

# 確率的勾配降下法の平滑化効果を利用した段階的最適化手法による ディープニューラルネットワークの大域的最適化

日本オペレーションズリサーチ2024年春季研究発表会

3/7(火) 10:00～10:20

05001736 明治大学 佐藤尚樹

01016200 明治大学 飯塚秀明



# 背景：機械学習（画像分類）

入力

$$z_i \in \mathbb{R}^d \quad (i = 1, 2, \dots, N)$$



$$d = 512 \times 512 \times 3$$

$$N = 3000000$$

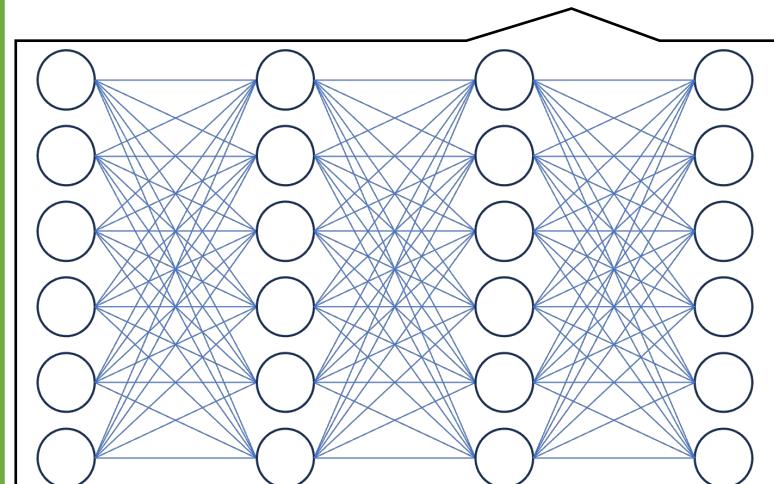
データセット

関数  $g: \mathbb{R}^d \rightarrow \mathbb{R}$

変数:  $x \in \mathbb{R}^D$ ,

$$z_i \in \mathbb{R}^d$$

Deep Neural Network



$D = 20\text{万} \sim 5\text{兆}$

出力

$$g(x, z_i) \in \mathbb{R}$$

犬:3

正解

$$y_i \in \mathbb{R}$$

猫:2

誤差  $|g(x, z_i) - y_i|$

誤差の平均

$$f(x) := \frac{1}{N} \sum_{i=1}^N |g(x, z_i) - y_i|$$

関数  $f(x)$  を最適化する

# 背景：連続最適化

## ▷ 最急降下法 (Gradient Descent)

$$\mathbf{d}_t := -\nabla f(\mathbf{x}_t)$$

全勾配

$$\mathbf{x}_{t+1} := \mathbf{x}_t + \alpha_t \mathbf{d}_t$$

## ▷ 確率的勾配降下法 (Stochastic Gradient Descent, SGD)

$$\mathbf{d}_t := -\nabla f_{S_t}(\mathbf{x}_t) = -\frac{1}{b} \sum_{i=1}^b \mathbf{G}_{\xi_{t,i}}(\mathbf{x}_t)$$

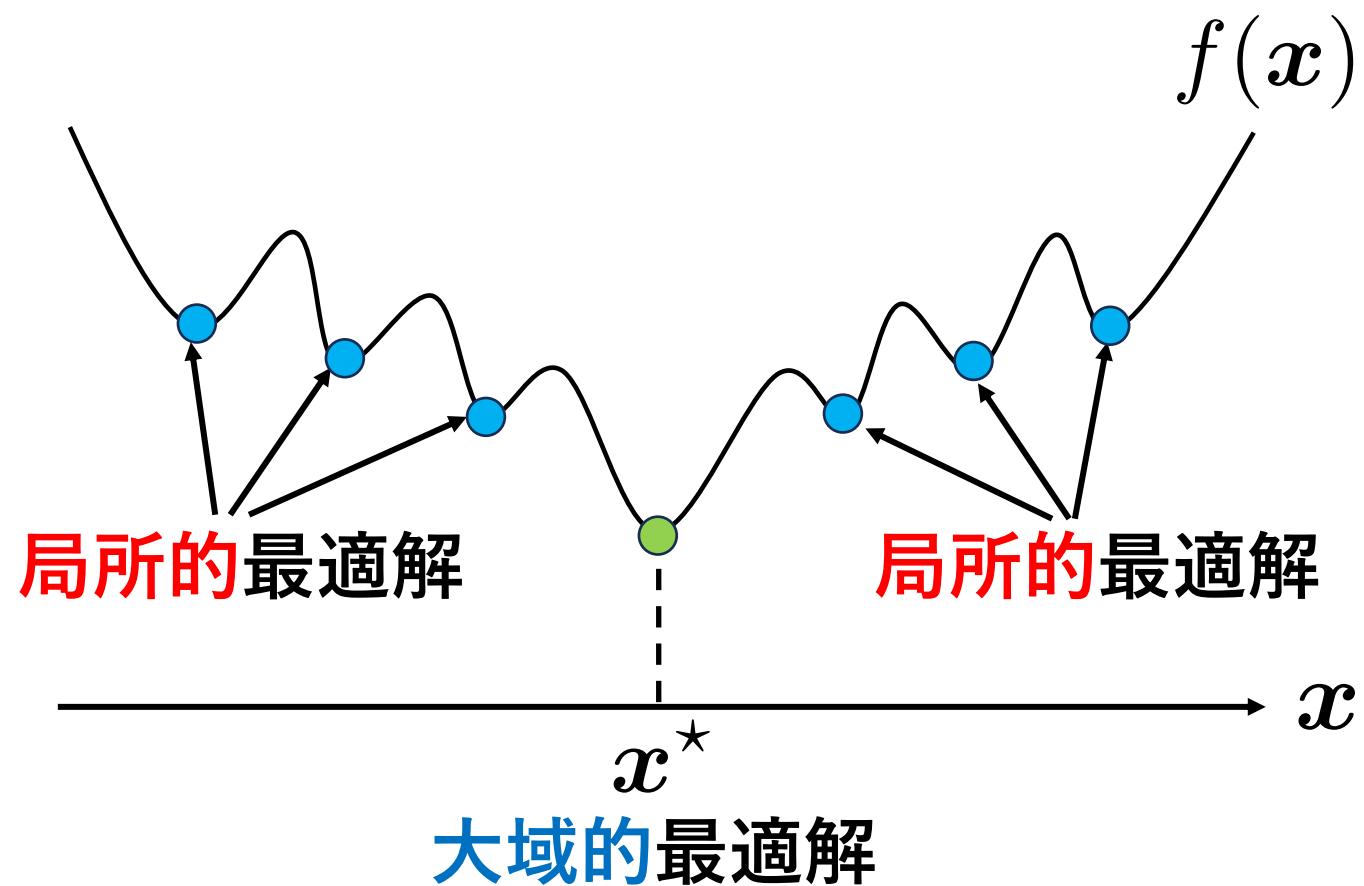
ミニバッチ  
確率的勾配

確率的勾配

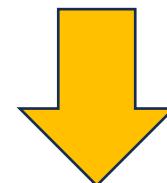
バッチサイズ

ランダムに選ばれた  $b$  個の確率的勾配の平均で代用(ex.  $b=32, 1024$ )

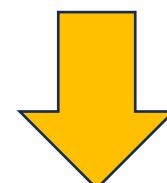
# 背景：大域的最適化の難しさ



機械学習に使用される最適化アルゴリズムは、全て勾配を利用する手法

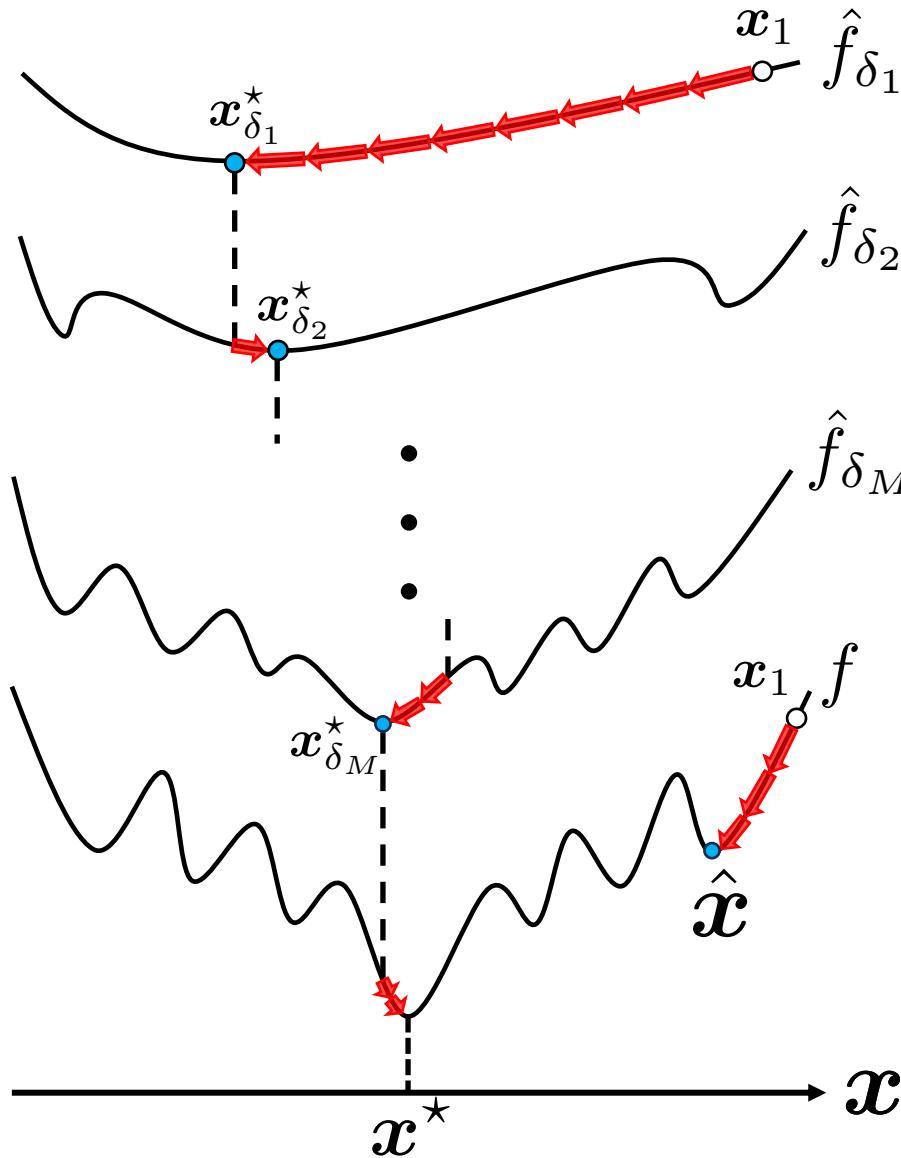


勾配に頼っている手法は、山を登れない



非凸関数の大域的最適解を見つけることは難しい

# 背景：段階的最適化 (Graduated Optimization)



▷ 1987年に提案された**大域的最適化**手法  
▷ 徐々に小さくなるノイズで**平滑化**された  
関数の列を順番に**最適化**する

▷ **目的関数の平滑化**

大きさ  $\delta$  のノイズで平滑化された関数は、

$$\hat{f}_\delta(x) := \mathbb{E}_{\mathbf{u} \sim B(\mathbf{0}; 1)} [f(x - \delta \mathbf{u})]$$

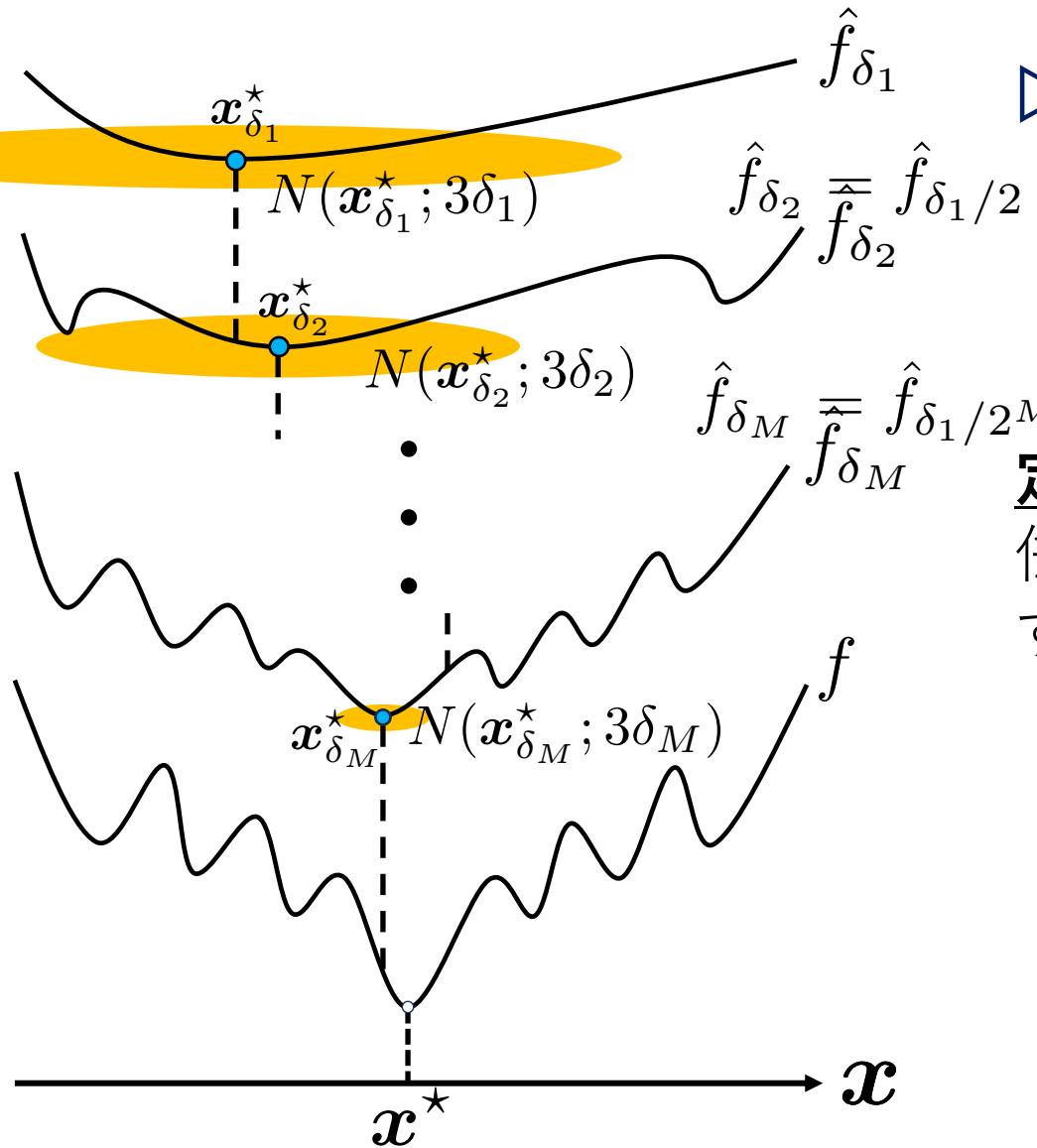
↑ 中心0, 半径1の閉球

$$= \int p(\mathbf{u}) f(x - \delta \mathbf{u}) d\mathbf{u}$$

— 計算可能 — 計算不可能 —

# 背景 : $\sigma$ -nice 関数[2]

[2] E. Hazan et al. “On Graduated Optimization for Stochastic Non-Convex Problems”  
In *Proceedings of the 33th International Conference on Machine Learning*, pp1833-1841, 2016



