

目的関数の平滑化とディープニューラルネットワークの汎化性能におけるモーメンタム法の慣性項の役割

日本オペレーションズ・リサーチ学会 2024年秋季研究発表会

9/11(水) 10:10~10:30

明治大学 佐藤尚樹

明治大学 飯塚秀明



背景：関数の平滑化

関数 $f: \mathbb{R}^d \rightarrow \mathbb{R}$ を、大きさ $\delta \in \mathbb{R}$ のノイズで平滑化した関数は、

$$\hat{f}_\delta(x) := \mathbb{E}_{\substack{\underline{u} \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{\sqrt{d}} I_d\right) \\ \text{確率変数} \\ \underline{u} \in \mathbb{R}^d}} [f(x - \underline{\delta u})]$$

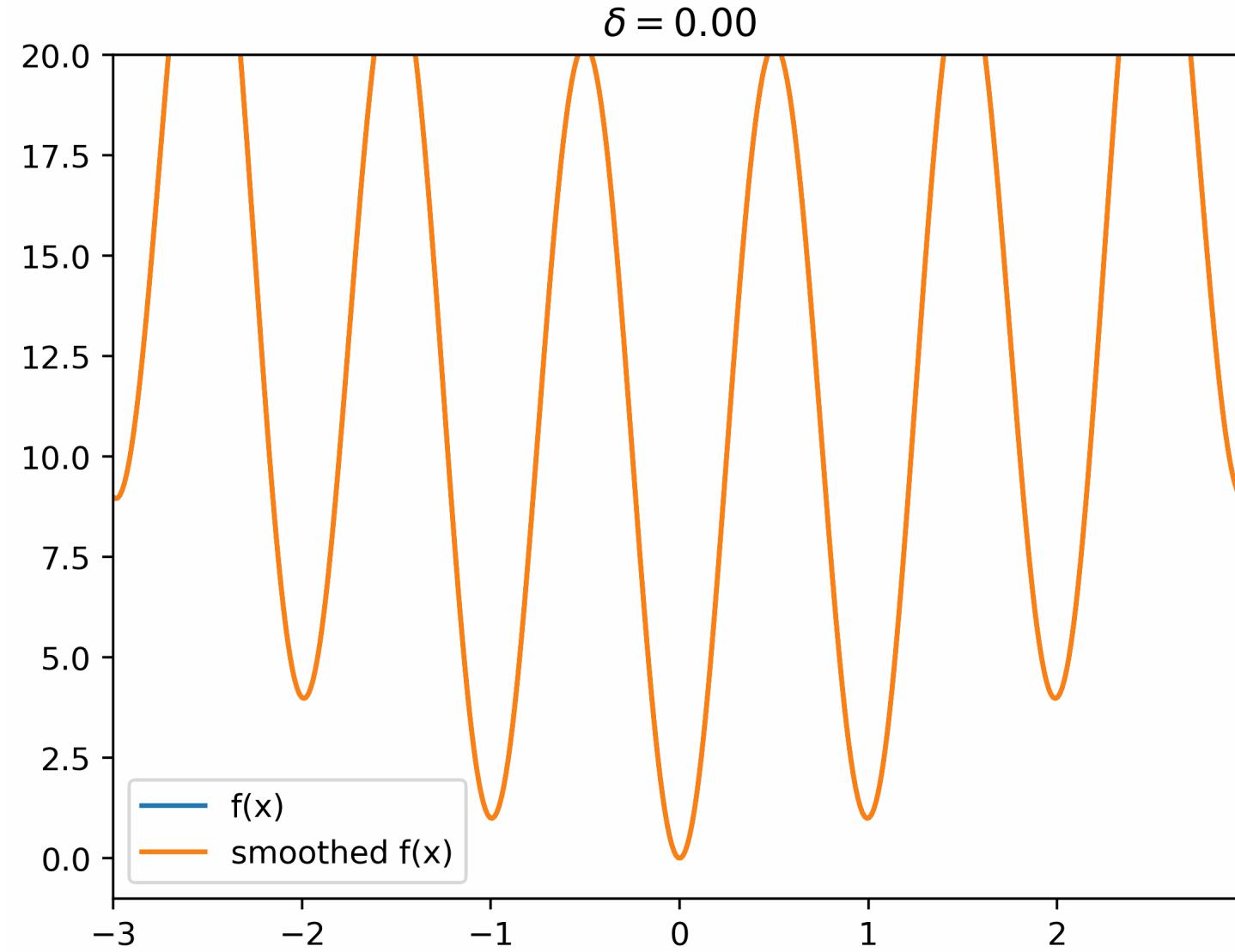
ノイズスケール
(平滑化の度合い)

平均 $\mathbf{0}$,
分散 $\frac{1}{\sqrt{d}} I_d$
の正規分布

背景：関数の平滑化

▷ 目的関数の平滑化

$$\hat{f}_\delta(x) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(x - \delta \mathbf{u})]$$



▷ 1変数のRastrigin's 関数

$$f(x) = x^2 - 10 \cos(2\pi x) + 10$$

背景：機械学習（画像分類）

入力

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



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

$$N = 3000000$$

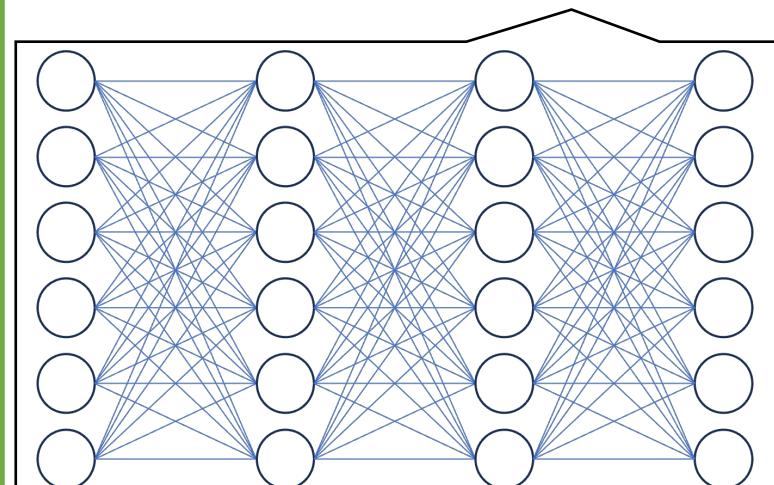
データセット

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)

$$d_t := -\nabla f(x_t)$$

全勾配

$$x_{t+1} := x_t + \eta_t d_t$$

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

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

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：連続最適化

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

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \eta_t \nabla f_{\mathcal{S}_t}(\mathbf{x}_t)$$

▷ 慣性項付き確率的勾配降下法 (SGD with momentum)

$$\mathbf{m}_t := \nabla f_{\mathcal{S}_t}(\mathbf{x}_t) + \beta \mathbf{m}_{t-1}$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \eta_t \mathbf{m}_t$$

慣性係数

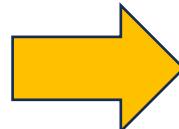
ただし、 $\beta \in [0, 1)$, $\mathbf{m}_{-1} = \mathbf{0}$ とする。

▷ SGDよりも収束が速く、しかも汎化性能も優れる。

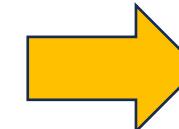
背景：機械学習における汎化能力

▷ 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



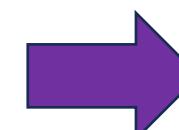
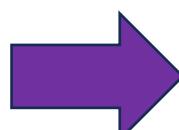
関数 $g: \mathbb{R}^d \rightarrow \mathbb{R}$
変数: $x \in \mathbb{R}^D$,
 $z_i \in \mathbb{R}^d$
Deep Neural Network
訓練済み
= パラメータ調整済



