2024/8/24-26

連続最適化および関連分野に関する夏季学校 2024 非凸最適化アルゴリズムと その計算量解析

丸茂 直貴

東京大学 大学院情報理工学系研究科 数理情報学専攻



扱う非凸最適化問題:無制約・平滑最適化 $\min_{x \in \mathbb{R}^d} f(x)$ $f: \mathbb{R}^d \to \mathbb{R}$ は十分滑らか、一般に非凸、 $\inf_{x \in \mathbb{R}^d} f(x) > -\infty$

- 目標: アルゴリズムとその計算量保証の理解
 - アルゴリズムの設計・解析における定石の習得
- ❶ 最急降下法
- **2** 3次正則化 Newton 法
- ❸ 再始動 Heavy-ball 法

例: Newton 法は (fの性質が良くても) 収束しないことがある C 演習2



例: 凸関数に対する加速勾配法 [Nesterov, 1983]

$$egin{aligned} x_{k+1} &= y_k - rac{1}{L}
abla f(y_k), \ y_{k+1} &= x_{k+1} + rac{k}{k+3} (x_{k+1} - x_k) \end{aligned}$$

- 謎の値 ^k/_{k+3} は 計算量解析 から導ける
- 加速勾配法は凸関数に対し最良の計算量を達成 😂
- 実用上も強い ☺



- どのような問題クラスに対して
- びのような点を求めるために
- C どのような計算を 最悪でも何回すればよいか? ※ 平均ケースなどを考えることもある

例:

- る 目的関数 f が線形,2次,凸,強凸,複数の関数の和,合成関数, ∇f **が Lipschitz 連続**, $\nabla^2 f$ が Lipschitz 連続, …
- **b** $||x x^*|| \leq \varepsilon$, $f(x) f(x^*) \leq \varepsilon$, $||\nabla f(x)|| \leq \varepsilon$, ...

c $f(x), \nabla f(x), \nabla^2 f(x), \nabla^2 f(x)v, \ldots$ *ε*-停留点



・3次正則化Newton法

・再始動 Heavy-ball 法

非凸最適化アルゴリズムの計算量解析の定石

暫定解の列 $(x_k)_{k\in\mathbb{N}}$

計算量解析の定石

- $||x_{k+1} x_k||$ が大 \implies **関数値減少量**が大 \bigcirc
- $\|x_{k+1}-x_k\|$ が小 \implies 勾配ノルム が小 \bigcirc

「どちらに転んでも嬉しい」という状況を作る

6/77

、 を両方示し, 組み合わせる

非凸最適化アルゴリズムの計算量解析の定石

暫定解の列 $(x_k)_{k\in\mathbb{N}}$

計算量解析の定石

- $||x_{k+1} x_k||$ が大 \implies **関数値減少量**が大 \bigcirc
- $||x_{k+1} x_k||$ が小 \implies 勾配ノルム が小 \bigcirc

ここからの話:最急降下法を例に

- 上をどう示す?
- 上をどう組み合わせる?
- 上が成り立つアルゴリズムをどう設計する? を見る

を両方示し, 組み合わせる

最も単純な例:最急降下法

7/77

仮定: $\nabla f o L$ -Lipschitz 連続性

$$\|\nabla f(x) - \nabla f(y)\| \le \mathbf{L} \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$$

更新式

$$x_{k+1} = x_k - \eta_k \nabla f(x_k) \qquad (\eta_k > 0)$$

計算量保証

$$\Delta \coloneqq f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)$$
 とおく、 $\eta_k = \frac{1}{L}$ とすると、 $\min_{0 \le i < k} \| \nabla f(x_i) \| \le \sqrt{\frac{2L\Delta}{k}}$

 \mathbf{C} ε -停留点を求めるための計算量: $\left\lceil \frac{2L\Delta}{\varepsilon^2} \right\rceil = \mathbf{O}(\varepsilon^{-2})$

仮定: $\nabla f o L$ -Lipschitz 連続性

$$\|\nabla f(x) - \nabla f(y)\| \le L \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$$

同値な条件 critic critic

一次近似しにくい点があると*L* が大きくなる

Taylor 展開:
$$f(x+u) \simeq f(x) + \langle \nabla f(x), u \rangle + \frac{1}{2} \langle \nabla^2 f(x) u, u \rangle$$

仮定: $\nabla f o L$ -Lipschitz 連続性

$$\|\nabla f(x) - \nabla f(y)\| \le L \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$$

同値な条件 ① 演習1(上の仮定) \iff $\|\nabla^2 f(x)\| \le L \quad (\forall x \in \mathbb{R}^d)$

仮定を満たさない例: $f(x) = |x|^{3/2}$



仮定: ▽f の *L*-Lipschitz 連続性

$$\|\nabla f(x) - \nabla f(y)\| \le L \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$$



Quiz: どちらの*L*が大きい?

仮定: ▽fのL-Lipschitz連続性

$$\|\nabla f(x) - \nabla f(y)\| \le L \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$$



Lipschitz 連続性から従う不等式

$$\nabla f \, \check{x} \, L$$
-Lipschitz 連続ならば,
$$\left| f(x) - \left(f(y) + \langle \nabla f(y), x - y \rangle \right) \right| \leq \frac{L}{2} \|x - y\|^2 \qquad (\forall x, y \in \mathbb{R}^d)$$
$$f(x) \, \mathcal{O}$$
一次近似



$$abla f$$
 が *L*-Lipschitz 連続ならば,
$$\left|f(x) - \left(f(y) + \langle \nabla f(y), x - y \rangle\right)\right| \leq \frac{L}{2} \|x - y\|^2 \qquad (\forall x, y \in \mathbb{R}^d)$$

10/77

Lipschitz 連続性から従う不等式: 上界補題

$$abla f$$
が *L*-Lipschitz 連続ならば,
 $f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} ||x - y||^2 \quad (\forall x, y \in \mathbb{R}^d)$

11/77



最急降下法の解析: 関数値減少量の評価

• 上界補題:

$$f(x_{k+1}) - f(x_k) \le \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

• 更新式:
$$x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$$

関数値減少量の評価

$$f(x_{k+1}) - f(x_k) \le -\frac{L}{2} \|x_{k+1} - x_k\|^2$$

最急降下法の解析: 勾配ノルムの評価

$$x_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$$
より

勾配ノルムの評価

$$\|
abla f(x_k)\| = L\|x_{k+1} - x_k\|$$

最急降下法の解析:2種類の評価を組み合わせる

A, **B** より
$$\frac{1}{2L} \|\nabla f(x_k)\|^2 \le f(x_k) - f(x_{k+1})$$

 $\Delta \coloneqq f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)$
 k について足す: $\frac{1}{2L} \sum_{i=0}^{k-1} \|\nabla f(x_i)\|^2 \le f(x_0) - f(x_k) \le \Delta$

$$\min_{0 \le i < k} \left\|
abla f(x_i)
ight\|^2 \le rac{1}{k} \sum_{i=0}^{\infty} \left\|
abla f(x_i)
ight\|^2$$
を使う: $\min_{0 \le i < k} \left\|
abla f(x_i)
ight\| \le \sqrt{rac{2L\Delta}{k}}$

14/77

最急降下法の解析:2種類の評価を組み合わせる

14/77

定石

・・・ ≤ f(x_k) - f(x_{k+1})を足して・・・ ≤ f(x₀) - f(x_k) ≤ ∆を得る
● 自明な関係 (最小値) ≤ (平均値)を使う

A, **B** より
$$\frac{1}{2L} \|\nabla f(x_k)\|^2 \le f(x_k) - f(x_{k+1})$$

 $\Delta \coloneqq f(x_0) - \inf_{x \in \mathbb{R}^d} f(x)$
 k について足す: $\frac{1}{2L} \sum_{i=0}^{k-1} \|\nabla f(x_i)\|^2 \le f(x_0) - f(x_k) \le \Delta$

$$\min_{0 \leq i < k} \left\|
abla f(x_i)
ight\|^2 \leq rac{1}{k} \sum\limits_{i=0}^{k-1} \left\|
abla f(x_i)
ight\|^2$$
を使う: $\min_{0 \leq i < k} \left\|
abla f(x_i)
ight\| \leq \sqrt{rac{2L\Delta}{k}}$

非凸最適化アルゴリズムの計算量解析の定石(再掲) 15/77

を両方示し, 組み合わせる

暫定解の列 $(x_k)_{k\in\mathbb{N}}$

計算量解析の定石

- $||x_{k+1} x_k||$ が大 \implies 関数値減少量が大 \bigcirc
- $||x_{k+1} x_k||$ が小 \implies 勾配ノルム が小 \bigcirc

ここからの話:最急降下法を例に

- 上をどう示す?
- 上をどう組み合わせる?
- 上が成り立つアルゴリズムをどう設計する?
 を見る

アルゴリズム設計の定石

① x_k 付近でfをよく近似する「シンプルな」関数 $\overline{f_k}$ を構成

2 \overline{f}_k を(近似的に)最小化する点を x_{k+1} とする



アルゴリズムの設計指針:近似関数の選び方

定石: 近似関数に課す条件

a
$$\bar{f}_k(x_k) = f(x_k), \ \nabla \bar{f}_k(x_k) = \nabla f(x_k)$$
 b $f(x_{k+1}) \le \bar{f}_k(x_{k+1})$

17/77

※ さらに $\overline{f_k}$ を強凸にすることも多い



アルゴリズムの設計指針: 近似関数の選び方

定石: 近似関数に課す条件

a
$$\bar{f}_k(x_k) = f(x_k), \ \nabla \bar{f}_k(x_k) = \nabla f(x_k)$$
 b $f(x_{k+1}) \le \bar{f}_k(x_{k+1})$

17/77

※ さらに $\overline{f_k}$ を強凸にすることも多い

- このように f_kをとると,
 関数値減少量 & 勾配ノルム の評価が容易 ②
- 最急降下法もこのような \bar{f}_k を使う手法と見なせる

アルゴリズムの設計指針:近似関数の選び方

定石: 近似関数に課す条件

$$\bullet \ \overline{f_k}(x_k) = f(x_k), \ \nabla \overline{f_k}(x_k) = \nabla f(x_k) \qquad \bullet \ f(x_{k+1}) \leq \overline{f_k}(x_{k+1})$$

17/77

※ さらに
$$\overline{f_k}$$
を強凸にすることも多い
 $\rightarrow \overline{f_k}(x_k) - \overline{f_k}(x_{k+1})$ が評価しやすい 〇



アルゴリズムの設計指針:最急降下法の場合

定石: 近似関数に課す条件

$$\begin{aligned} x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k) \\ &= \operatorname*{argmin}_{x \in \mathbb{R}^d} \left\{ \overline{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 \right\} \end{aligned}$$

• 🗛 は明らか

• $oldsymbol{B}$ は上界補題 $f(x) \leq \overline{f}_k(x)$ ($\forall x \in \mathbb{R}^d$) より成立

18/77

アルゴリズムの設計指針:最急降下法の場合

定石: 近似関数に課す条件

$$\bullet \ \overline{f_k}(x_k) = f(x_k), \ \nabla \overline{f_k}(x_k) = \nabla f(x_k) \qquad \bullet \ f(x_{k+1}) \leq \overline{f_k}(x_{k+1})$$

$$\begin{aligned} x_{k+1} &= x_k - \frac{1}{L} \nabla f(x_k) \\ &= \operatorname*{argmin}_{x \in \mathbb{R}^d} \left\{ \overline{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{L}{2} \|x - x_k\|^2 \right\} \end{aligned}$$

18/77

定石

- $\overline{f_k}$ を「Taylor 展開 + 正則化項」で構成
- Lipschitz 連続性を使って B を示す

ここまでのまとめ

計算量解析の定石

- $||x_{k+1} x_k||$ が大 \implies 関数値減少量が大 \bigcirc)を両方示し, ∫ 組み合わせる
- $||x_{k+1} x_k||$ が小 \implies 勾配ノルム が小 \bigcirc

アルゴリズム設計の定石

① x_k 付近でfをよく近似する「シンプルな」関数 $\frac{f_k}{f_k}$ を構成 **2** $\overline{f_k}$ を(近似的に)最小化する点を x_{k+1} とする

定石: 近似関数に課す条件

A $\overline{f_k}(x_k) = f(x_k), \ \nabla \overline{f_k}(x_k) = \nabla f(x_k)$ **B** $f(x_{k+1}) < f_k(x_{k+1})$ ここから少し脇道に逸れて 最急降下法のステップサイズ選択の話をします

実用上重要 かつ

最急降下法以外にも共通する 話題

Lipschitz 定数 L が未知の場合への対応

最急降下法の計算量(再掲)
最急降下法でステップサイズを
$$\eta_k = \frac{1}{L}$$
と定めると, $\min_{0 \le i < k} \| \nabla f(x_i) \| \le \sqrt{\frac{2L\Delta}{k}}$

Lが未知の場合どうする?