▷ 段階的最適化アルゴリズムが大域的最適解に収束するために必要な性質を満たす特別な**非凸**関数の族。

## 定義

任意の非凸関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は  $\sigma$ -nice 関数であるという。

(i) 任意の  $\delta > 0$  と  $x_\delta^*$  に対して、 $x_{\delta/2}^*$  が存在して、

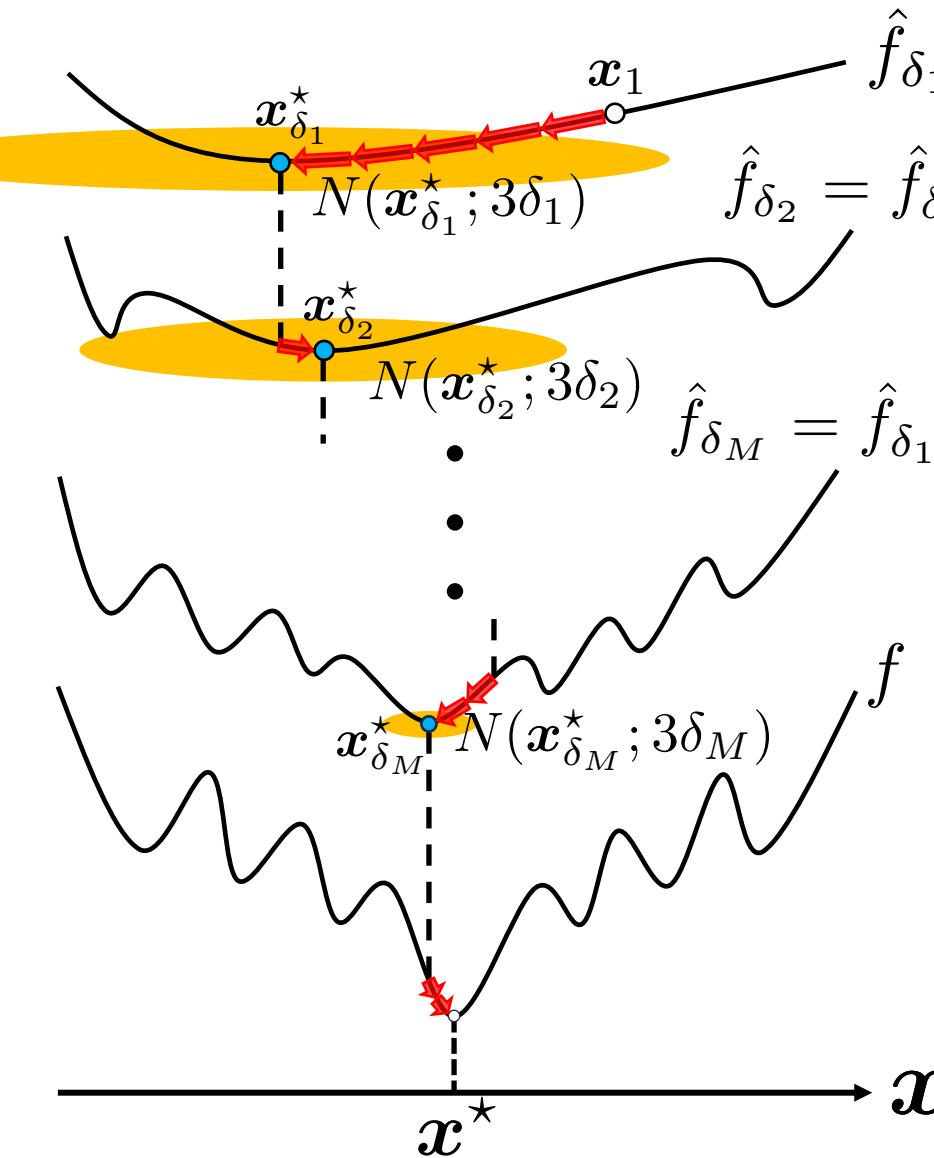
$$\|x_\delta^* - x_{\delta/2}^*\| \leq \frac{\delta}{2}$$

を満たす。

(ii) 任意の  $\delta > 0$  に対して、関数  $\hat{f}_\delta(x)$  は近傍  $N(x_\delta^*; 3\delta)$  で  $\sigma$ -強凸となる。

# 背景 : $\sigma$ -nice 関数 [2]

[2] E. Hazan et al. “On Graduated Optimization for Stochastic Non-Convex Problems”  
 In *Proceedings of the 33th International Conference on Machine Learning*, pp1833-1841, 2016



## Algorithm 1 Graduated optimization (実現不可能)

**Require:**  $\delta_1 > 0, x_1 \in N(x_{\delta_1}^*; 3\delta_1)$

**for**  $m = 1$  to  $M$  **do**

$T_m := \sigma\delta_m^2/32$  ← 反復回数を決定

$x_{m+1} := \text{SGD}(T_m, x_m, \hat{f}_{\delta_m})$  ← 平滑化された関数を SGD 等で最適化

$\delta_{m+1} := \delta_m/2$  ← ノイズレベルを更新

**end for**

**return**  $x_{M+1}$

## 定理5.1 [2]

$\sigma$ -nice 関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  に対して、Algorithm 1 を適用すると、 $\mathcal{O}(1/\epsilon^2)$  回の反復で  $f$  の大域的最適解  $x^*$  の  $\epsilon$ -近傍に到達する。

## 背景：SGDの確率的ノイズと平滑化

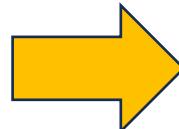
---

- ▷ 確率的勾配降下法は、なぜか汎化性の高い解に収束する。  
特に画像分類では、Adamよりも優れることがある。

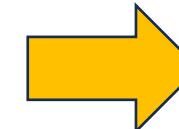
# 付録：機械学習における汎化能力

▷ 学習が完了したモデルが、訓練データ以外の入力に対しても正しく分類できることを『汎化性に優れる』という。

訓練に使われたデータ



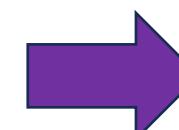
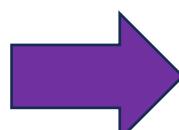
関数  $g: \mathbb{R}^d \rightarrow \mathbb{R}$   
変数:  $x \in \mathbb{R}^D$ ,  
 $z_i \in \mathbb{R}^d$   
Deep Neural Network  
訓練済み  
= パラメータ調整済



猫

訓練損失: 0.001  
訓練精度: 99.9%

訓練に使われていないデータ

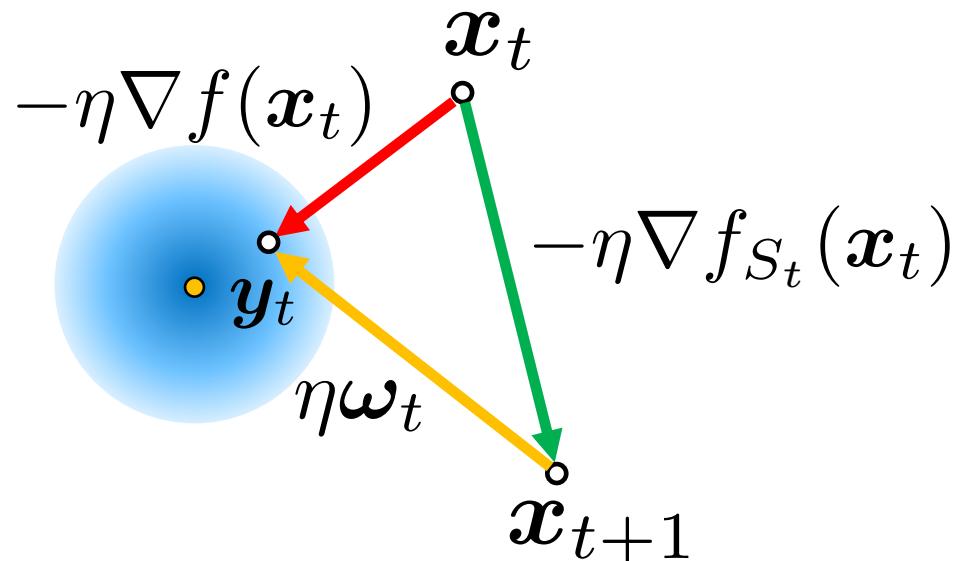


これも猫！

テスト精度: 75%

# 背景：SGDの確率的ノイズと平滑化

- ▷ 確率的勾配降下法(SGD)は、なぜか汎化性の高い解に収束する。特に画像分類では、Adamよりも優れることがある。
- ▷ SGDの確率的ノイズが役立っている？[1]



$$(\text{GD}) \quad y_t := x_t - \eta \nabla f(x_t)$$

$$(\text{SGD}) \quad x_{t+1} := x_t - \eta \nabla f_{S_t}(x_t)$$

$$(\text{確率的ノイズ}) \quad \omega_t := \frac{\nabla f_{S_t}(x_t) - \nabla f(x_t)}{\text{SGDの探索方向}} \quad \frac{\text{GDの探索方向}}{\text{(最急降下方向)}}$$

[1] R. Kleinberg et al. “An Alternative View: When Does SGD Escape Local Minima?”  
In *Proceedings of the 35th International Conference on Machine Learning*, pp2703-2712, 2018

# 背景 : SGDの確率的ノイズと平滑化

$$(\text{GD}) \quad \mathbf{y}_t := \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

$$(\text{SGD}) \quad \mathbf{x}_{t+1} := \mathbf{x}_t - \eta \nabla f_{S_t}(\mathbf{x}_t)$$

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim B(\mathbf{0}; 1)} [f(\mathbf{x} - \delta \mathbf{u})]$$

(確率的ノイズ)  $\boldsymbol{\omega}_t := \nabla f_{S_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)$  とすると、次の式が成り立つ。

$$\mathbb{E}_{\boldsymbol{\omega}_t} [\mathbf{y}_{t+1}] = \mathbf{y}_t - \eta \nabla \underbrace{\mathbb{E}_{\boldsymbol{\omega}_t} [f(\mathbf{y}_t - \eta \boldsymbol{\omega}_t)]}_{:= \hat{f}(\mathbf{y}_t)}$$

- ▷ ”関数  $f(\mathbf{x}_t)$  を SGD で最適化すること”と、”関数  $\hat{f}(\mathbf{y}_t)$  を 最急降下法 で最適化すること”は、期待値の意味では 等価 であると言える。
- ▷ ある程度 平滑化 された関数が、 最急降下法 で最適化されていると みなせる。

[1] R. Kleinberg et al. “An Alternative View: When Does SGD Escape Local Minima?”

In *Proceedings of the 35th International Conference on Machine Learning*, pp2703-2712, 2018

# 動機

---

- ▷  $\sigma$ -nice 関数がどれくらい特別な関数なのか分かっていない。  
そもそも存在するかどうか不明。
- ➡ 任意の非凸関数が  $\sigma$ -nice 関数になるための十分条件を明らかにしたい。
- ▷ SGDが持つ確率的ノイズには関数を平滑化する効果がある。
- ➡ 平滑化の度合い, すなわちノイズレベル  $\delta$  が何によって定まるのかを明らかにしたい。
- ➡ 訓練中にノイズレベル  $\delta$  が徐々に小さくなるようにすることで、SGDを利用した段階的最適化アルゴリズムを構築したい。
- ➡ DNNを含む非凸関数の大域的最適化を実現したい。

# 貢献

---

- ▷  $\sigma$ -nice 関数を拡張した new  $\sigma$ -nice 関数を定義し、任意の関数が new  $\sigma$ -nice 関数となるための十分条件を示した。
- ▷ SGD の確率的ノイズによる平滑化の度合いは、

$$\delta = \eta \sqrt{\frac{C^2}{b}}$$

で表せることを示した。ただし、 $\eta$  は学習率、 $b$  はバッチサイズ、 $C^2$  は確率的勾配の分散とする。

- ▷ SGD の平滑化特性を利用した暗黙的な段階的最適化アルゴリズムを提案し、それが  $\mathcal{O}\left(1/\epsilon^{\frac{1}{p}}\right)$  ( $p \in (0, 1]$ ) 回の反復で new  $\sigma$ -nice 関数の大域的最適解  $x^*$  の  $\epsilon$ -近傍に到達できることを示した。

# 準備：最適化問題

---

## 経験損失最小化問題

- ▷ 訓練データセット  $S := (z_1, z_2, \dots, z_n)$
- ▷ Deep Neural Network のパラメータ  $x \in \mathbb{R}^d$

$$\min_{x \in \mathbb{R}^d} f(x; S) = \frac{1}{n} \sum_{i=1}^n \underline{\underline{l(x; z_i)}} = \frac{1}{n} \sum_{i=1}^n \underline{\underline{l_i(x)}}$$

$i$  番目の訓練データ  $z_i$  に対する **損失関数**

- ▷ 非凸目的関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  を最適化する。

# 準備：仮定

---

## 仮定

(A1) 関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  は連続的微分可能で、 $L_g$ -平滑とする。

$$\forall x, y \in \mathbb{R}^d: \|\nabla f(x) - \nabla f(y)\| \leq L_g \|x - y\|$$

