

MATH466/MATH766

Math of machine learning

02/10-02/12 Lecture 9-10 Optimization Algorithm

References:

- Convex Optimization by Stephen Boyd and Lieven Vandenberghe
- Lecture notes by Prof. [Ryan Tibshirani](#)
<https://www.stat.cmu.edu/~ryantibs/convexopt-F15/>

Todays contents:

- strong convexity and Lipschitz smoothness
- convergence and convergence rate
- gradient descent
- subgradient method

Important concepts:

- strong convexity and Lipschitz smoothness
- sublinear convergence, linear convergence

Recommend reading:

-

0. Preparation

0.1. strongly convexity

- f is μ -strongly convex ($\mu > 0$) iff $f - \frac{\mu}{2} \| \cdot \|_2^2$ is convex

- Strongly convex \Rightarrow strictly convex \Rightarrow convex
 $\Leftrightarrow f(x) = x^*$, $f(x) = e^x$ $\Leftrightarrow f(x) = x$

- 0-th, 1-st and 2-nd equivalent conditions

$$(0) \quad f(\theta x + (1-\theta)y) \leq \theta f(x) + (1-\theta)f(y) - \frac{\theta(1-\theta)\mu}{2} \|x-y\|^2 \quad \forall x, y \\ \theta \in [0, 1]$$

$$\textcircled{1} \quad (1) \quad f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2, \quad \forall x, y$$

$$\textcircled{2} \quad \nabla^2 f(x) - \mu I \succeq 0, \quad \forall x$$

0.2. Lipschitz smoothness

- A differentiable function f is said to be L -Lipschitz differentiable or L -Lipschitz smooth if for some $L > 0$,

$$\| \nabla f(x) - \nabla f(y) \| \leq L \|x-y\|, \quad \forall x, y. \quad (*)$$

(This definition does not assume convexity of f .)

- 1-st and 2-nd order necessary conditions

$$(*) \Rightarrow (1) \quad f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|x-y\|^2$$

$$(*) \Rightarrow (2) \quad \nabla^2 f(x) \preceq L I$$

Interpretation :

Convexity \sim gradient is non-decreasing

Strong convexity \sim the growth of gradient has a low-bound

Lipschitz smoothness \sim the growth of gradient has an upper-bound

1. Convergence and Convergence Rate

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

The study of an optimization algorithm mainly focus on convergence rate
 local > convergence and complexity. # of iteration - per iteration
 global

- local convergence v.s. global convergence

Algorithm A has local convergence if there exists $r > 0$

such that when $d(x^0, x^*) < r$,

the $\{x^k\}$ generated by algorithm A satisfies

$$f(x^k) \rightarrow f(x^*), \text{ and } \text{dist}(x^k, x^*) \rightarrow 0.$$

Algorithm A has global convergence if

for any initialization x^0 ,

the $\{x^k\}$ generated by algorithm A satisfies

$$f(x^k) \rightarrow f(x^*), \text{ and } \text{dist}(x^k, x^*) \rightarrow 0.$$

- Convergence rate (assume $x^k \rightarrow \bar{x}$)

(1) Sublinear rate if $\lim_{k \rightarrow +\infty} \frac{\|x^{k+1} - \bar{x}\|}{\|x^k - \bar{x}\|} = 1$.

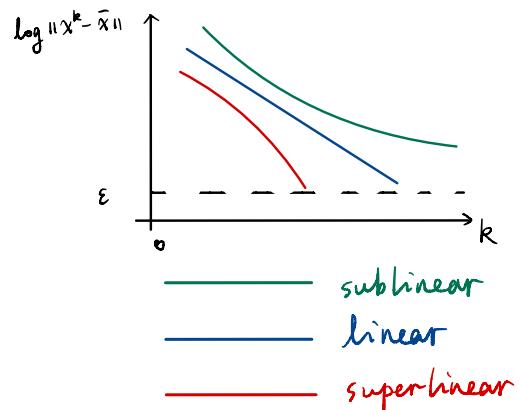
e.g. $x^k = \frac{1}{k+1}$, to have $\|x^k - \bar{x}\| < \varepsilon$, we need $k = O(\frac{1}{\varepsilon})$