Lipschitz 定数 *L* が未知の場合への対応

最急降下法の計算量(再掲)
最急降下法でステップサイズを
$$\eta_k = \frac{1}{L}$$
と定めると, $\min_{0 \le i < k} \| \nabla f(x_i) \| \le \sqrt{\frac{2L\Delta}{k}}$

定石

- アルゴリズムと解析に登場する真値 Lを推定値ℓで置換
- 解析で使った不等式が不成立なら, ℓを大きくしてやり直し

$$f(x_{k+1}) - f(x_k) \le \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{\ell}{2} ||x_{k+1} - x_k||^2$$

最急降下法 + Lipschitz 定数の推定

アルゴリズムから従う不等式:

•
$$f(x_{k+1}) - f(x_k) \le \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{\ell}{2} ||x_{k+1} - x_k||^2$$

• $\ell \leq \max\{\ell_{\text{init}}, \alpha L\} \eqqcolon \ell_{\max}$

最急降下法の解析: Lipschitz 定数を推定する場合

$$L$$
が既知の場合と同様に $\frac{1}{2\ell} \|\nabla f(x_k)\|^2 \le f(x_k) - f(x_{k+1})$ が成立

$$\frac{\ell \leq \max\{\ell_{\text{init}}, \alpha L\} \eqqcolon \ell_{\max}}{2\ell_{\max}} \|\nabla f(x_k)\|^2 \leq \frac{1}{2\ell} \|\nabla f(x_k)\|^2 \leq f(x_k) - f(x_{k+1})$$

23/77

を*k*について足すと,

$$\min_{0 \le i < k} \|\nabla f(x_i)\| \le \sqrt{\frac{2\ell_{\max}\Delta}{k}}$$

c kが更新される反復は高々 $\left\lceil \frac{2\ell_{\max}\Delta}{\varepsilon^2} \right
vert$ 回

※ kが更新されない反復もあることに注意

最急降下法の解析: Lipschitz 定数を推定する場合

24/77

kが更新されない反復:

7行目がn回実行されるとすると、 $\ell_{\text{init}}\alpha^n \leq \ell_{\text{max}}$ より

$$n \leq \log_{lpha} \left(rac{\ell_{\max}}{\ell_{\min}}
ight)$$

最急降下法の解析: Lipschitz 定数を推定する場合

上の最急降下法の反復数の上界

$$\left\lceil \frac{2\ell_{\max}\Delta}{\varepsilon^2} \right\rceil + \log_{\alpha} \left(\frac{\ell_{\max}}{\ell_{\min}} \right) \qquad (\ell_{\max} \coloneqq \max\{\ell_{\min}, \, \alpha L\})$$

24/77
最急降下法の解析: Lipschitz 定数を推定する場合

1:
$$\ell \leftarrow \ell_{\text{init}}, k \leftarrow 0$$
 $\triangleright \ell_{\text{init}} > 0$: 初期推定值
2: **loop**
3: $x \leftarrow x_k - \frac{1}{\ell} \nabla f(x_k)$
4: **if** $f(x) - f(x_k) \le \langle \nabla f(x_k), x - x_k \rangle + \frac{\ell}{2} ||x - x_k||^2$:
5: $x_{k+1} \leftarrow x, k \leftarrow k+1$
6: **else**
7: $\ell \leftarrow \alpha \ell$ $\triangleright \alpha > 1$: 定数

24/77

- ∇f の Lipschitz 定数 L が未知でも計算量保証が得られた 😂

最急降下法の解析: Lipschitz 定数を推定する場合

25/77

1:
$$\ell \leftarrow \ell_{init}, k \leftarrow 0$$
 $\triangleright \ell_{init} > 0$: 初期推定値

 2: loop

 3: $x \leftarrow x_k - \frac{1}{\ell} \nabla f(x_k)$

 4: if $f(x) - f(x_k) \le \langle \nabla f(x_k), x - x_k \rangle + \frac{\ell}{2} ||x - x_k||^2$:

 5: $x_{k+1} \leftarrow x, k \leftarrow k+1, \ell \leftarrow \beta \ell$
 $\triangleright 0 < \beta \le 1$: 定数

 6: else

 7: $\ell \leftarrow \alpha \ell$
 $\triangleright \alpha > 1$: 定数

上の最急降下法の反復数の上界 む 演習 3 $\left[\frac{2\ell_{\max}\Delta}{\varepsilon^2}\right]\log_{\alpha}\left(\frac{\alpha}{\beta}\right) + \log_{\alpha}\left(\frac{\ell_{\max}}{\ell_{\min}}\right)$

最急降下法の解析: Lipschitz 定数を推定する場合

1:
$$\ell \leftarrow \ell_{\text{init}}, k \leftarrow 0$$
 $\triangleright \ell_{\text{init}} > 0$: 初期推定值
2: loop
3: $x \leftarrow x_k - \frac{1}{\ell} \nabla f(x_k)$
4: if $f(x) - f(x_k) \le \langle \nabla f(x_k), x - x_k \rangle + \frac{\ell}{2} ||x - x_k||^2$:
5: $x_{k+1} \leftarrow x, k \leftarrow k+1, \ell \leftarrow \beta \ell$ $\triangleright 0 < \beta \le 1$: 定数
6: else

7: $\ell \leftarrow \alpha \ell$

1: $\ell \leftarrow$

3: 4: 5: 6:

▷ α > 1: 正釵

25/77

補足: Armijo 条件 [Armijo, 1966]

$$f(x) - f(x_k) \le c \langle \nabla f(x_k), x - x_k \rangle \qquad (0 < c < 1: \ \textbf{z} \texttt{D})$$

 $c = \frac{1}{2}$ とすると4行目の条件に一致

数值例: Rosenbrock 関数 [Rosenbrock, 1960]

$$\min_{(x,y)\in\mathbb{R}^2} \ (x-1)^2 + 100(y-x^2)^2$$

26/77



数値例: β < 1の利点



最急降下法の計算量 $O\left(\frac{L\Delta}{\varepsilon^2}\right)$ は良い? 悪い?

ー次法: $f(x) \geq \nabla f(x)$ の計算を通してのみ fの情報にアクセス

- 一次法の計算量下界 [Carmon et al., 2020]
- 任意の $L, \Delta, \varepsilon > 0$ と任意の一次法Aに対し、ある関数fが存在し、
- ∇*f*は*L*-Lipschitz連続
- $f(x_0) \inf_x f(x) \le \Delta$ (x_0 : みの初期点)
- Aは $c_{\varepsilon^2}^{\underline{L}\Delta}$ 回以下の関数値・勾配計算では $\|\nabla f(x)\| \leq \varepsilon$ なる x を 見つけられない (c: 定数)



ー次法: $f(x) \geq \nabla f(x)$ の計算を通してのみ fの情報にアクセス

- 一次法の計算量下界 [Carmon et al., 2020]
- 任意の $L, \Delta, \varepsilon > 0$ と任意の一次法Aに対し、ある関数fが存在し、
- ∇*f*は*L*-Lipschitz連続
- $f(x_0) \inf_x f(x) \le \Delta$ (x_0 : みの初期点)
- Aは $c_{\varepsilon^2}^{\underline{L}\Delta}$ 回以下の関数値・勾配計算では $\|\nabla f(x)\| \leq \varepsilon$ なる x を 見つけられない (c: 定数)

上の証明は難しいが,アルゴリズムを最急降下法に限れば 似たことが比較的容易に示せる **企**演習4

一次法にとって最悪の問題例



図は [Carmon et al., 2020] より

•最急降下法

• 3次正則化Newton法

•再始動Heavy-ball法

3次正則化 Newton 法: 導入

31/77

定石(再掲)

$$\mathbf{\Theta} \ \overline{f_k}(x_k) = f(x_k), \ \nabla \overline{f_k}(x_k) = \nabla f(x_k) \qquad \mathbf{\Theta} \ f(x_{k+1}) \leq \overline{f_k}(x_{k+1})$$

- $\overline{f_k}$ を「Taylor 展開 + 正則化項」で構成
- Lipschitz 連続性を使って ³ を示す

- 最急降下法では1次の Taylor 展開を使った
- 2次の Taylor 展開を使うと? **公 3 次正則化 Newton 法**

3次正則化 Newton 法: 導入

ここで用いる上界補題

 $abla^2 f$ がM-Lipschitz連続ならば,

$$f(x) \le f(y) + \langle \nabla f(y), x - y \rangle + \frac{1}{2} \|x - y\|_{\nabla^2 f(y)}^2 + \frac{M}{6} \|x - y\|^3 \\ \|x\|_A^2 \coloneqq \langle Ax, x \rangle$$

比較:最急降下法で用いた上界補題 ∇f が *L*-Lipschitz 連続ならば, $f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} ||x - y||^2$

3次正則化Newton法: 導入

ここで用いる上界補題

 $abla^2 f$ がM-Lipschitz連続ならば,

$$\begin{aligned} f(x) &\leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{1}{2} \|x - y\|_{\nabla^2 f(y)}^2 + \frac{M}{6} \|x - y\|^3 \\ \|x\|_A^2 &\coloneqq \langle Ax, x \rangle \end{aligned}$$

$$\bar{f}_{k}(x) \coloneqq f(x_{k}) + \langle \nabla f(x_{k}), x - x_{k} \rangle + \frac{1}{2} \|x - x_{k}\|_{\nabla^{2} f(x_{k})}^{2} + \frac{M}{6} \|x - x_{k}\|^{3}$$

a
$$\overline{f_k}(x_k) = f(x_k), \ \nabla \overline{f_k}(x_k) = \nabla f(x_k)$$
 b $f(x_{k+1}) \leq \overline{f_k}(x_{k+1})$

仮定: $\nabla^2 f \, \boldsymbol{O} \, M$ -Lipschitz 連続性 $\left\| \nabla^2 f(x) - \nabla^2 f(y) \right\| \le M \|x - y\| \qquad (\forall x, y \in \mathbb{R}^d)$

3次正則化 (Cubic Regularized) Newton 法 [Griewank, 1981]

$$x_{k+1} \in \operatorname*{argmin}_{x \in \mathbb{R}^d} \left\{ \overline{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{6} \|x - x_k\|^3 \right\}$$