(A2) 関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  は  $L_f$ -リップシツツ関数とする。

$$\forall x, y \in \mathbb{R}^d: |f(x) - f(y)| \leq L_f \|x - y\|$$

(A3)  $(x_t)_{t \in \mathbb{N}} \subset \mathbb{R}^d$  を最適化手法によって生成された点列とするとき、

(i) 任意の  $t \in \mathbb{N}$  に対して、次の式が成り立つとする。

$$\mathbb{E}_{\xi_t} [\mathbf{G}_{\xi_t}(x_t)] = \nabla f(x_t)$$

(ii) 次の式を満たす非負定数  $C^2$  が存在するとする。

$$\mathbb{E}_{\xi_t} \left[ \|\mathbf{G}_{\xi_t}(x_t) - \nabla f(x_t)\|^2 \right] \leq C^2 \quad \text{←確率的勾配の分散}$$

# 準備：仮定

## 仮定続き

(A4) 時刻  $t \in \mathbb{N}$  で、全勾配  $\nabla f$  は **ミニバッチ**  $\mathcal{S}_t$  で次のように近似されるとする。

$$\nabla f_{\mathcal{S}_t}(x_t) := \frac{1}{b} \sum_{i \in [b]} \mathbf{G}_{\xi_{t,i}}(x_t) = \frac{1}{b} \sum_{\{i: z_i \in \mathcal{S}_t\}} \nabla l_i(x_t)$$

ミニバッチ  
確率的勾配      確率的勾配

## 補題2.1

仮定(A3)(ii)と(A4)が成り立つとすると、任意の  $t \in \mathbb{N}$  に対して、

$$\mathbb{E}_{\xi_t} \left[ \left\| \frac{\nabla f_{\mathcal{S}_t}(x_t)}{\text{SGDの探索方向}} - \frac{\nabla f(x_t)}{\text{GDの探索方向 (最急降下方向)}} \right\|^2 \right] \leq \frac{C^2}{b}$$

が成り立つ。

# New $\sigma$ -nice 関数

## 定義 ( $\sigma$ -nice 関数)

任意の非凸関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は  $\sigma$ -nice 関数であるという。

(i) 任意の  $\underline{\delta > 0}$  と  $\mathbf{x}_\delta^*$  に対して、

$$\|\mathbf{x}_\delta^* - \mathbf{x}_{\delta/2}^*\| \leq \frac{\delta}{2} \quad \left( \gamma_m := \frac{1}{2} \right)$$

が成り立つ。

(ii) 任意の  $\underline{\delta > 0}$  に対して、関数  $\hat{f}_\delta(\mathbf{x})$  は近傍  $N(\mathbf{x}_\delta^*; \underline{3\delta})$  で  $\sigma$ -強凸となる。

## 定義 (New $\sigma$ -nice 関数)

任意の非凸リップシツツ関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は new  $\sigma$ -nice 関数であるという。

(i) 任意の  $\underline{\delta_m \in \mathbb{R}}$  と  $\mathbf{x}_{\delta_m}^*$  に対して、

$$\|\mathbf{x}_{\delta_m}^* - \mathbf{x}_{\delta_{m+1}}^*\| \leq |\delta_m| - |\delta_{m+1}|$$

が成り立つ。ただし、

$$\delta_{m+1} := \underline{\gamma_m \delta_m} \quad (\gamma_m \in (0, 1))$$

とする。

(ii) 任意の  $\underline{\delta_m \in \mathbb{R}}$  に対して、関数  $\hat{f}_{\delta_m}(\mathbf{x})$  は近傍  $N(\mathbf{x}_m^*; \underline{d_m \delta_m})$  で  $\sigma$ -強凸となる。

# New $\sigma$ -nice 関数

## 定義 ( $\sigma$ -nice 関数)

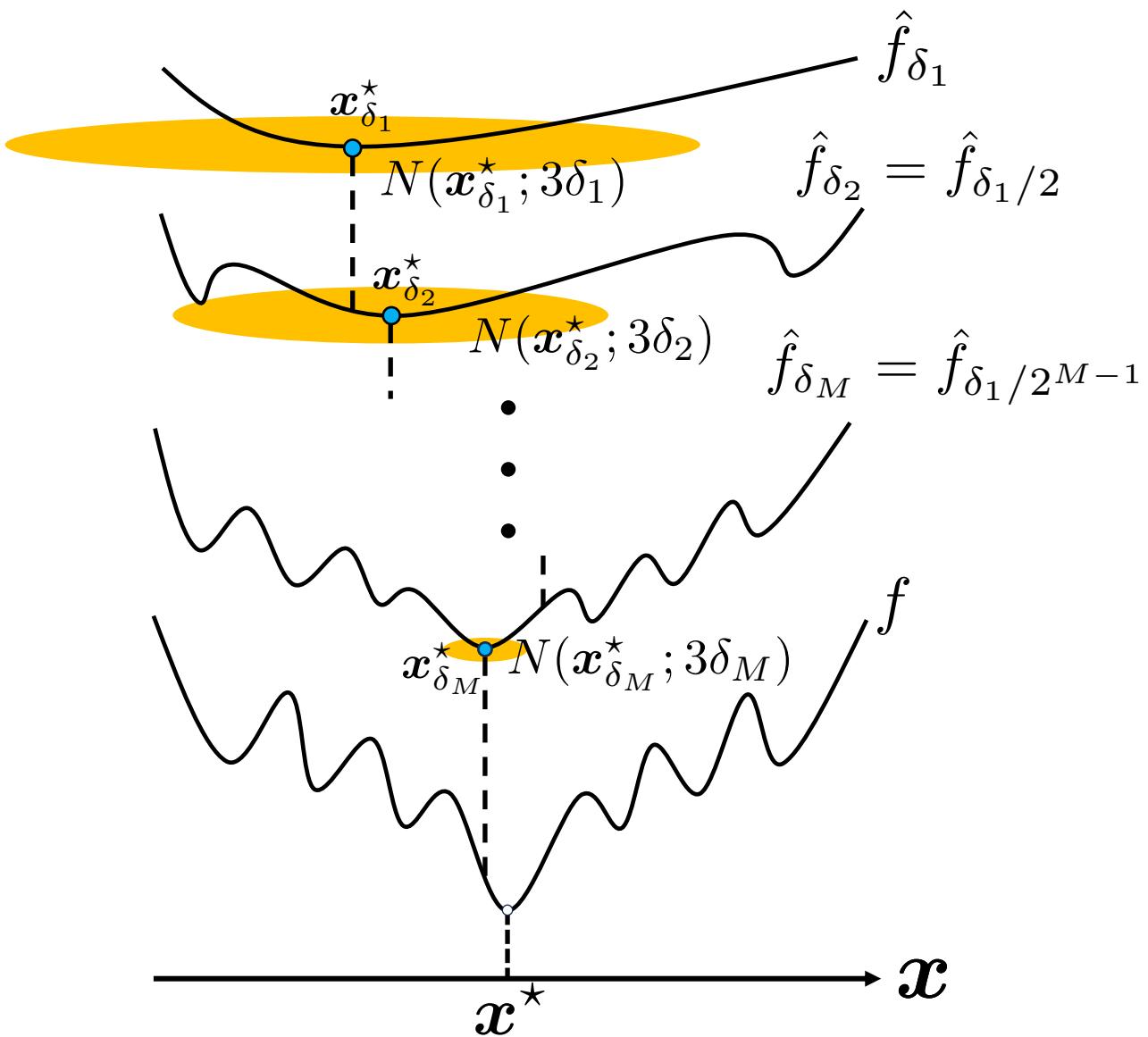
任意の非凸関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は  $\sigma$ -nice 関数であるという。

(i) 任意の  $\underline{\delta > 0}$  と  $x_\delta^*$  に対して、

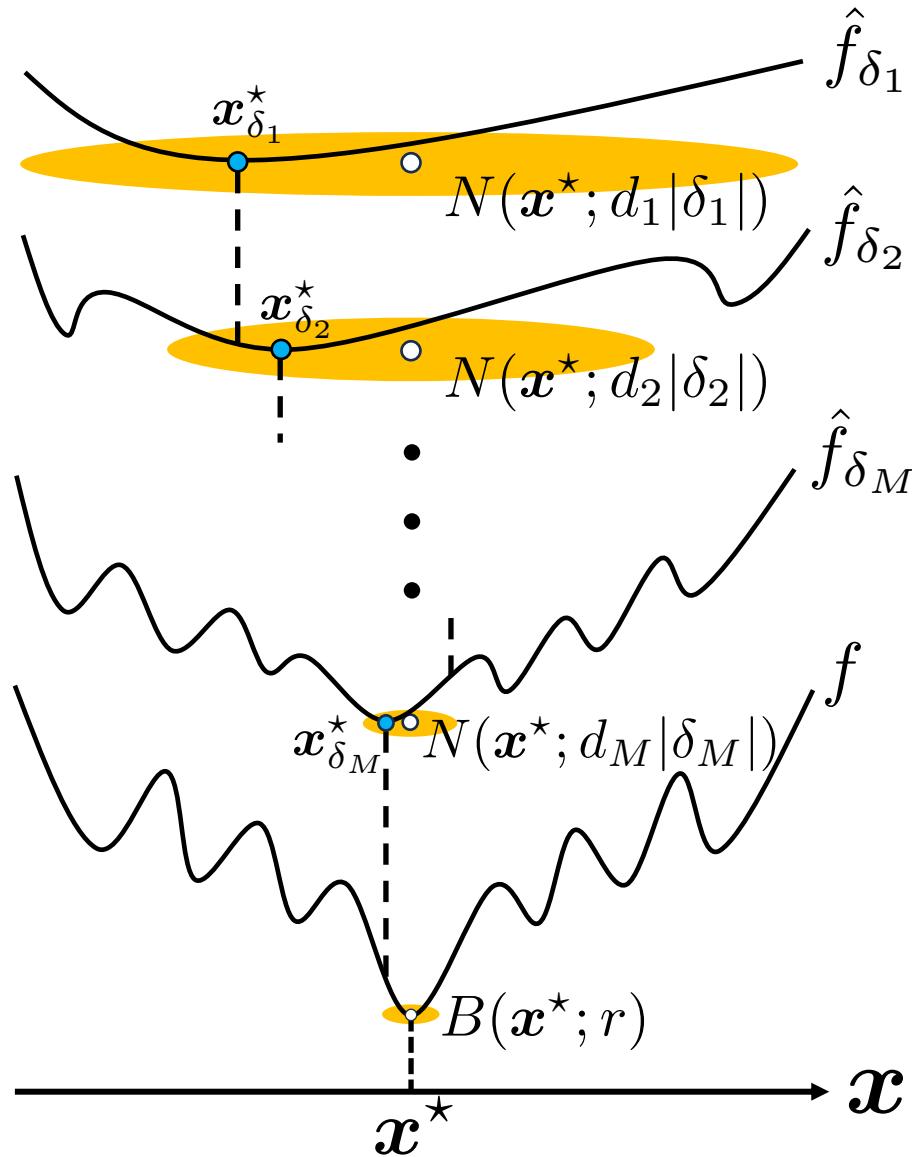
$$\|x_\delta^* - x_{\delta/2}^*\| \leq \frac{\delta}{2} \quad \left( \gamma_m := \frac{1}{2} \right)$$

が成り立つ。

(ii) 任意の  $\underline{\delta > 0}$  に対して、関数  $\hat{f}_\delta(x)$  は近傍  $N(x_\delta^*; 3\delta)$  で  $\sigma$ -強凸となる。



# New $\sigma$ -nice 関数



## 定義 (New $\sigma$ -nice 関数)

任意の非凸リップシツツ関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は new  $\sigma$ -nice 関数であるという。

(i) 任意の  $\underline{\delta_m \in \mathbb{R}}$  と  $x_{\delta_m}^*$  に対して、

$$\|x_{\delta_m}^* - x_{\delta_{m+1}}^*\| \leq |\delta_m| - |\delta_{m+1}|$$

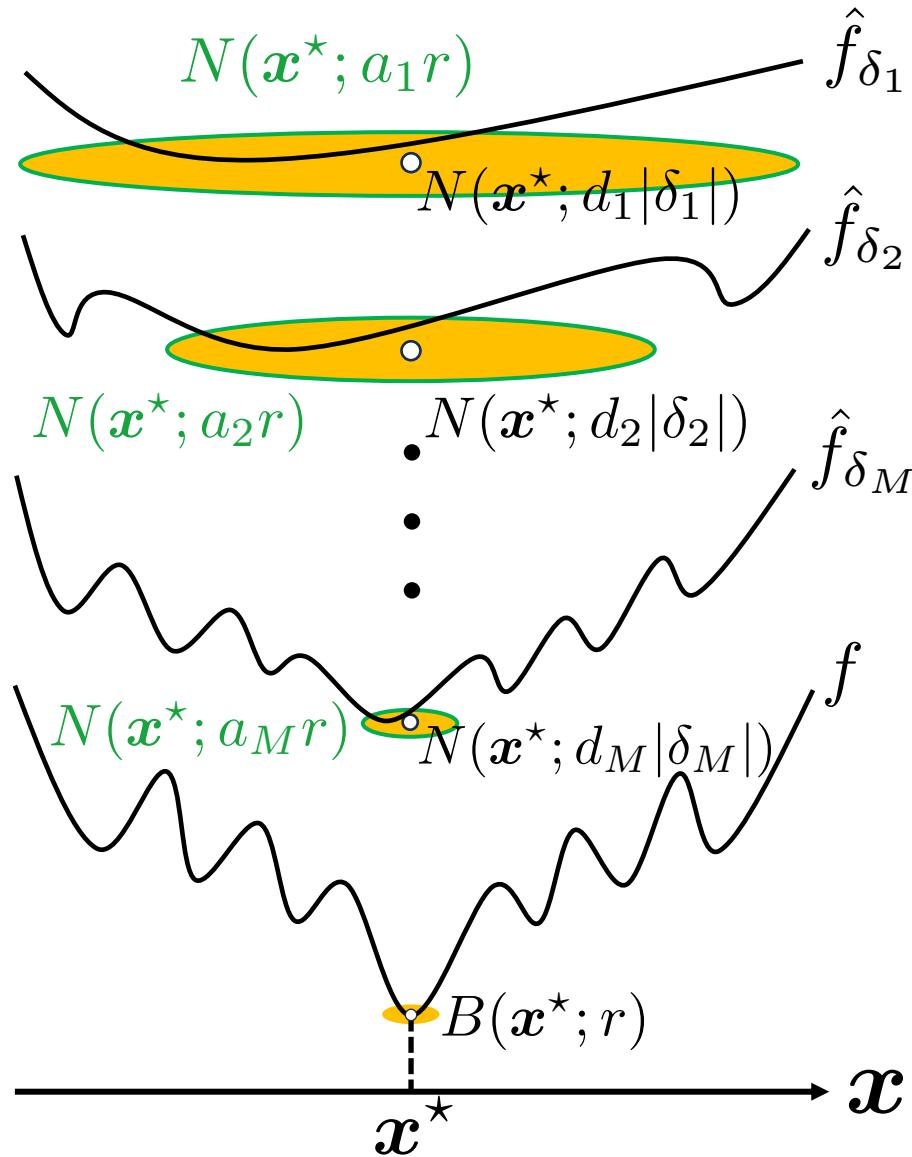
が成り立つ。ただし、

$$\delta_{m+1} := \underline{\gamma_m \delta_m} \quad (\gamma_m \in (0, 1))$$

とする。

(ii) 任意の  $\underline{\delta_m \in \mathbb{R}}$  に対して、関数  $\hat{f}_{\delta_m}(x)$  は近傍  $N(x^*; \underline{d_m} |\delta_m|)$  で  $\sigma$ -強凸となる。

