



EECS 559

Optimization Methods for SIPML

Lecture 4 – Smooth Unconstrained Optimization

Instructor: Prof. Qing Qu (qingqu@umich.edu)

Lecture Agenda

- Basics of Iterative Methods
- (Accelerated) Gradient Descent for Smooth Problems
- Linesearch Methods for Step Size
- Exploiting Function Properties for Faster Convergence
- Newton & Quasi-Newton Method

1/27/25

2



Smooth Unconstrained Minimization

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathcal{C}.$$

- Suppose $f(\mathbf{x})$ is continuously differentiable over \mathcal{C}
- Consider unconstrained problem $\mathcal{C} = \mathbb{R}^n$
- Solve the problem via *iterative optimization methods*, producing a sequence of points

$$\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k, \dots$$

starting from an initialization \mathbf{x}_0 .

1/27/25

3



Smooth Unconstrained Minimization

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathbb{R}^n.$$

Goal: generate a sequence $\{\mathbf{x}_k\}_{k \geq 1}$ from an initialization \mathbf{x}_0 , such that it quickly converges to a minimizer \mathbf{x}_* of f .

$$\text{complexity} = \text{per iter cost} \times \#\text{of iterations.}$$

1/27/25

4



Smooth Unconstrained Minimization

$$\text{complexity} = \text{per iter cost} \times \#\text{of iterations.}$$

- Per iteration cost.** How much computation it takes to generate the next point \mathbf{x}_{k+1} from \mathbf{x}_k ;
- Convergence rate.** How quickly the sequence $\{\mathbf{x}_k\}_{k \geq 1}$ converges, measured by

Distance to a minimizer	$\ \mathbf{x}_k - \mathbf{x}_*\ _2$
Sub-optimality in objective	$ f(\mathbf{x}_k) - f(\mathbf{x}_*) $
Gradient	$\ \nabla f(\mathbf{x}_k)\ _2$

1/27/25

5



Iterative Descent Algorithms

Starting with an initialization \mathbf{x}_0 , construct a sequence $\{\mathbf{x}_k\}_{k \geq 1}$ such that

$$f(\mathbf{x}_{k+1}) < f(\mathbf{x}_k), \quad k = 0, 1, \dots$$

- In each iteration, search in a **descent direction**:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \underbrace{\tau_k}_{\text{stepsize}} \cdot \underbrace{\mathbf{d}_k}_{\text{descent direction}}$$
- The direction \mathbf{d} is said to be a descent direction at \mathbf{x} if

$$f'(\mathbf{x}; \mathbf{d}) := \lim_{\tau \searrow 0} \frac{f(\mathbf{x} + \tau \mathbf{d}) - f(\mathbf{x})}{\tau} = \nabla f(\mathbf{x})^\top \mathbf{d} < 0.$$

1/27/25

6



Lecture Agenda

- Basics of Iterative Methods
- (Accelerated) Gradient Descent for Smooth Problems**
- Linesearch Methods for Step Size
- Exploiting Function Properties for Faster Convergence
- Newton & Quasi-Newton Method

1/27/25

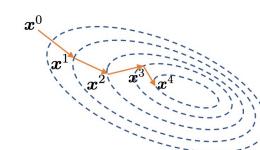
7



Gradient Descent (GD)

One of the most important examples: **gradient descent**

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k)$$



- can be traced back to Augustin Louis Cauchy's 1847 ...

1/27/25

8



Gradient Descent (GD)

One of the most important examples: **gradient descent**

$$\boxed{\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k).}$$

- Descent direction: $\mathbf{d}_k = -\nabla f(\mathbf{x}_k)$
- a.k.a., *steepest descent*, since from Cauchy-Schwarz inequality
 $\arg \min_{\|\mathbf{d}\|_2 \leq 1} f'(\mathbf{x}; \mathbf{d}) = \arg \min_{\|\mathbf{d}\|_2 \leq 1} \nabla f(\mathbf{x})^\top \mathbf{d} = -\nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|_2$.

1/27/25

9

 UNIVERSITY OF MICHIGAN

Gradient Descent (GD)

One of the most important examples: **gradient descent**

$$\boxed{\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k).}$$

- the scalar $\tau_k > 0$ is the *step size*:
 - can be determined *analytically* from the property of the function f
 - numerically* computed via a *line search*:

$$\tau_k = \arg \min_{t \geq 0} f(\mathbf{x}_k - t \cdot \nabla f(\mathbf{x}_k))$$

10

 UNIVERSITY OF MICHIGAN

Convergence of GD

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

Assumptions:

- $f : \mathbb{R}^n \mapsto \mathbb{R}$ is **convex** and **smooth**;
- The gradient $\nabla f(\mathbf{x})$ is **L -Lipschitz** with
 $\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}')\|_2 \leq L \cdot \|\mathbf{x} - \mathbf{x}'\|_2, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathbb{R}^n.$

1/27/25

11

 UNIVERSITY OF MICHIGAN

Convergence of GD

Lemma. (Sufficient function value decrease of GD)
 With step size $\tau_k = 1/L$, the gradient descent iterates $\{\mathbf{x}_k\}_{k \geq 1}$ generated by

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k)$$

satisfy the following:

$$\begin{aligned} f(\mathbf{x}_{k+1}) &\leq f(\mathbf{x}_k) - \frac{1}{2L} \|\nabla f(\mathbf{x}_k)\|_2^2 \\ &< f(\mathbf{x}_k) \quad (\|\nabla f(\mathbf{x}_k)\|_2 \neq 0) \end{aligned}$$

12

 UNIVERSITY OF MICHIGAN

Convergence of GD

Proof. Because ∇f is Lipschitz, by 1st-order Taylor approximation

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2, \quad \forall \mathbf{x}', \mathbf{x}$$