(2) Linear rate If $\lim_{k \rightarrow +\infty} \frac{\|x^{k+1} - \bar{x}\|}{\|x^k - \bar{x}\|} \leq r \in (0, 1)$.

e.g. $x^k = e^{-k}$, to have $\|x^k - \bar{x}\| < \varepsilon$, we need $k = O(\log \frac{1}{\varepsilon})$

(3) Superlinear rate if $\lim_{k \rightarrow +\infty} \frac{\|x^{k+1} - \bar{x}\|}{\|x^k - \bar{x}\|} = 0$

e.g. $x^k = e^{-2^k}$ to have $\|x^k - \bar{x}\| < \varepsilon$, we need $k = O(\log_2 \log \frac{1}{\varepsilon})$



2. Descent Method and Gradient Descent

Consider unconstrained smooth convex optimization

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

A natural idea of algorithm: $x^{(k+1)} = x^{(k)} + t^{(k)} p^{(k)}$

\uparrow \uparrow
 step size search direction

Q1. What is a good search direction?

Q2. What is a proper step-size?

Want: $f(x^{(k+1)}) \leq f(x^{(k)})$

$$f(x^{(k+1)}) = f(x^{(k)}) + t^{(k)} \nabla f(x^{(k)} + s p^{(k)})^T p^{(k)} \quad s \in (0, t^{(k)})$$

- Descent direction: if $\nabla f(x^{(k)})^T p^{(k)} < 0$, then $p^{(k)}$ is a descent direction
- Line-search for stepsize: if $p^{(k)}$ is a descent direction
then when $t^{(k)}$ is small enough, $f(x^{(k+1)}) \leq f(x^{(k)})$

But we don't want $t^{(k)}$ to be too small (why?)

for given $x^{(k)}$ and direction $p^{(k)}$

① exact line search: (in general not practical)

solve $\min_{t>0} f(x^{(k)} + t p^{(k)})$ for t ID opt

② backtracking line search: pick $\alpha \in (0, 0.5)$, $\beta \in (0, 1)$, $\gamma > 0$

$$t \leftarrow \gamma$$

while $f(x^{(k)} + t p^{(k)}) > f(x^{(k)}) + \alpha t \nabla f(x^{(k)})^T p^{(k)}$

(Armijo's cond.)

there are other

versions of line search

$$t \leftarrow \beta t$$

α: the fraction of decrease in linearized f we accept
typical choice: $\alpha \in (0.01, 0.3)$ less → more greedy

β: the amount of shrinkage when condition doesn't hold
typical choice: $\beta \in (0.1, 0.8)$ less → more crude

Gradient descent

$$p^{(k)} = -\nabla f(x^{(k)})$$

Alg. GD

Initialization $x^{(0)}$

while $\|\nabla f(x^{(k)})\|_2 > \epsilon$

$$x^{(k+1)} = x^{(k)} - t^{(k)} \nabla f(x^{(k)})$$

Prop. If f is L -smooth, $\nabla f(x^{(k)}) \neq 0$, then for any $t \in (0, \frac{2}{L})$

$$f(x^{(k+1)}) < f(x^{(k)}).$$

$$\text{Pf. } f(x^{(k+1)}) \leq f(x^{(k)}) + t \nabla f(x^{(k)})^T p^{(k)} + \frac{L}{2} \|t p^{(k)}\|^2$$

$$= f(x^{(k)}) - t \|\nabla f(x^{(k)})\|^2 + \frac{L}{2} t^2 \|\nabla f(x^{(k)})\|^2$$

$$= f(x^{(k)}) - (-\frac{L}{2} t^2 + t) \|\nabla f(x^{(k)})\|^2$$

$$\text{Since } \|\nabla f(x^{(k)})\|^2 < 0$$

And when $t \in (0, \frac{2}{L})$, $-\frac{L}{2} t^2 + t < 0$, $f(x^{(k+1)}) < f(x^{(k)})$ \square

If f is L -smooth, we usually pick $t = \frac{1}{L}$.