# New $\sigma$ -nice性の十分条件



## 命題3.2

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が十分小さい  $r > 0$  に対して閉球  $B(x^*; r)$  で  $\sigma$ -強凸だとすると、関数  $f$  が new  $\sigma$ -nice 関数となるための十分条件は、ノイズレベル  $\delta_m$  が次の条件を満たすことである。

$$\frac{2 \max \left\{ \left\| \nabla \hat{f}_{\delta_m}(x^*) \right\|, \left\| \nabla \hat{f}_{\delta_{m+1}}(x^*) \right\| \right\}}{\sigma(1 - \gamma_m)} \leq |\delta_m| = |\delta^-|$$

## $\delta^-$ が存在する確率 $p(a_m)$

$$0 < p(a_m) := \frac{2 \arccos \left( \frac{r \sqrt{a_m^2 - 1}}{\|x^* - x\| \|u_m\|} \right)}{\pi} < 1$$

→  $a_m$  が大きいとき、 $p(a_m)$  はほぼ0

# 背景 : SGDの確率的ノイズと平滑化(再掲)

$$(\text{GD}) \quad \mathbf{y}_t := \mathbf{x}_t - \eta \nabla f(\mathbf{x}_t)$$

$$(\text{SGD}) \quad \mathbf{x}_{t+1} := \mathbf{x}_t - \eta \nabla f_{S_t}(\mathbf{x}_t)$$

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim B(\mathbf{0}; 1)} [f(\mathbf{x} - \delta \mathbf{u})]$$

(確率的ノイズ)  $\boldsymbol{\omega}_t := \nabla f_{S_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)$  とすると、次の式が成り立つ。

$$\mathbb{E}_{\boldsymbol{\omega}_t} [\mathbf{y}_{t+1}] = \mathbf{y}_t - \eta \nabla \underbrace{\mathbb{E}_{\boldsymbol{\omega}_t} [f(\mathbf{y}_t - \eta \boldsymbol{\omega}_t)]}_{:= \hat{f}(\mathbf{y}_t)}$$

- ▷ ”関数  $f(\mathbf{x}_t)$  を SGD で最適化すること”と、”関数  $\hat{f}(\mathbf{y}_t)$  を 最急降下法 で最適化すること”は、期待値の意味では 等価 であると言える。
- ▷ ある程度 平滑化 された関数が、 最急降下法 で最適化されていると みなせる。

[1] R. Kleinberg et al. “An Alternative View: When Does SGD Escape Local Minima?”

In *Proceedings of the 35th International Conference on Machine Learning*, pp2703-2712, 2018

# SGDの平滑化特性

## 補題2.1

$$\mathbb{E}_{\xi_t} \left[ \|\nabla f_{S_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)\|^2 \right] \leq \frac{C^2}{b}$$

## ▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim B(\mathbf{0}; 1)} [f(\mathbf{x} - \delta \mathbf{u})]$$

▷  $\omega_t := \nabla f_{S_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)$  と補題2.1から、

$$\omega_t = \frac{C}{\sqrt{b}} \mathbf{u}_t \quad (\mathbf{u}_t \sim B(\mathbf{0}; 1))$$

↑ 中心0, 半径1の閉球

が成り立つ。したがって、

$$\begin{aligned} \mathbb{E}_{\omega_t} [\mathbf{y}_{t+1}] &= \mathbf{y}_t - \eta \nabla \mathbb{E}_{\omega_t} [f(\mathbf{y}_t - \underline{\eta \omega_t})] \\ &= \mathbf{y}_t - \eta \nabla \mathbb{E}_{\mathbf{u}_t \sim B(\mathbf{0}; 1)} \left[ f \left( \mathbf{y}_t - \frac{\eta C}{\sqrt{b}} \mathbf{u}_t \right) \right] \end{aligned}$$

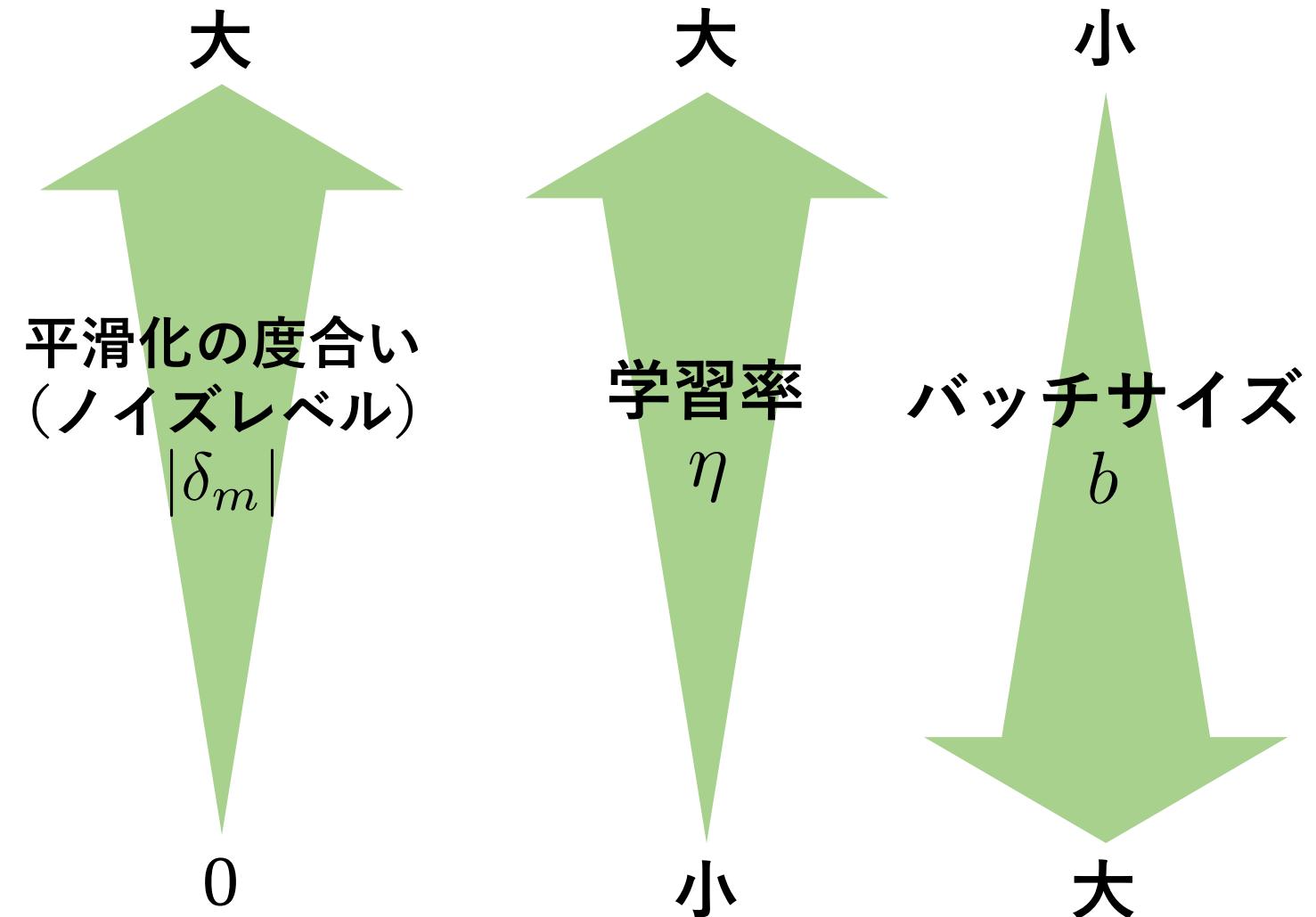
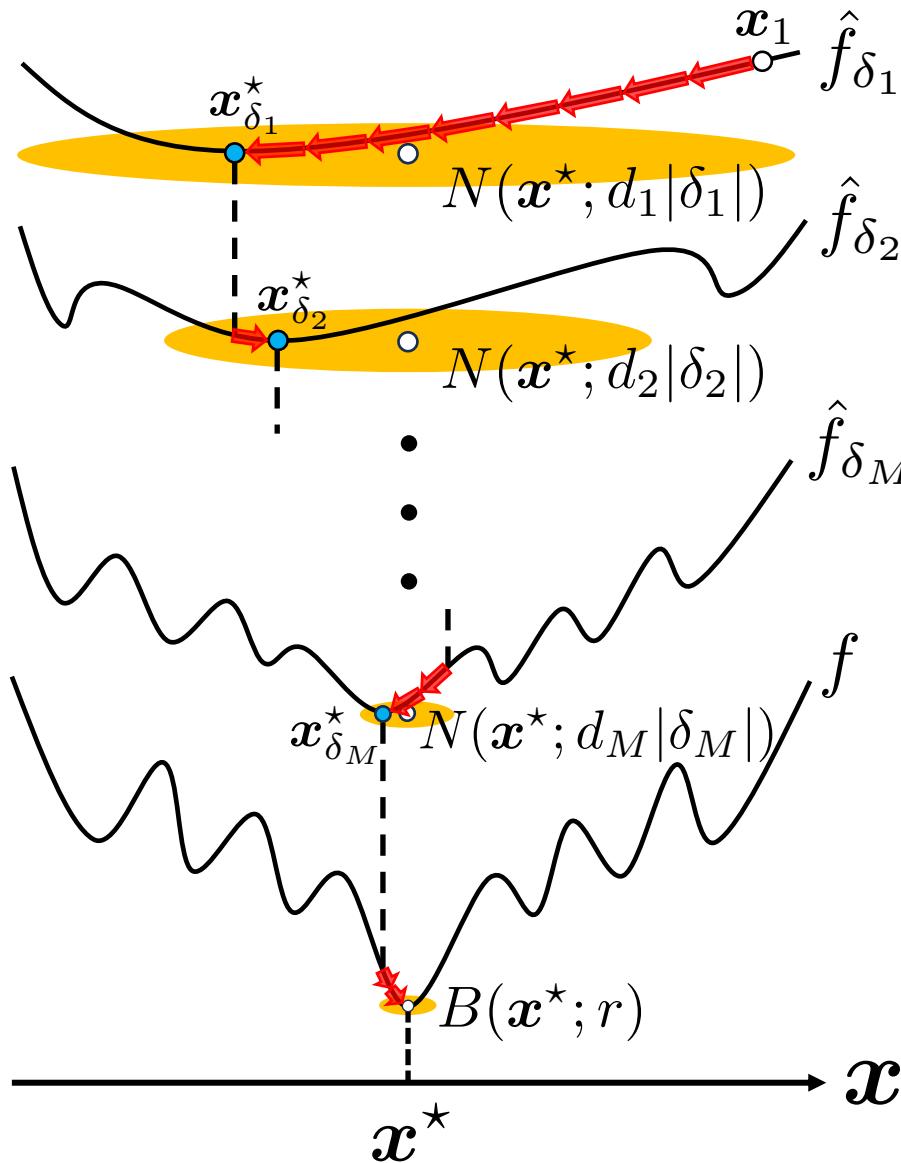
$$= \mathbf{y}_t - \eta \nabla \hat{f}_{\frac{\eta C}{\sqrt{b}}}(\mathbf{y}_t)$$

大きさ  $\eta C / \sqrt{b}$  のノイズで平滑化された関数

が成り立つ。

$$\delta = \eta \sqrt{\frac{C^2}{b}}$$

# SGDの平滑化特性



$$\delta = \eta \sqrt{\frac{C^2}{b}}$$

# 暗黙的な段階的最適化アルゴリズム

## Algorithm 3 Implicit Graduated Optimization

**Require:**  $\epsilon > 0, p \in (0, 1], \bar{d} > 0, \mathbf{x}_1, \eta_1, b_1$

$$\delta_1 := \frac{\eta_1 C}{\sqrt{b_1}}$$

$$\alpha_0 := \min \left\{ \frac{\sqrt{b_1}}{4L_f \eta_1 C (1+\bar{d})}, \frac{\sqrt{b_1}}{\sqrt{2\sigma} \eta_1 C} \right\}, M^p := \frac{1}{\alpha_0 \epsilon}$$

**for**  $m = 1$  to  $M + 1$  **do**

**if**  $m \neq M + 1$  **then**

$$\epsilon_m := \sigma^2 \delta_m^2, T_F := H_4 / (\epsilon_m - H_3 \eta_m)$$

$$\gamma_m := \frac{(M-m)^p}{\{M-(m-1)\}^p}$$

$$\kappa_m / \sqrt{\lambda_m} = \gamma_m \quad (\kappa_m \in (0, 1], \lambda_m \geq 1)$$

**end if**

$$\mathbf{x}_{m+1} := \text{GD}(T_F, \mathbf{x}_m, \hat{f}_{\delta_m}, \eta_m, b_m)$$

$$\eta_{m+1} := \kappa_m \eta_m, b_{m+1} := \lambda_m b_m$$

$$\delta_{m+1} := \frac{\eta_{m+1} C}{\sqrt{b_{m+1}}}$$

**end for**

**return**  $\mathbf{x}_{M+2}$

## 定理3.4

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が  $L_f$ -リップシツツ new  $\sigma$ -nice 関数だとすると、アルゴリズム3は、 $\mathcal{O}\left(1/\epsilon^{\frac{1}{p}}\right)$  ( $p \in (0, 1]$ ) 回の反復で、関数  $f$  の**大域的最適解**  $\mathbf{x}^*$  の  $\epsilon$ -近傍に到達する。

