

Heavy-Ball Method

$$\min_{x \in \mathbb{R}^n} f(x)$$

f -convex, Lipschitz gradient.

Algorithm (Heavy-Ball Method)

Init: $x_0 \in \mathbb{R}^n$, $L > 0$, momentum parameter $0 \leq \beta \leq 1$
Fix $K \geq 1$ the number of iters. $S_0 = 0 \in \mathbb{R}^n$

For $k=0 \dots K-1$:

1. Aggregate: $S_{k+1} = \beta S_k + f'(x_k)$

2. Perform a step: $x_{k+1} = x_k - \frac{1}{L} S_{k+1}$.

Return a point \bar{x} with the best desired acc. measure.

The Gradient Method $\beta := 0$.

• $f(x_k) - f^* \leq \frac{2LR^2}{k+4}$, $R = \|x_0 - x^*\|$.

• $\min_{0 \leq i \leq k-1} \|f'(x_i)\| \leq \frac{4LR}{k}$, $k \geq 1$.

• Strongly convex, smooth: $\mu I \preceq f''(x) \preceq LI \quad \forall x$

Linear rate:

$$f(x_k) - f^* \leq \exp\left(-k \cdot \frac{\mu}{L}\right) \cdot (f(x_0) - f^*)$$

To find: $f(x_k) - f^* \leq \varepsilon \Rightarrow K = \frac{L}{\mu} \ln \frac{f(x_0) - f^*}{\varepsilon}$.

• $\beta = 1$.

• Form average points $\bar{x}_k = \frac{1}{k} \sum_{i=0}^{k-1} x_i$

• Take the minimal gradient:

$$g_k = \min \{ \|f'(\bar{x}_1)\|, \dots, \|f'(\bar{x}_k)\| \}.$$

• The rate of HB (f is convex quadratic):

$$g_k^2 \leq \frac{12L(f(x_0) - f^*)}{k(k+1)} \quad g_k = O\left(\frac{1}{k}\right)$$

• Assume $f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$, $\mu = \lambda_{\min}(A) > 0$.

$$f(x) - f^* \leq \frac{1}{2\mu} \|f'(x)\|^2 \quad \forall x$$

$$f(\bar{x}_k^*) - f^* \leq \frac{6L}{\mu} \cdot \frac{1}{k^2} (f(x_0) - f^*)$$

We can use restarts

$$K = \sqrt{\frac{12L}{\mu}}. \text{ Then}$$

$$f(\bar{x}_K^*) - f^* \leq \frac{1}{2} (f(x_0) - f^*)$$

\Rightarrow We need only $T = \log_2 \frac{f(x_0) - f^*}{\varepsilon}$ to
get $f(\bar{x}) - f^* \leq \varepsilon$.

Theorem The complexity of the HB method
with restarts is

$$\sqrt{\frac{12L}{\mu}} \cdot \log_2 \frac{f(x_0) - f^*}{\varepsilon}$$

oracle calls (matrix-vector products).

- $\beta = 0$ Gradient Method $\hat{O}\left(\frac{L}{\mu}\right)$
- $\beta = 1$ with restarts : $\hat{O}\left(\sqrt{\frac{L}{\mu}}\right) \rightarrow$ optimal
- $\beta \approx 1 - \sqrt{\frac{\mu}{L}}$

• In practice : $\beta \approx 0.99$

• Only for Quadratic Functions.

Best Method: Conjugate Gradient Method.

• Nesterov's Fast Gradient Method

$$f(x) = \frac{1}{2} \langle Ax, x \rangle - \langle b, x \rangle$$

$$f'(x) = Ax - b, \quad f'(x^*) = Ax^* - b = 0$$

$$f'(x) = A(x - x^*)$$

$$x_{u+1} = x_u - \alpha f'(x_u) + \beta (x_u - x_{u-1})$$

$$\underbrace{x_{u+1} - x^*}_{r_{u+1}} = \underbrace{x_u - x^*}_{r_u} - \alpha A \underbrace{(x_u - x^*)}_{r_u} + \beta \left(\underbrace{x_u - x^*}_{r_u} - \underbrace{x_{u-1} + x^*}_{r_{u-1}} \right)$$

$$\mathbb{R}^{2u} \ni \begin{bmatrix} r_{u+1} \\ r_u \end{bmatrix} = \begin{bmatrix} I(1+\beta) - \alpha A & -\beta I \\ I & 0 \end{bmatrix} \begin{bmatrix} r_u \\ r_{u-1} \end{bmatrix}$$

$$\left\| \begin{bmatrix} r_{u+1} \\ r_u \end{bmatrix} \right\| \leq \left\| C \cdot \begin{bmatrix} r_u \\ r_{u-1} \end{bmatrix} \right\|$$

$$\boxed{\beta = 0 \quad \alpha = \frac{1}{L} \quad \|C\| \leq \left(1 - \frac{\mu}{L}\right)}$$

\rightarrow min, $\alpha, \beta \Rightarrow$ get the optimal rate.