計算量保証 [Nesterov and Polyak, 2006]
$$\min_{1 \le i \le k} \|
abla f(x_i) \| \le \left(\frac{12\sqrt{M}\Delta}{k} \right)^{2/3}$$

c ε -停留点を求めるための計算量: $\left[\frac{12\sqrt{M}\Delta}{\varepsilon^{3/2}}\right] = O(\varepsilon^{-3/2})$

CRN法:子問題の非凸性

$$\begin{split} \min_{x \in \mathbb{R}^d} \quad f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{6} \|x - x_k\|^3 \\ \downarrow \quad \mathbf{u} \coloneqq x - x_k, \quad \mathbf{g} \coloneqq \nabla f(x_k), \quad \mathbf{H} \coloneqq \nabla^2 f(x_k) \\ \min_{u \in \mathbb{R}^d} \quad \langle g, u \rangle + \frac{1}{2} \|u\|_H^2 + \frac{M}{6} \|u\|^3 \end{split}$$

CRN 子問題は一般に非凸
☆ 問題の特殊性を利用して解く



$$\min_{u\in\mathbb{R}^d} \hspace{0.1 cm} \langle g,u
angle + rac{1}{2} \|u\|_{H}^2 + rac{M}{6} \|u\|^3$$

補題:子問題の最適解の性質 心 演習 5 $u \in \mathbb{R}^d$ が上の問題の最適解ならば、 $\sigma \coloneqq \frac{M}{2} ||u||$ とおくと $(H + \sigma I)u = -g, \quad H + \sigma I \succeq O$

※ 実は ← も成立

3 4

$$\min_{u\in\mathbb{R}^d} \hspace{0.1 cm} \langle g,u
angle + rac{1}{2} \|u\|_{H}^2 + rac{M}{6} \|u\|^3$$

補題:子問題の最適解の性質 心 演習 5 $u \in \mathbb{R}^d$ が上の問題の最適解ならば,

$$\sigma \coloneqq \frac{M}{2} \|u\| \qquad \texttt{とおくと} \qquad (H + \sigma I)u = -g, \qquad H + \sigma I \succeq O$$

直交行列で対角化: H = VΛV^T, 固有値 λ₁ ≤ ··· ≤ λ_d
ũ := V^Tu, ğ := V^Tg

$$\sigma = \frac{M}{2} \|\tilde{u}\|, \qquad (\Lambda + \sigma I)\tilde{u} = -\tilde{g}, \qquad \sigma \ge -\lambda_1$$

CRN法:子問題の最適解の性質(続き)

$$\sigma = \frac{M}{2} \|\tilde{u}\|, \qquad (\Lambda + \sigma I)\tilde{u} = -\tilde{g}, \qquad \sigma \ge -\lambda_1$$

$$ilde{g}_1
eq 0$$
 を仮定すると $\sigma>-\lambda_1$
このとき $ilde{u}_i=-rac{ ilde{g}_i}{\lambda_i+\sigma}$

$$ightarrow \sigma = rac{M}{2} \sqrt{\sum\limits_{i=1}^d rac{ ilde{g}_i^2}{(\lambda_i + \sigma)^2}}$$



右辺は $\sigma > -\lambda_1$ で単調減少 \mathbf{C} 二分法で解ける \mathbf{C}

※ $\tilde{g}_1 = 0$ の場合は Hard case と呼ばれ、少々厄介 (割愛)

CRN法:子問題の最適解の性質(続き)

$$\sigma = \frac{M}{2} \|\tilde{u}\|, \qquad (\Lambda + \sigma I)\tilde{u} = -\tilde{g}, \qquad \sigma \ge -\lambda_1$$

$$ilde{g}_1
eq 0$$
 を仮定すると $\sigma > -\lambda_1$
このとき $ilde{u}_i = -rac{ ilde{g}_i}{\lambda_i + \sigma}$

$$ightarrow \sigma = rac{M}{2} \sqrt{\sum\limits_{i=1}^d rac{ ilde{g}_i^2}{(\lambda_i + \sigma)^2}}$$



- σ が十分大ならば, \tilde{u} は — \tilde{g} 方向に近い
- $\sigma m \lambda_1$ に近づくにつれ, $\tilde{u} m e_1$ 成分が増えていく

36/77

- CRN 法は,解の更新方向 $u \coloneqq x x_k$ として
- $-\nabla f(x_k)$: 最急降下方向
- ∇²f(x_k)の最小固有値に対応する固有ベクトル方向
 を上手く組み合わせている!

固有ベクトル方向を活用することで<mark>鞍点</mark>を回避しやすい ☺

- σ が十分大ならば, \tilde{u} は — \tilde{g} 方向に近い
- $\sigma \dot{m} \lambda_1$ に近づくにつれ, $\tilde{u} \sigma e_1$ 成分が増えていく

数値例: CRN法は鞍点を回避しやすい



ここからCRN法の計算量解析をします

最急降下法で学んだ定石通り

計算量解析の定石(再掲)

- $||x_{k+1} x_k||$ が大 \implies 関数値減少量が大 \bigcirc を両方示し, 組み合わせる
- $||x_{k+1} x_k||$ が小 \implies 勾配ノルム が小 \bigcirc

CRN 法の解析: 関数値減少量の評価

既に示したこと $u_k \coloneqq x_{k+1} - x_k, \ \sigma_k \coloneqq \frac{M}{2} \|u_k\|$ とおくと, $(H + \sigma_k I)u_k = -g, \qquad H + \sigma_k I \succeq O$

CRN 法の解析: 関数値減少量の評価

既に示したこと $(H + \sigma_k I)u_k = -g, \qquad H + \sigma_k I \succeq O$ 上界補題 $f(x_{k+1}) - f(x_k) \leq \langle g, u_k \rangle + rac{1}{2} \|u_k\|_H^2 + rac{M}{6} \|u_k\|^3$ $(H + \sigma_k I)u_k = -g = -\|u_k\|_{H + \sigma_k I}^2 + \frac{1}{2}\|u_k\|_H^2 + \frac{M}{6}\|u_k\|^3$ $= -\frac{1}{2} \|u_k\|_{H+\sigma_k I}^2 - \frac{\sigma_k}{2} \|u_k\|^2 + \frac{M}{6} \|u_k\|^3$ $H + \sigma_k I \succeq O \qquad \leq -\frac{\sigma_k}{2} \|u_k\|^2 + \frac{M}{\epsilon} \|u_k\|^3$ $\sigma_k = \frac{M}{2} \|u_k\| = -\frac{M}{12} \|u_k\|^3$

CRN 法の解析: 勾配ノルムの評価

ここで使う補題 $\nabla^2 f \, t^{m} M$ -Lipschitz 連続ならば, $\left\| \nabla f(x) - \left(\nabla f(y) + \nabla^2 f(y)(x-y) \right) \right\| \le \frac{M}{2} \|x-y\|^2 \qquad (\forall x, y \in \mathbb{R}^d)$

最急降下法で用いた補題(再掲) $\nabla f \, \check{x} L$ -Lipschitz 連続ならば, $\left| f(x) - \left(f(y) + \langle \nabla f(y), x - y \rangle \right) \right| \leq \frac{L}{2} \|x - y\|^2 \quad (\forall x, y \in \mathbb{R}^d)$

と全く同様に示せる

ここで使う補題 $abla^2 f \, t^{\prime} M$ -Lipschitz 連続ならば, $\left\| \nabla f(x) - \left(\nabla f(y) + \nabla^2 f(y)(x-y) \right) \right\| \leq \frac{M}{2} \|x-y\|^2 \qquad (\forall x, y \in \mathbb{R}^d)$

上の補題と三角不等式
$$\|\nabla f(x_{k+1})\| \le \|\nabla f(x_k) + \nabla^2 f(x_k)u_k\| + \frac{M}{2} \|u_k\|^2$$
$$= \|g + Hu_k\| + \frac{M}{2} \|u_k\|^2$$
$$(H + \sigma_k I)u_k = -g = \sigma_k \|u_k\| + \frac{M}{2} \|u_k\|^2$$
$$\sigma_k = \frac{M}{2} \|u_k\| = M \|u_k\|^2$$

CRN 法の解析: 2種類の評価を組み合わせる

Δ, **B** より
$$\frac{\|
abla f(x_{k+1})\|^{3/2}}{12\sqrt{M}} \leq f(x_k) - f(x_{k+1})$$

これを k について足し、最急降下法のときと同様に評価すると、

$$\min_{1 \le i \le k} \|\nabla f(x_i)\| \le \left(\frac{12\sqrt{M}\Delta}{k}\right)^{2/3}$$

ゆ 計算量は
$$\left\lceil \frac{12\sqrt{M}\Delta}{\varepsilon^{3/2}} \right\rceil = O(\varepsilon^{-3/2})$$

CRN 法の解析: 2種類の評価を組み合わせる

これを k について足し、最急降下法のときと同様に評価すると、

$$\min_{1 \le i \le k} \|
abla f(x_i)\| \le \left(rac{12\sqrt{M}\Delta}{k}
ight)^{2/k}$$

ゆ計算量は
$$\left[\frac{12\sqrt{M}\Delta}{\varepsilon^{3/2}}\right] = O(\varepsilon^{-3/2})$$

3

41/77

CRN 法の解析: 2次の停留点への収束

C
$$\nabla^2 f(x_{k+1}) \succeq \nabla^2 f(x_k) - M \|u_k\| I \succeq -\frac{3}{2} M \|u_k\| I$$

Lipschitz 連続性 $\nabla^2 f(x_k) \succeq -\sigma_k I = -\frac{M}{2} \|u_k\| I$

(A, B, C) を組み合わせると、 $\|\nabla f(x)\| \leq \varepsilon$ に加えて

$$\nabla^2 f(x) \succeq -\frac{3}{2}\sqrt{M\varepsilon}I$$

も満たす点xが $\left[\frac{12\sqrt{M}\Delta}{\varepsilon^{3/2}}\right] = O(\varepsilon^{-3/2})$ 反復で見つかることがわかる

CRN 法の解析: 2次の停留点への収束

 $\varepsilon \rightarrow 0$ とすると、2次の停留点:

$$abla f(x) = \mathbf{0}, \qquad
abla^2 f(x) \succeq O$$

♪ CRN 法は鞍点を回避しやすい

Δ, **B**, **C** を組み合わせると,
$$\|\nabla f(x)\| \leq \varepsilon$$
に加えて

$$\nabla^2 f(x) \succeq -\frac{3}{2} \sqrt{M\varepsilon} I$$

も満たす点xが $\left[\frac{12\sqrt{M}\Delta}{\varepsilon^{3/2}}\right] = O(\varepsilon^{-3/2})$ 反復で見つかることがわかる

ここから CRN 法のその他の話題をいくつか紹介

- 子問題を近似的に解く
- Lipschitz 定数 M の推定
- Lazy CRN法
- Hesse 行列・ベクトル積の利用

子問題を近似的に解く [Cartis et al., 2012]

• 正則化を強めに設定:

$$\bar{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{3} \|x - x_k\|^3$$

• 近似解の条件: **④** $\|\nabla \bar{f}_k(x_{k+1})\| \le \frac{M}{2} \|u_k\|^2$, **⑤** $\bar{f}_k(x_{k+1}) \le \bar{f}_k(x_k)$

$$u_k \coloneqq x_{k+1} - x_k$$

正則化を強めると子問題の1次の停留点を求めるだけで十分に 😂

※ 元問題の2次の停留点を求めたい場合はもう少し頑張る必要あり

44/77

子問題を近似的に解く: 関数値減少量の評価

• 正則化を強めに設定:

$$\bar{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{3} \|x - x_k\|^3$$

• 近似解の条件: **③** $\|\nabla \bar{f}_k(x_{k+1})\| \le \frac{M}{2} \|u_k\|^2$, **⑤** $\bar{f}_k(x_{k+1}) \le \bar{f}_k(x_k)$