## Algorithm 4 Gradient Descent (GD)

**Require:**  $T_F, \hat{\mathbf{x}}_1, F, \eta, b$

**for**  $t = 1$  to  $T_F$  **do**

$$\hat{\mathbf{x}}_{t+1} := \hat{\mathbf{x}}_t - \eta \nabla F(\mathbf{x}_t)$$

**end for**

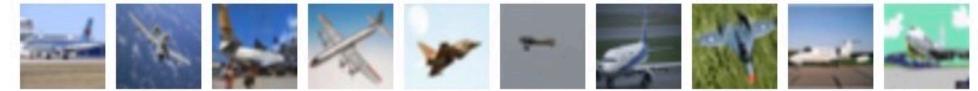
**return**  $\hat{\mathbf{x}}_{T_F+1} = \text{GD}(T_F, \hat{\mathbf{x}}_1, F, \eta)$

# 数値実験

---

- ▷ モデル: ResNet18
- ▷ データセット: CIFAR100
  - ▷ 100クラスの画像分類
  - ▷ 5万枚の訓練データ
  - ▷ 1万枚のテストデータ
  - ▷ 画像サイズは $32 \times 32 \times 3$
- ▷ 交差エントロピー誤差
- ▷ NVIDIA GeForce RTX4090

airplane



automobile



bird



cat



deer



dog



frog



horse



ship



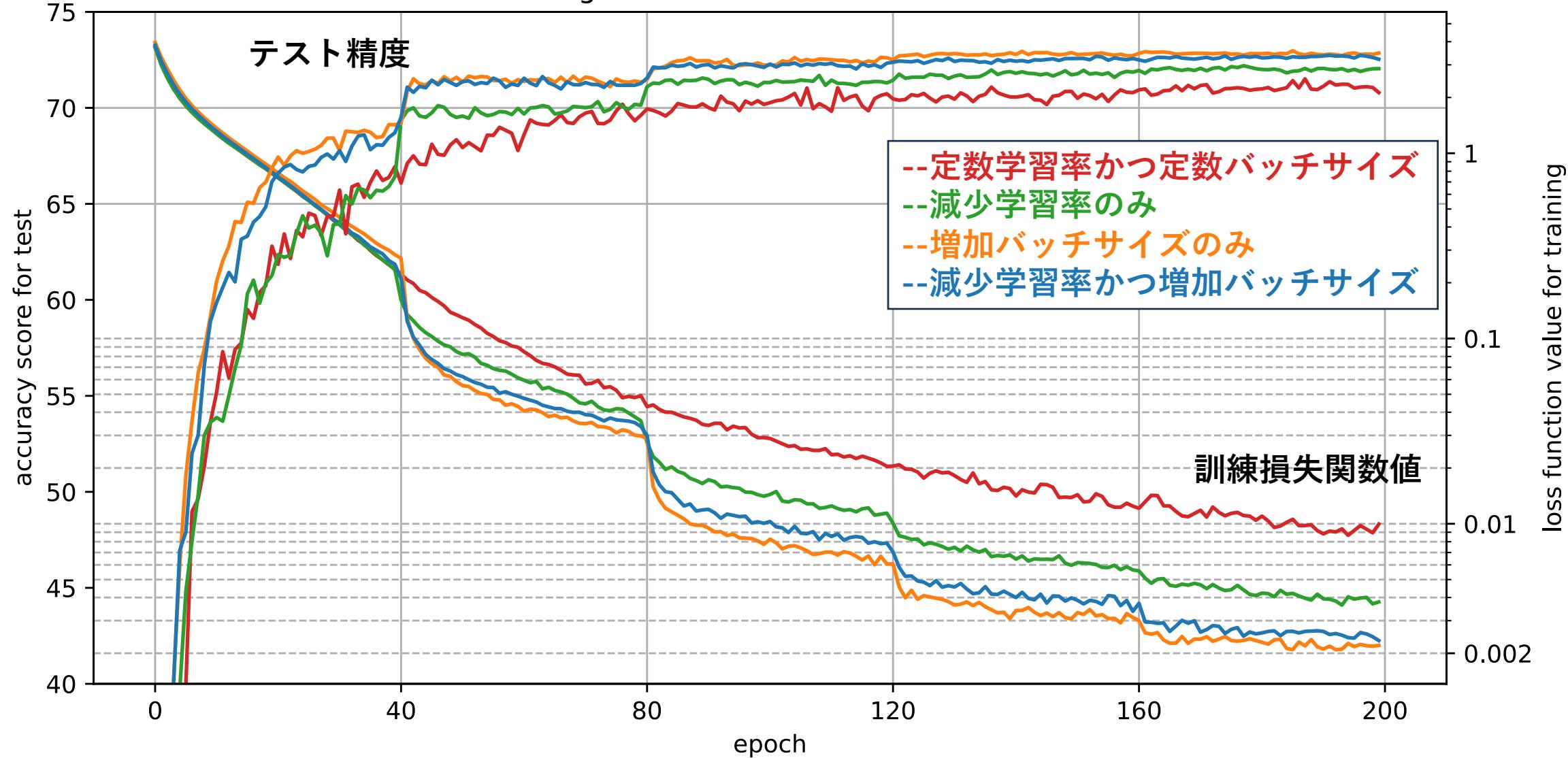
truck



# 数值実験—暗黙的な段階的最適化

$$\delta = \eta \sqrt{\frac{C^2}{b}}$$

Training ResNet18 on CIFAR100 dataset



# 結論

---

- ▷  $\sigma$ -nice 関数を定義できない場合があることを示し、new  $\sigma$ -nice 関数を定義した。また、任意の関数がnew  $\sigma$ -nice 関数となるための十分条件を示した。
- ▷ SGDには関数を平滑化する効果があり、  
その度合いは学習率とバッチサイズによって決まる。
$$\delta = \eta \sqrt{\frac{C^2}{b}}$$
- ▷ この性質を利用した暗黙的な段階的最適化アルゴリズムアルゴリズムを提案し、このアルゴリズムが  $\mathcal{O}\left(1/\epsilon^{\frac{1}{p}}\right)$  ( $p \in (0, 1]$ ) 回の反復で大域的最適解  $x^*$  の  $\epsilon$ -近傍に到達することを示した。
- ▷ したがって、計算資源が十分にあれば、シンプルなSGDでもDNNを含む非凸関数の大域的最適化が可能となる。

# 付録：なぜ段階的最適化なのか

- ▷ 機械学習分野で既に成功している手法だからです。
- ▷ スコアベースモデル（2019）と拡散モデル（2020）
  - ▷ 敵対的生成ネットワーク（GANs）を上回る性能を誇る、最先端の画像生成モデルで、暗黙的に段階的最適化のアプローチを利用しています。

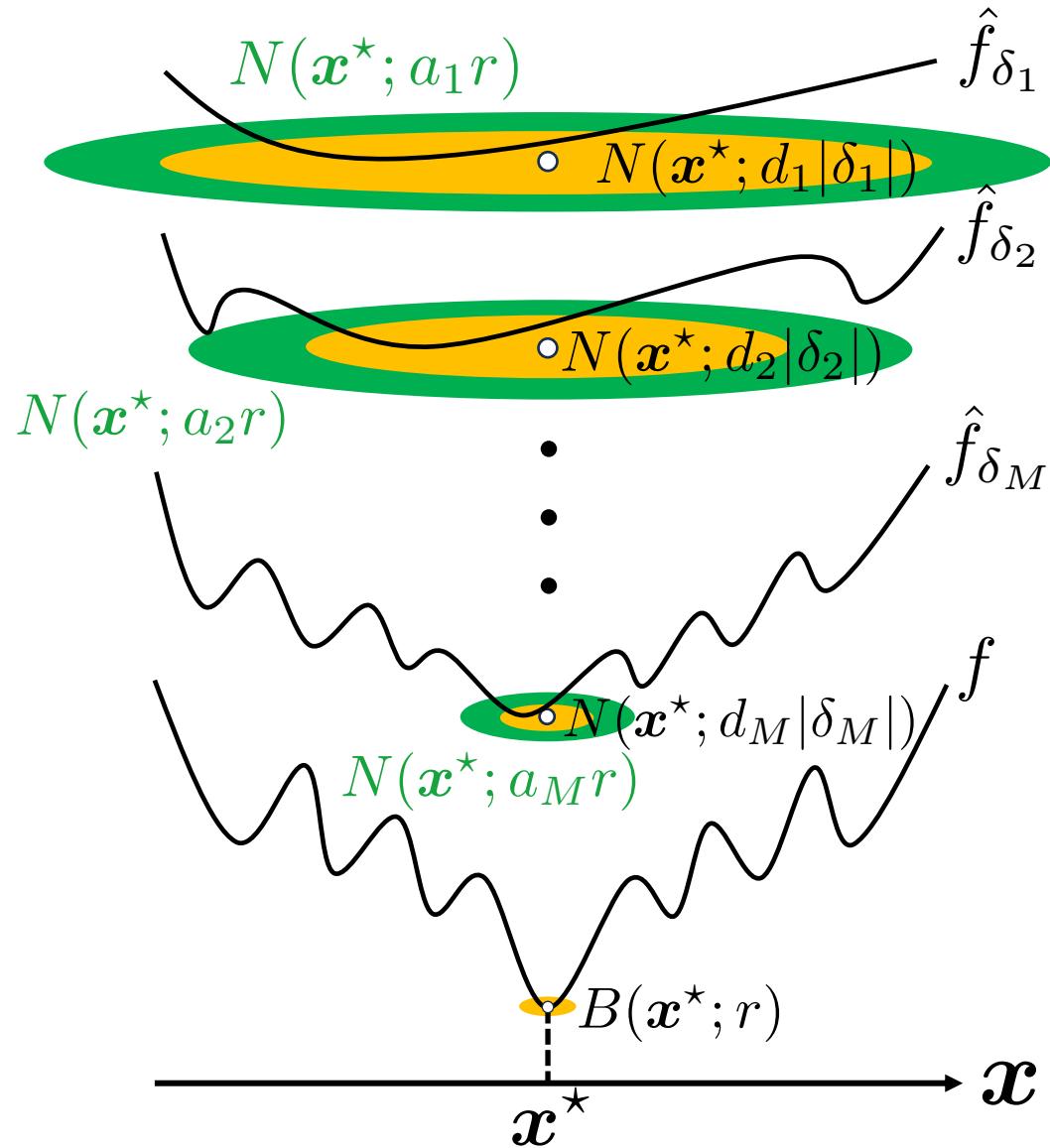


小さいノイズ

大きいノイズ

→ 段階的最適化の視点から収束解析を提供できる可能性があります。

# 付録 : New $\sigma$ -nice 性の十分条件



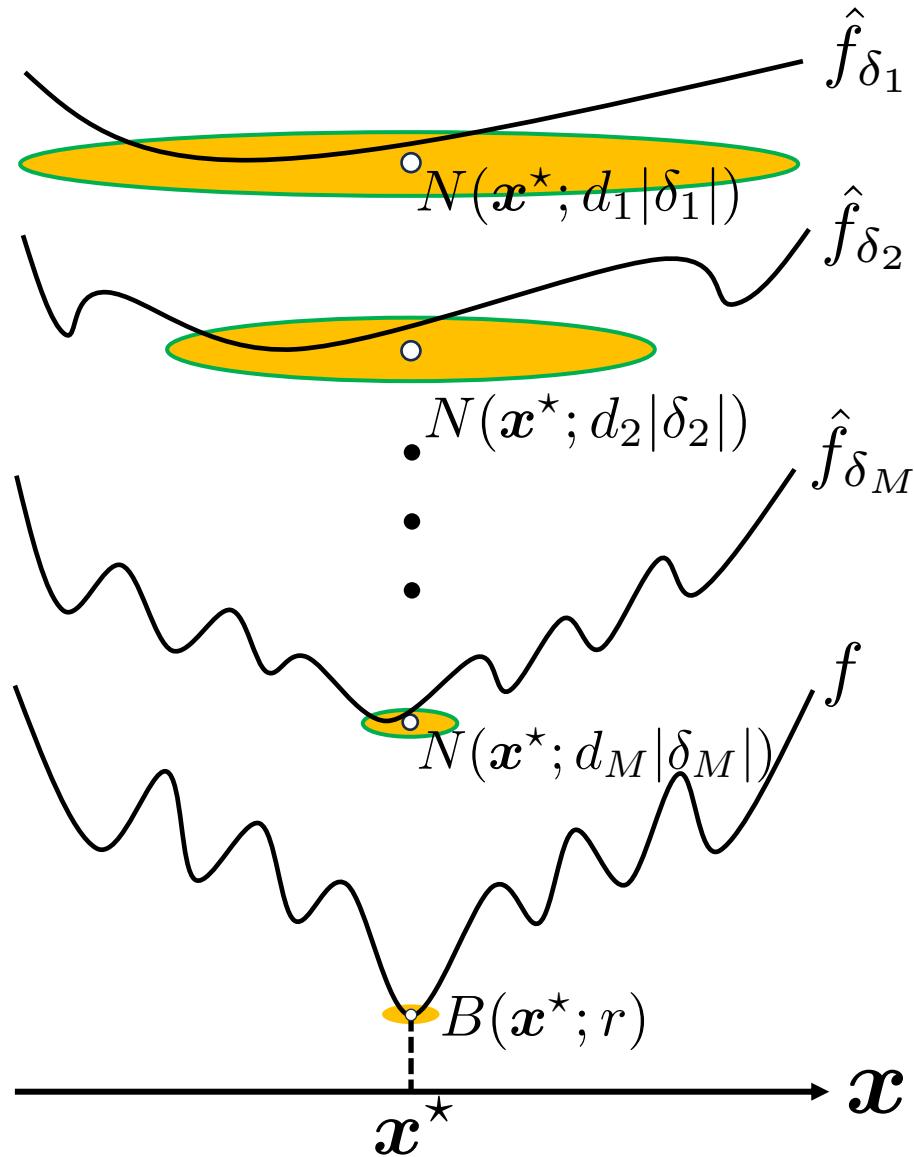
## 命題3.1

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が十分小さい  $r > 0$  に対して閉球  $B(x^*; r)$  で  $\sigma$ -強凸だとし、ノイズレベル  $\delta_m$  が  $|\delta_m| = |\delta^-|$  を満たすならば、平滑化された関数  $\hat{f}_{\delta_m}$  は近傍  $\underline{N(x^*; a_m r)} (a_m > \sqrt{2})$  で  $\sigma$ -強凸となる。ただし、

$$|\delta^-| := \left| \|x^* - x\| \|\mathbf{u}_m\| \cos \theta - \sqrt{\|x^* - x\|^2 \|\mathbf{u}_m\|^2 \cos^2 \theta - r^2 (a_m^2 - 1)} \right|$$

