

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アルゴリズム による経験損失最小化問題のための大域的最適化

佐藤 尚樹, 飯塚 秀明 (明治大学)



要約：暗黙的な段階的最適化で経験損失関数(非凸関数)の大域的最適化を実現できる。

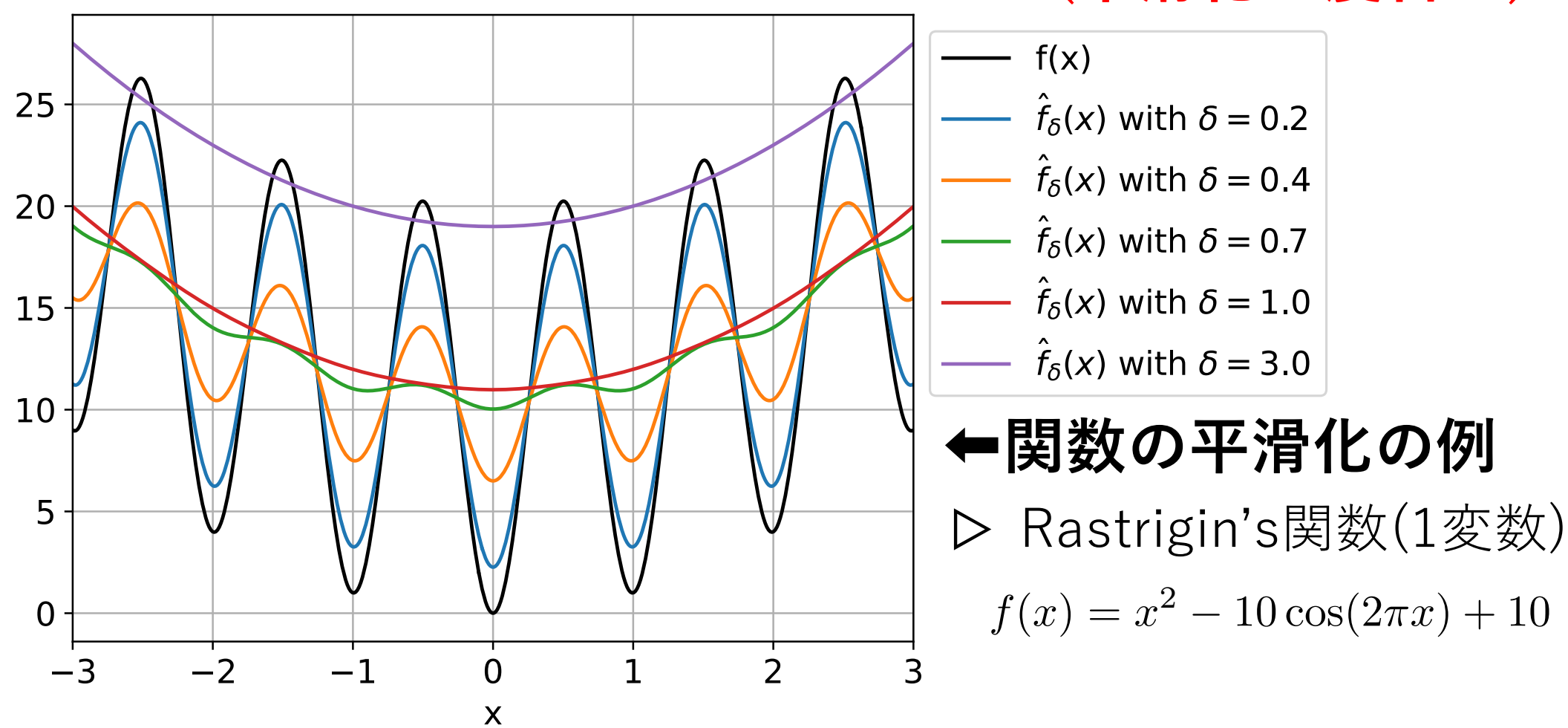


I. 段階的最適化 (Graduated Optimization)

- 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大域的最適化手法。
- 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, 次のように定義される。