Lower Bounds for Smooth Convex

Optimization

$$GM: f(x_k) - f^* \leq O\left(\frac{LR^2}{k}\right)$$

Theorem let $L > 0$. For any first-order optimization algorithm running for $k \geq 1$, \exists convex function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ with $n \geq 2k+1$, it has Lipschitz grad. with const $L > 0$ such that:

$$f(x_k) - f^* \geq \frac{3LR^2}{16(k+1)^2}.$$

Simplify everything:

- $x_0 = 0$, $p(x) = f(x - x_0)$
- $L = \text{const.}$ ($L = 4$)
- Linear Span Methods

Assume $x_{k+1} \in L_{k+1} = \text{span}\{f'(x_0), \dots, f'(x_k)\}$.

$$L_0 \subseteq L_1 \subseteq L_2 \subseteq L_3 \subseteq \dots \quad \text{Krylov Subspaces}$$

$$L_0 = \{0\}$$

$$L_1 = \text{span}\{f'(x_0)\}$$

$$L_2 = \text{span}\{f'(x_0), f'(x_1)\}, \dots$$

Fix $k \geq 1$.

$$f_k(x) = \frac{1}{2} \left[\sum_{i=1}^{k-1} (x^{(i)} - x^{(i+1)})^2 + \sum_{i=k}^n (x^{(i)})^2 \right] - \langle b, x \rangle$$

$$= \frac{1}{2} \left[(x^{(1)} - x^{(2)})^2 + (x^{(2)} - x^{(3)})^2 + \dots + (x^{(k-1)} - x^{(k)})^2 + (x^{(k)})^2 + \dots + (x^{(n)})^2 \right] - \langle b, x \rangle$$

$$= \frac{1}{2} \|C_k x\|^2 - \langle b, x \rangle = \frac{1}{2} \langle \underbrace{C_k^T C_k}_{A_k} x, x \rangle - \langle b, x \rangle$$

$$C_k = \begin{bmatrix} \mathbb{1}_n & 0 \\ 0 & I_{n-k} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$D_k = \begin{bmatrix} 1 & -1 & & & 0 \\ & 1 & -1 & & \\ & & 1 & -1 & \\ & & & 1 & -1 \\ 0 & & & & 1 \end{bmatrix} \in \mathbb{R}^{k \times k}$$

$$D_k^T D_k = \begin{bmatrix} 1 & & & & \\ -1 & 1 & & & \\ & -1 & 1 & & \\ & & -1 & 1 & \\ & & & -1 & 1 \\ 0 & & & & 0 \end{bmatrix}$$

$$A_k = \left[\begin{array}{ccccc|c} 1 & -1 & 0 & 0 & 0 & 0 \\ -1 & 2 & -1 & & & \\ 0 & -1 & 2 & -1 & & \\ 0 & & -1 & 2 & -1 & \\ 0 & 0 & & -1 & 2 & \\ \hline & & & & & I \end{array} \right]$$

$$f_k(x) = \frac{1}{2} \langle A_k x, x \rangle - \langle b, x \rangle$$

• Simple enough, $b = e_1$

• Graphs



Laplacian of this graph

$$L = D - C$$

$$\begin{aligned} \langle f_n''(x)h, h \rangle &= \underbrace{\sum_{i=1}^{k-1} (h^{(i)} - h^{(i+1)})^2 + \sum_{i=k}^n (h^{(i)})^2}_{\geq 0} \\ &\leq \sum_{i=1}^{k-1} 2(h^{(i)})^2 + 2(h^{(i+1)})^2 + \sum_{i=k}^n (h^{(i)})^2 \leq \\ &\leq 4 \cdot \|h\|_2^2 \Rightarrow \boxed{L=4} \end{aligned}$$

$$x_n^* = \underset{x \in \mathbb{Q}^n}{\text{argmin}} f_n(x)$$

$$\left\{ \begin{array}{ll} x^{(1)} - x^{(2)} - \mathbf{1} = 0 & i=1 \\ 2x^{(i)} - x^{(i-1)} - x^{(i+1)} = 0 & 2 \leq i \leq k \\ 2x^{(k)} - x^{(k-1)} = 0 & i=k \\ x^{(i)} = 0 & k < i \leq n \end{array} \right. \quad b = \mathbf{e}_1$$

$$\boxed{\begin{aligned} (x_n^*)^{(1)} &= k \\ (x_n^*)^{(2)} &= k-1 \\ &\dots \\ (x_n^*)^{(k)} &= 1 \\ (x_n^*)^{(i)} &= 0 \end{aligned}}$$

for $k < i \leq n$.

$$f_n^* = \frac{1}{2} \langle A_n x_n^*, x_n^* \rangle - \langle \mathbf{e}_1, x_n^* \rangle = -\frac{1}{2} \langle \mathbf{e}_1, x_n^* \rangle = -\frac{k}{2}.$$

$$\|x_0 - x_n^*\|^2 = \|x_n^*\|^2 = \sum_{i=1}^k i^2 = \frac{k(k+1)(2k+1)}{6} \leq \frac{(k+1)^3}{6}.$$

Proposition Assume $\{x_u\}$ that satisfies: $x_u \in L_k$. Then

$$L_k \subseteq \mathbb{R}^{n,k} = \{x \in \mathbb{R}^n \mid x^{(i)} = 0 \quad i \geq k+1\}.$$

Proof

$$L_0 = \{0\} \in \mathbb{R}^{n,0}$$

$$L_1 = \text{span}\{f'(x_0)\} = \text{span}\{-e_1\} \in \mathbb{R}^{n,1}$$

By induction, $f'(x) = A_k x - b$, $x \in \mathbb{R}^{n,k}$
 $\mathbb{R}^{n,k+1}$ for tridiag. matrix.

Proposition: $f_u(x)$, $f_{u+p}(x)$, $p \geq 0$

$$\underline{f_u(x)} \equiv \underline{f_{u+p}(x)} \quad \forall x \in \mathbb{R}^{n,k}.$$

To prove the lower bound: $f(x) \equiv f_{2k+1}(x)$

$$f(x_u) = f_{2k+1}(x_u) = f_u(x_u) \geq f_u^* = -\frac{k}{2}.$$

But

$$f^* = f_{2k+1}^*(x_{2k+1}^*) = -\frac{2k+1}{2}. \quad \blacksquare$$