とし、 $x \in \underline{N(x^*; a_m r)}$ ,  $\mathbf{u}_m \sim B(\mathbf{0}; 1)$  で、 $\theta$  は  $\mathbf{u}_m$  と  $x^* - x$  のなす角とする。また、  
 $d_m := a_m r / |\delta^-|$  とすると、関数  $\hat{f}_{\delta_m}$  は近傍  $\underline{N(x^*; d_m |\delta_m|)}$  でも  $\sigma$ -強凸となる。

# 付録 : New $\sigma$ -nice 性の十分条件



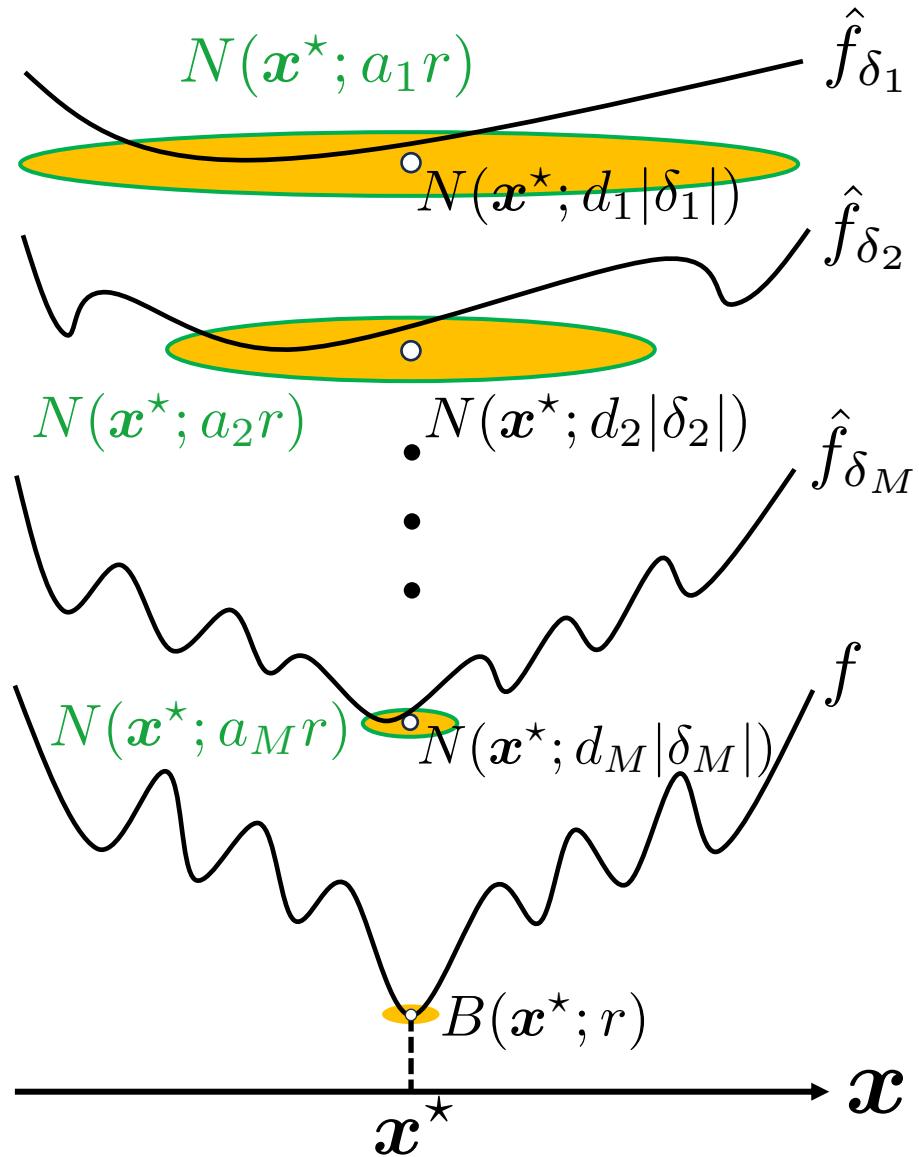
## 命題3.1

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が十分小さい  $r > 0$  に対して閉球  $B(x^*; r)$  で  $\sigma$ -強凸だとし、ノイズレベル  $\delta_m$  が  $|\delta_m| = |\delta^-|$  を満たすならば、平滑化された関数  $\hat{f}_{\delta_m}$  は近傍  $N(x^*; a_m r)$  ( $a_m > \sqrt{2}$ ) で  $\sigma$ -強凸となる。ただし、

$$|\delta^-| := \left| \|x^* - x\| \|\mathbf{u}_m\| \cos \theta - \sqrt{\|x^* - x\|^2 \|\mathbf{u}_m\|^2 \cos^2 \theta - r^2 (a_m^2 - 1)} \right|$$

とし、 $x \in N(x^*; a_m r)$ ,  $\mathbf{u}_m \sim B(\mathbf{0}; 1)$  で、 $\theta$  は  $\mathbf{u}_m$  と  $x^* - x$  のなす角とする。また、  
 $d_m := a_m r / |\delta^-|$  とすると、関数  $\hat{f}_{\delta_m}$  は近傍  $N(x^*; d_m |\delta_m|)$  でも  $\sigma$ -強凸となる。

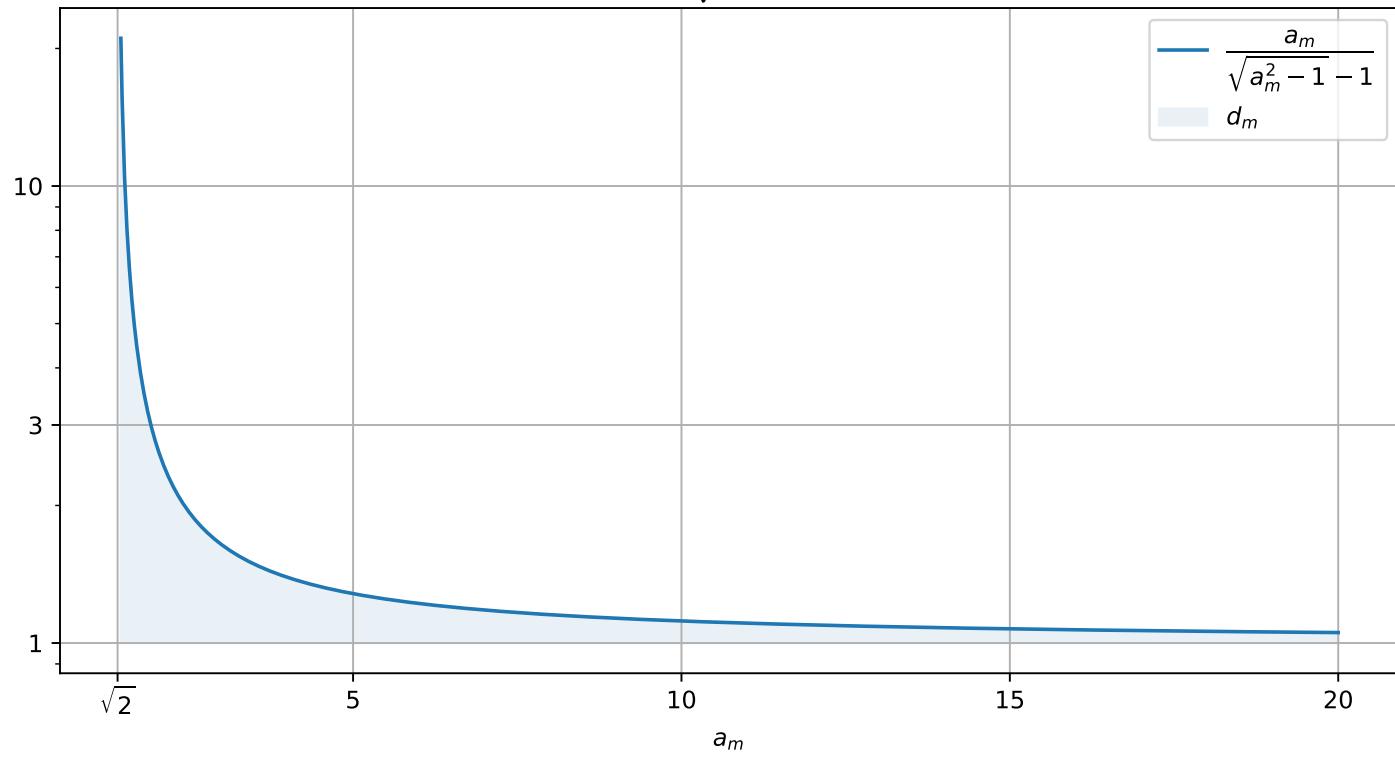
# 付録 : New $\sigma$ -nice 性の十分条件



$d_m$  がとるべき値の範囲

$$d_m := a_m r / |\delta^-|$$

$$\rightarrow 1 \leq d_m \leq \frac{a_m}{\sqrt{a_m^2 - 1} - 1}$$



# 付録 : $\sigma$ -nice 関数 (再掲)

## 定義 ( $\sigma$ -nice 関数)

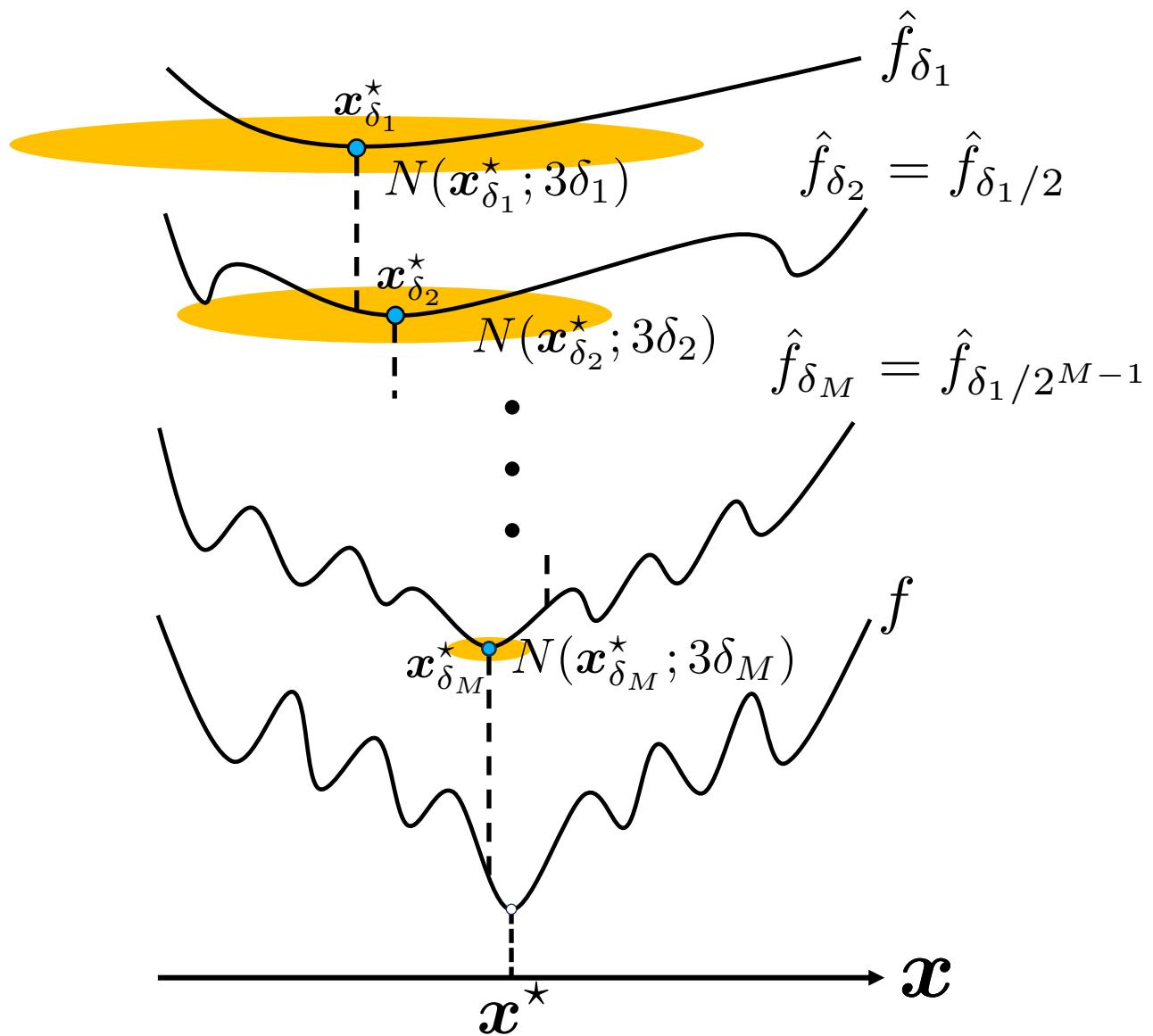
任意の非凸関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は  $\sigma$ -nice 関数であるという。

(i) 任意の  $\delta > 0$  と  $x_\delta^*$  に対して、

$$\|x_\delta^* - x_{\delta/2}^*\| \leq \frac{\delta}{2}$$

が成り立つ。

(ii) 任意の  $\delta > 0$  に対して、関数  $\hat{f}_\delta(x)$  は近傍  $N(x_\delta^*; 3\delta)$  で  $\sigma$ -強凸となる。

