

Recap:

$$\textcircled{1} \hat{w} = \underset{w}{\operatorname{argmin}} f(w) \quad f(w) := \frac{1}{N} \sum_{i=1}^N L(y_i, h_w(x_i))$$

$\textcircled{2}$ we need optimization to solve $\textcircled{1}$ (approximately)

$$\text{stationary point} \leftarrow \nabla f(\hat{w}) = 0 \quad (\text{necessary condition})$$

$\textcircled{3}$ convexity \rightarrow guarantee stationary point \Leftrightarrow minimizer
 (approximately \Rightarrow exactly)

characterization of convexity:

$$\begin{cases} f(\lambda x + (1-\lambda)y) \leq \lambda f(x) + (1-\lambda)f(y) & \forall x, y, \lambda \in [0, 1] \\ f(y) \geq f(x) + \nabla f(x)^T (y-x) & \forall x, y \end{cases}$$

$\textcircled{4}$ C-strongly convexity \rightarrow guarantee the unique minimizer
 \Rightarrow stationary point \rightarrow unique minimizer

characterization: $\begin{cases} f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{C}{2} \|y-x\|_2^2 \\ \langle \nabla f(x) - \nabla f(y), x-y \rangle \geq C \|y-x\|_2^2 \end{cases}$

$f(x)$ is C-strongly convex $\Leftrightarrow g(x) = f(x) - \frac{C}{2} \|x\|_2^2$ is convex

Today's lecture

1. Positive Definite / Semi-Definite

smallest eigenvalue

↑

$$\textcircled{1} A \text{ is PD} \Leftrightarrow x^T A x > 0 \text{ for } \forall x \neq 0 \Leftrightarrow \underline{\lambda_n > 0}$$

$$\textcircled{2} A \text{ is PSD} \Leftrightarrow x^T A x \geq 0 \text{ for } \forall x \in \mathbb{R}^n \Leftrightarrow \underline{\lambda_n \geq 0}$$

$$\textcircled{3} A \geq cI \Leftrightarrow x^T A x \geq c \|x\|_2^2 \text{ for } \forall x \in \mathbb{R}^n \\ \Leftrightarrow \lambda_n \geq c$$

2. Hessian and convexity

if f is well-conditioned, then H_f is symmetric

$$[H_f(x)]_{ij} := \frac{\partial^2 f}{\partial x_i \partial x_j}(x)$$

$$H_f : x \in \mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$$

can characterize convexity & c-strongly convexity

a) $f \in C^2$ is convex on D

$\Leftrightarrow H_f(x)$ is Positive Semi-Definite on all $x \in D$

$$[\text{e.g.}] f(x) = \frac{1}{2} x^T A x \Rightarrow f(x) = \|x\|_2^2 \text{ is convex}$$

$$\rightarrow \underline{\nabla f(x)} = \frac{1}{2} A x + \frac{1}{2} A^T x \\ = A x \quad (\text{if } A \text{ is symmetric})$$

$$\rightarrow \underline{\nabla^2 f(x)} = H_f(x) = A$$

Therefore, $f(\cdot)$ is convex \Leftrightarrow A is PSD

Proof: " \Rightarrow " f is convex

$$\Rightarrow \forall x, y, f(y) \geq f(x) + \nabla f(x)^T (y - x)$$

(Analyze: we want to show for $\forall x, \underline{z^T H_f(x) z} \geq 0, \forall z$)

\Rightarrow choose $y = x + \lambda z$: (fix x and randomly choose z)

then we have $f(x + \lambda z) \geq f(x) + \lambda \nabla f(x)^T z$

Also, from Taylor Expansion:

$$f(x + \lambda z) = f(x) + \nabla f(x)^T (\lambda z) + \frac{1}{2} (\lambda z)^T H_f(x) \cdot (\lambda z) + o(\lambda^2)$$

$$\geq f(x) + \nabla f(x)^T (\lambda z)$$

$$\Rightarrow \lambda^2 \cdot z^T H_f(x) z + o(\lambda^2) \geq 0$$

$$\Rightarrow z^T H_f(x) z + o(1) \geq 0 \quad \forall \lambda \geq 0 \quad (\text{let } \lambda \rightarrow 0)$$

$$\Rightarrow z^T H_f(x) z \geq 0 \text{ holds for } \forall z \in \mathbb{R}^n$$

