f(x)$ ×××××××× x ×××××××× s.t. ××××××××××



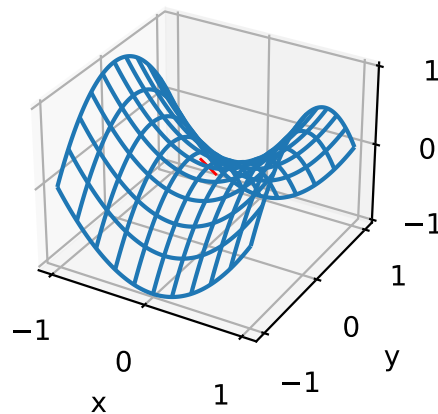
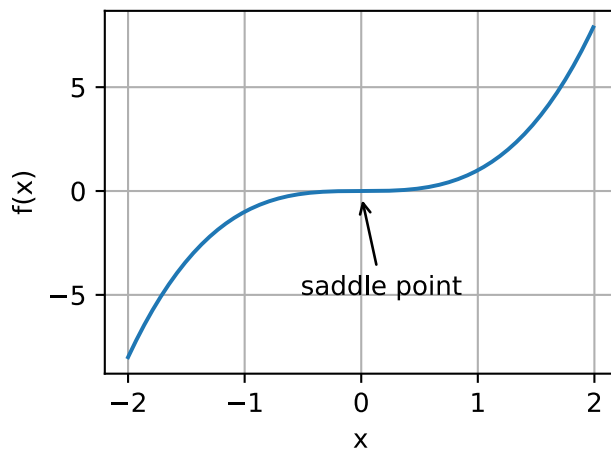
Challenges in Deep Learning Optimization

1. **XXXXXX**(local minimum)

$$\begin{array}{l} \text{f(x) x f(x) x f(x) f(x) f(x) f(x) f(x) f(x)} \\ \text{x f(x) f(x)} \end{array}$$


Challenges in Deep Learning Optimization

2. $\times \times$ (saddle point)

[illegible]

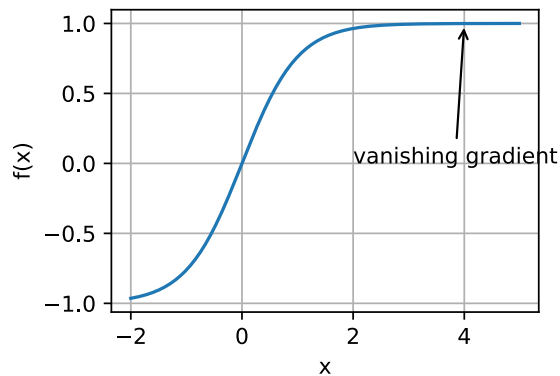
Challenges in Deep Learning Optimization

3. vanishing

vanishing gradient problem

vanilla

$f(x) = \tanh(x)$ $x = 4$ $f'(x) = 1 - \tanh^2(x)$ $f'(4) = 0.0013$
 $x = 4$



Challenges in Deep Learning Optimization

XXXX XXXX X XX XXXXXXXXXXXXXXXXXXXXXXXX

X XXXX XX

XX

XXXXXXXXXXXXXXXXXXXXXXXXXXXX XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

XXXXXXXXXXXXXXXXXXXX

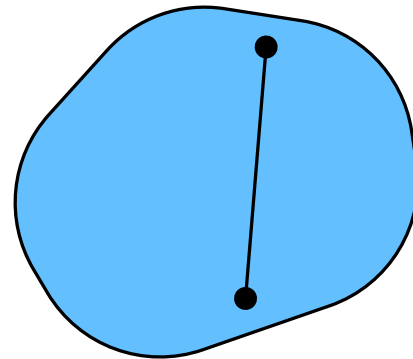
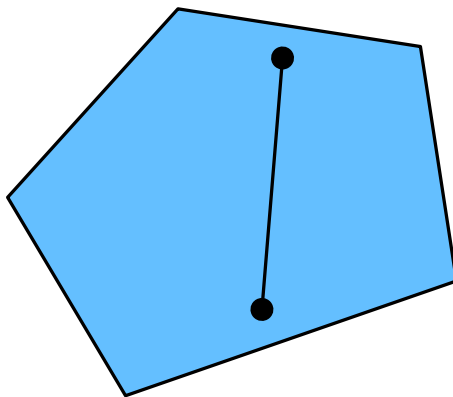
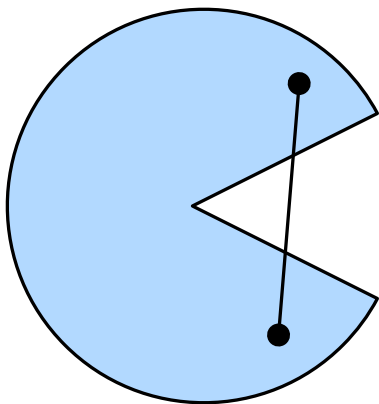
Basic Definition