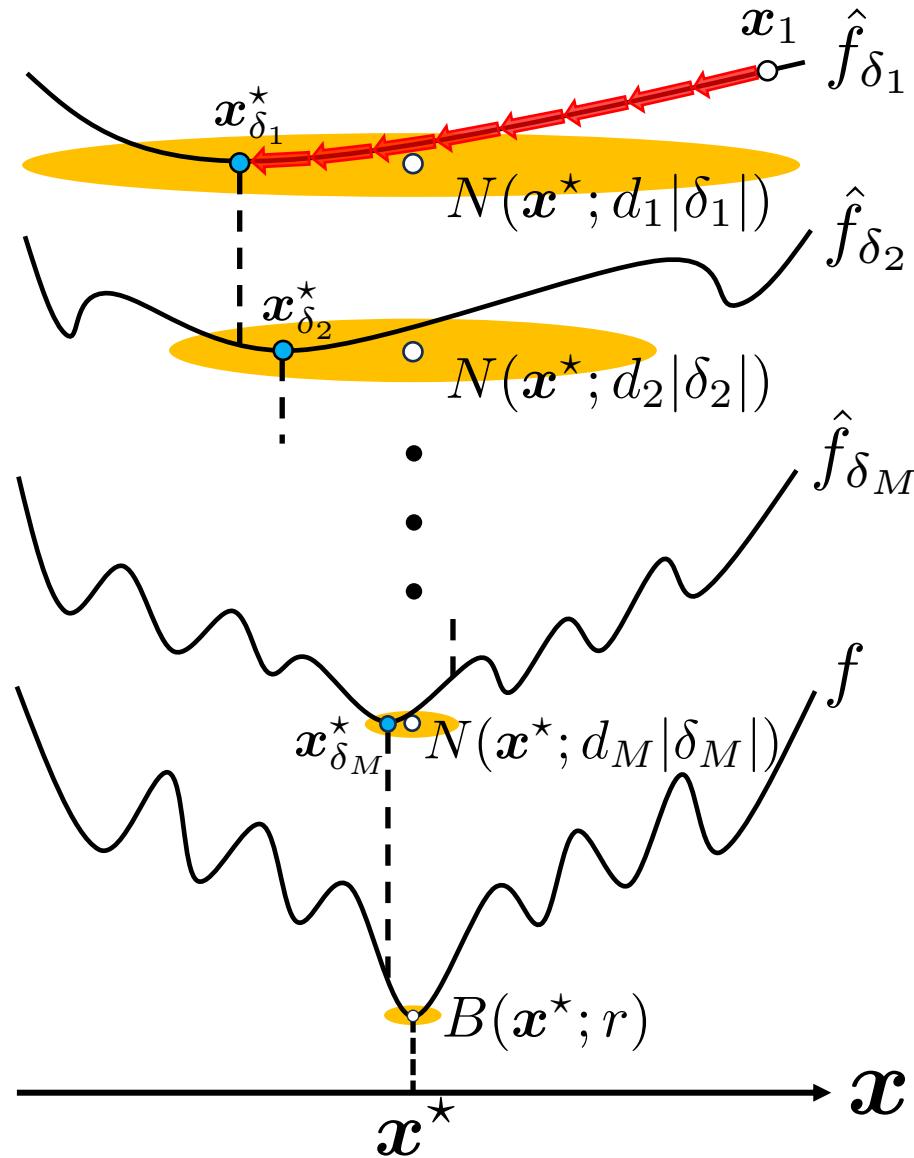


## 付録：SGDの挙動についての理論的洞察

---

- ▷なぜ大きいバッチサイズは**汎化性**を悪化させるのか。  
→ バッチサイズが大きすぎると、最適化される関数は元の非凸関数に近付き、悪い局所的最適解に陥りやすくなるため。
  
- ▷なぜ**減少学習率**または**増加バッチサイズ**は定数よりも優れるのか。  
→ 訓練中に学習率を減少させる、あるいはバッチサイズを増加させることは、ノイズレベルを小さくすることと等価。  
→ 減少学習率または増加バッチサイズを使うことは、まさに**段階的最適化**のアプローチになっているため、定数よりも優れる。

# 付録：最適なノイズスケジューリング



## 命題3.3

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が new  $\sigma$ -nice 関数ならば、任意の  $m \in [M-1]$  に対して、

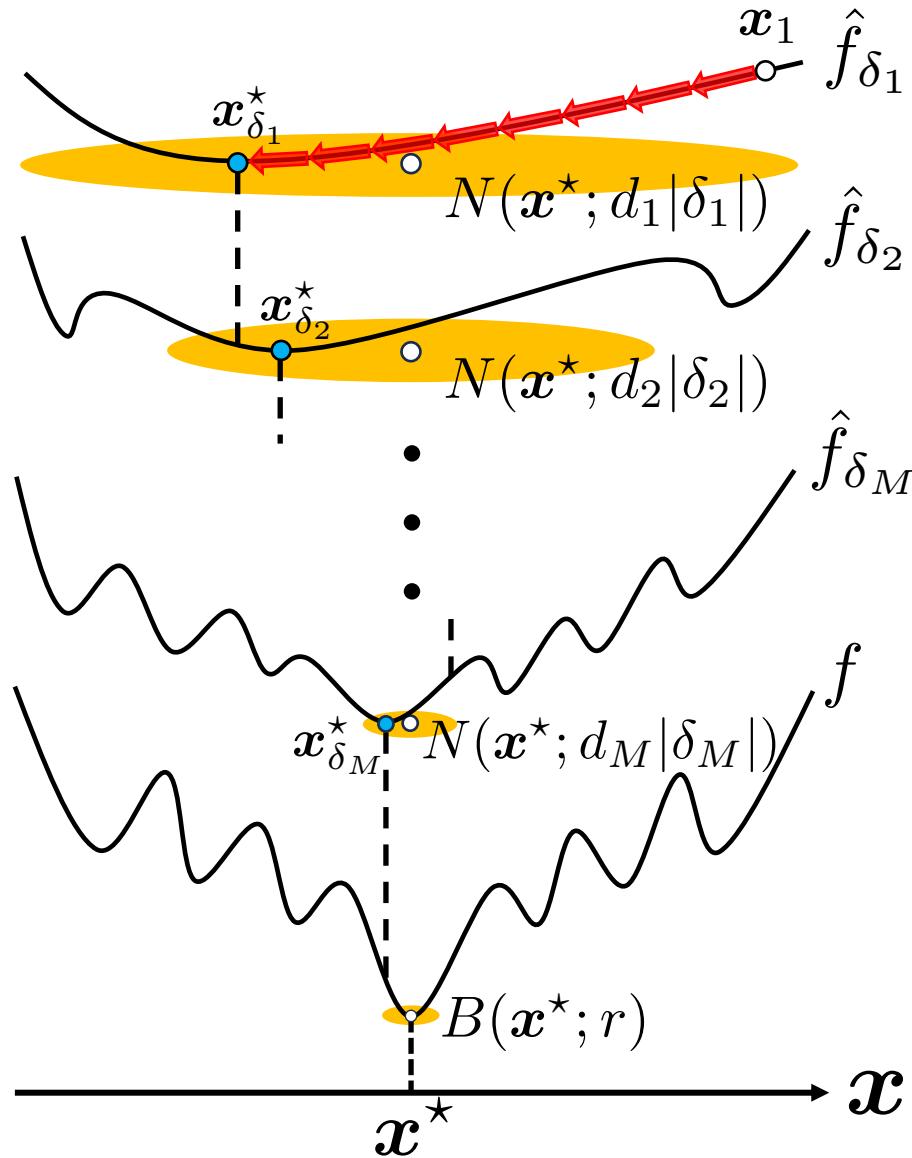
$\|x_{\delta_m}^* - x^*\| < d_{m+1} |\delta_{m+1}|$

が成り立つ。

- ▷ 次の関数の最適化の初期点が、必ず **強凸領域**に含まれる。
- ▷ 最初の初期点  $x_1$  が関数  $\hat{f}_{\delta_1}$  の **強凸領域**  $N(x^*; d_1 |\delta_1|)$  に含まれていれば、大域的最適解  $x^*$  に到達できる。
- ▷ 命題3.3が成り立つためには、

$\gamma_m \in \left( \frac{1}{d_{m+1}}, 1 \right)$  が必要。

# 付録：最適なノイズスケジューリング

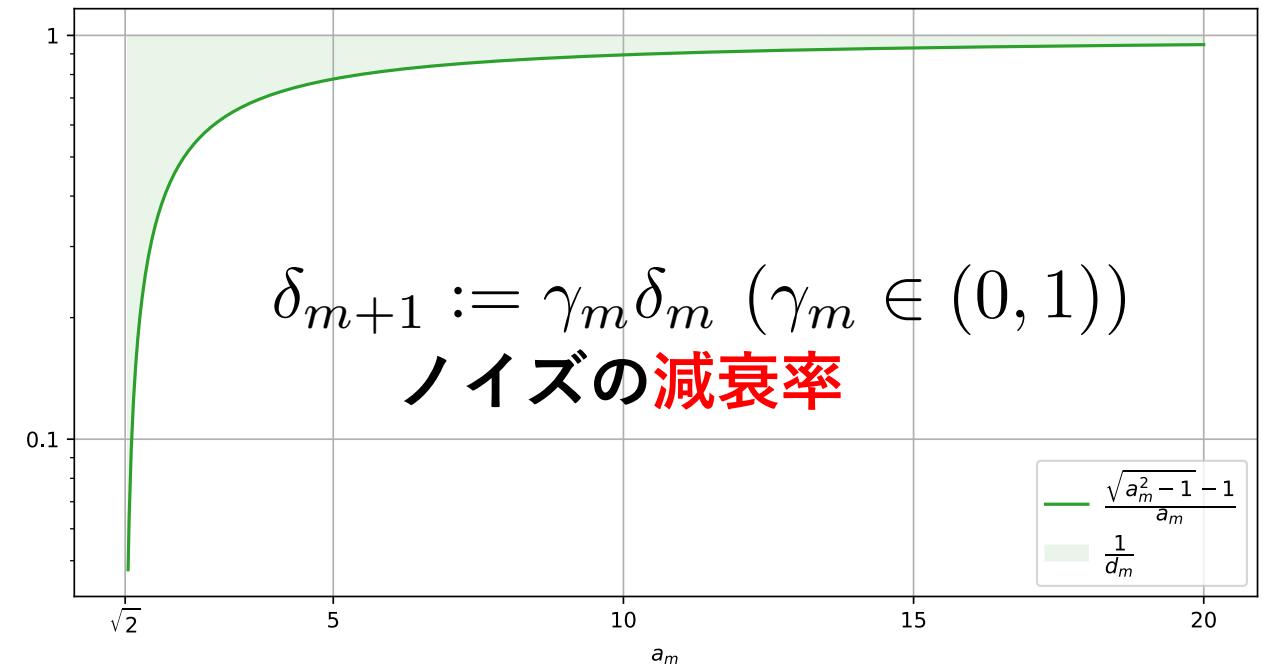


▷ 命題3.3が成り立つためには、

$$\gamma_m \in \left( \frac{1}{d_{m+1}}, 1 \right) \text{ が必要。}$$

$$d_m := a_m r / |\delta^-|$$

$$\rightarrow \frac{\sqrt{a_m^2 - 1} - 1}{a_m} \leq \frac{1}{d_m} \leq 1$$



# 付録 : $\sigma$ -nice 関数 (再掲)

## 定義 ( $\sigma$ -nice 関数)

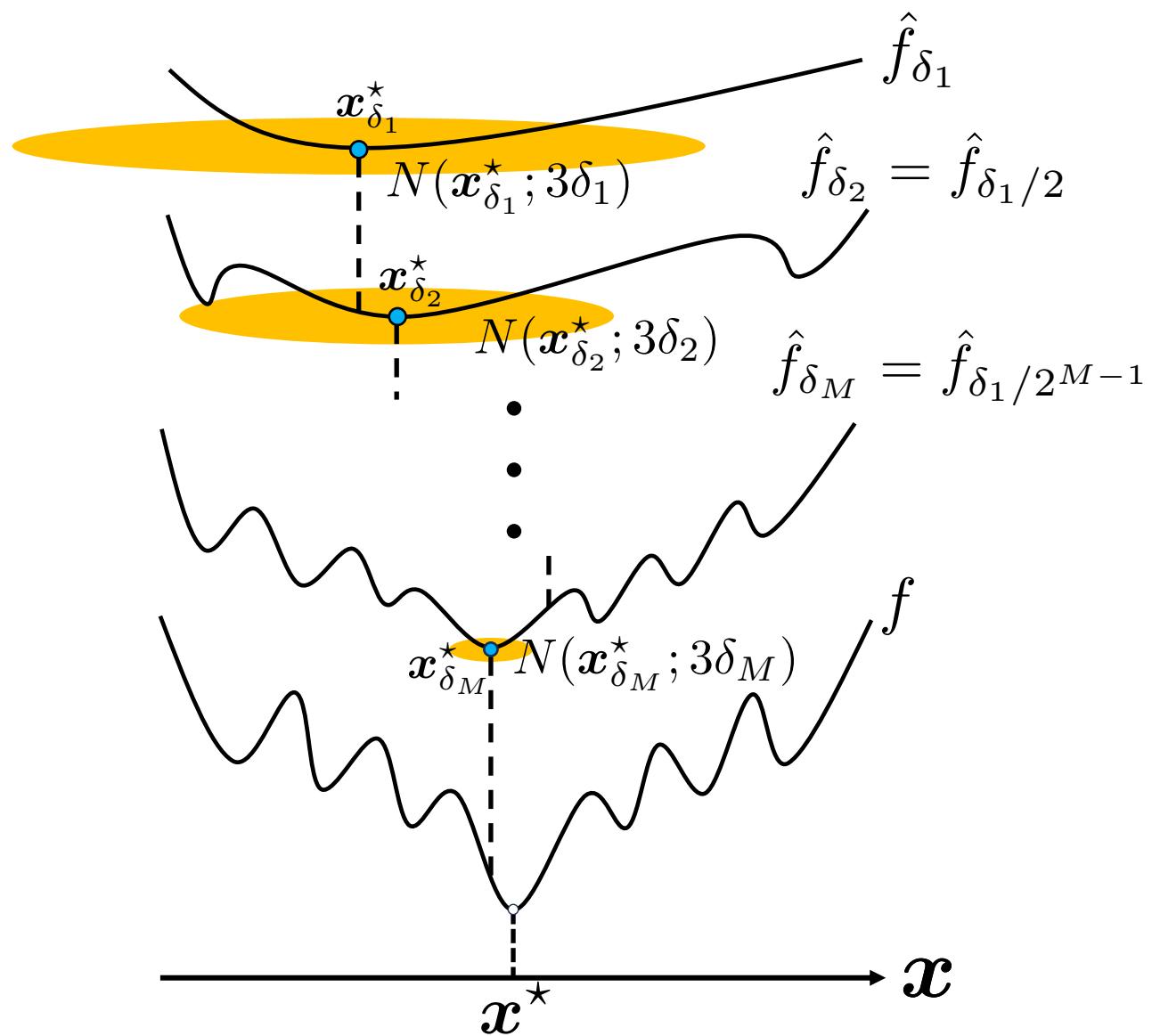
任意の非凸関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が次の2つの条件を満たすとき、関数  $f$  は  $\sigma$ -nice 関数であるという。

(i) 任意の  $\underline{\delta > 0}$  と  $x_\delta^*$  に対して、

$$\|x_\delta^* - x_{\delta/2}^*\| \leq \frac{\delta}{2} \quad \left( \gamma_m := \frac{1}{2} \right)$$