 $u_k \coloneqq x_{k+1} - x_k$

44/77

$$f(x_{k+1}) - f(x_k) = f(x_{k+1}) - \bar{f}_k(x_k)$$

上界補題 $\leq \left(\bar{f}_k(x_{k+1}) - \frac{M}{6} \|u_k\|^3\right) - \bar{f}_k(x_k)$
B $\leq -\frac{M}{6} \|u_k\|^3$

|子問題を近似的に解く: 勾配ノルムの評価

• 正則化を強めに設定:

$$\bar{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{3} \|x - x_k\|^3$$

44/77

• 近似解の条件: **④** $\|\nabla \bar{f}_k(x_{k+1})\| \le \frac{M}{2} \|u_k\|^2$, **⑤** $\bar{f}_k(x_{k+1}) \le \bar{f}_k(x_k)$

$$\begin{aligned} \|\nabla f(x_{k+1})\| &\leq \left\|\nabla \bar{f}_{k}(x_{k+1})\right\| + \left\|\nabla f(x_{k+1}) - \nabla \bar{f}_{k}(x_{k+1})\right\| \\ &\triangleq \frac{M}{2} \|u_{k}\|^{2} + \left\|\nabla f(x_{k+1}) - \nabla \bar{f}_{k}(x_{k+1})\right\| \\ &= \frac{M}{2} \|u_{k}\|^{2} + \left\|\nabla f(x_{k+1}) - \nabla f(x_{k}) - \nabla^{2} f(x_{k})u_{k} - M\|u_{k}\|u_{k}\right\| \\ &\triangleq \frac{M}{2} \|u_{k}\|^{2} \leq \frac{M}{2} \|u_{k}\|^{2} + \frac{M}{2} \|u_{k}\|^{2} + M\|u_{k}\|^{2} = 2M \|u_{k}\|^{2} \end{aligned}$$

• 正則化を強めに設定:

$$\bar{f}_k(x) \coloneqq f(x_k) + \langle \nabla f(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|_{\nabla^2 f(x_k)}^2 + \frac{M}{3} \|x - x_k\|^3$$

• 近似解の条件: **④** $\|\nabla \bar{f}_k(x_{k+1})\| \leq \frac{M}{2} \|u_k\|^2$, **B** $\bar{f}_k(x_{k+1}) \leq \bar{f}_k(x_k)$

- 正則化を強めたおかげで証明がシンプルに 😂
- 「子問題の最適解の性質」の補題も不使用

最急降下法と同様に Lipschitz 定数 *M* の推定が可能 😂

Mの推定値mに対し,解析で用いた不等式

• $f(x_{k+1}) - f(x_k) \le \langle \nabla f(x_k), u_k \rangle + \frac{1}{2} \|u_k\|_{\nabla^2 f(x_k)}^2 + \frac{m}{6} \|u_k\|^3$

•
$$\|\nabla f(x_{k+1}) - \nabla f(x_k) - \nabla^2 f(x_k) u_k\| \le \frac{m}{2} \|u_k\|^2$$

の成立をチェック

成立していなければ <mark>m</mark>を大きくしてやり直す
Lazy CRN法 [Doikov et al., 2023]

$ abla^2 f(\mathbf{x}_0)$	reuse Hessian		$ abla^2 f(\mathbf{x}_m)$		
$ abla f(\mathrm{x}_0)$	$ abla f(\mathrm{x}_1)$		$ abla f(\mathrm{x}_{m-1})$	$ abla f(\mathrm{x}_m)$	(図は [Doikov et al., 2023] より)

定理: Lazy CRN 法の反復数 [Doikov et al., 2023]

- Hesse 行列を *m* 反復再利用
- 正則化項 $\frac{M}{6} \|u\|^3 \ge mM \|u\|^3$ に

計算量:
$$O\left(\frac{\sqrt{M}\Delta}{\varepsilon^{3/2}}\right)(\mathsf{G}+\mathsf{H}) \longrightarrow O\left(\frac{\sqrt{M}\Delta}{\varepsilon^{3/2}}\right)\left(\sqrt{m}\mathsf{G}+\frac{\mathsf{H}}{\sqrt{m}}\right)$$

$$oldsymbol{C} oldsymbol{m} \simeq rac{\mathsf{H}}{\mathsf{G}}$$
と設定すれば,計算 $\mathbb{E} \operatorname{O} \left(rac{\sqrt{M}\Delta}{arepsilon^{3/2}}
ight) \sqrt{\mathsf{G}\cdot\mathsf{H}}$ 🕲

$$\min_{u \in \mathbb{R}^d} \ \langle g, u
angle + rac{1}{2} \| u \|_H^2 + rac{M}{6} \| u \|^3$$

dが大きいとき, $H \coloneqq
abla^2 f(x_k) \in \mathbb{R}^{d imes d}$ を分解するのは非現実的… 😕

Hesse 行列・ベクトル積 (HVP) $\nabla^2 f(x) v$ をオラクルとして利用

- $x, v \in \mathbb{R}^d$ に対し、 $\nabla^2 f(x) v$ の計算が容易な場合がある 例: $\nabla^2 f(x)$ が疎 or 低ランク、高速自動微分
- 勾配で近似も可能: $abla^2 f(x) v \simeq rac{
 abla f(x+\delta v)
 abla f(x)}{\delta}$ (非推奨)
- メモリが節約でき、d ~ 10⁶ 程度でも動く ☺

HVP が使える CRN 子問題の解法いろいろ

• Lanczos法 [Carmon and Duchi, 2018, 2020; Cartis et al., 2011; Gould and Simoncini, 2020] 等

48/77

[Royer et al., 2020]

- 共役勾配法 + 最小固有値計算
- 2(d+1)次の一般化固有値問題に帰着 [Lieder, 2020]
 上の基となった信頼領域法版: [Adachi et al., 2017]
- (*d*+1)次の2次固有値問題に帰着 [Jia et al., 2022]
- 制約付き凸最適化に帰着 → 加速射影勾配法 [Jiang et al., 2021]
- 局所最適化とジャンプを繰り返す [Cristofari et al., 2019]
- 小さい固有値のみで方程式を近似 [Gao et al., 2022]

CRN子問題のためのLanczos法

$$\min_{u\in\mathbb{R}^d} \hspace{0.1 cm} \langle g,u
angle + rac{1}{2} \|u\|_{H}^2 + rac{M}{6} \|u\|^3$$

Krylov 部分空間に制限: u ∈ span{g, Hg,..., H^{m-1}g}

•
$$[g, Hg, \dots, H^{m-1}g] = QR と \mathsf{QR} 分解し, u = Qv とおく$$

$$\min_{\boldsymbol{v}\in\mathbb{R}^{m}} \langle Q^{\top}g, \boldsymbol{v}\rangle + \frac{1}{2} \|\boldsymbol{v}\|_{\boldsymbol{T}}^{2} + \frac{M}{6} \|\boldsymbol{v}\|^{3}$$

- $T \coloneqq Q^\top HQ$ は**三重対角** 😂
- HVPの計算回数はO(m)
- Hard case 対策に乱数を使うことも (例: [Carmon and Duchi, 2018])

CRN 子問題のための Lanczos 法: m の設定

部分空間の次元 m の選び方 [Carmon and Duchi, 2020, Lemma 6.1]

$$m \ge \left\lceil 24 \cdot 2^{3/4} rac{\sqrt{L}}{(Marepsilon)^{1/4}} \left(1 + rac{1}{8} \log^2 \left(rac{4d}{\delta^2}
ight)
ight)
ight
ceil = \Theta ig(arepsilon^{-1/4}ig)$$

ならば,確率 $1 - \delta$ で,十分な精度の子問題の解が得られる

CRN 子問題のための Lanczos 法: mの設定

部分空間の次元 mの選び方 [Carmon and Duchi, 2020, Lemma 6.1]

$$m \ge \left\lceil 24 \cdot 2^{3/4} \frac{\sqrt{L}}{(M\varepsilon)^{1/4}} \left(1 + \frac{1}{8} \log^2 \left(\frac{4d}{\delta^2} \right) \right) \right\rceil = \Theta\left(\varepsilon^{-1/4}\right)$$

50/77

ならば,確率 $1-\delta$ で,十分な精度の子問題の解が得られる

× O
$$\left(\frac{\sqrt{M}\Delta}{\varepsilon^{3/2}}\right)$$
 (CRN 法の反復数)

CRN 法全体での HVP 計算回数 [Carmon and Duchi, 2020, Prop 6.3] $O(1) \cdot \frac{\sqrt{L}M^{1/4}\Delta}{\varepsilon^{7/4}} \left(1 + \log^2\left(\frac{d}{\delta^2}\right) + \log\left(\frac{\sqrt{L}M^{1/4}\Delta}{\varepsilon^{7/4}}\right)\right) = \tilde{O}(\varepsilon^{-7/4})$

	Lipschitz 仮定	計算量 / Δ	2次停留点
1 最急降下法	∇f	$O(L\varepsilon^{-2})$ G	
2 CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})(G+H)$	\checkmark
3 Lazy CRN	$ abla^2 f$	$\mathrm{O}(\sqrt{M}\varepsilon^{-\frac{3}{2}})\sqrt{G\cdotH}$	\checkmark
• CRN + Lanczo	os $\nabla f \And \nabla^2 f$	$ ilde{\mathbf{O}}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})HVP$	\checkmark
		$+ O(\sqrt{M}\varepsilon^{-\frac{3}{2}})G$	

	Lipschitz 仮定	計算量 / Δ	2次停留点	
1 最急降下法	∇f	$O(L\varepsilon^{-2})$ G		
2 CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})(G+H)$	\checkmark	
3 Lazy CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})\sqrt{G\cdotH}$	\checkmark	
CRN + Lanczos	$ abla f$ & $ abla^2 f$	$rac{ ilde{O}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})HVP}{+\operatorname{O}(\sqrt{M}arepsilon^{-rac{3}{2}})G}$	\checkmark	
• Lipschitz 定数 <i>L</i> , <i>M</i>				
● 精度 ε	L 29	なを考慮し(アルコリスム選択		
● オラクルコスト G, H, HVP (♪ 演習 1(c), (d)				
 子問題の計算コスト 				

	Lipschitz 仮定	計算量 / Δ	2次停留点
1 最急降下法	∇f	$\mathrm{O}(Larepsilon^{-2})G$	
2 CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})(G+H)$	\checkmark
3 Lazy CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})\sqrt{G\cdotH}$	\checkmark
CRN + Lanczos	$ abla f$ & $ abla^2 f$	$ ilde{\mathrm{O}}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})HVP$	\checkmark
		$+ \operatorname{O}(\sqrt{M}\varepsilon^{-\frac{3}{2}})G$	

例:

- 「*d* が大」のとき, **1** or **4** が有力
- さらに「 ε が小」or「 $M \ll L$ 」で,HVPが使えるなら 4 が有力

	Lipschitz 仮定	計算量 /Δ	2次停留点
1 最急降下法	abla f	${ m O}(Larepsilon^{-2}){ m G}$	
2 CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})(G+H)$	\checkmark
3 Lazy CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})\sqrt{G\cdotH}$	\checkmark
CRN + Lanczos	$ abla f$ & $ abla^2 f$	$ ilde{ extsf{O}}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}}) extsf{HVP}$	\checkmark
		$+ \operatorname{O}(\sqrt{M}\varepsilon^{-\frac{3}{2}})G$	