□□□□□□□□□□Lipschitz □□□□□□

Convex Set

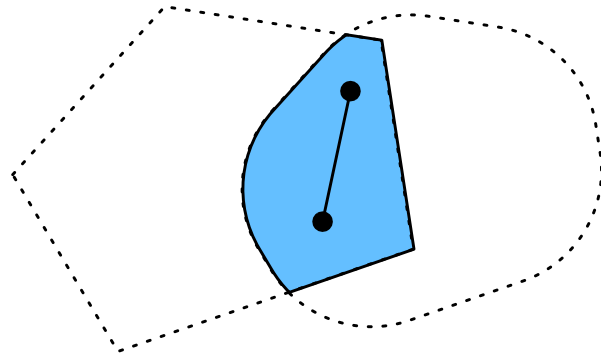
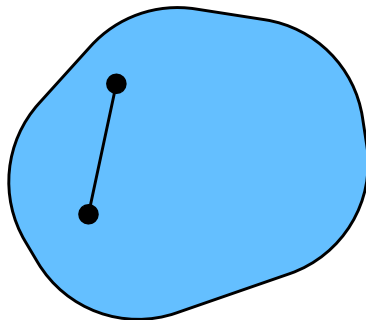
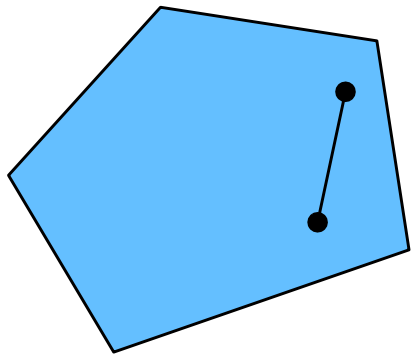
Let \mathcal{X} be a set. \mathcal{X} is convex if for any $x, y \in \mathcal{X}$ and $\alpha \in [0, 1]$,

$$\alpha x + (1 - \alpha)y \in \mathcal{X}$$



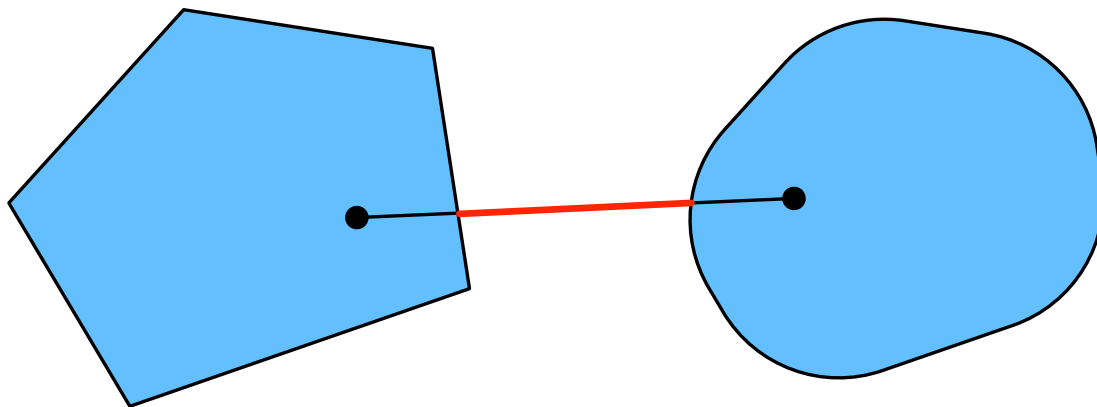
Convex Set

$\text{Convex Set } \mathcal{X} \cap \mathcal{Y} \text{ is Convex}$



Convex Set

xxxxxxxxxx \mathcal{X} x \mathcal{Y} xxxxxxxxxx $\mathcal{X} \cup \mathcal{Y}$ xxxxxxxx