Now, let $\mathbf{x}' = \mathbf{x} - \tau \nabla f(\mathbf{x})$

$$\begin{aligned} f(\mathbf{x}') &\leq f(\mathbf{x}) - \tau \|\nabla f(\mathbf{x})\|_2^2 + \frac{\tau^2 L}{2} \|\nabla f(\mathbf{x})\|_2^2 \\ &= f(\mathbf{x}) - \tau \left(1 - \frac{\tau L}{2}\right) \|\nabla f(\mathbf{x})\|_2^2 \end{aligned}$$

1/27/25

13

 UNIVERSITY OF MICHIGAN

Convergence of GD

Proof. Because ∇f is Lipschitz, by 1st-order Taylor approximation

$$f(\mathbf{x}') \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{x}' - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x}' - \mathbf{x}\|_2^2, \quad \forall \mathbf{x}', \mathbf{x}$$

Now, let $\mathbf{x}' = \mathbf{x} - \tau \nabla f(\mathbf{x})$

$$\begin{aligned} f(\mathbf{x}') &\leq f(\mathbf{x}) - \tau \|\nabla f(\mathbf{x})\|_2^2 + \frac{\tau^2 L}{2} \|\nabla f(\mathbf{x})\|_2^2 \\ &= f(\mathbf{x}) - \tau \left(1 - \frac{\tau L}{2}\right) \|\nabla f(\mathbf{x})\|_2^2 \\ &\leq f(\mathbf{x}) - \frac{1}{L} (1 - 1/2) \|\nabla f(\mathbf{x})\|_2^2 = f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned}$$

14

 UNIVERSITY OF MICHIGAN

Convergence of GD

Theorem. ($O(1/k)$ sublinear convergence of GD)

Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex and smooth with its gradient L -Lipschitz. If we run GD for k iterations with a fixed step size $\tau = 1/L$, it will yield a sequence $\{\mathbf{x}_k\}_{k \geq 1}$ such that

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{L}{2k} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 = O\left(\frac{1}{k}\right),$$

Moreover, as $k \rightarrow +\infty$, then $\mathbf{x}_k \rightarrow \mathbf{x}_*$.

1/27/25

Theorem D.1 in High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. John Wright, Yi Ma.

15

 UNIVERSITY OF MICHIGAN

Convergence of GD

Proof. Given f is convex & smooth, and \mathbf{x}_* is an optimal solution

$$f(\mathbf{x}_*) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{x}_* - \mathbf{x})$$

$$\iff f(\mathbf{x}) \leq f(\mathbf{x}_*) + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{x}_*)$$

By our previous lemma $\mathbf{x}' = \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) \implies f(\mathbf{x}') \leq f(\mathbf{x}) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2$

$$f(\mathbf{x}') \leq f(\mathbf{x}_*) + \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{x}_*) - \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2$$

$$\iff f(\mathbf{x}') - f(\mathbf{x}_*) \leq \frac{L}{2} \left[\frac{2}{L} \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{x}_*) - \frac{1}{L^2} \|\nabla f(\mathbf{x})\|_2^2 \right]$$

16

 UNIVERSITY OF MICHIGAN

Convergence of GD

Therefore, we have

$$\begin{aligned} f(\mathbf{x}') - f(\mathbf{x}_*) &\leq \frac{L}{2} \left[\frac{2}{L} \nabla f(\mathbf{x})^\top (\mathbf{x} - \mathbf{x}_*) - \frac{1}{L^2} \|\nabla f(\mathbf{x})\|_2^2 \right] \\ &= \frac{L}{2} \left[\|\mathbf{x} - \mathbf{x}_*\|_2^2 - \left\| \mathbf{x} - \frac{1}{L} \nabla f(\mathbf{x}) - \mathbf{x}_* \right\|_2^2 \right] \\ &= \frac{L}{2} \left[\|\mathbf{x} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}' - \mathbf{x}_*\|_2^2 \right] \end{aligned}$$

1/27/25

17

UNIVERSITY OF MICHIGAN

Convergence of GD

Let $\mathbf{x}_k = \mathbf{x}_{k-1} - \frac{1}{L} \nabla f(\mathbf{x}_{k-1})$, then we have

$$\begin{aligned} \sum_{i=1}^k (f(\mathbf{x}_i) - f(\mathbf{x}_*)) &\leq \frac{L}{2} \sum_{i=1}^k \left[\|\mathbf{x}_{i-1} - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_i - \mathbf{x}_*\|_2^2 \right] \\ &= \frac{L}{2} \left[\|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 - \|\mathbf{x}_k - \mathbf{x}_0\|_2^2 \right] \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 \end{aligned}$$

Therefore, given $f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq f(\mathbf{x}_i) - f(\mathbf{x}_*)$ for $\forall i \leq k$,

$$k(f(\mathbf{x}_k) - f(\mathbf{x}_*)) \leq \sum_{i=1}^k (f(\mathbf{x}_i) - f(\mathbf{x}_*)) \leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2$$

18

UNIVERSITY OF MICHIGAN

Suboptimal Convergence Rate of GD

Suppose the iterates $\{\mathbf{x}_i\}_{i \geq 0}$ is generated from a *black-box model*:

$$\mathbf{x}_{k+1} = \mathcal{F}_k \left(\{\mathbf{x}_i\}_{i=0}^k, \{f(\mathbf{x}_i)\}_{i=0}^k, \{\nabla f(\mathbf{x}_i)\}_{i=0}^k \right)$$

Theorem (Nesterov'03). For every positive L and R , there exists a convex differentiable f with ∇f L -Lipschitz, and an initial point \mathbf{x}_0 satisfying $\|\mathbf{x}_0 - \mathbf{x}_*\|_2 \leq R$, such that

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \geq c \frac{LR}{k^2} = \Omega\left(\frac{1}{k^2}\right),$$

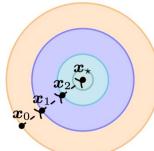
where $c > 0$ is some numerical constant.