例:

- 「*d* が小」で H がそれほど大きくないなら, 2 or 3 が有力
- さらに「H ~ G」なら 2 が有力 (定数倍の都合)

• 最急降下法

・3次正則化Newton法

• 再始動 Heavy-ball 法

より良い計算量の追求

	Lipschitz 仮定	計算量 / Δ	2次停留点
1 最急降下法	∇f	$O(L\varepsilon^{-2})$ G	
2 CRN	$ abla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})(G+H)$	\checkmark
3 Lazy CRN	$\nabla^2 f$	$O(\sqrt{M}\varepsilon^{-\frac{3}{2}})\sqrt{G\cdotH}$	\checkmark
• CRN + Lanczos	$ abla f$ & $ abla^2 f$	$\tilde{O}(\sqrt{L}M^{\frac{1}{4}}\varepsilon^{-\frac{7}{4}})$ HVP	\checkmark
		$+ O(\sqrt{M\epsilon^2})$	
❻ これから説明	∇f & $ abla^2 f$	$\mathrm{O}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})G$	

abla f と $abla^2 f$ の Lipschitz 連続性を仮定する手法(一部) 53/77

		計算量 / Δ	2次停留点	決定的	
	4	$\tilde{O}(\sqrt{L}M^{\frac{1}{4}}\varepsilon^{-\frac{7}{4}})HVP+O(\sqrt{M}\varepsilon^{-\frac{3}{2}})G$	\checkmark		
	а	$ ilde{\mathrm{O}}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})G$	\checkmark		
	Ь	$ ilde{\mathrm{O}}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})G$		\checkmark	
	5	$\mathrm{O}(\sqrt{L}M^{rac{1}{4}}arepsilon^{-rac{7}{4}})G$		\checkmark	
4	CRN 系	[Agarwa	al et al., 2017; Ca	rmon and Du	chi, 2020]
	加速勾配	法 + 特殊な正則化 + 最小固有値近似詞	†算	[Carmon et	al., 2018]
	Damped	Newton-CG + 最小固有值近似計算		[Royer et	al., 2020]
a) HVP を使わず最小固有値近似計算 [Allen-Zhu and Li, 2018; Xu et :			al., 2017]	
	加速勾配	法 + 摂動		[Jin et	al., 2018]
b	Convex until proven guilty [Carmon et al., 20			al., 2017]	
6	• Heavy-ball 法 or 加速勾配法 + 再始動 [Li and Lin, 2022, 2023; Marumo and Takeda, 2024a,b]				

余談: 初期のアルゴリズムの例 [Agarwal et al., 2017]

54/77

Algorithm 1 FastCubic($f, x_0, \varepsilon, L, L_2$) Algorithm 2 FastCubicMin $(q, \mathbf{H}, L, L_2, \kappa)$ (main algorithm for cubic minimization) **Input:** f(x) that satisfies (2.1) with L_2 and L_3 a starting vector x_0 ; a target accuracy ε . **Input:** q a vector, **H** a symmetric matrix, parameters κ , L and L₂ which satisfies $-L_2 \mathbf{I} \prec \mathbf{H} \prec L_2 \mathbf{I}$. 1: $\kappa \leftarrow \left(\frac{900}{\varepsilon L}\right)^{1/2}$. **Output:** (λ, v, v_{\min}) 1: $B \leftarrow L_2 + \sqrt{L \|g\|} + \frac{1}{r}$. 2: for t = 0 to ∞ do $\frac{1}{2} \sum_{\tilde{\epsilon} \leftarrow 1/(1000 \,(\max\{L, \|g\|, \frac{3\kappa}{10}, B, 1\})^{20})} 10000 \,(\cdots)^{20}$ $m_t(h) \triangleq \nabla f(x_t)^\top h + \frac{h^\top \nabla^2 f(x_t)h}{2} + \frac{L}{6} \|h\|^3$ 3: $(\lambda, v, v_{\min}) \leftarrow \mathsf{FastCubicMin} \left(\nabla f(x_t), \nabla^2 f(x_t), L, L_2, \kappa \right)$ 3: $\lambda_0 \leftarrow 2B$. 4: 4: for i = 0 to ∞ do $h' \leftarrow \text{either } v \text{ or } \frac{\lambda v_{\min}}{2I} \text{ whichever gives smaller value for } m_t(h);$ 5: Compute v such that $||v + (\mathbf{H} + \lambda_i \mathbf{I})^{-1}g|| \leq \tilde{\varepsilon}$. Set $x_{t+1} \triangleq x_t + h'$ 6: if $L||v|| \in [2\lambda_i - L\tilde{\varepsilon}, 2\lambda_i + L\tilde{\varepsilon}]$ then if $m_t(h') > -\frac{\varepsilon^{3/2}}{\alpha/L}$ then return x_{t+1} . \diamond c is a constant; we proved $c = 2.4 * 10^6$ works return $(\lambda_i, v, \emptyset)$. 8 end for else if $L||v|| > 2\lambda_i + L\tilde{\varepsilon}$ then $c = 2.4 \cdot 10^{6}$ return BinarySearch($\lambda_1 = \lambda_{i-1}, \lambda_2 = \lambda_i, \tilde{\varepsilon}$). else if $L \|v\| < 2\lambda_i - L\tilde{\varepsilon}$ then Let Power Method find vector w that is 9/10-appx leading eigenvector of $(\mathbf{H} + \lambda_i \mathbf{I})^{-1}$: Algorithm 3 BinarySearch($\lambda_1, \lambda_2, \tilde{\varepsilon}$) (binary search subroutine) $\frac{9}{10}\lambda_{\max}((\mathbf{H}+\lambda_i\mathbf{I})^{-1}) \le w^{\top}(\mathbf{H}+\lambda_i\mathbf{I})^{-1}w \le \lambda_{\max}((\mathbf{H}+\lambda_i\mathbf{I})^{-1}) .$ Input: $\lambda_1 \geq \lambda_2$, $L \| (\mathbf{H} + \lambda_1 \mathbf{I})^{-1} g \| \leq 2\lambda_1$, $L \| (\mathbf{H} + \lambda_2 \mathbf{I})^{-1} g \| \geq 2\lambda_2$, $\lambda_2 + \lambda_{\min}(\mathbf{H}) > 0$ **Output:** (λ, v, \emptyset) Compute a vector \tilde{w} such that $\|\tilde{w} - (\mathbf{H} + \lambda_i \mathbf{I})^{-1} w\| \le \hat{\varepsilon} \triangleq \frac{1}{60R}$. 12: 1: for t = 1 to ∞ do $\Delta \leftarrow \frac{1}{2} \frac{1}{\tilde{w}^{\top} w - \tilde{\epsilon}}$. 13: $\lambda_{\text{mid}} \leftarrow \frac{\lambda_1 + \lambda_2}{2}$ if $\Delta > \frac{1}{2n}$ then 14: Compute vector v such that $||v + (\mathbf{H} + \lambda_{\text{mid}}\mathbf{I})^{-1}g|| \leq \tilde{\varepsilon}/2$ 3. $\tilde{\lambda}_{i+1} \leftarrow \lambda_i - \frac{\Delta}{2}$. 15: if $L \|v\| \in [2\lambda_{\text{mid}} - L\tilde{\varepsilon}, 2\lambda_{\text{mid}} + L\tilde{\varepsilon}]$ then 4: if $\tilde{\lambda}_{i+1} > 0$ then $\lambda_{i+1} \leftarrow \tilde{\lambda}_{i+1}$ else $\lambda_{i+1} \leftarrow 0$ 16. return $(\lambda_{\text{mid}}, v, \emptyset)$ 5: else 17: 6: else if $L||v|| + L\tilde{\varepsilon} \leq 2\lambda_{\text{mid}}$ then 18. Use AppxPCA to find any unit vector v_{\min} such that $v_{\min}^{\top} \mathbf{H} v_{\min} \leq \lambda_{\min}(\mathbf{H}) + \frac{1}{10r}$. 7: $\lambda_1 \leftarrow \lambda_{mid}$ Flip the sign of v_{\min} so that $q^{\top}v_{\min} < 0$. else if $L||v|| - L\tilde{\varepsilon} > 2\lambda_{mid}$ then 19: return (λ_i, v, v_{\min}) . 20. $\lambda_2 \leftarrow \lambda_{\text{mid}}$ Q٠ end if 21: end if end if 99. 11: end for 23: end for

最新のアルゴリズムはよりシンプル & 実用的 😂

- どちらも慣性を利用する手法
- 近似関数 \bar{f}_k の最小化としては解釈しにくい(?)

$$x_{k+1} = x_k - \eta_k \nabla f(x_k) + \theta_k (x_k - x_{k-1})$$

加速勾配法 (Accelerated gradient descent; AGD) [Nesterov, 1983] $x_{k+1} = x_k - \eta_k \nabla f(x_k + \theta_k(x_k - x_{k-1})) + \theta_k(x_k - x_{k-1})$



$$x_{k+1} = x_k - \eta_k \nabla f(x_k) + \theta_k (x_k - x_{k-1})$$

加速勾配法 (Accelerated gradient descent; AGD) [Nesterov, 1983] $x_{k+1} = x_k - \eta_k \nabla f (x_k + \theta_k (x_k - x_{k-1})) + \theta_k (x_k - x_{k-1})$ $\eta_k = \eta, \ \theta_k = 1 - \gamma \sqrt{\eta} \ \varepsilon$ 定め, $\eta \to 0$ ($\gamma \ge 0$: 定数)

HB 微分方程式

$$\ddot{x}(t) = -\gamma \dot{x}(t) - \nabla f(x(t))$$

$$x_{k+1} = x_k - \eta_k \nabla f(x_k) + \theta_k (x_k - x_{k-1})$$

加速勾配法 (Accelerated gradient descent; AGD) [Nesterov, 1983]

$$x_{k+1} = x_k - \eta_k \nabla f \left(x_k + \theta_k (x_k - x_{k-1}) \right) + \theta_k (x_k - x_{k-1})$$

凸の場合:

- HBは強凸2次関数に対し最良の計算量 [Polyak, 1964]
 η_k, θ_k が k に非依存だと一般の強凸関数に対し非最適 [Goujaud et al., 2023]
- AGD は一般の凸関数 / 強凸関数に対し最適 [Nesterov, 1983, 2004]

$$x_{k+1} = x_k - \eta_k \nabla f(x_k) + \theta_k (x_k - x_{k-1})$$

加速勾配法 (Accelerated gradient descent; AGD) [Nesterov, 1983]

$$x_{k+1} = x_k - \eta_k \nabla f \left(x_k + \theta_k (x_k - x_{k-1}) \right) + \frac{\theta_k (x_k - x_{k-1})}{\theta_k (x_k - x_{k-1})}$$

非凸の場合:

∇f と ∇²f が Lipschitz 連続ならば、HB も
 AGD も 再始動
 と組み合わせると強い

途中で
$$x_{k-1} = x_k$$
と再設定



図は [O'Donoghue and Candès, 2015] より

最新のHB/AGD [Marumo and Takeda, 2024a,b] の利点

Lipschitz 定数 *L*, *M* の自動推定が可能で,真の値は入力不要

57/77

最新のHB/AGD [Marumo and Takeda, 2024a,b] の利点

Lipschitz 定数 *L*, *M* の自動推定が可能で,真の値は入力不要

57/77

参考: それ以前の一次法は*L*,*M*,*ε*の入力が必要

- 例: [Li and Lin, 2022] では慣性係数 $heta_k = 1 rac{2(Marepsilon)^{1/4}}{\sqrt{L}}$
- •計算量解析で以下の不等式を使用. Mの推定は難しそう 😣