猫

経験損失: 0.001
訓練精度: 99.9%

訓練に使われていないデータ

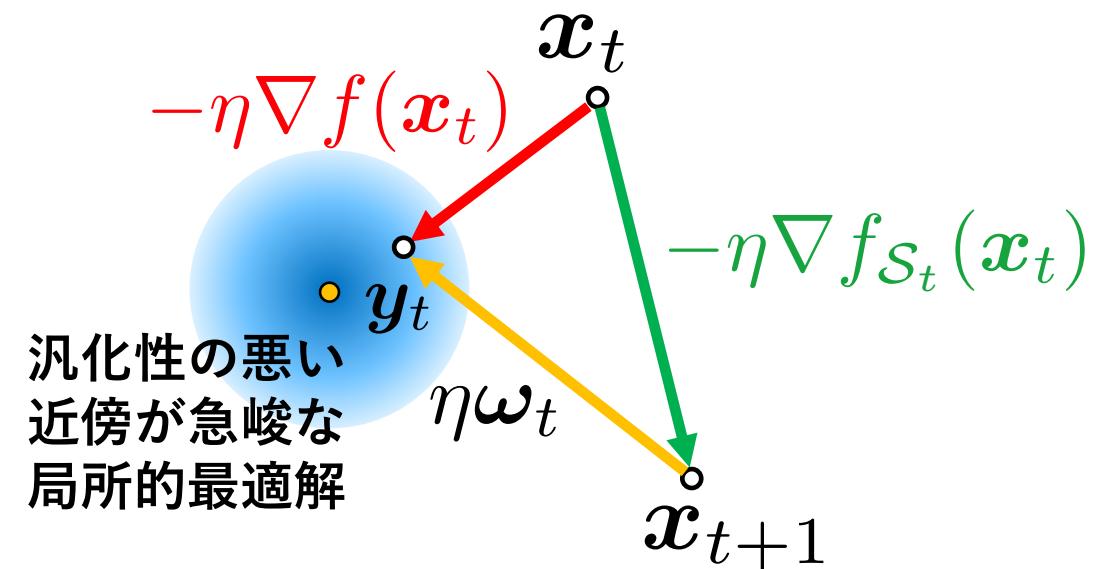


これも猫！

テスト精度: 75%
↑ これが汎化性能！

背景：SGDの確率的ノイズと平滑化

- ▶ 非凸関数の最適化において、確率的勾配降下法(SGD)は、なぜか最急降下法(GD)よりも汎化性の高い解に収束する。
 - ▶ SGDの確率的ノイズが役立っている？[1]



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

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

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

[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

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

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

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(\mathbf{x} - \delta \mathbf{u})]$$

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

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

- ▷ “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

背景 : SGDの平滑化特性

[2] N.Sato and H. Iiduka “Using Stochastic Gradient Descent to Smooth Nonconvex Functions: Analysis of Implicit Graduated Optimization with Optimal Noise Scheduling”, arXiv2311.08745 (2023)

補題A.1

$$\mathbb{E}_{\xi_t} \left[\|\nabla f_{\mathcal{S}_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)\|^2 \right] \leq \frac{C_{\text{opt}}^2}{b}$$

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(\mathbf{x} - \delta \mathbf{u})]$$

▷ $\omega_t^{\text{SGD}} := \nabla f_{\mathcal{S}_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)$ と補題A.1から、

$$\omega_t^{\text{SGD}} = \sqrt{\frac{C_{\text{SGD}}^2}{b}} \mathbf{u}_t \quad \left(\mathbf{u}_t \sim \mathcal{N}\left(\mathbf{0}; \frac{1}{\sqrt{d}} I_d\right) \right)$$

が成り立つ。したがって、

$$\mathbb{E}_{\omega_t^{\text{SGD}}} [\mathbf{y}_{t+1}] = \mathbb{E}_{\omega_t^{\text{SGD}}} [\mathbf{y}_t] - \eta \nabla \mathbb{E}_{\omega_t^{\text{SGD}}} [f(\mathbf{y}_t - \eta \underline{\omega_t^{\text{SGD}}})]$$

$$= \mathbb{E}_{\omega_t^{\text{SGD}}} [\mathbf{y}_t] - \eta \nabla \mathbb{E}_{\mathbf{u}_t \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} \left[f \left(\mathbf{y}_t - \eta \sqrt{\frac{C_{\text{SGD}}^2}{b}} \mathbf{u}_t \right) \right]$$

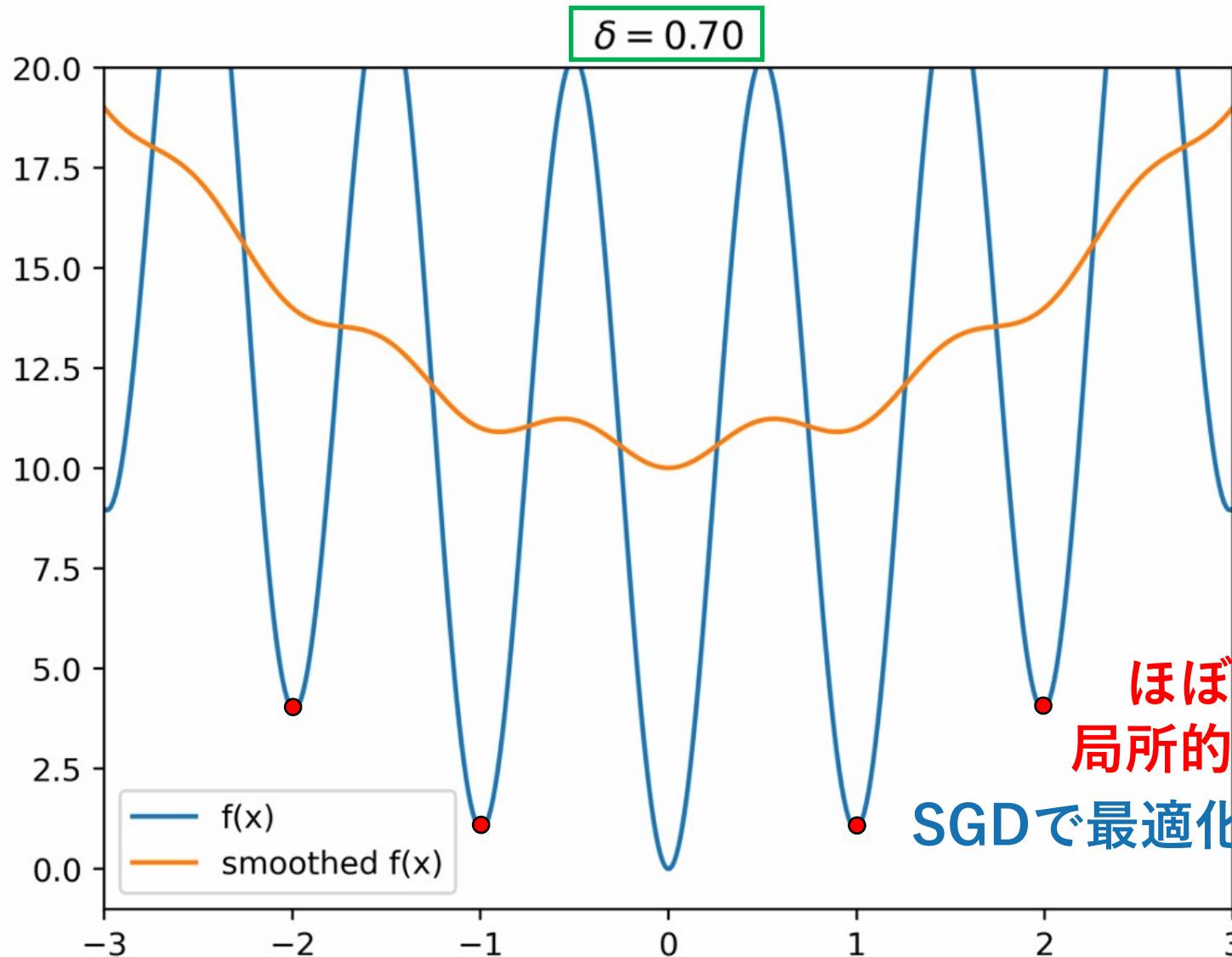
$$= \mathbb{E}_{\omega_t^{\text{SGD}}} [\mathbf{y}_t] - \eta \nabla \hat{f}_{\frac{\eta \sqrt{\frac{C_{\text{SGD}}^2}{b}}}{\eta \sqrt{\frac{C_{\text{SGD}}^2}{b}}}} (\mathbf{y}_t)$$