1/27/25 Theorem D.2 in High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. John Wright, Yi Ma.

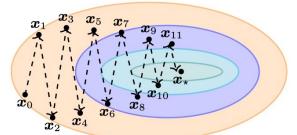
19

UNIVERSITY OF MICHIGAN

Slow Convergence of GD



well-conditioned
(fast convergence)



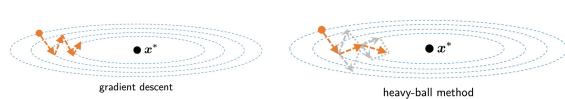
Ill-conditioned
(slow convergence)

1/27/25 Theorem D.2 in High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. John Wright, Yi Ma.

19

UNIVERSITY OF MICHIGAN

Heavy Ball Methods (Polyak'64)



$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \nabla f(\mathbf{x}_k) + \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1})$$

momentum term

- This is also called **momentum method**.
- Basis for the popular ADAM for training modern neural networks.



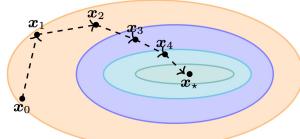
B. Polyak

1/27/25

21

UNIVERSITY OF MICHIGAN

Heavy Ball Methods (Polyak'64)



$$\begin{aligned} \mathbf{x}_{k+1} &= \mathbf{x}_k - \tau_k \nabla f(\mathbf{x}_k) \\ &+ \beta_k (\mathbf{x}_k - \mathbf{x}_{k-1}) \end{aligned}$$

momentum term

- This is also called **momentum method**.
- Basis for the popular ADAM for training modern neural networks.
- The worst-case convergence is still $O(1/k)$.

1/27/25

22

UNIVERSITY OF MICHIGAN

Nesterov's Method (Nesterov'83)

Generate two sequences $\{\mathbf{x}_k\}_{k \geq 1}$ and $\{\mathbf{p}_k\}_{k \geq 1}$

$$\begin{aligned} \mathbf{p}_{k+1} &= \mathbf{x}_k + \beta_k \cdot (\mathbf{x}_k - \mathbf{x}_{k-1}), \\ \mathbf{x}_{k+1} &= \mathbf{p}_{k+1} - \alpha \nabla f(\mathbf{p}_{k+1}) \end{aligned}$$

- Alternates between gradient updates \mathbf{x}_k and proper extrapolation \mathbf{p}_k .
- Not a descent method, i.e., we may not have $f(\mathbf{x}_{k+1}) \leq f(\mathbf{x}_k)$.
- With properly chosen α and β_k , the method can achieve optimal convergence rate $O(1/k^2)$



1/27/25

23

UNIVERSITY OF MICHIGAN

Nesterov's Method (Nesterov'83)

Theorem. (Convergence of accelerated GD)
Suppose $f : \mathbb{R}^n \mapsto \mathbb{R}$ is convex and smooth with its gradient L -Lipschitz. The iterates $\{\mathbf{x}_k\}_{k \geq 1}$ generated by the accelerated GD method satisfy

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{L}{2(k+1)^2} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2 = O\left(\frac{1}{k^2}\right).$$

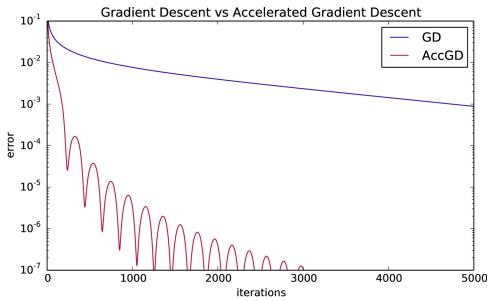
Moreover, as $k \rightarrow +\infty$, then $\mathbf{x}_k \rightarrow \mathbf{x}_*$.

1/27/25 Theorem D.3 in High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. John Wright, Yi Ma.

24

UNIVERSITY OF MICHIGAN

Nesterov's Method (Nesterov'83)



1/27/25

25

UNIVERSITY OF MICHIGAN

Lecture Agenda

- Basics of Iterative Methods
- (Accelerated) Gradient Descent for Smooth Problems
- Linesearch Methods for Step Size**
- Exploiting Function Properties for Faster Convergence
- Newton & Quasi-Newton Method

1/27/25

26

UNIVERSITY OF MICHIGAN

Exact Linesearch

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \tau_k \cdot \nabla f(\mathbf{x}_k)$$

- The choice $\tau = 1/L$ relies on the knowledge of the function property of $f(\mathbf{x})$

- Another more practical strategy is using linesearch

$$\tau_k = \arg \min_{t \geq 0} f(\mathbf{x}_k - t \cdot \nabla f(\mathbf{x}_k))$$

which could be expensive to solve...