" \Leftarrow " Now, $H_f(x)$ is PSD for $\forall x \in D$

$$\Rightarrow z^T H_f(x) z \geq 0 \text{ for } \forall z \in \mathbb{R}^n, \forall x \in D$$

$$\Rightarrow \text{for arbitrary } y \in \mathbb{R}^n, y = x + z \quad (z = y - x)$$

$$\begin{aligned} \text{then } f(y) &= f(x) + \nabla f(x)^T (y - x) \\ &\quad + \frac{1}{2} z^T H_f(\hat{x}) z \\ &\geq f(x) + \nabla f(x)^T (y - x) \end{aligned}$$

$\Rightarrow f(\cdot)$ is a convex function

#

b) $f(x)$ is a c -strongly convex function on D

$$\Leftrightarrow H_f(x) \succeq cI \text{ at each } x \in D$$

Pf Sketch. $f(x)$ is c -strongly convex

$$\Leftrightarrow g(x) = f(x) - \frac{c}{2} \|x\|_2^2 \text{ is convex}$$

$$\Leftrightarrow H_g(x) \succeq 0 \text{ for } \forall x \in D$$

$$\Leftrightarrow H_f(x) - cI \succeq 0 \text{ for } \forall x \in D$$

$$\Leftrightarrow H_f(x) \succeq cI \text{ for } \forall x \in D$$

[e.g.] ① $f(x) = e^x$ $x \in \mathbb{R}$

↓
Compute Hessian !!!

↓
convex but not c -strongly convex for arbitrary c

$f''(x) = e^x \rightarrow 0$ when $x \rightarrow -\infty$

\Rightarrow there does not exist $c > 0$ s.t. $e^x \geq c$ for all $x \in \mathbb{R}$!

$\Rightarrow f(\cdot)$ is not c -strongly convex on \mathbb{R}

② $f(x) = e^x$ $x \in \mathbb{R}^+$

→ c -strongly convex for $c \leq 1$

③ $f(x) = \|x\|_2^2$ → c -strongly convex for $c \leq 2$

→ $H_f(x) = 2I$

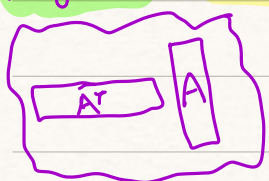
④ $f(x) = \|Ax - b\|_2^2$ → convex

→
$$\begin{cases} \nabla f(x) = 2A^T(Ax - b) \\ H_f(x) = 2A^T A \geq 0 \end{cases}$$

Rmk: 1. in e.g. 4, $f(x)$ is at least convex.

2. in order to further determine whether $f(x)$ is c -strongly convex or not, check the smallest eigenvalue of

matrix $A^T A$ ($\lambda_n \geq c$ or not)



↕
whether A is full column rank matrix

Implementation

3. Algos to find w^* such that $\nabla f(w^*) = 0$

Rmk: actually there are 2 perspectives

a) w^* s.t. $\nabla f(w^*) = 0$

b) $w^* = \underset{w}{\operatorname{argmin}} f(w)$

idea: use $w^{(k+1)}$ to approximate w^*

here, we focus on this iterative algorithm

a) stop criterion

$$w_1, w_2, \dots, w_k \rightarrow w^*$$

① stop at $k=T$, T is pre-determined \Rightarrow deep learning

② stop at $\|\nabla f(w_k)\| \leq \text{tolerance}$

tolerance \Rightarrow small precision requirement

③ stop at $\|w_{k+1} - w_k\| \leq \text{tolerance}$

b) Question: How many iterations are required to satisfy

$$\|w_k - w^*\| \leq \epsilon$$

convergence rate

c) update cost

$$w_{k+1} = w_k + \alpha_k p_k$$

p_k $\begin{cases} \text{first order information} \rightarrow \nabla f(w_k) \\ \text{second order information} \rightarrow H_f(w_k) \end{cases}$

d) scalability & parallelization

formulation

4. Newton's Method

$$\begin{cases} \textcircled{1} \nabla f(w) = 0 \approx \nabla f(w_k) + H_f(w_k) \cdot (w - w_k) = 0 \\ \textcircled{2} w^* = \underset{w}{\operatorname{argmin}} f(w) \approx w = \underset{w}{\operatorname{argmin}} \left\{ \nabla f(w_k)^T (w - w_k) + \frac{1}{2} (w - w_k)^T H_f(w_k) (w - w_k) \right\} \end{cases}$$

$$\Rightarrow \boxed{\omega_{k+1} = \omega_k - H_f(\omega_k)^{-1} \nabla f(\omega_k)} \rightarrow \boxed{\text{Newton's method}}$$

→ Theorem: (local quadratic convergence rate)
(no need of convexity)

Suppose $\|\omega_0 - \omega^*\|$ is sufficiently small, $H_f(\omega)^{-1}$ & $\nabla f(\omega)$ are L -Lipschitz, then there exists M s.t.