が成り立つ。

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

$$\delta^{\text{SGD}} = \eta \sqrt{\frac{C_{\text{SGD}}^2}{b}}$$

背景：SGDの確率的ノイズによる暗黙的な目的関数の平滑化



$$\delta^{\text{SGD}} = \eta \sqrt{\frac{C_{\text{SGD}}^2}{b}}$$

確率的ノイズの大きさ
= 平滑化の度合い

背景：慣性項と確率的ノイズの矛盾

▷ 慎性項は確率的ノイズを削減しているはず。

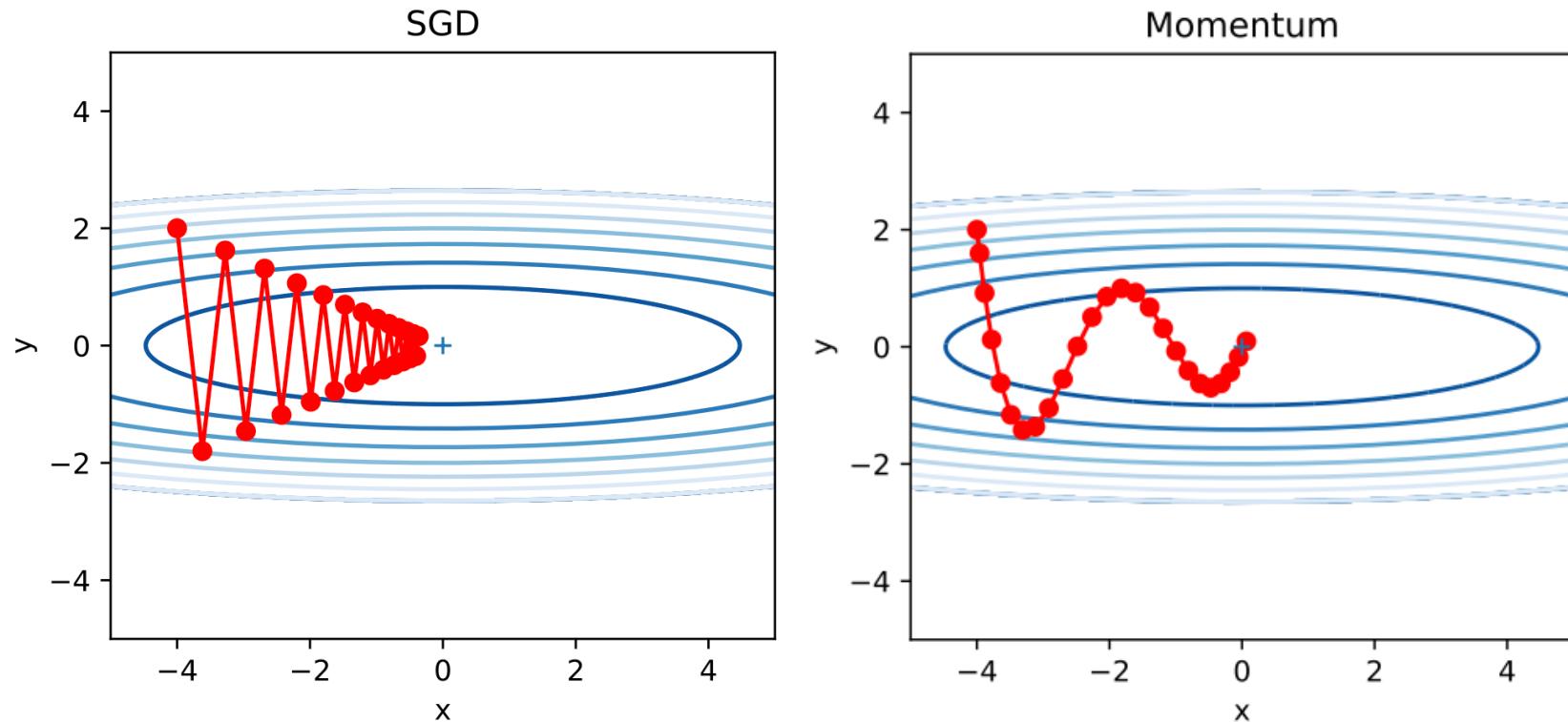
▷ SGD

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \eta_t \nabla f_{\mathcal{S}_t}(\mathbf{x}_t)$$

▷ Momentum

$$\mathbf{m}_t := \nabla f_{\mathcal{S}_t}(\mathbf{x}_t) + \beta \mathbf{m}_{t-1}$$

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \eta_t \mathbf{m}_t$$



▷ 十分な大きさの確率的ノイズが高い汎化性をもたらすはず。

▷ しかし、Momentumの方がSGDよりも高い汎化性をもたらす。

動機

- ▷ 慣性項と確率的ノイズの間にある矛盾を解決したい。

<方針>

- ▷ SGDの確率的ノイズによる目的関数の平滑化の議論を SGD with momentum に拡張する。
- ▷ 慣性係数 β が平滑化の度合いにどのように現れるかを確認する。
- ▷ 数値実験で理論を検証する。

貢献

- ▷ SGD with momentumの確率的ノイズによる平滑化の度合いは、

$$\delta^{\text{SHB}} = \eta \sqrt{\left(1 + \hat{\beta}\right) \frac{C_{\text{SHB}}^2}{b} + \hat{\beta} K_{\text{SHB}}^2}, \quad \delta^{\text{NSHB}} = \eta \sqrt{\frac{1}{1 - \beta} \frac{C_{\text{NSHB}}^2}{b}}$$

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

- ▷ なぜ、いつ、慣性項が汎化性能を向上させるか明らかにした。
- ▷ 慣性項と確率的ノイズの間にある矛盾を解決した。
すなわち、「慣性項によって削減される確率的ノイズ」と、「汎化性能に寄与する確率的ノイズ」が別物であることを示した。

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_f -リップシツツ関数とする。

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

(A2) $(x_t)_{t \in \mathbb{N}} \subset \mathbb{R}^d$ 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。

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

(ii) 次の式を満たす非負定数 C_{opt}^2 が存在するとする。

$$\mathbb{E}_{\xi_t} \left[\|\mathbf{G}_{\xi_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)\|^2 \right] \leq C_{\text{opt}}^2$$