1/27/25

27

UNIVERSITY OF MICHIGAN

Backtracking Linesearch

- Armijo condition:** for some $c_1 \in (0, 1)$

$$f(\mathbf{x}_k + \tau_k \cdot \mathbf{d}_k) < f(\mathbf{x}_k) + c_1 \cdot \tau_k \cdot \mathbf{d}_k^\top \nabla f(\mathbf{x}_k)$$

Purpose: ensure *sufficient decrease* of objective values

- Wolfe condition:** for some $0 < c_1 < c_2 < 1$

$$f(\mathbf{x}_k + \tau_k \cdot \mathbf{d}_k) < f(\mathbf{x}_k) + c_1 \cdot \tau_k \cdot \mathbf{d}_k^\top \nabla f(\mathbf{x}_k)$$

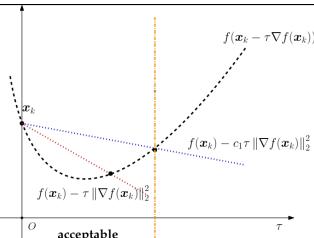
$$\mathbf{d}_k^\top \nabla f(\mathbf{x}_k + \tau_k \cdot \mathbf{d}_k) \geq c_2 \mathbf{d}_k^\top \nabla f(\mathbf{x}_k) \quad (\text{curvature condition})$$

Purpose: rule out *unacceptably small steps*

1/27/25

28

UNIVERSITY OF MICHIGAN



Algorithm 1 Backtracking Linesearch for Gradient Descent

```

Initialize  $\tau = 1$ ,  $c_1 \in (0, 1)$ ,  $\alpha \in (0, 1)$ 
while  $f(\mathbf{x}_k - \tau \nabla f(\mathbf{x}_k)) > f(\mathbf{x}_k) - c_1 \tau \|\nabla f(\mathbf{x}_k)\|_2^2$  do
     $\tau \leftarrow \alpha \cdot \tau$ 
end while

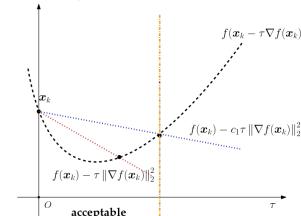
```

1/27/25

29

UNIVERSITY OF MICHIGAN

Backtracking Linesearch for GD



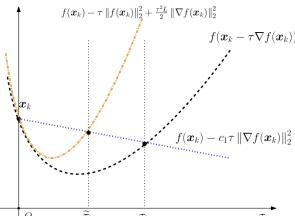
- $f(\mathbf{x}_k) - c_1 \tau \|\nabla f(\mathbf{x}_k)\|_2^2$ lies above $f(\mathbf{x}_k - \tau \nabla f(\mathbf{x}_k))$ for small τ
- ensures sufficient decrease of objective values.

1/27/25

30

UNIVERSITY OF MICHIGAN

Backtracking Linesearch for GD



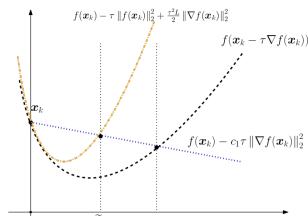
In practice, backtracking linesearch often (but not always) provides good estimates on the local Lipschitz constants of the gradient $\nabla f(\mathbf{x})$.

1/27/25

31

UNIVERSITY OF MICHIGAN

Backtracking Linesearch for GD



Question. In training deep network, why we do *not* use linesearch for adapting the learning rate?

1/27/25

32

UNIVERSITY OF MICHIGAN

Backtracking Linesearch for GD

Theorem. (Boyd, Vandenberghe'04) Let the function $f(\mathbf{x})$ be μ -strongly convex and L -smooth. With backtracking linesearch, the GD method converges with

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \left(1 - \min\left\{2c_1\mu, \frac{2\alpha c_1\mu}{L}\right\}\right)^k (f(\mathbf{x}_0) - f(\mathbf{x}_*))$$

where \mathbf{x}_* is the unique minimizer.

1/27/25

33

UNIVERSITY OF MICHIGAN

Lecture Agenda

- Basics of Iterative Methods
- (Accelerated) Gradient Descent for Smooth Problems
- Linesearch Methods for Step Size
- Exploiting Function Properties for Faster Convergence**
- Newton & Quasi-Newton Method

1/27/25

34

UNIVERSITY OF MICHIGAN

Assumptions of Strong Convexity (SC)

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathbb{R}^n.$$

Assumption 1. The function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is *smooth* with gradient $\nabla f(\mathbf{x})$ being L -Lipschitz.

Assumption 2. The function is μ -strongly convex. That is, for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

• Furthermore, if $f \in \mathcal{C}^2$, then $\mathbf{0} \prec \mu\mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L\mathbf{I}$.

1/27/25

35

UNIVERSITY OF MICHIGAN

Linear Convergence of GD

Theorem. (Linear convergence for strong convexity) Let $f : \mathbb{R}^n \mapsto \mathbb{R}$ be smooth and μ -strongly convex, and suppose its gradient $\nabla f(\mathbf{x})$ is L -Lipschitz. Choose a fixed step size $\tau_k \equiv \tau = \frac{2}{\mu+L}$, then

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2,$$

where $\kappa = \frac{L}{\mu}$ is the condition number, and \mathbf{x}_* is the minimizer.

- A direct consequence of Lipschitzness of $\nabla f(\mathbf{x})$

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \frac{L}{2} \left(\frac{\kappa - 1}{\kappa + 1}\right)^{2k} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2$$

1/27/25 Theorem D.4 in High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications. John Wright, Yi Ma.

36

UNIVERSITY OF MICHIGAN

Linear Convergence of GD

Proof. First, by the fundamental theorem of calculus:

$$\nabla f(\mathbf{x}_k) = \nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_*) = \left[\int_0^1 \nabla^2 f(\mathbf{x}(t)) dt \right] (\mathbf{x}_k - \mathbf{x}_*)$$

where $\mathbf{x}(t) = (1-t)\mathbf{x}_k + t\mathbf{x}_*$.