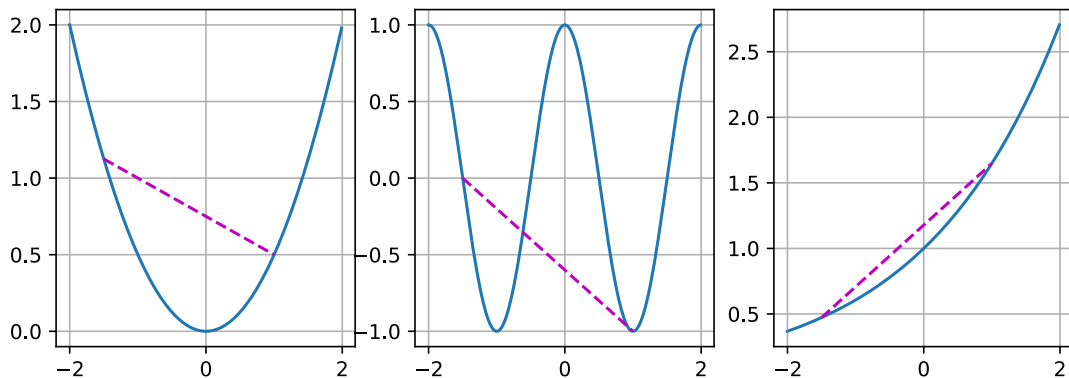


Convex Function

Let \mathcal{X} be a convex set, $x, y \in \mathcal{X}$ and $\alpha \in [0, 1]$

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y)$$

Then f is convex



Properties of Convex Function

1. Subdifferential

Let f be a convex function on \mathcal{X} and $x^* \in \mathcal{X}$. The subdifferential of f at x^* is the set of all subgradients of f at x^* .

Definition:

$x^* \in \mathcal{X}$ is a subgradient of f at x^* if for all $x \in \mathcal{X}$, $0 < |x - x^*| < p$, $f(x) \geq f(x^*) + p(x - x^*)$.

Let $x^* \in \mathcal{X}$ and $f(x) < f(x^*)$. Then $\alpha = 1 - \frac{p}{|x - x^*|}$ is a subgradient of f at x^* .

$$\begin{aligned} f(\alpha x + (1 - \alpha)x^*) &\leq \alpha f(x) + (1 - \alpha)f(x^*) \\ &< \alpha f(x^*) + (1 - \alpha)f(x^*) \\ &= f(x^*) \end{aligned}$$

$x^* \in \mathcal{X}$ is a subgradient of f at x^* if and only if $f(x) \geq f(x^*) + p(x - x^*)$ for all $x \in \mathcal{X}$.

Properties of Convex Function

2. $\square\square\square\square\square\square\square\square\square\square$

3. $\square\square\square\square\square\square\square\square$

...

Strong Convex

Definition 2.1 (Strongly μ -Convex Function)

A function $f: \mathcal{D} \rightarrow \mathbb{R}$ is called μ -strongly convex if

$f: \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex if for all $x, y \in \mathcal{D}$

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|^2, \quad \forall x, y \in \mathcal{D}$$

where $\mu > 0$ is the strong convexity parameter.

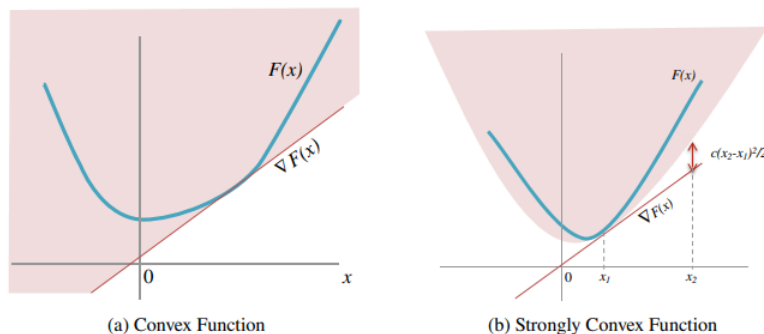


Fig. 2.3 Illustration of a convex function and a c -strongly convex function

Strong Convex

[illegible]

$$\nabla f(y) = \nabla f(x^*) = 0$$

$$\langle \nabla f(x), x - x^* \rangle \geq \frac{\mu}{2} \|x - x^*\|^2$$

$$\|f - \mu - \frac{\mu}{2} \| \cdot \| ^2 \|$$

Lipschitz Continuity

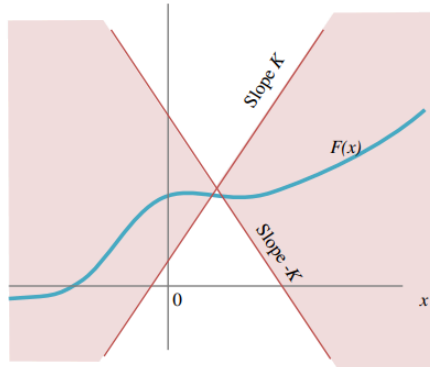
L -Lipschitz Lipschitz

$f : R^d \rightarrow R$ $\| \cdot \|$ $L > 0$ $x, y \in R^d$

$$|f(x) - f(y)| \leq L \|x - y\|$$

f $\| \cdot \|$ L -Lipschitz

Fig. 2.1 Illustration of Lipschitz continuity. If a scalar function $F(x)$ is K -Lipschitz continuous, then for any point x , the function lies inside the region bounded by lines of slope K and $-K$ that pass through the point $(x, F(x))$