最適化手法に依存する

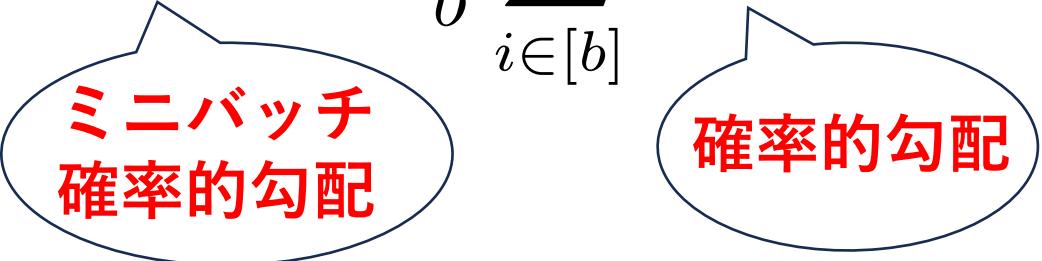
←各最適化手法の
確率的勾配の分散

準備：仮定

仮定続き

(A3) 時刻 $t \in \mathbb{N}$ で、全勾配 ∇f は **ミニバッチ** \mathcal{S}_t で次のように近似されるとする。

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



(A4) 任意の $t \in \mathbb{N}$ で、それぞれの最適化手法に対して、次の式を満たす正の定数 K_{opt} が存在するとする。

$$\mathbb{E} \left[\|\nabla f(x_t)\|^2 \right] \leq K_{\text{opt}}^2$$

準備：アルゴリズム

Algorithm 1 Stochastic Heavy Ball (SHB)

```
Require:  $x_0 \in \mathbb{R}^d$ ,  $\eta > 0$ ,  $\beta \in [0, 1)$ ,  $\mathbf{m}_{-1} := \mathbf{0}$ 
for  $t = 0$  to  $T - 1$  do
     $\mathbf{m}_t := \nabla f_{\mathcal{S}_t}(\mathbf{x}_t) + \beta \mathbf{m}_{t-1}$ 
     $\mathbf{x}_{t+1} := \mathbf{x}_t - \eta \mathbf{m}_t$ 
end for
return  $\mathbf{x}_T$ 
```

[3] I. Gitman, H. Lang, P. Zhang, and L. Xiao,
“Understanding the role of momentum in stochastic
gradient methods,” in *Advances in Neural Information
Processing Systems*, vol. 32, 2019, pp. 9630–9640.

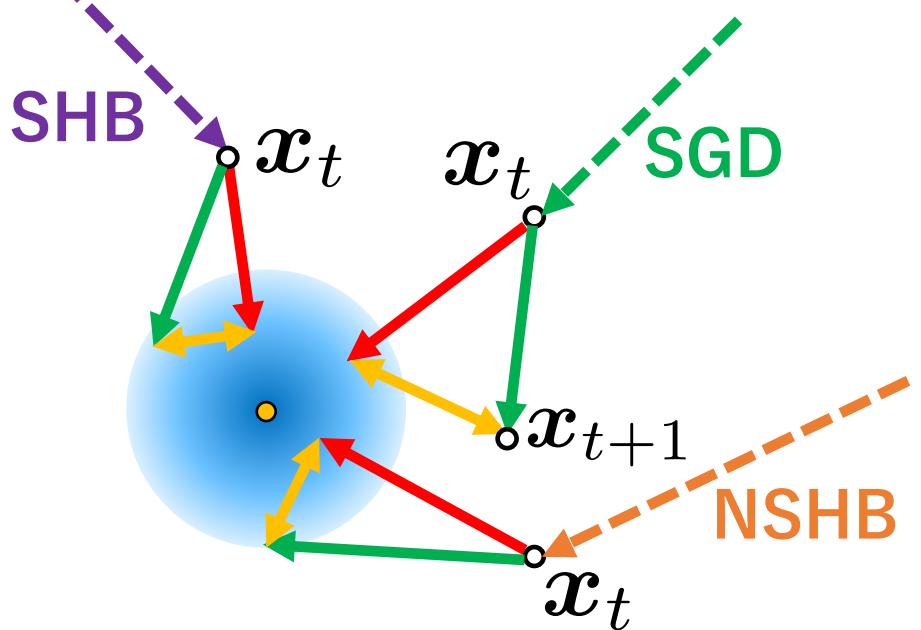
- ▷ 凸結合をとらない
- ▷ あまり理論解析されない
- ▷ 数値実験でよく利用される
 - ▷ PyTorchやTensorFlowが
提供するSGD with momentum

Algorithm 2 Normalized Stochastic Heavy Ball (NSHB)

```
Require:  $x_0 \in \mathbb{R}^d$ ,  $\eta > 0$ ,  $\beta \in [0, 1)$ ,  $\mathbf{d}_{-1} := \mathbf{0}$ 
for  $t = 0$  to  $T - 1$  do
     $\mathbf{d}_t := (1 - \beta) \nabla f_{\mathcal{S}_t}(\mathbf{x}_t) + \beta \mathbf{d}_{t-1}$ 
     $\mathbf{x}_{t+1} := \mathbf{x}_t - \eta \mathbf{d}_t$ 
end for
return  $\mathbf{x}_T$ 
```

- ▷ 凸結合をとる
- ▷ よく理論解析される
- ▷ 数値実験で利用されない

勾配ノイズ（従来の確率的ノイズ）



SGDの
勾配ノイズ $\frac{\nabla f_{S_t}(x_t)}{\text{green}} - \frac{\nabla f(x_t)}{\text{red}}$

SHBの
勾配ノイズ $\frac{\nabla f_{S_t}(x_t)}{\text{green}} - \frac{\nabla f(x_t)}{\text{red}}$

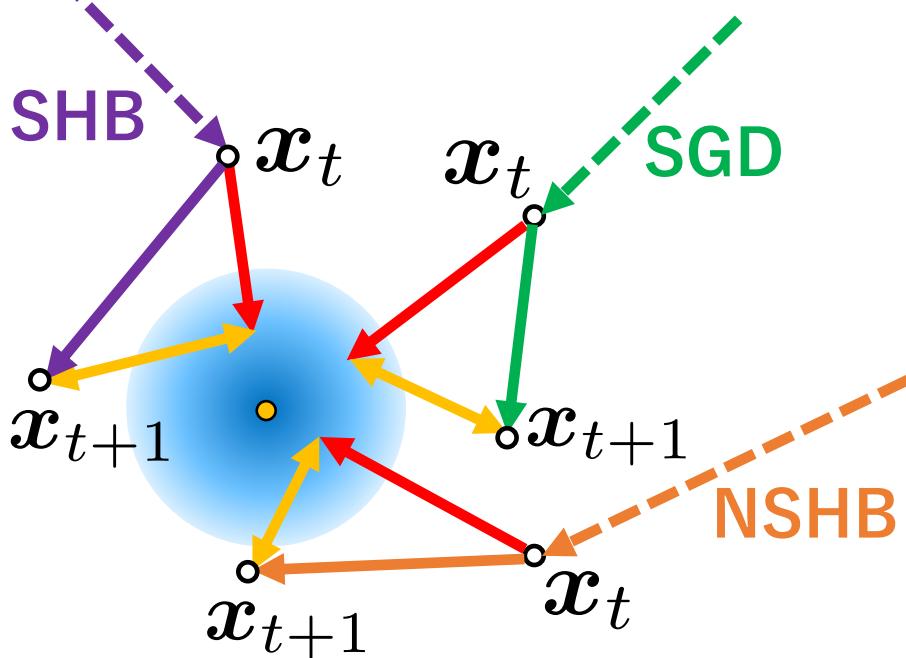
NSHBの
勾配ノイズ $\frac{\nabla f_{S_t}(x_t)}{\text{green}} - \frac{\nabla f(x_t)}{\text{red}}$