Then, we have

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &= \|\mathbf{x}_k - \tau \nabla f(\mathbf{x}_k) - \mathbf{x}_*\|_2 \\ &= \|(\mathbf{x}_k - \mathbf{x}_*) - \tau (\nabla f(\mathbf{x}_k) - \nabla f(\mathbf{x}_*))\|_2 \\ &= \left\| \left(\mathbf{x}_k - \tau \int_0^1 \nabla^2 f(\mathbf{x}(t)) dt \right) (\mathbf{x}_k - \mathbf{x}_*) \right\|_2 \end{aligned}$$

1/27/25

37

UNIVERSITY OF MICHIGAN

Linear Convergence of GD

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &\leq \left\| \mathbf{x}_k - \tau \int_0^1 \nabla^2 f(\mathbf{x}(t)) dt \right\|_2 \|\mathbf{x}_k - \mathbf{x}_*\|_2 \\ &\leq \left[\sup_{0 \leq t \leq 1} \|\mathbf{x}_k - \tau \nabla^2 f(\mathbf{x}(t))\|_2 \right] \|\mathbf{x}_k - \mathbf{x}_*\|_2 \end{aligned}$$

By the fact that f is μ -strongly convex, $\nabla^2 f(\mathbf{x}(t)) \geq \mu\mathbf{I}$, $\forall t$, we have

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &\leq (1 - \tau\mu) \|\mathbf{x}_k - \mathbf{x}_*\|_2 \\ &= \left(1 - \frac{2\mu}{\mu+L}\right) \|\mathbf{x}_k - \mathbf{x}_*\|_2 = \frac{L-\mu}{L+\mu} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \\ \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &\leq \frac{\kappa - 1}{\kappa + 1} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \leq \left(\frac{\kappa - 1}{\kappa + 1}\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2 \end{aligned}$$

38

UNIVERSITY OF MICHIGAN

Linear Convergence of GD

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &\leq \left\| \mathbf{x}_k - \tau \int_0^1 \nabla^2 f(\mathbf{x}(t)) dt \right\|_2 \|\mathbf{x}_k - \mathbf{x}_*\|_2 \\ &\leq \left[\sup_{0 \leq t \leq 1} \|\mathbf{x}_k - \tau \nabla^2 f(\mathbf{x}(t))\|_2 \right] \|\mathbf{x}_k - \mathbf{x}_*\|_2 \end{aligned}$$

By the fact that f is μ -strongly convex, $\nabla^2 f(\mathbf{x}(t)) \geq \mu\mathbf{I}$, $\forall t$, we have

$$\begin{aligned} \|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 &\leq (1 - \tau\mu) \|\mathbf{x}_k - \mathbf{x}_*\|_2 \\ &= \left(1 - \frac{2\mu}{\mu+L}\right) \|\mathbf{x}_k - \mathbf{x}_*\|_2 = \frac{L-\mu}{L+\mu} \|\mathbf{x}_k - \mathbf{x}_*\|_2 \end{aligned}$$

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \langle \nabla f(\mathbf{x}_*), \mathbf{x}_k - \mathbf{x}_* \rangle + \frac{L}{2} \|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \leq \frac{L}{2} \left(\frac{\kappa - 1}{\kappa + 1}\right)^{2k} \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2$$

1/27/25

39

UNIVERSITY OF MICHIGAN

Is SC Necessary for Linear Convergence?

The strong convexity requirement can often be relaxed:

- Local (restricted) strong convexity;
- Regularity condition;
- Polyak-Lojasiewicz condition, etc.

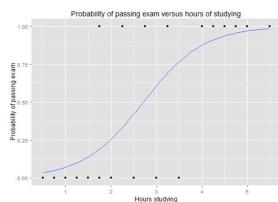
40

UNIVERSITY OF MICHIGAN

Example: Logistic Regression

For binary classification, we want to learn a linear classifier to classify the data.

Hours z_i	Pass y_i	Hours z_i	Pass y_i
0.5	0	2.75	1
0.75	0	3.00	0
1.00	0	3.25	1
1.25	0	3.50	0
1.5	0	4.00	1
1.75	0	4.25	1
1.75	1	4.50	1
2.00	0	4.75	1
2.25	1	5.00	1
2.50	0	5.50	1



1/27/25

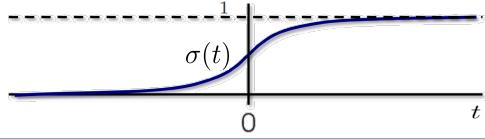
41

UNIVERSITY OF MICHIGAN

Example: Logistic Regression

- Given a bunch of training data $\{\mathbf{z}_i, y_i\}_{i=1}^m$ with $\mathbf{z} \in \mathbb{R}^n$, we want to learn the linear classifier $h(\mathbf{z}) = \mathbf{w}^\top \mathbf{z} + b$, so that we can make a decision by

$$y_i = \begin{cases} 1 & \text{if } \sigma(h(\mathbf{z}_i)) \geq 1/2, \\ 0 & \text{if } \sigma(h(\mathbf{z}_i)) < 1/2. \end{cases} \quad \text{with } \sigma(h(\mathbf{z})) = \frac{1}{1+\exp(-h(\mathbf{z}))}.$$



42

UNIVERSITY OF MICHIGAN

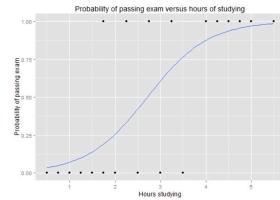
Example: Logistic Regression

Probability of passing an exam versus hours of study

Hours z_i	Pass y_i	Hours z_i	Pass y_i
0.5	0	2.75	1
0.75	0	3.00	0
1.00	0	3.25	1
1.25	0	3.50	0
1.5	0	4.00	1
1.75	0	4.25	1
1.75	1	4.50	1
2.00	0	4.75	1
2.25	1	5.00	1
2.50	0	5.50	1

$$\text{Prob.} = \frac{1}{1+\exp(-\mathbf{w}^\top \mathbf{z} - b)}$$

$$w = 1.5046, b = -4.0777.$$



1/27/25

43

UNIVERSITY OF MICHIGAN

Example: Logistic Regression

For simplicity, let us assume $y_i \in \{-1, 1\}$ instead of $y_i \in \{0, 1\}$. To learn the parameter \mathbf{w} , by maximum likelihood estimation (MLE), we have