Smoothness

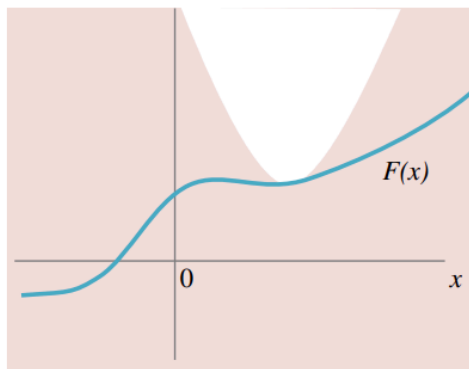
Definition 2.1 (Lipschitz smoothness)

A function $f : R^d \rightarrow R$ is L -Lipschitz smooth if $L > 0$ and for all $x, y \in R^d$

$$f(x) - f(y) \leq \nabla f(y)^\top (x - y) + \frac{L}{2} \|x - y\|^2$$

where $\|\cdot\|$ is the L_2 -norm.

Fig. 2.2 Illustration of Lipschitz smoothness. If a scalar function $F(x)$ is L -Lipschitz smooth, then for any point x , the function lies inside the shaded region shown in the picture, as specified by (2.2)



Smoothness

□□□□□□

$$|\nabla f(x) - \nabla f(y)| \leq L \|x - y\|$$

□□□□□ f L -~~□□□□□□□□□□□□□□□□~~ ∇f L -Lipschitz □□□□

Gradient Decent & Mini-batch Stochastic Gradient Decent

□□□□□□□□□□□□□□□□□□

Gradient Decent

□□□□□□□□

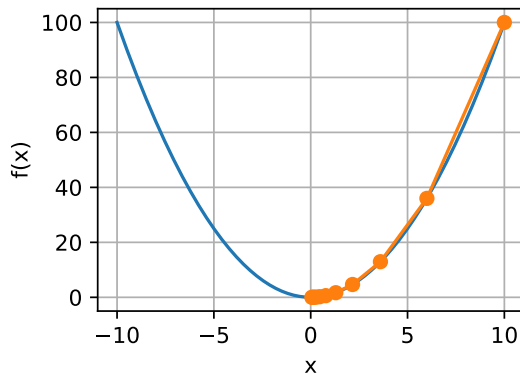
$$w_{t+1} = w_t - \eta_t \nabla f(w_t)$$

□□□ η_t □□□□□ w_t □□□□ $f(w_t)$ □□□□□□□ $\eta_t \leq \eta_{t-1} \leq \dots \leq \eta_1$ □□□□□□

□□□□□□□□□□□□□□□□

$$\|w_{t+1} - w^*\|^2 = \|w_t - \eta_t \nabla f(w_t) - w^*\|^2 \tag{1}$$

$$= \|w_t - w^*\|^2 - 2\eta_t \langle \nabla f(w_t), w_t - w^* \rangle + \eta_t^2 \|\nabla f(w_t)\|^2 \tag{2}$$



Gradient Decent

□□□□□□□□□□□□□□

▣▣ $f(x)$ ▣ x ▣▣▣▣▣▣▣▣▣▣▣▣▣▣▣▣

$$f(x + \epsilon) = f(x) + \epsilon f'(x) + o(\epsilon^2)$$

▣▣▣▣▣▣▣▣▣ $\eta > 0$ ▣▣ $\epsilon = -\eta f'(x)$ □□□□

$$f(x - \eta f'(x)) = f(x) - \eta f'(x)^2 + o(\eta^2 f'(x)^2)$$

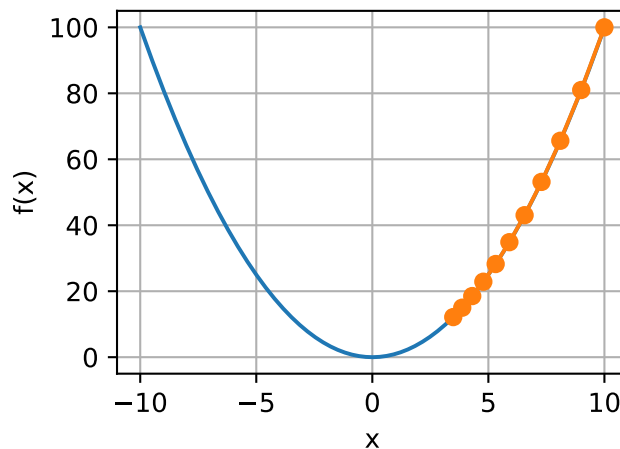
▣▣ $f'(x) \neq 0$ ▣▣▣▣▣▣▣▣▣▣ η ▣▣▣▣▣▣▣ $o(\eta^2 f'(x)^2)$ ▣▣▣▣▣▣▣▣▣▣

$$f(x - \eta f'(x)) < f(x)$$

Gradient Decent (Learning Rate)