補題A.1. 仮定(A2)(ii)と(A3)の下で、任意の時刻 $t \in \mathbb{N}$ に対して、次が成り立つ。

$$\mathbb{E}_{\xi_t} \left[\left\| \nabla f_{S_t}(x_t) - \nabla f(x_t) \right\|^2 \right] \leq \frac{C_{\text{opt}}^2}{b}$$

∴ 勾配ノイズ = 確率的勾配の分散

探索方向ノイズ



SGDの
探索方向ノイズ $\omega_t^{\text{SGD}} := \frac{\nabla f_{S_t}(x_t)}{\text{SGDの探索方向}} - \frac{\nabla f(x_t)}{\text{GDの探索方向}}$

SHBの
探索方向ノイズ $\omega_t^{\text{SHB}} := \frac{m_t}{\substack{\text{SHB} \\ \text{の探索方向}}} - \frac{\nabla f(x_t)}{\text{GDの探索方向}}$

NSHBの
探索方向ノイズ $\omega_t^{\text{NSHB}} := \frac{d_t}{\substack{\text{NSHB} \\ \text{の探索方向}}} - \frac{\nabla f(x_t)}{\text{GDの探索方向}}$

定理3.1. 仮定(A2)(ii),(A3),(A4)の下で、任意の時刻 $t \in \mathbb{N}$ に対して、次が成り立つ。

$$\mathbb{E} [\|\omega_t^{\text{SHB}}\|] \leq \sqrt{\frac{C_{\text{SHB}}^2}{b} + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2} \left(\frac{C_{\text{SHB}}^2}{b} + K_{\text{SHB}}^2 \right)}, \quad \mathbb{E} [\|\omega_t^{\text{NSHB}}\|] \leq \sqrt{\frac{1}{1 - \beta} \frac{C_{\text{NSHB}}^2}{b}}$$

SHBの探索方向ノイズと平滑化

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

$$(\text{SHB}) \quad \mathbf{x}_{t+1} := \mathbf{x}_t - \eta \mathbf{m}_t$$

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(\mathbf{x} - \delta \mathbf{u})]$$

探索方向ノイズ $\omega_t^{\text{SHB}} := \mathbf{m}_t - \nabla f(\mathbf{x}_t)$ とすると、次の式が成り立つ。

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

- ▷ “関数 $f(\mathbf{x}_t)$ をSHB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。

SHBの平滑化特性

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(\mathbf{x} - \delta \mathbf{u})]$$

定理3.1

$$\mathbb{E} [\|\boldsymbol{\omega}_t^{\text{SHB}}\|] \leq \sqrt{\left(1 + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2}\right) \frac{C_{\text{SHB}}^2}{b} + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2} K_{\text{SHB}}^2} =: A$$

▷ 定理3.1から、

$$\boldsymbol{\omega}_t^{\text{SHB}} = \sqrt{\left(1 + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2}\right) \frac{C_{\text{SHB}}^2}{b} + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2} K_{\text{SHB}}^2} \mathbf{u}_t = A \mathbf{u}_t \quad \left(\mathbf{u}_t \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{\sqrt{d}} I_d\right) \right)$$

が成り立つ。したがって、

$$\begin{aligned} \mathbb{E}_{\boldsymbol{\omega}_t^{\text{SHB}}} [\mathbf{y}_{t+1}] &= \mathbb{E}_{\boldsymbol{\omega}_t^{\text{SHB}}} [\mathbf{y}_t] - \eta \nabla \mathbb{E}_{\boldsymbol{\omega}_t^{\text{SHB}}} [f(\mathbf{y}_t - \eta \underline{\boldsymbol{\omega}_t^{\text{SHB}}})] \\ &= \mathbb{E}_{\boldsymbol{\omega}_t^{\text{SHB}}} [\mathbf{y}_t] - \eta \nabla \mathbb{E}_{\mathbf{u}_t \sim \mathcal{N}\left(\mathbf{0}, \frac{1}{\sqrt{d}} I_d\right)} [f(\mathbf{y}_t - \eta \underline{A \mathbf{u}_t})] \\ &= \mathbb{E}_{\boldsymbol{\omega}_t^{\text{SHB}}} [\mathbf{y}_t] - \eta \nabla \underline{\hat{f}_{\eta A}(\mathbf{y}_t)} \quad \text{大きさ } \eta A \text{ のノイズで} \\ &\quad \text{平滑化された関数} \end{aligned}$$

が成り立つ。

最適化手法ごとの平滑化の度合い

$$\delta^{\text{SGD}} = \eta \sqrt{\frac{C_{\text{SGD}}^2}{b}},$$

∴ 探索方向ノイズ = 平滑化の度合い

$$\delta^{\text{SHB}} = \eta \sqrt{\left(1 + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2}\right) \frac{C_{\text{SHB}}^2}{b} + \frac{\beta(\beta^2 - \beta + 1)}{(1 - \beta)^2} K_{\text{SHB}}^2},$$

≈ 3.13

$$\delta^{\text{NSHB}} = \eta \sqrt{\frac{1}{1 - \beta} \frac{C_{\text{NSHB}}^2}{b}}$$

- ▷ どれくらいの大きさか？
- ▷ 未知数の推定が必要

仮定(A2)(ii) [確率的勾配の分散]

$$\mathbb{E}_{\xi_t} \left[\|\mathbf{G}_{\xi_t}(\mathbf{x}_t) - \nabla f(\mathbf{x}_t)\|^2 \right] \leq C_{\text{opt}}^2$$

仮定(A4) [勾配ノルムの上界]

$$\mathbb{E} \left[\|\nabla f(\mathbf{x}_t)\|^2 \right] \leq K_{\text{opt}}^2$$

確率的勾配の分散 C_{opt}^2 の推定

命題4.1. 仮定(A1)-(A4)の下で、任意の $\epsilon > 0$ に対して、次が成り立つ。

$$b_{\text{SGD}}^* > \frac{\eta C_{\text{SGD}}^2}{\epsilon^2}, \quad b_{\text{SHB}}^* > \frac{\eta(\beta^2 - \beta + 1)C_{\text{SHB}}^2}{\beta(1 - \beta)^2\epsilon^2}, \quad b_{\text{NSHB}}^* > \frac{\eta C_{\text{NSHB}}^2}{(1 - \beta)\epsilon^2}$$

▷ b_{opt}^* は、 $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} [\|\nabla f(x_t)\|^2] \leq \epsilon^2$ を達成するために必要な

計算量を最小にするバッチサイズ(数値実験で測定可能)

▷ ϵ は訓練の停止条件 (例えば、 $\epsilon = 0.5$)

停止条件

$$\triangleright C_{\text{SGD}}^2 < \frac{b_{\text{SGD}}^* \epsilon^2}{\eta}, \quad C_{\text{SHB}}^2 < \frac{b_{\text{SHB}}^* \beta(1 - \beta)^2 \epsilon^2}{\eta(\beta^2 - \beta + 1)}, \quad C_{\text{NSHB}}^2 < \frac{b_{\text{NSHB}}^* (1 - \beta) \epsilon^2}{\eta}$$