$$\|\omega_{k+1} - \omega^*\| \leq M \|\omega_k - \omega^*\|^2$$

$$\boxed{\nabla f(\omega^*) = 0}$$

stationary point

[e.g.] $\omega_0 = 1.1 \quad \omega_1 = 1.01 \quad \omega_2 = 1.0001 \dots$

[Pf Sketch]:
$$\begin{cases} \omega_{k+1} = \omega_k - H_f(\omega_k)^{-1} \cdot \nabla f(\omega_k) \\ \omega^* \text{ satisfy that } \nabla f(\omega^*) = 0 \end{cases}$$

$$\Leftrightarrow \nabla f(\omega_k) + H_f(\bar{\omega})(\omega^* - \omega_k) = 0$$

$$\boxed{\bar{\omega} \in (\omega_k, \omega^*)} \Rightarrow \omega^* = \omega_k - H_f(\bar{\omega})^{-1} \nabla f(\omega_k)$$

Thus,
$$\omega_{k+1} - \omega^* = (H_f(\bar{\omega})^{-1} - H_f(\omega_k)^{-1}) \cdot \nabla f(\omega_k)$$

$$= \underbrace{(H_f(\bar{\omega})^{-1} - H_f(\omega_k)^{-1})}_{L\text{-lip. of } H_f(\omega)} \cdot \underbrace{(\nabla f(\omega_k) - \nabla f(\omega^*))}_{L\text{-lip. of } \nabla f(\omega)}$$

Rmk: ① quadratic local convergence rate

② converge to ω^* , which is the stationary point

(not minimizer)

step size issue

③ no global convergence guarantee

→ line search rule modification (globalize)

Hessian matrix issue

④ if $H_f(\omega_k)$ is not PD, then Newton's direction

$$p_k = -H_f(\omega_k)^{-1} \nabla f(\omega_k) \text{ may not be descent direction}$$

⑤ invert / record Hessian (attain $H_f(\cdot)^{-1}$) is expensive
(n^2) space

Possible Solution:

quasi-newton

- a) approximate $H_f(w_k)^T$ via $\nabla f(w_k)$ [BFGS]
- b) solve linear equation $H_f(w_k) p_k = -\nabla f(w_k)$ via iterative methods

5. Gradient Descent

① idea: $w_{k+1} = w_k + \alpha_k p_k$

\Downarrow

$$f(w_{k+1}) \approx f(w_k) + \alpha_k \nabla f(w_k)^T p_k + o(\alpha_k)$$

$\nabla f(w_k)^T p_k$ is the steepest descent direc.

$p_k^{GD} = -\nabla f(w_k)$

Rmk: ① if α_k is sufficiently small, $f(w_{k+1}) \leq f(w_k)$

\Downarrow

descent direction

② Another formulation:

$$p_k^{GD} = \underset{p}{\operatorname{argmin}} \underbrace{f(w_k) + \nabla f(w_k)^T p}_{\text{approximation term}} + \underbrace{\frac{1}{2} \|p\|_2^2}_{\text{proximal term}}$$

② Determine step length α_k

$$\begin{cases} \text{exact line search: } \alpha_k = \underset{\alpha}{\operatorname{argmin}} f(w_k + \alpha p_k) \\ \text{inexact line search: } \underline{\text{backtracking}} \xrightarrow{\text{similar}} \begin{cases} \text{exponential decay} \\ \text{cosine schedule} \end{cases} \\ \text{fixed step size } \alpha_k \equiv \alpha \rightarrow \text{most popular} \end{cases}$$

③ Thm:

\rightarrow convex

if $f(w) \geq 0$ & ∇f is L-Lipschitz, then fixed step size GD with $\alpha_k \equiv \alpha \leq \frac{1}{L}$ satisfy:

\rightarrow sub-linear

$$0 \leq f(w_n) - f(w^*) \leq \frac{1}{n} \frac{\|w_0 - w^*\|_2^2}{2\alpha}$$

$w^* = \operatorname{argmin} f(w)$

④ Thm: (much better than previous result due to strongly convexity)
 if f is c -strongly convex and $\alpha \leq \frac{c}{L^2}$, then Linear converg
 $w^* = \operatorname{argmin} f(w)$
 $0 \leq f(w_n) - f(w^*) \leq (1 - c\alpha)^n L \|w_0 - w^*\|_2^2$

exponential convergence

Rmk: in application, if our $f(\cdot)$ is strongly convex, then we can choose our learning rate small but not too small to achieve a good convergence rate.

satisfy $\alpha \leq \frac{c}{L^2}$

make $1 - c\alpha \ll 1$