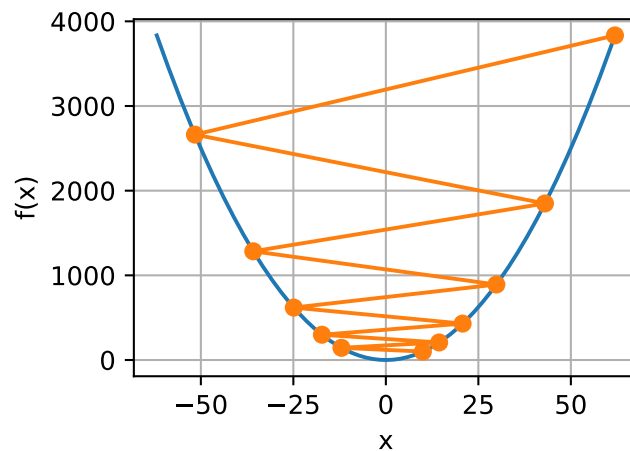
□ □ □ □ □ □ □

1.



Gradient Decent (Learning Rate)

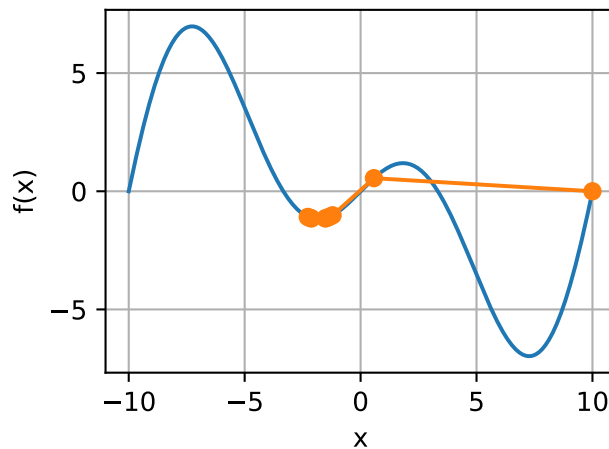
2.



Gradient Decent (Learning Rate)

3.

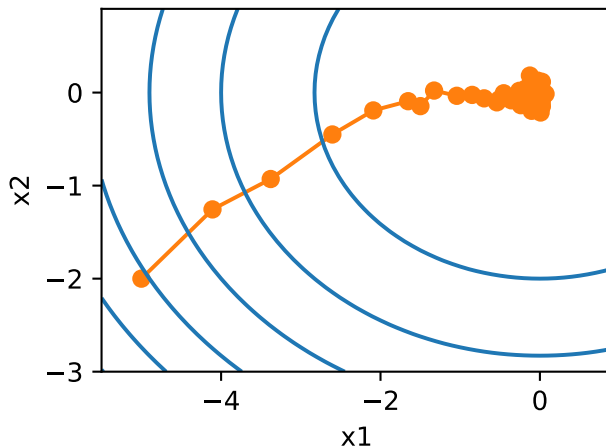


Stochastic Gradient Decent

sample (\mathbf{x}_i, y_i)

$$w_{t+1} = w_t - \eta_t \nabla f(w_t; \xi_t)$$

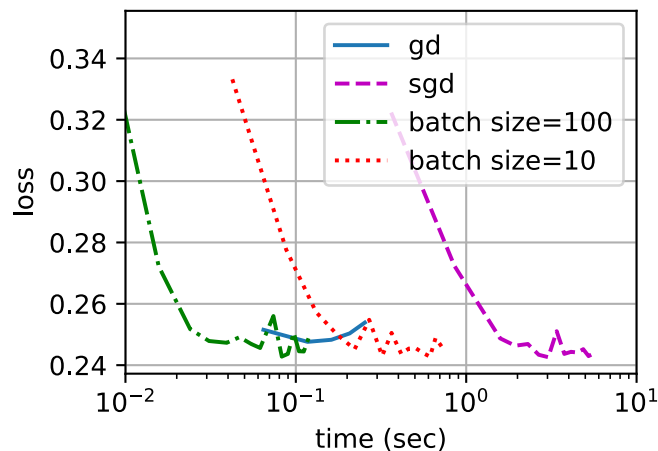
$$g(w_t) = \nabla f(w_t; \xi_t)$$



Mini-batch Stochastic Gradient Decent

Mini-batch Stochastic Gradient Descent $\frac{1}{b} \sum_{i=1}^b \nabla f(w_t; \xi_{t,i})$ Mini-batch Stochastic Gradient Descent $\square\square\square\square$