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i \mathbf{w}^\top \mathbf{z}_i))$$

$$\nabla^2 f(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \underbrace{\frac{\exp(-y_i \mathbf{w}^\top \mathbf{z}_i)}{(1 + \exp(-y_i \mathbf{w}^\top \mathbf{z}_i))^2} \mathbf{z}_i \mathbf{z}_i^\top}_{\rightarrow 0 \text{ as } \mathbf{w} \rightarrow \infty} \xrightarrow{\mathbf{w} \rightarrow \infty} 0$$

$\Rightarrow f$ is 0-strongly convex

Does it mean that we no longer have linear convergence?

44

UNIVERSITY OF MICHIGAN

Local Strong Convexity (LSC)

Definition. (Local Strong Convexity) A function $f(\mathbf{x})$ is locally μ -strongly convex if it is μ -strongly convex within a ball of the global minimizer \mathbf{x}_* : $\mathcal{B}(\delta) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{x}_*\|_2 \leq \delta\}$

Theorem. Let f be locally μ -strongly convex and its gradient is L -Lipschitz with

$$\mu \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq L \mathbf{I}, \quad \forall \mathbf{x} \in \mathcal{B}_0,$$

where $\mathcal{B}_0 := \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{x}_*\|_2 \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2\}$. Then the linear convergence of GD continues to hold within \mathcal{B}_0 .

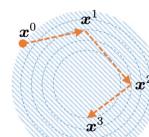
1/27/25

45

UNIVERSITY OF MICHIGAN

Local Strong Convexity

$$\mathcal{B}_0 = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{x}_*\|_2 \leq \|\mathbf{x}_0 - \mathbf{x}_*\|_2\}$$



- Suppose $\mathbf{x}_k \in \mathcal{B}_0$. From our previous analysis, $\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2 \leq \frac{\kappa-1}{\kappa+1} \|\mathbf{x}_k - \mathbf{x}_*\|_2$
- Thus, we have $\mathbf{x}_{k+1} \in \mathcal{B}_0$, so that we repeat the analysis for the next iteration ...

46

UNIVERSITY OF MICHIGAN

Local Strong Convexity

Back to the LR example, the local strong convexity parameter is controlled by

$$\mu = \inf_{\mathbf{w} \in \mathcal{B}_0(\mathbf{w}_*)} \lambda_{\min} \left(\frac{1}{m} \sum_{i=1}^m \frac{\exp(-y_i \mathbf{w}^\top \mathbf{z}_i)}{(1 + \exp(-y_i \mathbf{w}^\top \mathbf{z}_i))^2} \mathbf{z}_i \mathbf{z}_i^\top \right),$$

which can be strictly bounded away from 0.

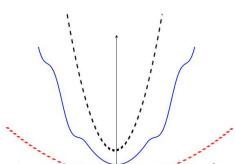
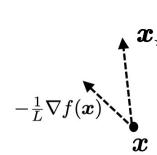
For instance, when $\mathbf{z}_i \sim i.i.d. \mathcal{N}(\mathbf{0}, \mathbf{I})$, we can show that $\mu \geq c_0$ for some universal constant $c_0 > 0$ if $m/n > 2$ (Sur et al.'17)

1/27/25

47

UNIVERSITY OF MICHIGAN

(Local) Regularity Condition

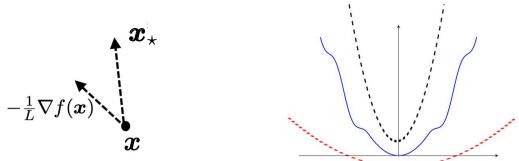


Definition. A function satisfies the regularity condition, if $\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_* \rangle \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}_*\|_2^2 + \frac{1}{2L} \|\nabla f(\mathbf{x})\|_2^2$ for any $\mathbf{x} \in \mathbb{R}^n$ or $\mathbf{x} \in \mathcal{B}(\mathbf{x}_*, \delta)$.

48

UNIVERSITY OF MICHIGAN

(Local) Regularity Condition



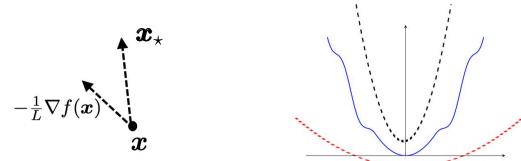
- An alternative for strong convexity and smoothness.
- Compared to strong convexity (which involves $\forall(\mathbf{x}, \mathbf{x}')$), regularity condition restrict to $(\mathbf{x}, \mathbf{x}_*)$ for any \mathbf{x} .

1/27/25

49

UNIVERSITY OF MICHIGAN

(Local) Regularity Condition



Theorem. Suppose f satisfies the regularity condition. If $\tau_k = \frac{1}{L}$, then the GD method converges *linearly*

$$\|\mathbf{x}_k - \mathbf{x}_*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^k \|\mathbf{x}_0 - \mathbf{x}_*\|_2^2.$$

Polyak-Lojasiewicz (PL) Condition

$$\|\nabla f(\mathbf{x})\|_2^2 \geq 2\mu(f(\mathbf{x}) - f(\mathbf{x}_*)), \quad \forall \mathbf{x}.$$

- It guarantees that gradient $\|\nabla f(\mathbf{x})\|_2$ grows fast as we move away from the minimizer $\nabla f(\mathbf{x}_*) = \mathbf{0}$;
- It guarantees that every stationary point is a global minimizer.