学習率

慣性係数

この式で推定可能！

確率的勾配の分散 C_{opt}^2 の推定

▷ CIFAR100 データセットで ResNet18 を訓練するとき、

停止条件を $\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E} \left[\|\nabla f(x_t)\|^2 \right] \leq \epsilon^2 = (0.5)^2$ とすると、

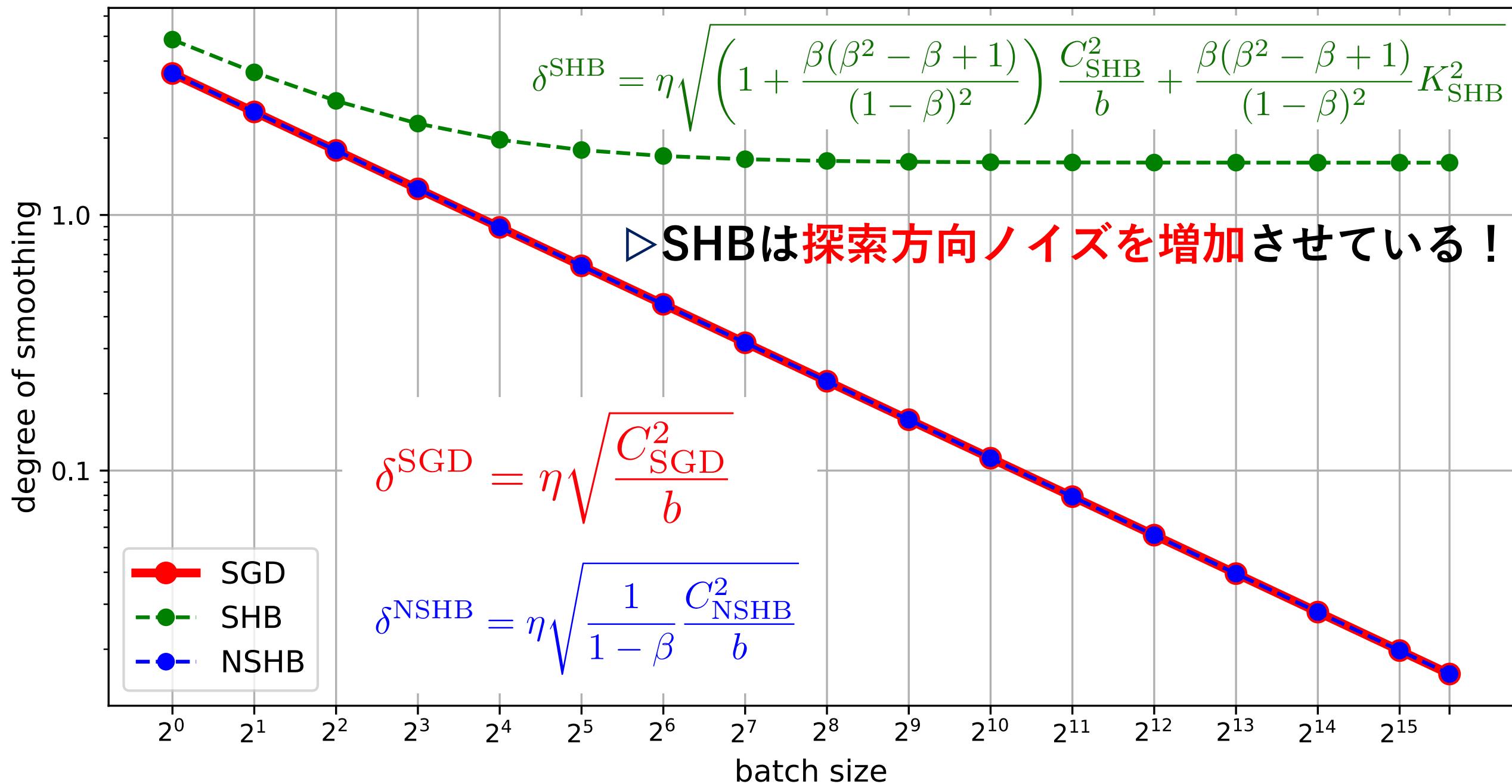
$$b_{\text{SGD}}^* = 2^9, \quad b_{\text{SHB}}^* = 2^{10}, \quad b_{\text{NSHB}}^* = 2^9$$

となった。なお、訓練に使用された学習率は $\eta = 0.1$ で、
慣性係数は $\beta = 0.9$ だった。これらの値を代入すると、

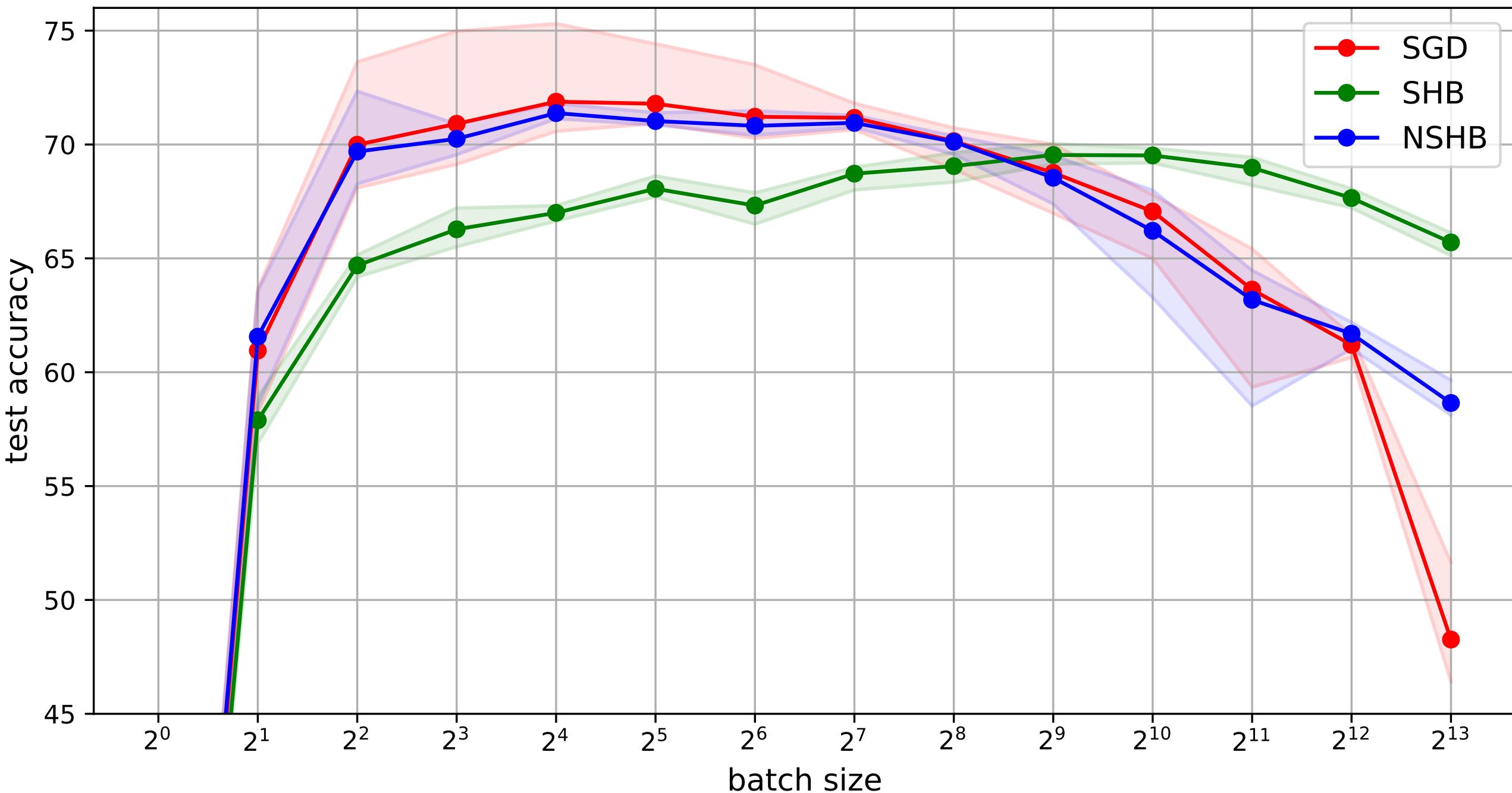
$$C_{\text{SGD}}^2 < \frac{b_{\text{SGD}}^* \epsilon^2}{\eta}, \quad C_{\text{SHB}}^2 < \frac{b_{\text{SHB}}^* \beta (1 - \beta)^2 \epsilon^2}{\eta (\beta^2 - \beta + 1)}, \quad C_{\text{NSHB}}^2 < \frac{b_{\text{NSHB}}^* (1 - \beta) \epsilon^2}{\eta}$$
$$= 1280 \quad = 25.3 \quad = 128$$