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

ノイズスケール
(平滑化の度合い)



関数の平滑化の例

- Rastrigin's関数(1変数)
- $f(x) = x^2 - 10 \cos(2\pi x) + 10$

σ -nice関数 [Hazan et al., 2016]

- (i) 任意の $\delta_m > 0$ と $x_{\delta_m}^*$ に対して, $\|x_{\delta_m}^* - x_{\delta_{m+1}}^*\| \leq \delta_{m+1} := \frac{\delta_m}{2}$.
- (ii) 任意の $\delta_m > 0$ に対して, 関数 \hat{f}_{δ_m} は近傍 $N(x_{\delta_m}^*; 3\delta_m)$ で σ -強凸である。

II. 確率的勾配降下法の平滑化特性

- 時刻 t のパラメータを x_t とする。最急降下法で更新した先を y_t , SGDで更新した先を x_{t+1} とすると,

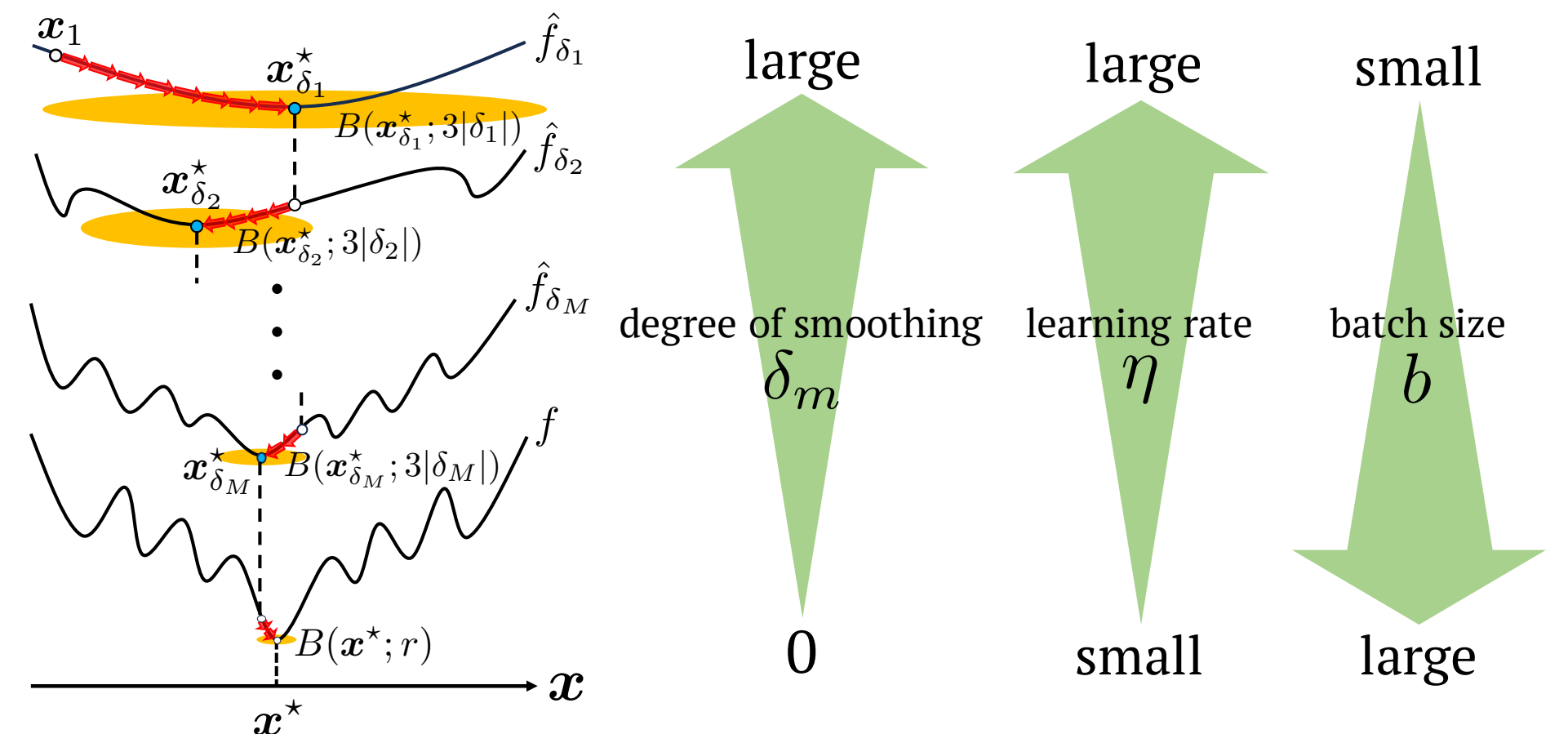
$$\begin{aligned} (GD) \quad y_t &:= x_t - \eta \nabla f(x_t) \\ (SGD) \quad x_{t+1} &:= x_t - \eta \nabla f_{S_t}(x_t) \end{aligned}$$