•
$$f(x) - f(y) \le \langle \nabla f(y), x - y \rangle + \frac{1}{2} ||x - y||_{\nabla^2 f(y)}^2 + \frac{M}{6} ||x - y||^3$$

• $\left\| \nabla f(x) - \nabla f(y) - \nabla^2 f(y)(x - y) \right\| \le \frac{M}{2} ||x - y||^2$

最新のHB/AGD [Marumo and Takeda, 2024a,b] の利点

Lipschitz 定数 *L*, *M* の自動推定が可能で,真の値は入力不要

57/77

参え
$$\nabla^2 f(y)$$
を含まない不等式のみを使って解析する
• 例: [Li and Lin, 2022] では慣性係数 $\theta_k = 1 - \frac{2(M\varepsilon)^{1/4}}{\sqrt{L}}$
• 計算量解析で以下の不等式を使用. *M* の推定は難しそう ②
• $f(x) - f(y) \leq \langle \nabla f(y), x - y \rangle + \frac{1}{2} ||x - y||_{\nabla^2 f(y)}^2 + \frac{M}{6} ||x - y||^3$

•
$$\left\| \nabla f(x) - \nabla f(y) - \nabla^2 f(y)(x-y) \right\| \le \frac{M}{2} \|x-y\|^2$$

計算量解析に用いる Hessian-free 不等式1

$$f(x) - f(y) \le \langle \nabla f(y), x - y \rangle + \frac{1}{2} \|x - y\|_{\nabla^2 f(y)}^2 + \frac{M}{6} \|x - y\|^3$$
の代わりに

$$f(x) - f(y) \le \frac{1}{2} \langle \nabla f(x) + \nabla f(y), x - y \rangle + \frac{M}{12} ||x - y||^{3}$$

●証明は演習6

- 台形則の誤差評価の多次元版
- $\nabla^2 f(y)$ なしでも $f(x) - f(y) \leq \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2$ より良いオーダー f(x) - f(y)

計算量解析に用いる Hessian-free 不等式2

$$\left\|\nabla f(x) - \nabla f(y) - \nabla^2 f(y)(x-y)\right\| \leq \frac{M}{2} \|x-y\|^2$$
の代わりに

$$\left\|\nabla f(\bar{x}) - \sum_{i=1}^{n} \lambda_i \nabla f(x_i)\right\| \le \frac{M}{2} \sum_{i=1}^{n} \lambda_i \|x_i - \bar{x}\|^2$$

•
$$\bar{x} \coloneqq \sum_{i=1}^{n} \lambda_i x_i$$
: x_1, \ldots, x_n の重みつき平均

• 直観: M が小 $\implies \nabla f$ は affine に近い \implies 左辺が小

• Jensenの不等式
$$f(\bar{x}) - \sum_{i=1}^{n} \lambda_i f(x_i) \le 0$$
 に類似

凸最適化で $f(\bar{x})$ が小さいことを示すのによく使われる

♪ 証明は演習6

ここからHB法の**計算量解析**をします

- 簡単のため, *L*, *M* は既知と仮定
- ステップサイズ $\eta_k = \frac{1}{L}$,慣性係数 $\theta_k = 1$ と設定

計算量解析の定石(再掲)

• $||x_{k+1} - x_k||$ が大 \implies **関数値減少量**が大 \bigcirc • $||x_{k+1} - x_k||$ が小 \implies **勾配ノルム** が小 \bigcirc 組み合わせる

HB法の解析: 関数値減少量 & 勾配ノルム の評価 61/77

これ以降
$$S_k \coloneqq \sum_{i=0}^{k-1} \|x_{i+1} - x_i\|^2$$
, $\bar{x}_k \coloneqq \frac{1}{k} \sum_{i=0}^{k-1} x_i$, $x_k^\star \coloneqq \operatorname*{argmin}_{x \in \{x_1, \dots, x_k\}} f(x)$

謎の仮定
$$\sqrt{(k-1)^5S_{k-1}} \leq rac{3L}{4M}$$
の下,

•
$$S_k$$
が大 ⇒ 関数値減少量が大 ②: $f(\mathbf{x}_k^{\star}) - f(x_0) \le -\frac{LS_k}{4k}$
• S_k が小 ⇒ 勾配ノルム が小 ③: $\min_{1 \le i \le k} \|\nabla f(\bar{x}_i)\| \le 3L\sqrt{\frac{S_k}{k^3}}$

これまでとの違い:

●証明は演習7,8

TO

- k反復分をまとめて評価. $||x_{k+1} x_k||$ の代わりに S_k が登場
- 平均解 x̄_k で勾配ノルムを評価

HB法の解析: 関数値減少量 & 勾配ノルム の評価 61/77

これ以降
$$S_k \coloneqq \sum_{i=0}^{k-1} \|x_{i+1} - x_i\|^2$$
, $\bar{x}_k \coloneqq \frac{1}{k} \sum_{i=0}^{k-1} x_i$, $x_k^\star \coloneqq \operatorname*{argmin}_{x \in \{x_1, \dots, x_k\}} f(x)$

謎の仮定
$$\sqrt{(k-1)^5 S_{k-1}} \leq \frac{3L}{4M}$$
の下,
• S_k が大 \implies 関数値減少量が大 〇: $f(\mathbf{x}_k^{\star}) - f(x_0) \leq -\frac{LS_k}{4k}$
• S_k が小 \implies 勾配ノルム が小 〇: $\min_{1 \leq i \leq k} \|\nabla f(\bar{x}_i)\| \leq 3L\sqrt{\frac{S_k}{k^3}}$

●証明は演習7,8

- 謎の仮定の動機は? → 次ページ
- 謎の仮定の正当化は? → 再始動の活用

HB法の解析: 関数値減少量の評価方針 (演習 7 のヒント) 62/77

2種類の不等式

$$f(x_k^*) - f(x_0) \le \frac{1}{2k - 1} \left(-\frac{LS_k}{2} + \frac{M}{6} (k - 1) \sum_{i=0}^{k-2} \|x_{i+1} - x_i\|^3 \right)$$

�� 厄介な
$$\sum_{i} \|x_{i+1} - x_{i}\|^{3}$$
を仮定 $\sqrt{(k-1)^{5}S_{k-1}} \leq \frac{3L}{4M}$ で処理

$$f(x_k^\star) - f(x_0) \le -\frac{LS_k}{4k}$$

HB法の解析: 勾配ノルムの評価方針 (演習8のヒント) 63/77

もう一つの Hessian-free 不等式

$$\left\|\nabla f(\bar{x}_k) - \frac{1}{k} \sum_{i=0}^{k-1} \nabla f(x_i)\right\| \le \frac{M}{2k} \sum_{i=0}^{k-1} \|x_i - \bar{x}_k\|^2$$

V

を使って頑張ると、

.

$$\begin{split} \min_{1 \le i \le k} \|\nabla f(\bar{x}_i)\| \le 2L\sqrt{\frac{S_k}{k^3}} + \frac{M}{6}(k-1)S_{k-1} \\ & \roppose \frac{M}{6}(k-1)S_{k-1} \, \varepsilon \, \text{仮定} \, \sqrt{(k-1)^5 S_{k-1}} \le \frac{3L}{4M} \, \mathfrak{C} \mathfrak{U} \mathfrak{U} \\ & \min_{1 \le i \le k} \|\nabla f(\bar{x}_i)\| \le 3L\sqrt{\frac{S_k}{k^3}} \end{split}$$

再始動 Heavy-ball 法

1:
$$(x_0, x_{-1}) \leftarrow (x_{\text{init}}, x_{\text{init}}), \ k \leftarrow 0$$

2: **loop**

ĸ

3:
$$k \leftarrow k+1$$

4:
$$x_k \leftarrow 2x_{k-1} - x_{k-2} - \frac{1}{L} \nabla f(x_{k-1})$$

5: **if** $\sqrt{k^5 S_k} > \frac{3L}{4M}$:

$$\triangleright \theta_k = 1 \text{ O HB }$$
法

6:
$$(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), \ k \leftarrow 0$$

再始動により常に
$$\sqrt{(k-1)^5 S_{k-1}} \leq rac{3L}{4M}$$
を保証

い
$$\begin{cases} f(x_k^{\star}) - f(x_0) \leq -\frac{LS_k}{4k} \\ \min_{1 \leq i \leq k} \|\nabla f(\bar{x}_i)\| \leq 3L\sqrt{\frac{S_k}{k^3}} \end{cases} & \\ \text{が常に成立 } \mathfrak{S}_k \end{cases}$$

再始動 Heavy-ball 法

1:
$$(x_0, x_{-1}) \leftarrow (x_{\text{init}}, x_{\text{init}}), \ k \leftarrow 0$$

2: **loop**

 $k \leftarrow k + 1$ 3:

4:
$$x_k \leftarrow 2x_{k-1} - x_{k-2} - \frac{1}{L}\nabla f(x_{k-1})$$

$$\triangleright \theta_k = 1 \mathcal{O} HB 法$$

5: If
$$\sqrt{k^\circ} S_k > \frac{1}{4M}$$
:
6: $(x_0, x_{-1}) \leftarrow (x_1^\star, x_1^\star)$,

$$(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), \ k \leftarrow 0$$

- HB 法は、 → → 凸 2 次 関 数 に 対 し て は 最 滴
- 再始動 HB は, 2次関数でないことの悪影響を再始動でリセット

再始動 HB 法の解析: 2種類の評価を組み合わせる

65/77

$$\begin{array}{l} \bullet \quad f(x_k^{\star}) - f(x_0) \leq -\frac{LS_k}{4k} \\ \bullet \quad \min_{1 \leq i \leq k} \|\nabla f(\bar{x}_i)\| \leq 3L \sqrt{\frac{S_k}{k^3}} \end{array} \end{array}$$

•
$$\varepsilon < \min_{1 \le i \le k} \| \nabla f(\bar{x}_i) \|$$
 と B より, $S_k > \left(\frac{\varepsilon}{3L}\right)^2 k^3$

• 第k反復で再始動時,条件 $\sqrt{k^5 S_k} > \frac{3L}{4M}$ より $S_k > \left(\frac{3L}{4M}\right)^2 k^{-5}$

よって,

$$S_k = S_k^{7/8} S_k^{1/8} > \left(\left(\frac{\varepsilon}{3L}\right)^2 k^3 \right)^{7/8} \left(\left(\frac{3L}{4M}\right)^2 k^{-5} \right)^{1/8} = \frac{\varepsilon^{7/4}}{3\sqrt{6L^3}M^{1/4}} k^2$$

再始動 HB 法の解析: 2種類の評価を組み合わせる

65/77

③, **③** を組み合わせると, $\frac{arepsilon^{7/4}}{12\sqrt{6L}M^{1/4}} k \leq f(x_0) - f(x_k^{\star})$

(最後の再始動までの合計反復数) = $\sum k \leq \frac{12\sqrt{6L}M^{1/4}}{\epsilon^{7/4}}\Delta$ 最後の再始動以降の反復数は $\frac{6\sqrt{L\Delta}}{\epsilon} = O(\epsilon^{-1})$ (容易,詳細割愛)

$$\bar{x}_k \coloneqq \frac{1}{k} \sum_{i=0}^{k-1} x_i \, \check{m} \, \| \nabla f(\bar{x}_k) \| \le \varepsilon \, \varepsilon \, \check{n} \, \mathsf{MOT}$$
 満たすまでの合計反復数は
$$\frac{12\sqrt{6L} M^{1/4} \Delta}{\varepsilon^{7/4}} + \mathcal{O}(\varepsilon^{-1}) \qquad \text{以下}$$

66/77

アルゴリズム (再掲)
1:
$$(x_0, x_{-1}) \leftarrow (x_{init}, x_{init}), k \leftarrow 0$$

2: **loop**
3: $k \leftarrow k + 1$
4: $x_k \leftarrow 2x_{k-1} - x_{k-2} - \frac{1}{L} \nabla f(x_{k-1})$ $\triangleright \theta_k = 1 \text{ O HB 法}$
5: **if** $\sqrt{k^5 S_k} > \frac{3L}{4M}$: $\triangleright S_k \coloneqq \sum_{i=0}^{k-1} ||x_{i+1} - x_i||^2$
6: $(x_0, x_{-1}) \leftarrow (x_k^*, x_k^*), k \leftarrow 0$

ここからHB法の発展的話題を紹介

- *L*, *M* が未知の場合
- Universal HB法
- 再始動なしHB法?