▷ 確かに SGD with momentum は 勾配ノイズを削減 している！

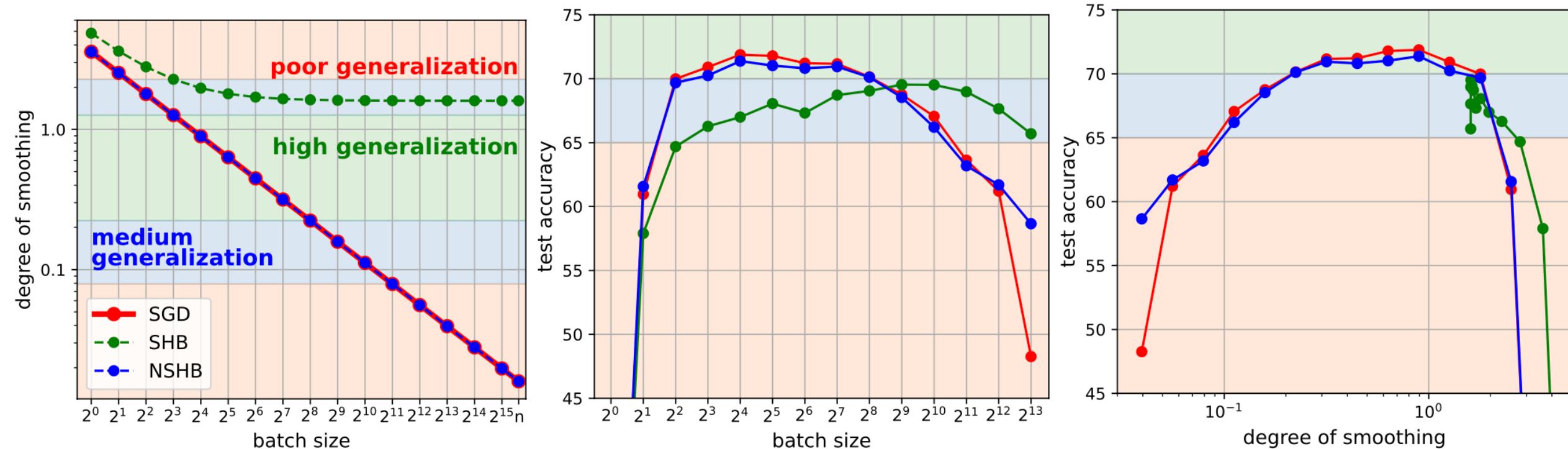
Training ResNet18 on CIFAR100 dataset



Training ResNet18 on CIFAR100 dataset



平滑化の度合いと汎化性能の相関



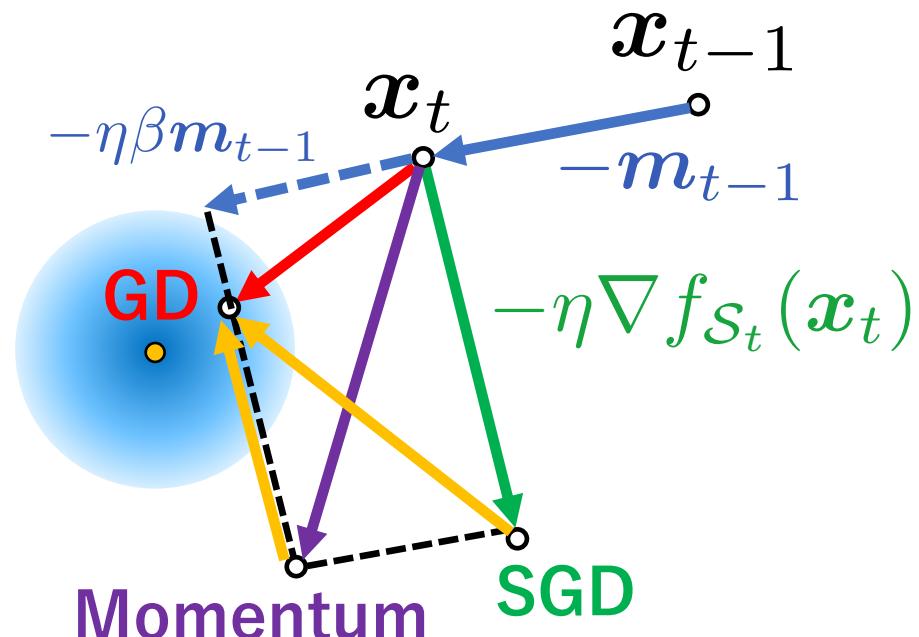
- ▷ バッチサイズが大きくなっても、SHBだけは平滑化の度合いが十分大きいので、汎化性能が悪化しない。
- ▷ NSHBの平滑化の度合いはSGDと同等なので、SGDとNSHBの汎化性能はほぼ一致している。

結論

- ▷ SGD with momentumには目的関数を平滑化する効果があり、その度合いは学習率、バッチサイズ、慣性係数などで決まる。
- ▷ 次の2つの主張が競合しないことを示した。
 - ▷ 慣性項は確率的ノイズを削減しているはず。
 - ▷ 十分な大きさの確率的ノイズが高い汎化性をもたらすはず。
- ▷ 慣性項は勾配ノイズを削減するが、逆に探索方向ノイズを増加させる。そして探索方向ノイズが汎化性能に寄与する。
- ▷ 適度な大きさの平滑化の度合いが高い汎化性能をもたらす。
- ▷ SHBはバッチサイズが大きいときに有効で、逆にバッチサイズが小さいときは慣性項がある意義がなくなる。
- ▷ SHBとNSHBは全くの別物で、SHBの方が優れている。

没・背景：慣性項と確率的ノイズの矛盾

- ▷ 慣性項は確率的ノイズを削減しているはず。
- ▷ 十分な大きさの確率的ノイズが高い汎化性をもたらすはず。
- ▷ しかし、Momentumの方がSGDよりも高い汎化性をもたらす。



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

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

$$(\text{Momentum}) \quad m_t := \nabla f_{\mathcal{S}_t}(x_t) + \beta m_{t-1}$$

$$\hat{y}_t := x_t - \eta m_t$$

SGDの
確率的ノイズ

$$\omega_t^{\text{SGD}} := \frac{\nabla f_{\mathcal{S}_t}(x_t) - \nabla f(x_t)}{\text{SGDの探索方向}}$$