確率的ノイズ $\omega_t := \nabla f_{S_t}(x_t) - \nabla f(x_t)$

SGDの探索方向 (最急降下方向)

$$\begin{aligned} \mathbb{E}_{\omega_t}[y_{t+1}] &= \mathbb{E}_{\omega_t}[y_t] - \eta \nabla \mathbb{E}_{\omega_t}[f(y_t - \eta \omega_t)] \\ &= \mathbb{E}_{\omega_t}[y_t] - \eta \nabla \mathbb{E}_{u \sim \mathcal{N}(0, \frac{1}{\sqrt{\delta}} I_d)} \left[f\left(y_t - \frac{\eta C}{\sqrt{\delta}} u\right) \right] \\ &= \mathbb{E}_{\omega_t}[y_t] - \eta \nabla \hat{f}_{\frac{\eta C}{\sqrt{\delta}}}(y_t) \end{aligned}$$

が成り立つから, 関数 f をSGDで最適化することと, 関数 $\hat{f}_{\frac{\eta C}{\sqrt{\delta}}}$ をGDで最適化することは期待値の意味で等価。



III. 暗黙的な段階的最適化 (Implicit Graduated Optimization)

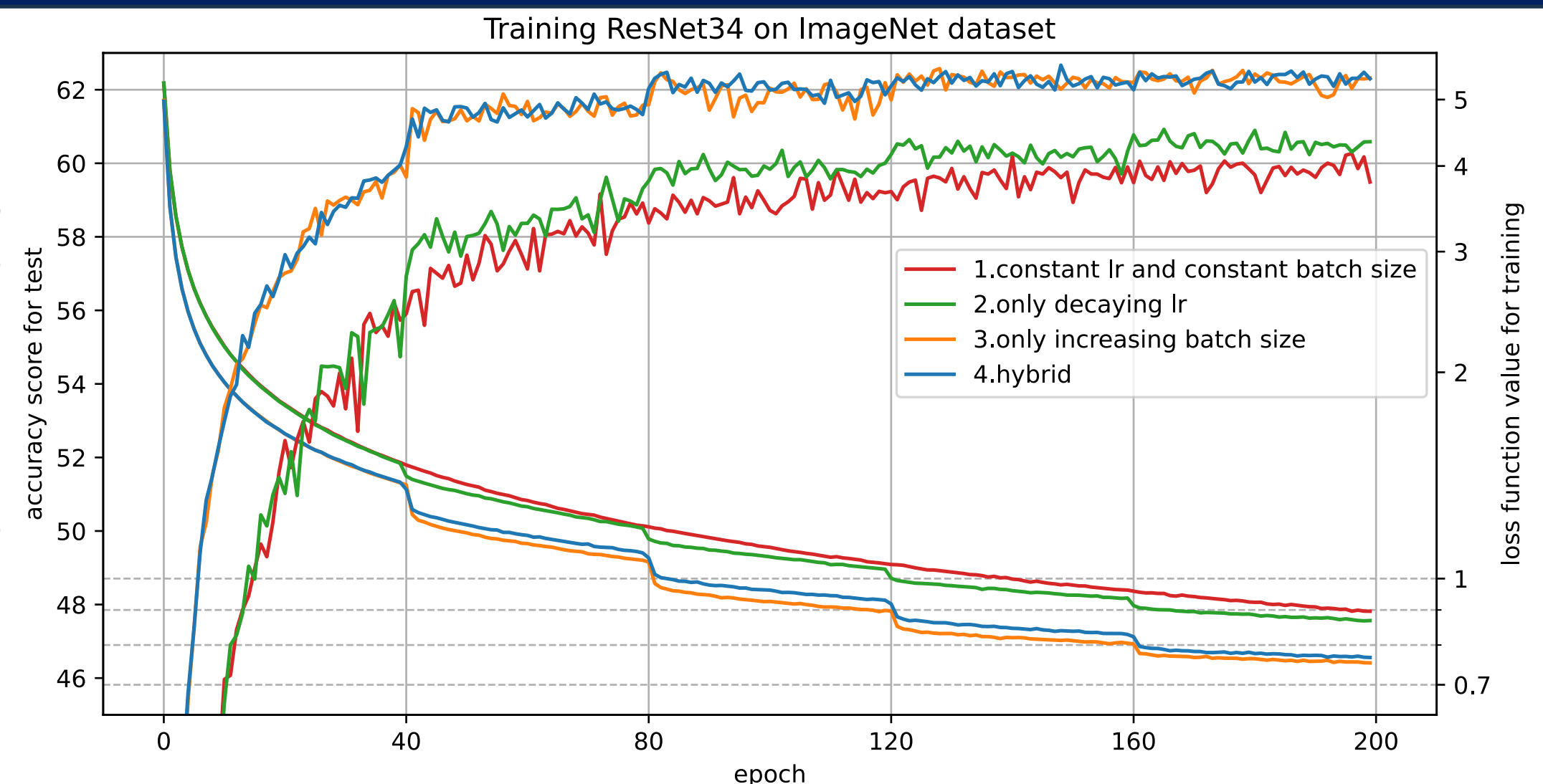
Algorithm 1 Implicit Graduated Optimization