1/27/25

51

UNIVERSITY OF MICHIGAN

Polyak-Lojasiewicz (PL) Condition

Theorem. Suppose f satisfies the PL condition and its gradient ∇f is L -Lipschitz. Then if $\tau_k = \frac{1}{L}$, we have

$$f(\mathbf{x}_k) - f(\mathbf{x}_*) \leq \left(1 - \frac{\mu}{L}\right)^k (f(\mathbf{x}_0) - f(\mathbf{x}_*)).$$

- It guarantees linear convergence to the optimal objective value;
- It does *not* imply the uniqueness of global minimizer.

1/27/25

1/27/25

https://www.cs.ubc.ca/~jmutini/documents/PL_talk.pdf

52

UNIVERSITY OF MICHIGAN

Example: Overparameterized Linear Regression

Given data points $\{y_i, \mathbf{z}_i\}_{i=1}^m$ with $\mathbf{z}_i \in \mathbb{R}^n$, find a linear model that best fits the data

$$\min_{\mathbf{w} \in \mathbb{R}^n} f(\mathbf{w}) = \frac{1}{2m} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{z}_i)^2.$$

Over-parameterization: model dimension $n >$ sample size m
— a regime of particular interest in deep learning

1/27/25

53

UNIVERSITY OF MICHIGAN

Example: Overparameterized Linear Regression

The problem is convex but *not* strongly convex, because

$$\nabla^2 f(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \mathbf{z}_i \mathbf{z}_i^\top \text{ is rank-deficient if } n > m$$

But for the most “non-degenerate” cases, one has $f(\mathbf{w}_*) = 0$ and the PL condition is met, and hence GD converges linearly.

1/27/25

1/27/25

54

UNIVERSITY OF MICHIGAN

Example: Overparameterized Linear Regression

Fact. Suppose that $\mathbf{Z} = [\mathbf{z}_1 \cdots \mathbf{z}_m]^\top \in \mathbb{R}^{m \times n}$ has rank m , and that $\tau_k = \lambda_{\max}^{-1}(\mathbf{Z}\mathbf{Z}^\top)$. Then GD obeys

$$f(\mathbf{w}_k) - f(\mathbf{w}_*) \leq \left(1 - \frac{\lambda_{\min}(\mathbf{Z}\mathbf{Z}^\top)}{\lambda_{\max}(\mathbf{Z}\mathbf{Z}^\top)}\right)^k (f(\mathbf{w}_0) - f(\mathbf{w}_*)), \quad \forall k.$$

- Very mild assumption on $\{\mathbf{z}_i\}_{i=1}^m$ and no assumptions on $\{y_i\}_{i=1}^m$
- There are many global minimizers, but GD has *implicit bias* to a minimizer closest to the initialization \mathbf{w}_0 .

1/27/25

55

UNIVERSITY OF MICHIGAN

Linear Convergence without Strong Convexity

Many (local) conditions give linear rates that are weaker than strong-convexity (SC):

- 1963: Polyak-Lojasiewicz (PL)
- 1993: Error bounds
- 2000: Quadratic growth (QG)
- 2013-2015: essential SC, weak SC, restricted secant inequality, restricted SC, optimal SC, semi-SC, etc.

All of the above imply PL except QG.

1/27/25

https://www.cs.ubc.ca/~jmutini/documents/PL_talk.pdf

56

UNIVERSITY OF MICHIGAN

Lecture Agenda

- Basics of Iterative Methods
- (Accelerated) Gradient Descent for Smooth Problems
- Linesearch Methods for Step Size
- Exploiting Function Properties for Faster Convergence
- **Newton & Quasi-Newton Method**

1/27/25

57

UNIVERSITY OF MICHIGAN

Newton's Direction

$$\min_{\mathbf{x}} f(\mathbf{x}), \quad \text{s.t. } \mathbf{x} \in \mathbb{R}^n.$$

At the current iterate \mathbf{x}_k , consider the quadratic approximation

$$f(\mathbf{x}_k + \mathbf{d}) \approx \underbrace{f(\mathbf{x}_k) + \mathbf{d}^\top \nabla f(\mathbf{x}_k) + \frac{1}{2} \mathbf{d}^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d}}_{Q_k(\mathbf{d})}$$

Newton direction \mathbf{d}_k^N : minimizes $Q_k(\mathbf{d})$

$$\mathbf{d}_k^N = -(\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k).$$

1/27/25

58

UNIVERSITY OF MICHIGAN

Newton's Method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{\tau_k}{\text{stepsize}} \cdot \frac{\mathbf{d}_k^N}{\text{descent direction}}.$$

- **Descent direction:** when $\nabla^2 f(\mathbf{x}_k) \succ \mathbf{0}$
- $\nabla f(\mathbf{x}_k)^\top \mathbf{d}_k^N = -(\mathbf{d}_k^N)^\top \nabla^2 f(\mathbf{x}_k) \mathbf{d}_k^N < -c_k \|\mathbf{d}_k^N\|_2^2 < 0$.
- **Stepsize:** natural choice is $\tau_k = 1$, but can be chosen based on Armijo linesearch

$$f(\mathbf{x}_k + \alpha \mathbf{d}_k^N) \leq f(\mathbf{x}_k) + \alpha \nabla f(\mathbf{x}_k)^\top \mathbf{d}_k^N.$$

1/27/25

59

UNIVERSITY OF MICHIGAN

Newton's Method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{\tau_k}{\text{stepsize}} \cdot \frac{\mathbf{d}_k^N}{\text{descent direction}}.$$

- **Fast Convergence:** Newton's method usually has *local quadratic convergence*.
- **Expensive:** it requires storing and inverting $\nabla^2 f(\mathbf{x})$.
- **Unstable:** $\nabla^2 f \succ \mathbf{0}$ might not hold for general nonlinear problems.