$$w_{t+1} = w_t - \frac{\eta_t}{b} \sum_{i=1}^b \nabla f(w_t; \xi_{t,i})$$



Mini-batch Stochastic Gradient Decent

□□□□□□□□□□□□□□□□□□□□

Unbiased Estimate

mini-batch SGD □□□□□□□□□□

$$\mathbb{E}_{\xi}[g(w; \xi)] = \nabla f(w)$$

Gradient Bounded Variance Assumption

×××××××××××× mini-batch sample ×××××× $\nabla f(w; \xi)$ □□□□□

$$\text{Var}(\nabla f(w; \xi)) \leq \sigma^2$$

Mini-batch Stochastic Gradient Decent

□□□□□□□□□□□□□□□□□□□□□□

$$\text{Var}(g(w; \xi)) = \mathbb{E}_{\xi}[\|g(w; \xi)\|^2] - \|\mathbb{E}_{\xi}[g(w; \xi)]\|^2 \leq \frac{\sigma^2}{b}$$

□□

$$\mathbb{E}_{\xi}[\|g(w; \xi)\|^2] \leq \|\nabla f(w)\|^2 + \frac{\sigma^2}{b}$$

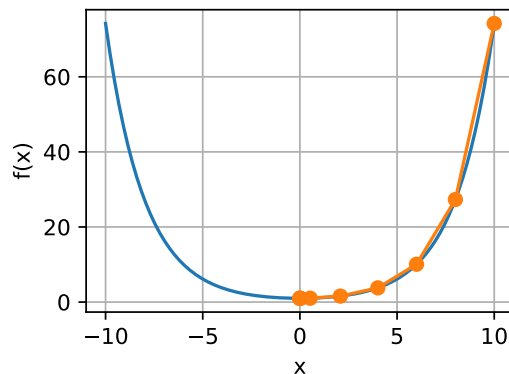
□ □ □ □ □ □ □ □ □

$$f(x + \epsilon) = f(x) + \epsilon \nabla f(x) + \frac{1}{2} \epsilon^2 \nabla^2 f(x) + o(\epsilon^3)$$

$$\mathbf{H} = \nabla^2 f(x) \text{ Hessian } n \times n \quad i, j \quad \frac{\partial^2 f}{\partial x_i \partial x_j}$$

$$\nabla f(x) = \frac{f(x+\epsilon) - f(x)}{\epsilon} \quad 0 \leq \epsilon \leq 1$$

$$\nabla f(x) + \mathbf{H}\epsilon = 0 \Rightarrow \epsilon = -\mathbf{H}^{-1} \nabla f(x)$$



Newton's Method

☒☒ Hessian

□ □ □ □ □ □ □

Convergence Analysis

XXXXXXXXXXXXXXXXXXXX & GD mini-batch SGD XXXXX & mini-batch SGD

□ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □ □

- $$\epsilon(T) \rightarrow 0 \quad \text{as } T \rightarrow \infty \quad \log \epsilon(T) \rightarrow -\infty$$

- XXX 0 XXXXXX

- $$3. \quad \min_{t=1, \dots, T} \mathbb{E} \|\nabla f(w_t)\|^2 \leq \epsilon(T) \quad \square \square \frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(w_t)\|^2 \quad \square \square 0$$

Convergence Analysis (μ -strongly convex and L -smooth & GD)

Assume $f: \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex and L -smooth. Let $\eta \leq \frac{1}{L}$. Assume w_0 and $f(w_T)$ are bounded. Then

$$f(w_t) - f^* \leq (1 - \eta\mu)^t (f(w_0) - f^*)$$

Proof:

□□□□□□

$$\begin{aligned} f(w_{t+1}) - f(w_t) &= f(w_t - \eta \nabla f(w_t)) - f(w_t) \\ &\leq \nabla f(w_t)^\top (w_t - \eta \nabla f(w_t) - w_t) + \frac{L}{2} \|w_t - \eta \nabla f(w_t) - w_t\|^2 \\ &= -\eta \|\nabla f(w_t)\|^2 + \frac{L}{2} \eta^2 \|\nabla f(w_t)\|^2 \\ &= \left(\frac{L}{2} \eta^2 - \eta \right) \|\nabla f(w_t)\|^2 \end{aligned}$$

$$\eta \|\nabla f(x)\|^2$$

$$2\mu(f(x) - f^*) \leq \|\nabla f(x)\|^2$$

□□□

$$\begin{aligned} f(w_{t+1}) - f(w_t) &\leq \eta \left(1 - \frac{L}{2}\eta\right) \left(-\|\nabla f(w_t)\|^2\right) \\ &\leq \eta \left(1 - \frac{L}{2}\eta\right) (2\mu(f(w_t) - f^*)) \end{aligned}$$

$$\eta \leq \frac{1}{L} \implies \left(1 - \frac{L}{2}\eta\right) \geq \frac{1}{2}$$

$$f(w_{t+1}) - f(w_t) \leq -\eta\mu(f(w_t) - f^*)$$

$$f(x) - f^* + f^* - f(w_t) \leq -\eta\mu(f(w_t) - f^*)$$

$$\begin{aligned} f(w_{t+1}) - f^* + f^* - f(w_t) &\leq -\eta\mu(f(w_t) - f^*) \\ \implies f(w_{t+1}) - f^* &\leq (1 - \eta\mu)(f(w_t) - f^*) \end{aligned}$$

□□□□□□□□

$$\begin{aligned} f(w_{t+1}) - f^* &\leq (1 - \eta\mu)(f(w_t) - f^*) \\ &\leq (1 - \eta\mu)(1 - \eta\mu)(f(w_{t-1}) - f^*) \\ &\leq \dots \\ &\leq (1 - \eta\mu) \dots (1 - \eta\mu)(f(w_1) - f^*) \\ &\leq (1 - \eta\mu)^{t+1}(f(w_0) - f^*) \end{aligned}$$

xxxxxxxx boundary □□□□□□□□

xxxxxxxx $f(w_t) - f^* \leq \epsilon$ □□□

$$\begin{aligned} (1 - \eta\mu)^t(f(w_0) - f^*) &\leq \epsilon \\ \Rightarrow t \log(1 - \eta\mu) + \log(f(w_0) - f^*) &\leq \log \epsilon \\ \Rightarrow t \log \frac{1}{1 - \eta\mu} - \log(f(w_0) - f^*) &\geq \log \frac{1}{\epsilon} \\ \Rightarrow t &= O(\log \frac{1}{\epsilon}) \end{aligned}$$

xxxxxxxx GD □□□□□□□□

Convergence Analysis (μ -strongly convex and L -smooth & mini-batch SGD)

Assume $f \in R^d$ μ -strongly convex L -smooth $\eta \leq \frac{1}{L}$ w_0 mini-batch SGD $\mathbb{E}[f(w_t)]$ bounded

$$\mathbb{E}[f(w_t)] - f(w^*) - \frac{\eta L \sigma^2}{2\mu b} \leq (1 - \eta\mu)^t \left(\mathbb{E}[f(w_0)] - f(w^*) - \frac{\eta L \sigma^2}{2\mu b} \right)$$

Proof:

□□□□

$$f(w_{t+1}) - f(w_t) \leq \nabla f(w_t)^\top (w_{t+1} - w_t) + \frac{L}{2} \|w_{t+1} - w_t\|^2 \quad (3)$$

$$\leq -\eta \nabla f(w_t)^\top g(w_t; \xi_t) + \frac{L}{2} \|\eta g(w_t; \xi_t)\|^2 \quad (4)$$

$$\leq -\eta \nabla f(w_t)^\top g(w_t; \xi_t) + \frac{L\eta^2}{2} \|g(w_t; \xi_t)\|^2 \quad (5)$$

□□□□□□□□

$$\mathbb{E}[f(w_{t+1}) - f(w_t)] \leq -\eta \mathbb{E}[\nabla f(w_t)^\top g(w_t; \xi_t)] + \frac{L\eta^2}{2} \mathbb{E}[||g(w_t; \xi_t)||^2] \tag{6}$$

$$\mathbb{E}[f(w_{t+1})] - f(w_t) \leq -\eta ||\nabla f(w_t)||^2 + \frac{L\eta^2}{2} ||\nabla f(w_t)||^2 + \frac{L\eta^2 \sigma^2}{2b} \tag{7}$$

$$\leq (\eta - \frac{L\eta^2}{2})(-||\nabla f(w_t)||^2) + \frac{L\eta^2 \sigma^2}{2b} \tag{8}$$