Require: $\epsilon, x_1 \in B(x_{\delta_1}^*; 3\delta_1), \eta_1 > 0, b_1 \in [n], \gamma \geq 0.5$
 $\delta_1 := \frac{\eta_1 C}{\sqrt{b_1}}, \alpha_0 := \min\left\{\frac{1}{16L_f\delta_1}, \frac{1}{\sqrt{2\sigma\delta_1}}\right\}, M := \log_\gamma \alpha_0 \epsilon$
for $m = 1$ to $M + 1$ **do**
 if $m \neq M + 1$ **then**
 $\epsilon_m := \sigma\delta_m^2/2, T_m := H_m/\epsilon_m$
 $\kappa_m/\sqrt{\lambda_m} = \gamma (\kappa_m \in (0, 1], \lambda_m \geq 1)$
 end if
 $x_{m+1} := \text{GD}(T_m, x_m, \hat{f}_{\delta_m}, \eta_m)$
 $\eta_{m+1} := \kappa_m \eta_m, b_{m+1} := \lambda_m b_m$
 $\delta_{m+1} := \frac{\eta_{m+1} C}{\sqrt{b_{m+1}}}$
end for
return x_{M+2}

M個の平滑化された関数の最適化に使われる。

Algorithm 2 Gradient Descent

Require: $T_m, \hat{x}_1^{(m)}, \hat{f}_{\delta_m}, \eta > 0$
for $t = 1$ to T_m **do**
 $\hat{x}_{t+1}^{(m)} := \hat{x}_t^{(m)} - \eta \nabla \hat{f}_{\delta_m}(x_t)$
end for
return $\hat{x}_{T_m+1}^{(m)}$



- SGDよりも高いテスト精度と低い損失を達成。

IV. 平滑化の度合い, Sharpness, 汎化性能の関係

- CIFAR100データセットによるResNet18の訓練でそれぞれを計測した。
- 学習率 $\eta \in \{0.01, 0.05, 0.1, 0.5\}$, バッチサイズ $b \in \{2, 4, \dots, 8192\}$, 200epochs