With this choice, we have the following convergence results.

Thm 1. Assume f is L -smooth, and $f(x) \geq f^*$ for any x

$$\text{then } \min \left\{ \|\nabla f(x^{(k)})\|^2 \right\}_{k=0}^{K-1} \leq \frac{2L(f(x^{(0)}) - f^*)}{K}$$

Thm 2. In addition to Thm 1 assumptions,

if f is convex and $f(x) \geq f(x^*)$ for any x .

$$\text{then } f(x^{(k)}) - f(x^*) \leq \frac{L \|x^{(0)} - x^*\|^2}{2k}$$

Thm 3. In addition to Thm 1, 2 assumptions,

if f is μ -strongly convex.

$$\text{then } f(x^{(k)}) - f(x^*) \leq (1 - \frac{\mu}{L})^k (f(x^{(0)}) - f(x^*))$$

$$\text{and } \|x^{(k)} - x^*\|^2 \leq (1 - \frac{\mu}{L})^k \|x^{(0)} - x^*\|^2$$

3. Subgradient Method

$$\text{direction } p^{(k)} = g^{(k)} \in \partial f(x^{(k)})$$

Alg. subGD

Initialization $x^{(0)}$

while $\|g^{(k)}\|_2 > \varepsilon$

$$x^{(k+1)} = x^{(k)} - t^{(k)} g^{(k)}$$

Subgradient is NOT Necessarily a descent method

i.e. $f(x^{(k+1)}) \leq f(x^{(k)})$ may not hold

$$\text{e.g. } f(x) = |x|, x^{(k)} = 0, g^{(k)} = 1$$

as long as $t > 0$, $f(x^{(k+1)}) > f(x^{(k)})$

so we keep track of best iterate $x_{\text{best}}^{(k)}$ among $x^{(0)}, \dots, x^{(k)}$ so far.

$$\text{i.e. } f(x_{\text{best}}^{(k)}) = \min_{k=0, \dots, K} f(x^{(k)})$$

or consider an average output $\bar{x}^{(K)} := \frac{1}{K} \sum_{k=1}^K x^{(k)}$

Thm. Assume that $f(x) \geq f(x^*)$ for any x

and f is L -Lipschitz continuous i.e. $|f(x) - f(y)| \leq L \|x - y\|_2$

$$\text{Then } f(x_{\text{best}}^{(k)}) - f(x^*) \leq \frac{R^2 + L^2 \sum_{k=1}^K t_k^2}{2 \sum_{k=1}^K t_k}, \quad (R = \|x^{(1)} - x^*\|)$$

If we pick a fixed step-size t , then $\lim_{k \rightarrow \infty} f(x_{\text{best}}^{(k)}) \leq f(x^*) + \frac{L^2 t}{2}$

The optimal value may NOT be achieved in the limit.

Smaller step-size may reduce the gap but requires more iterations.

Solution : diminishing step-size e.g. $t_k = \frac{1}{k}$

Convergence rate : To have $\frac{R^2 + L^2 \sum_{k=1}^K t_k^2}{2 \sum_{k=1}^K t_k} \leq \varepsilon$, need at least $(\frac{R L}{\varepsilon})^2$ iterations

$$\frac{R^2 + L^2 \sum_{k=1}^K t_k^2}{2 \sum_{k=1}^K t_k} \geq \frac{R L}{\sqrt{K}}, \quad "=" \text{ when } t_k = \frac{R/L}{\sqrt{k}}$$