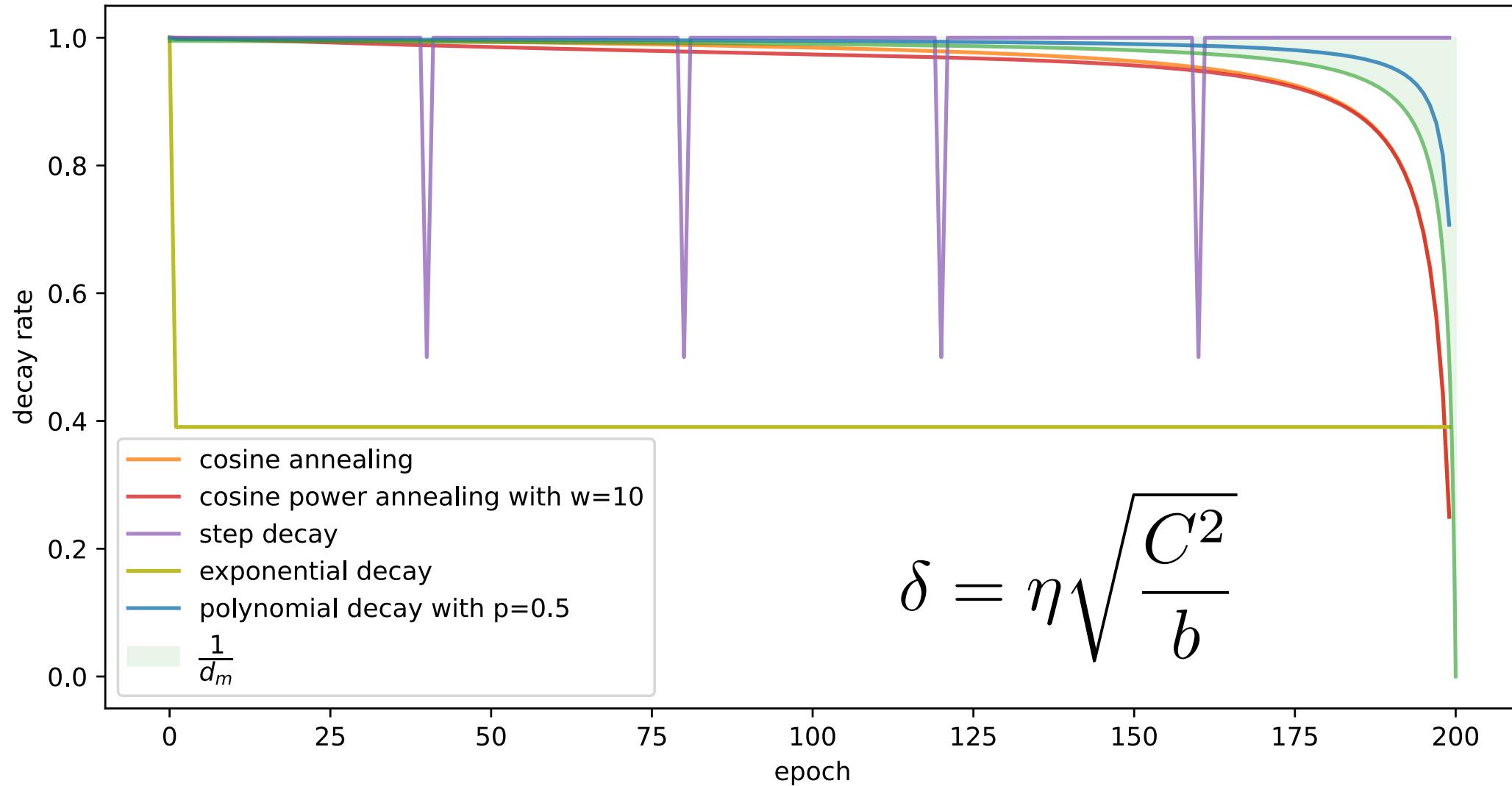
が成り立つ。

(ii) 任意の  $\underline{\delta > 0}$  に対して、関数  $\hat{f}_\delta(x)$  は近傍  $N(x_\delta^*; 3\delta)$  で  $\sigma$ -強凸となる。

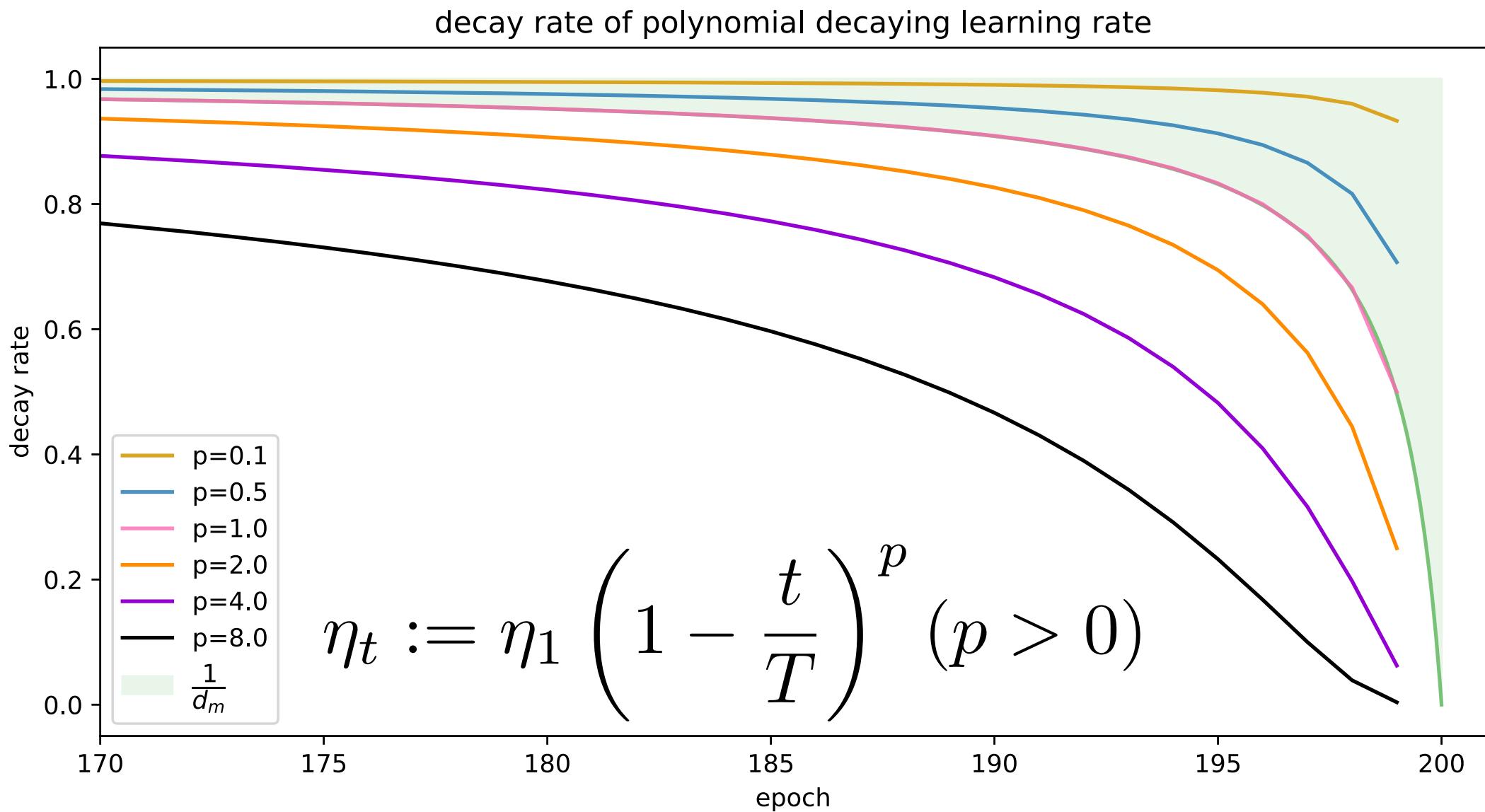


# 付録：最適な学習率スケジューリング

▷最適なノイズの減衰率は、直ちに最適な学習率の減衰率となる。



# 付録：最適な学習率スケジューリング



# 付録：暗黙的な段階的最適化アルゴリズム

## Algorithm 3 Implicit Graduated Optimization

**Require:**  $\epsilon > 0, p \in (0, 1), \bar{d} > 0, \mathbf{x}_1, \eta_1, b_1$

$$\delta_1 := \frac{\eta_1 C}{\sqrt{b_1}}$$

$$\alpha_0 := \min \left\{ \frac{\sqrt{b_1}}{4L_f \eta_1 C (1+\bar{d})}, \frac{\sqrt{b_1}}{\sqrt{2\sigma} \eta_1 C} \right\}, M^p := \frac{1}{\alpha_0 \epsilon}$$

**for**  $m = 1$  to  $M + 1$  **do**

**if**  $m \neq M + 1$  **then**

$$\epsilon_m := \sigma^2 \delta_m^2, T_F := H_4 / (\epsilon_m - H_3 \eta_m)$$

$$\gamma_m := \frac{(M-m)^p}{\{M-(m-1)\}^p}$$

$$\kappa_m / \sqrt{\lambda_m} = \gamma_m \quad (\kappa_m \in (0, 1], \lambda_m \geq 1)$$

**end if**

$$\mathbf{x}_{m+1} := \text{GD}(T_F, \mathbf{x}_m, \hat{f}_{\delta_m}, \eta_m, b_m)$$

$$\eta_{m+1} := \kappa_m \eta_m, b_{m+1} := \lambda_m b_m$$

$$\delta_{m+1} := \frac{\eta_{m+1} C}{\sqrt{b_{m+1}}}$$

**end for**

**return**  $\mathbf{x}_{M+2}$

## 定理3.4

関数  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  が  $L_f$ -リップシツツ new  $\sigma$ -nice 関数だとすると、アルゴリズム3は、 $\mathcal{O}\left(1/\epsilon^{\frac{1}{p}}\right)$  ( $p \in (0, 1]$ ) 回の反復で、関数  $f$  の**大域的最適解**  $\mathbf{x}^*$  の  $\epsilon$ -近傍に到達する。

## Algorithm 4 Gradient Descent (GD)

**Require:**  $T_F, \hat{\mathbf{x}}_1, F, \eta, b$

**for**  $t = 1$  to  $T_F$  **do**

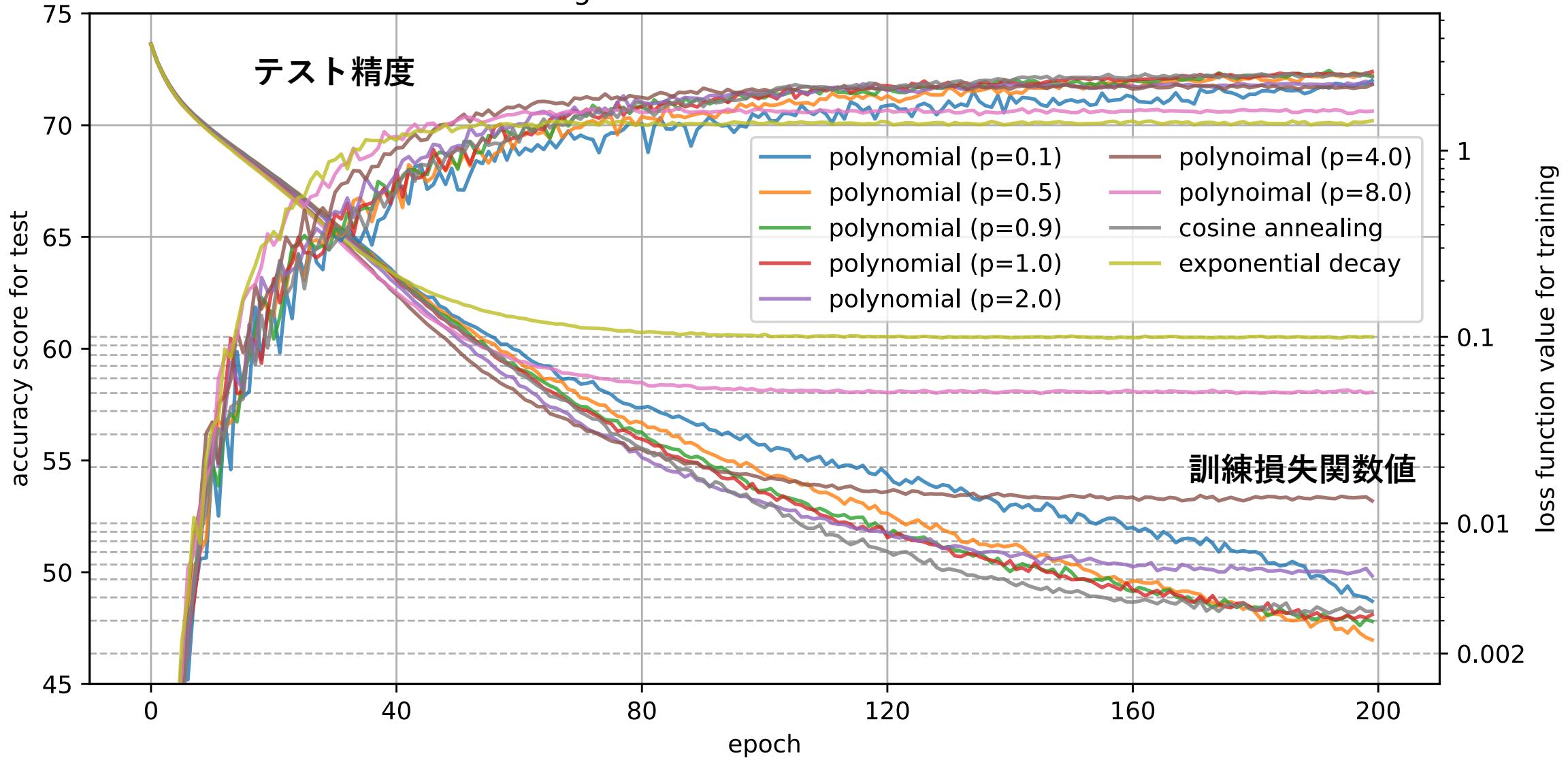
$$\hat{\mathbf{x}}_{t+1} := \hat{\mathbf{x}}_t - \eta \nabla F(\mathbf{x}_t)$$

**end for**

**return**  $\hat{\mathbf{x}}_{T_F+1} = \text{GD}(T_F, \hat{\mathbf{x}}_1, F, \eta)$

# 付録：数値実験—最適な学習率の減衰率

Training ResNet18 on CIFAR100 dataset



## 付録：数値実験—最適な学習率の減衰率