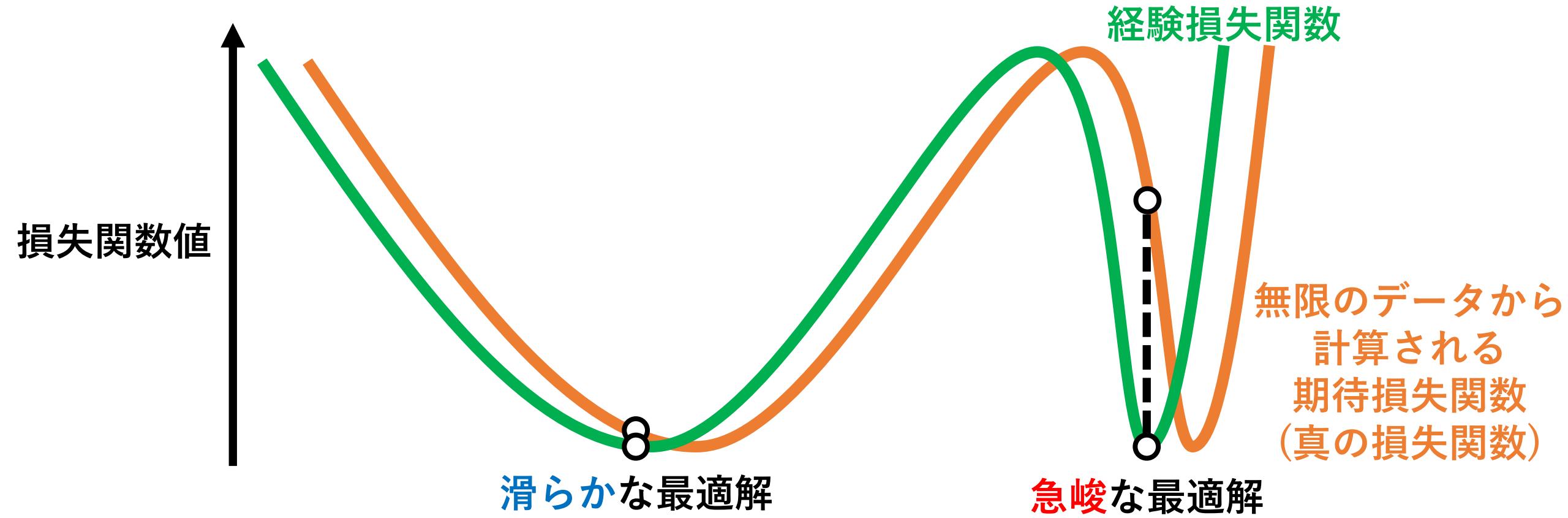
GDの探索方向

$$\omega_t^{\text{Momentum}} := \frac{m_t - \nabla f(x_t)}{\text{Momentumの探索方向}}$$

GDの探索方向

付録：汎化性能と経験損失関数の形状

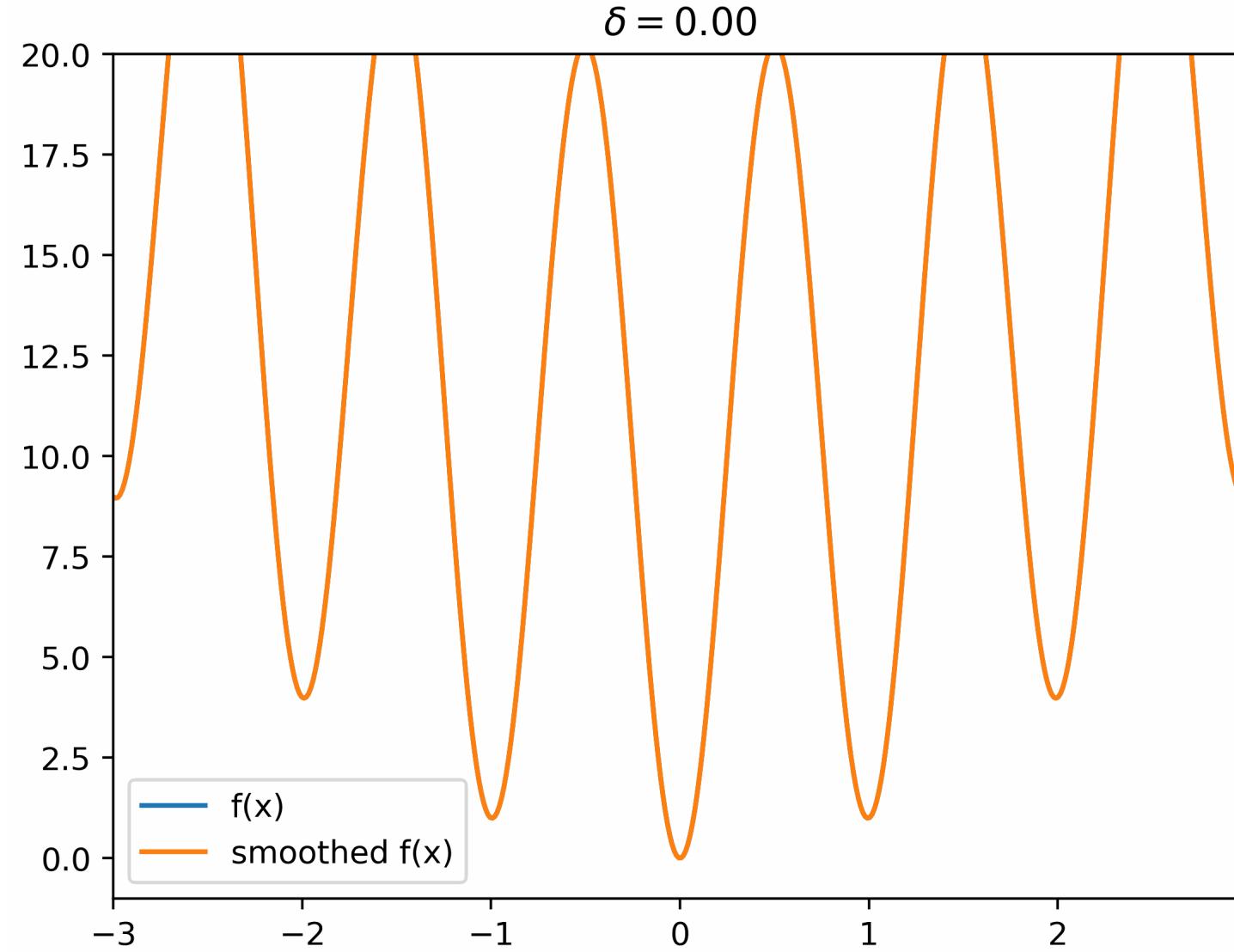
- ▷ 「その近傍が滑らかな最適解」は、「その近傍が急峻な最適解」よりも優れた汎化性能をもたらす。



付録：関数の平滑化

▷ 目的関数の平滑化

$$\hat{f}_\delta(\mathbf{x}) := \mathbb{E}_{\mathbf{u} \sim \mathcal{N}(\mathbf{0}, \frac{1}{\sqrt{d}} I_d)} [f(\mathbf{x} - \delta \mathbf{u})]$$



▷ 1変数のRastrigin's 関数

$$f(x) = x^2 - 10 \cos(2\pi x) + 10$$

$$|\hat{f}_\delta(\mathbf{x}) - f(\mathbf{x})| \leq |\delta| L_f$$