□□□ GD □□□□□□□□ $\eta < \frac{1}{L}$ □□□□□□□□□□

$$\mathbb{E}[f(w_{t+1})] - f(w_t) \leq \frac{\eta}{2} (-||\nabla f(w_t)||^2) + \frac{L\eta^2 \sigma^2}{2b} \tag{9}$$

$$\leq \frac{\eta}{2} (-2\mu(f(w_t) - f^*)) + \frac{L\eta^2 \sigma^2}{2b} \tag{10}$$

$$\leq -\eta\mu(\mathbb{E}[f(w_t)] - f^*) + \frac{L\eta^2 \sigma^2}{2b} \tag{11}$$

$$\mathbb{E}[f(w_t)] - f^* - f^* \leq 0$$

$$\mathbb{E}[f(w_{t+1})] - f^* \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^*) + \frac{L\eta^2\sigma^2}{2b}$$

$$\mathbb{E}[f(w_{t+1})] - f^* + x \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^* + x)$$

$$\mathbb{E}[f(w_{t+1})] - f^* + x \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^* + x)$$

$$\mathbb{E}[f(w_{t+1})] - f^* + x \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^* + x)$$

$$(1 - \eta\mu)x - x = \frac{L\eta^2\sigma^2}{2b}$$

$$x = -\frac{L\eta\sigma^2}{2\mu b}$$

$$\mathbb{E}[f(w_{t+1})] - f^* - \frac{L\eta\sigma^2}{2\mu b} \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^* - \frac{L\eta\sigma^2}{2\mu b})$$

$$\mathbb{E}[f(w_{t+1})] - f^* - \frac{L\eta\sigma^2}{2\mu b} \leq (1 - \eta\mu)(\mathbb{E}[f(w_t)] - f^* - \frac{L\eta\sigma^2}{2\mu b})$$

11

$$\mathbb{E}[f(w_t)] - f(w^*) - \frac{\eta L \sigma^2}{2\mu b} \leq (1 - \eta\mu)^t \left(\mathbb{E}[f(w_0)] - f(w^*) - \frac{\eta L \sigma^2}{2\mu b} \right)$$

111

[illegible]

Convergence Analysis (non-convex & mini-batch SGD)

Assume $f \in R^d$ with L -Lipschitz continuous and $\eta \leq \frac{1}{L}$ w_0 t mini-batch SGD f bounded

$$\mathbb{E}\left[\frac{1}{t} \sum_{i=1}^t \|\nabla f(w_i)\|^2\right] \leq \frac{L\sigma^2}{b} + \frac{2(f(w_0) - f(w_{\inf}))}{t\eta}$$

Proof:

$$L\eta < \frac{1}{L}$$

$$\mathbb{E}[f(w_{t+1})] - f(w_t) \leq -\eta \|\nabla f(w_t)\|^2 + \frac{L\eta^2 \sigma^2}{2b}$$

□□□□□□□□□□ t □□□

$$\frac{1}{t} \sum_{i=1}^t \mathbb{E}[f(w_{i+1})] - f(w_i) \leq -\frac{\eta}{2t} \sum_{i=1}^t \|\nabla f(w_i)\|^2 + \frac{L\eta^2\sigma^2}{2b} \tag{12}$$

$$\frac{1}{t} \sum_{i=1}^t \|\nabla f(w_i)\|^2 \leq -\frac{2\mathbb{E}[f(w_{t+1}) - f(w_0)]}{\eta t} + \frac{L\eta\sigma^2}{2b} \tag{13}$$

$$\leq \frac{2(f(w_0) - f(w_{\text{inf}}))}{\eta t} + \frac{L\eta\sigma^2}{2b} \tag{14}$$

□□□

□□□ $t = O(\frac{1}{\epsilon})$ □□□□□□□□□□□□□□□□

Other Convergence Analysis

Other Convergence Analysis

Condition	GD	SGD
Convex	$O(\frac{1}{\sqrt{T}})$	$O(\frac{1}{\sqrt{T}})$
+ Lipschitz	$O(\frac{1}{T})$	$O(\frac{1}{\sqrt{T}})$
+ Strongly Convex	$O(c^T)$	$O(\frac{1}{T})$

10-725/36-725: Convex Optimization(Fall 2018), Lecture 24: November 26, Ryan Tibshirani

Distributed Optimization

- Distributed Synchronous SGD
- Distributed Asynchronous SGD
- Federated Learning

Reference

1. [Convex Optimization], Stephen Boyd
2. [Optimization Algorithm for Distributed Machine Learning], Gauri Joshi
3. [Dive into Deep Learning], D2L.ai
4. [Optimization Methods for Large-Scale Machine Learning], Léon Bottou

Learning more about optimization and convergence analysis:

<https://blog.bj-yan.top/tags/convergence-analysis/>