1/27/25

60

UNIVERSITY OF MICHIGAN

Quasi-Newton's Method

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \frac{\tau_k}{\text{stepsize}} \cdot \frac{\mathbf{d}_k^{QN}}{\text{descent direction}}$$

$$\mathbf{d}_k^{QN} = -\underbrace{\mathbf{H}_k}_{\text{surrogate of } (\nabla^2 f(\mathbf{x}_k))^{-1}} \nabla f(\mathbf{x}_k)$$

Challenges:

- \mathbf{H}_k needs to be a good approximation of $(\nabla^2 f(\mathbf{x}_k))^{-1}$
- \mathbf{H}_k must be a descent direction with $\mathbf{H}_k \succ \mathbf{0}$
- It can be more efficiently computed than $(\nabla^2 f(\mathbf{x}_k))^{-1}$

1/27/25

61

UNIVERSITY OF MICHIGAN

Quasi-Newton's Method

Idea: consider the following quadratic approximation of f

$$Q_H^k(\mathbf{x}) := f(\mathbf{x}_{k+1}) + (\mathbf{x} - \mathbf{x}_{k+1})^\top \nabla f(\mathbf{x}_{k+1}) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_{k+1})^\top \mathbf{H}_k^{-1} (\mathbf{x} - \mathbf{x}_{k+1}).$$

Gradient matching for the last two iterates:

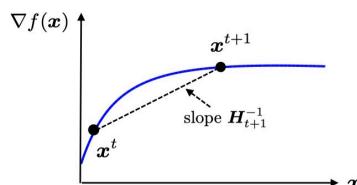
$$\nabla Q_H^k(\mathbf{x}_k) = \nabla f(\mathbf{x}_k), \quad \nabla Q_H^k(\mathbf{x}_{k+1}) = \nabla f(\mathbf{x}_{k+1})$$

1/27/25

62

UNIVERSITY OF MICHIGAN

Secant Equation



$$\mathbf{H}_{k+1}^{-1} \cdot \underbrace{(\mathbf{x}_{k+1} - \mathbf{x}_k)}_{\mathbf{s}_k} = \underbrace{\nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)}_{\mathbf{y}_k}$$

1/27/25

63

UNIVERSITY OF MICHIGAN

Quasi-Newton Direction

$$\mathbf{H}_{k+1}^{-1} \cdot \mathbf{s}_k = \mathbf{y}_k$$

- admit an infinite number of solutions;
- \mathbf{H}_k and \mathbf{H}_{k+1} are close and the difference is low-rank;
- \mathbf{H}_{k+1} needs to be positive definite

$$\mathbf{H}_0 \succ \mathbf{0}, \quad \mathbf{s}_k^\top \mathbf{y}_k > 0.$$

1/27/25

64

UNIVERSITY OF MICHIGAN

Quasi-Newton Direction

- Symmetric rank-1 update (SR1):

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^\top}{(\mathbf{s}_k - \mathbf{H}_k \mathbf{y}_k)^\top \mathbf{y}_k}$$

NOT preserving positive definiteness.

- BFGS rank-2 update:

$$\mathbf{H}_{k+1} = \mathbf{H}_k + \frac{(\mathbf{s}_k^\top \mathbf{y}_k + \mathbf{y}_k^\top \mathbf{H}_k \mathbf{y}_k) \mathbf{s}_k \mathbf{s}_k^\top}{(\mathbf{s}_k^\top \mathbf{y}_k)^2} - \frac{\mathbf{H}_k \mathbf{y}_k \mathbf{s}_k^\top + \mathbf{s}_k \mathbf{y}_k^\top \mathbf{H}_k}{\mathbf{s}_k^\top \mathbf{y}_k}$$

BFGS Method (1970)



BFGS Method (1970)

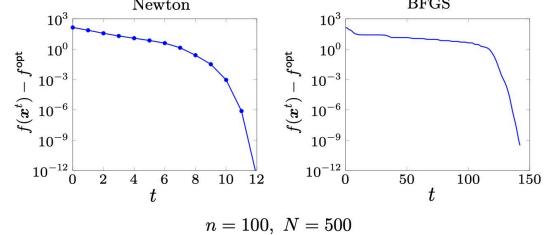
Theorem (Informal). Suppose $f \in \mathcal{C}^2$ is strongly convex, and has Lipschitz-continuous Hessian. Under mild conditions, BFGS achieves superlinear convergence

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}_*\|_2}{\|\mathbf{x}_k - \mathbf{x}_*\|_2} = 0.$$

- Iteration complexity: from $O(n^3)$ to $O(n^2)$;
- Asymptotic result: holds when $k \rightarrow \infty$;
- Limited memory version reduces to $O(n)$, but only linear convergence.

BFGS Method

$$\min_{\mathbf{x} \in \mathbb{R}^n} \mathbf{c}^\top \mathbf{x} - \sum_{i=1}^N \log(b_i - \mathbf{a}_i^\top \mathbf{x})$$



$n = 100, N = 500$

References & Further Readings

- *High-Dimensional Data Analysis with Low-Dimensional Models: Principles, Computation, and Applications*. John Wright, Yi Ma. ([Appendix D](#))
- *Numerical Optimization*, Jorge Nocedal, and Stephen Wright, Springer. ([Chapter 2, 3, & 6](#))