最急降下法のときと似た方法で推定したいが, 各反復で*L*の推定値ℓが変わると解析が困難に 😕

方針:

- 不等式チェックによりℓが小さすぎると判明したら再始動
- 再始動時以外には ℓを更新しない
L, M が未知の場合: L の推定

1:
$$(x_0, x_{-1}) \leftarrow (x_{\text{init}}, x_{\text{init}}), \ \ell \leftarrow \ell_{\text{init}}, \ k \leftarrow 0$$

2: **loop**

- 3: $k \leftarrow k+1$
- 4: $x_k \leftarrow 2x_{k-1} x_{k-2} \frac{1}{\ell} \nabla f(x_{k-1})$
- 5: **if** $f(x_k) f(x_{k-1}) > \langle \nabla f(x_{k-1}), x_k x_{k-1} \rangle + \frac{\ell}{2} ||x_k x_{k-1}||^2$:
- 6: $(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), k \leftarrow 0, \ell \leftarrow \alpha \ell$ ▷ $\alpha > 1$: 定数
- 7: else if $\sqrt{k^5 S_k} > rac{3\ell}{4M}$:
- 8: $(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), \quad k \leftarrow 0$
- *L*が既知の場合と同様の計算量保証が成立
- 8行目にℓを減少させるステップを入れても OK ☺

L, M が未知の場合: M の推定

Mの推定値mは毎反復更新しても解析に影響なし

解析で使った不等式

• $f(x_i) - f(x_{i-1}) \le \frac{1}{2} \langle \nabla f(x_i) + \nabla f(x_{i-1}), x_i - x_{i-1} \rangle + \frac{m}{12} \|x_i - x_{i-1}\|^3 \ (1 \le i \le k)$ • $\|\nabla f(\bar{x}_{i+1})\| \le \frac{\ell}{i+1} \|x_{i+1} - x_i\| + \frac{m}{4} iS_i \qquad (1 \le i \le k)$

を満たす最小の*m*を第*k*反復の推定値*m_k*とする

$$m_{k} \coloneqq \max\left\{ \max_{1 \le i \le k} \frac{12}{\|x_{i} - x_{i-1}\|^{3}} \left(f(x_{i}) - f(x_{i-1}) - \frac{1}{2} \langle \nabla f(x_{i}) + \nabla f(x_{i-1}), x_{i} - x_{i-1} \rangle \right), \\ \max_{1 \le i \le k} \frac{4}{iS_{i}} \left(\|\nabla f(\bar{x}_{i+1})\| - \frac{\ell}{i+1} \left\| x_{i} - x_{i-1} - \frac{1}{\ell} \nabla f(x_{i}) \right\| \right) \right\}$$

L, M が未知の場合: M の推定

Mの推定値mは毎反復更新しても解析に影響なし

解析で使った不等式

- $f(x_i) f(x_{i-1}) \leq \frac{1}{2} \langle \nabla f(x_i) + \nabla f(x_{i-1}), x_i x_{i-1} \rangle + \frac{m}{12} \|x_i x_{i-1}\|^3 \ (1 \leq i \leq k)$
- $\|\nabla f(\bar{x}_{i+1})\| \le \frac{\ell}{i+1} \|x_{i+1} x_i\| + \frac{m}{4} iS_i$ $(1 \le i \le k)$

を満たす最小の*m*を第*k*反復の推定値*m_k*とする

Mの推定を容易にするための鍵:

- Hessian-free 不等式の利用
- *x_k*の計算に*m_k*不使用

L, M が未知の場合:完全版アルゴリズム

1:
$$(x_0, x_{-1}) \leftarrow (x_{\text{init}}, x_{\text{init}}), \ \ell \leftarrow \ell_{\text{init}}, \ k \leftarrow 0$$

2: **loop**

- 3: $k \leftarrow k+1$
- 4: $x_k \leftarrow 2x_{k-1} x_{k-2} \frac{1}{\ell} \nabla f(x_{k-1})$
- 5: *M*の推定値*m_k*の計算
- 6: **if** $f(x_k) f(x_{k-1}) > \langle \nabla f(x_{k-1}), x_k x_{k-1} \rangle + \frac{\ell}{2} ||x_k x_{k-1}||^2$:

7:
$$(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), \quad k \leftarrow 0, \quad \ell \leftarrow \alpha \ell$$

8: else if
$$\sqrt{k^5 S_k} > rac{3\ell}{4m_k}$$
 :

9:
$$(x_0, x_{-1}) \leftarrow (x_k^\star, x_k^\star), \ k \leftarrow 0, \ \ell \leftarrow \beta \ell$$

- ∇f が Lipschitz 連続なら,最急降下法で $O(\varepsilon^{-2})$
- $\nabla f, \nabla^2 f$ が Lipschitz 連続なら,再始動 HB 法で $O(\epsilon^{-7/4})$

仮定によってアルゴリズムを使い分けるのは大変 😕

疑問:単一のアルゴリズムで両方の保証を達成可能か?

- ∇f が Lipschitz 連続なら,最急降下法で $O(\varepsilon^{-2})$
- $\nabla f, \nabla^2 f$ が Lipschitz 連続なら,再始動 HB 法で $O(\epsilon^{-7/4})$

仮定によってアルゴリズムを使い分けるのは大変 😣

疑問:単一のアルゴリズムで両方の保証を達成可能か?

- 可能! [Marumo and Takeda, 2024b]
- しかも「中間」の仮定にも対応可能

仮定: $\nabla^2 f$ の Hölder 連続性 $\left\|\nabla^2 f(x) - \nabla^2 f(y)\right\| \le H_{\nu} \|x - y\|^{\nu} \qquad (\forall x, y \in \mathbb{R}^d)$ ※ ∇f が *L*-Lipschitz 連続 \implies $H_0 < 2L$ 0 [0, 1]ν $4 + 3\nu$ p in $O(\varepsilon^{-p})$ $\mathbf{2}$ $\overline{2+2\nu}$ 最急降下法 [M-Takeda, 2024b] [Li and Lin, 2022] より強い仮定、より良い計算量

Hölder 連続性を使う計算量解析の例: [Cartis et al., 2017, 2019; Devolder et al., 2014; Dvurechensky, 2017; Ghadimi et al., 2019; Grapiglia and Nesterov, 2017, 2019, 2020; Lan, 2015; Nesterov, 2015]...

Universal 再始動 HB 法の反復数 [Marumo and Takeda, 2024b] $\inf_{\nu \in [0,1]} \left\{ 91\sqrt{L}H_{\nu}^{\frac{1}{2+2\nu}} \varepsilon^{-\frac{4+3\nu}{2+2\nu}} \Delta \right\} + O(\varepsilon^{-1})$

ポイント:

- $H_{\nu} = +\infty \ x \le \nu \ x$ or V
- $H_0 \leq 2L$ より,計算量は最急降下法の $O\left(\frac{L\Delta}{\varepsilon^2}\right)$ 以下
- $\nu = 1$ とすると $O\left(\frac{\sqrt{L}M^{1/4}\Delta}{\varepsilon^{7/4}}\right)$

再始動なしHB法?

75/77

- HB 法の O(*ε*^{-7/4}) 収束を示すのに再始動を用いた
- O(ε^{-7/4}) やÕ(ε^{-7/4})の他手法は再始動以上に複雑なことをする
- 凸の場合,加速勾配法は再始動なしで最良の計算量を達成

疑問:再始動なしのシンプルなHB法でも $O(\epsilon^{-7/4})$ は達成可能か?

再始動なしHB法?

- HB 法の $O(\epsilon^{-7/4})$ 収束を示すのに再始動を用いた
- O(ε^{-7/4}) やÕ(ε^{-7/4})の他手法は再始動以上に複雑なことをする
- 凸の場合,加速勾配法は再始動なしで最良の計算量を達成

疑問:再始動なしのシンプルなHB法でも $O(\epsilon^{-7/4})$ は達成可能か?

- HB 微分方程式は再始動なしで O(ε^{-7/4}) 収束 [Okamura et al., 2024]
- 微分方程式を離散化して得られる HB 法の解析は
 完全解決していない

HB 微分方程式の収束解析 [Okamura et al., 2024]

T > 0: 定数 • $\boldsymbol{\gamma} \coloneqq (3M)^{\frac{2}{7}} \left(\frac{\Delta}{T}\right)^{\frac{1}{7}}$ • $\ddot{x}(t) = -\gamma \dot{x}(t) - \nabla f(x(t)),$ $x(0) = x_0,$ $\dot{x}(0) = 0$ • $\bar{x}(t) \coloneqq \int_{0}^{t} \frac{\gamma e^{\gamma s}}{e^{\gamma t} - 1} x(s) ds$ $\implies \qquad \min_{0 \le t \le T} \|\nabla f(\bar{x}(t))\| \le \frac{7}{6} (3M)^{\frac{1}{7}} \left(\frac{\Delta}{T}\right)^{\frac{4}{7}} + \mathcal{O}\left(T^{-\frac{10}{7}}\right)$

• 再始動 HB では $\gamma = 0$ ($\theta_k = 1$).上は摩擦 ($\gamma > 0$) で再始動を代替

76/77

離散化して O(ε^{-7/4})の HB 法を得るには現状追加の仮定が必要

最急降下法,CRN,再始動HB,Levenberg-Marquardt法 企演習9

アルゴリズム設計の定石

- fの近似 $\overline{f_k}$ を最小化する点を x_{k+1} とする
- **A** $\bar{f}_k(x_k) = f(x_k), \ \nabla \bar{f}_k(x_k) = \nabla f(x_k)$ **B** $f(x_{k+1}) \leq \bar{f}_k(x_{k+1})$

他: 強凸な \bar{f}_k ,強めの正則化,慣性,Lipschitz 定数の推定,…

計算量解析の定石

- $||x_{k+1} x_k||$ が大 \implies **関数値減少量**が大
- $\|x_{k+1} x_k\|$ が小 \implies 勾配ノルム が小

他: Lipschitz 連続性から ³ を示す,(最小値) ≤ (平均値),…



- S. Adachi, S. Iwata, Y. Nakatsukasa, and A. Takeda. Solving the trust-region subproblem by a generalized eigenvalue problem. <u>SIAM Journal on Optimization</u>, 27(1):269–291, 2017. URL https://doi.org/10.1137/16M1058200.
- N. Agarwal, Z. Allen-Zhu, B. Bullins, E. Hazan, and T. Ma. Finding approximate local minima faster than gradient descent. In <u>Proceedings of the 49th Annual ACM SIGACT Symposium on Theory of Computing</u>, STOC 2017, pages 1195–1199, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450345286. URL https://doi.org/10.1145/3055399.3055464.
- Z. Allen-Zhu and Y. Li. NEON2: Finding local minima via first-order oracles. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, <u>Advances in Neural Information Processing Systems</u>, volume 31. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/file/d4b2aeb2453bdadaa45cbe9882ffefcf-Paper.pdf.
- L. Armijo. Minimization of functions having Lipschitz continuous first partial derivatives. <u>Pacific Journal of</u> <u>Mathematics</u>, 16(1):1-3, 1966. URL https://dx.doi.org/10.2140/pjm.1966.16.1.
- Y. Carmon and J. C. Duchi. Analysis of Krylov subspace solutions of regularized non-convex quadratic problems. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, <u>Advances in Neural</u> <u>Information Processing Systems</u>, volume 31. Curran Associates, Inc., 2018. URL https: //proceedings.neurips.cc/paper_files/paper/2018/file/349f36aa789af083b8e26839bd498af9-Paper.pdf.
- Y. Carmon and J. C. Duchi. First-order methods for nonconvex quadratic minimization. <u>SIAM Review</u>, 62(2):395–436, 2020. URL https://doi.org/10.1137/20M1321759.



- Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. "Convex until proven guilty": Dimension-free acceleration of gradient descent on non-convex functions. In D. Precup and Y. W. Teh, editors, <u>Proceedings of the 34th</u> <u>International Conference on Machine Learning</u>, volume 70 of <u>Proceedings of Machine Learning Research</u>, pages 654–663. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/carmon17a.html.
- Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Accelerated methods for nonconvex optimization. <u>SIAM Journal on</u> Optimization, 28(2):1751–1772, 2018. URL https://doi.org/10.1137/17M1114296.
- Y. Carmon, J. C. Duchi, O. Hinder, and A. Sidford. Lower bounds for finding stationary points I. <u>Mathematical</u> <u>Programming</u>, 184(1):71–120, 2020. URL https://doi.org/10.1007/s10107-019-01406-y.
- C. Cartis, N. I. M. Gould, and P. L. Toint. Adaptive cubic regularisation methods for unconstrained optimization. Part I: Motivation, convergence and numerical results. <u>Mathematical Programming</u>, 127(2):245–295, 2011. URL https://doi.org/10.1007/s10107-009-0286-5.
- C. Cartis, N. I. M. Gould, and P. L. Toint. On the oracle complexity of first-order and derivative-free algorithms for smooth nonconvex minimization. <u>SIAM Journal on Optimization</u>, 22(1):66–86, 2012. URL https://doi.org/10.1137/100812276.
- C. Cartis, N. I. M. Gould, and P. L. Toint. Worst-case evaluation complexity of regularization methods for smooth unconstrained optimization using Hölder continuous gradients. <u>Optimization Methods and Software</u>, 32(6): 1273–1298, 2017. URL https://doi.org/10.1080/10556788.2016.1268136.

参考文献 III

- C. Cartis, N. I. M. Gould, and P. L. Toint. Universal regularization methods: Varying the power, the smoothness and the accuracy. SIAM Journal on Optimization, 29(1):595–615, 2019. URL https://doi.org/10.1137/16M1106316.
- A. Cristofari, T. Dehghan Niri, and S. Lucidi. On global minimizers of quadratic functions with cubic regularization. Optimization Letters, 13(6):1269–1283, 2019. URL https://doi.org/10.1007/s11590-018-1316-0.
- O. Devolder, F. Glineur, and Y. Nesterov. First-order methods of smooth convex optimization with inexact oracle. Mathematical Programming, 146(1):37–75, 2014. URL https://doi.org/10.1007/s10107-013-0677-5.
- N. Doikov, E. M. Chayti, and M. Jaggi. Second-order optimization with lazy Hessians. In A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, <u>Proceedings of the 40th International Conference on</u> <u>Machine Learning</u>, volume 202 of <u>Proceedings of Machine Learning Research</u>, pages 8138–8161. PMLR, 23–29 Jul 2023. URL https://proceedings.mlr.press/v202/doikov23a.html.
- P. Dvurechensky. Gradient method with inexact oracle for composite non-convex optimization. <u>arXiv preprint</u>, 2017. URL https://arxiv.org/abs/1703.09180.
- Y. Gao, M.-C. Yue, and M. Ng. Approximate secular equations for the cubic regularization subproblem. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh, editors, <u>Advances in Neural Information Processing</u> <u>Systems</u>, volume 35, pages 14250–14260. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/ paper_files/paper/2022/file/5be69a584901a26c521c2b51e40a4c20-Paper-Conference.pdf.
- S. Ghadimi, G. Lan, and H. Zhang. Generalized uniformly optimal methods for nonlinear programming. Journal of Scientific Computing, 79(3):1854–1881, 2019. URL https://doi.org/10.1007/s10915-019-00915-4.



- B. Goujaud, A. Taylor, and A. Dieuleveut. Provable non-accelerations of the heavy-ball method. <u>arXiv preprint</u> <u>arXiv:2307.11291</u>, 2023.
- N. I. M. Gould and V. Simoncini. Error estimates for iterative algorithms for minimizing regularized quadratic subproblems. <u>Optimization Methods and Software</u>, 35(2):304–328, 2020. URL https://doi.org/10.1080/10556788.2019.1670177.
- G. N. Grapiglia and Y. Nesterov. Regularized Newton methods for minimizing functions with Hölder continuous Hessians. SIAM Journal on Optimization, 27(1):478–506, 2017. URL https://doi.org/10.1137/16M1087801.
- G. N. Grapiglia and Y. Nesterov. Accelerated regularized Newton methods for minimizing composite convex functions. SIAM Journal on Optimization, 29(1):77–99, 2019. URL https://doi.org/10.1137/17M1142077.
- G. N. Grapiglia and Y. Nesterov. Tensor methods for minimizing convex functions with Hölder continuous higher-order derivatives. SIAM Journal on Optimization, 30(4):2750–2779, 2020. URL https://doi.org/10.1137/19M1259432.
- A. Griewank. The modification of Newton's method for unconstrained optimization by bounding cubic terms. Technical report, Technical report NA/12, Department of Applied Mathematics and Theoretical Physics, University of Cambridge, 1981.
- X. Jia, X. Liang, C. Shen, and L.-H. Zhang. Solving the cubic regularization model by a nested restarting Lanczos method. <u>SIAM Journal on Matrix Analysis and Applications</u>, 43(2):812–839, 2022. URL https://doi.org/10.1137/21M1436324.

参考文献 V

- R. Jiang, M.-C. Yue, and Z. Zhou. An accelerated first-order method with complexity analysis for solving cubic regularization subproblems. <u>Computational Optimization and Applications</u>, 79(2):471–506, 2021. URL https://doi.org/10.1007/s10589-021-00274-7.
- C. Jin, P. Netrapalli, and M. I. Jordan. Accelerated gradient descent escapes saddle points faster than gradient descent. In S. Bubeck, V. Perchet, and P. Rigollet, editors, <u>Proceedings of the 31st Conference On Learning Theory</u>, volume 75 of <u>Proceedings of Machine Learning Research</u>, pages 1042–1085. PMLR, 06–09 Jul 2018. URL https://proceedings.mlr.press/v75/jin18a.html.
- G. Lan. Bundle-level type methods uniformly optimal for smooth and nonsmooth convex optimization. <u>Mathematical</u> <u>Programming</u>, 149(1):1–45, 2015. URL https://doi.org/10.1007/s10107-013-0737-x.
- H. Li and Z. Lin. Restarted nonconvex accelerated gradient descent: No more polylogarithmic factor in the O(ε^{-7/4}) complexity. In K. Chaudhuri, S. Jegelka, L. Song, C. Szepesvari, G. Niu, and S. Sabato, editors, Proceedings of the <u>39th International Conference on Machine Learning</u>, volume 162 of Proceedings of Machine Learning Research, pages 12901–12916. PMLR, 17–23 Jul 2022. URL https://proceedings.mlr.press/v162/1i22o.html.
- H. Li and Z. Lin. Restarted nonconvex accelerated gradient descent: No more polylogarithmic factor in the O(ε^{-7/4}) complexity. Journal of Machine Learning Research, 24(157):1–37, 2023. URL http://imlr.org/papers/v24/22-0522.html.
- F. Lieder. Solving large-scale cubic regularization by a generalized eigenvalue problem. <u>SIAM Journal on Optimization</u>, 30(4):3345–3358, 2020. URL https://doi.org/10.1137/19M1291388.



- 83/77
- N. Marumo and A. Takeda. Parameter-free accelerated gradient descent for nonconvex minimization. <u>SIAM Journal on</u> <u>Optimization</u>, 34(2):2093–2120, 2024a. URL https://doi.org/10.1137/22M1540934.
- N. Marumo and A. Takeda. Universal heavy-ball method for nonconvex optimization under Hölder continuous Hessians. Mathematical Programming, 2024b. URL https://doi.org/10.1007/s10107-024-02100-4.
- Y. Nesterov. A method for solving a convex programming problem with convergence rate $O(1/k^2)$. Soviet Mathematics Doklady, 269(3):372–376, 1983.
- Y. Nesterov. <u>Introductory Lectures on Convex Optimization: A Basic Course</u>. Springer, New York, 2004. URL https://doi.org/10.1007/978-1-4419-8853-9.
- Y. Nesterov. Universal gradient methods for convex optimization problems. <u>Mathematical Programming</u>, 152(1): 381–404, 2015. URL https://doi.org/10.1007/s10107-014-0790-0.
- Y. Nesterov and B. T. Polyak. Cubic regularization of Newton method and its global performance. <u>Mathematical</u> Programming, 108(1):177–205, 2006. URL https://doi.org/10.1007/s10107-006-0706-8.
- B. O'Donoghue and E. Candès. Adaptive restart for accelerated gradient schemes. <u>Foundations of Computational</u> <u>Mathematics</u>, 15(3):715–732, 2015. URL https://doi.org/10.1007/s10208-013-9150-3.
- K. Okamura, N. Marumo, and A. Takeda. Primitive heavy-ball dynamics achieves $O(\varepsilon^{-7/4})$ convergence for nonconvex optimization. arXiv preprint arXiv:2406.06100, 2024.

参考文献 VII

- B. T. Polyak. Some methods of speeding up the convergence of iteration methods. <u>USSR Computational Mathematics</u> <u>and Mathematical Physics</u>, 4(5):1–17, 1964. ISSN 0041-5553. URL <u>https://doi.org/10.1016/0041-5553(64)90137-5</u>.
- H. H. Rosenbrock. An automatic method for finding the greatest or least value of a function. <u>The Computer Journal</u>, 3 (3):175–184, 01 1960. ISSN 0010-4620. URL https://doi.org/10.1093/comjnl/3.3.175.
- C. W. Royer, M. O'Neill, and S. J. Wright. A Newton-CG algorithm with complexity guarantees for smooth unconstrained optimization. <u>Mathematical Programming</u>, 180(1):451–488, 2020. URL https://doi.org/10.1007/s10107-019-01362-7.
- Y. Xu, R. Jin, and T. Yang. NEON+: Accelerated gradient methods for extracting negative curvature for non-convex optimization. arXiv preprint, 2017. URL https://arxiv.org/abs/1712.01033.