# LEC 3    DSA 5202

## Recap:

1. convexity $\longleftrightarrow$ $H_f(x) \succcurlyeq 0$ for $\forall x \in D$

   c-strongly convexity $\longleftrightarrow$ $H_f(x) \succcurlyeq cI$ for $\forall x \in D$

   (variants) $\Downarrow$

2. Newton's method (pure & line search)

   local convergence

   global convergence

3. Gradient Descent method $\longrightarrow$ some convergence result (fixed step length)

   (require convexity condition)

## IDEA of these algos:

$$w^* = \underset{w}{\arg\min} \, f(w) \xrightarrow[\text{condition}]{\text{necessary}} \nabla f(w^*) = 0$$

$\longrightarrow$ only for linear $f(\cdot)$, $\nabla f(w^*) = X^T(Xw^* - y) = 0$ is tractable

$$\Downarrow$$

$$w^* = (X^TX)^{-1}X^T y$$

$\longrightarrow$ for most $f(\cdot)$. WE WANT TO FIND A SET OF $\{w_i\}_{i=1}^{K}$

   s.t $\quad w_i \longrightarrow w^*$ $(i \to \infty)$   [iterative approach]

$$\begin{cases} \text{Newton Update}: \quad w_{k+1} = w_k - H_f(w_k)^{-1} \nabla f(w_k) \\ \text{GD Update}: \quad w_{k+1} = w_k - \alpha_k \nabla f(w_k) \end{cases}$$

# Today's lecture

## 1. Issue of Gradient Descent

① Consider our high-level problem :

Population Risk Minimization

$z = (x, y)$

$f(w)$

$F(w, z)$

intractable $\Leftarrow$ $w^* = \underset{w}{\arg\min}$ $\boxed{\mathbb{E}_{(x,y) \sim \mathcal{D}} [ \ell(h_w(x), y)]}$
due to $\mathcal{D}$

e.g. $\begin{cases} \ell(z, z') = \|z - z'\|_2^2 \longrightarrow \boxed{\text{regression}} \\ \ell(z, z') = D_{KL}(z', z) = \sum_i z'(i) \log \dfrac{z'(i)}{z(i)} \longrightarrow \boxed{\text{classification}} \end{cases}$

Now, $\nabla f(w) = \nabla \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)] \longrightarrow$ intractable

$\neq \mathbb{E}_{z \sim \mathcal{D}} [\nabla F(w, z)]$

② Surrogate : $\boxed{\text{PRM} \longrightarrow \text{ERM}}$ $\boxed{\text{LLN}}$

$\longrightarrow$ Empirical Risk Minimization

$\hat{f}_n(w)$

$\hat{w} = \underset{w}{\arg\min} \boxed{\dfrac{1}{n} \sum_{i=1}^{n} F(w, z_i)} \xrightarrow{p} \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)]$

$\underset{f(w)}{\|}$

Recap: $\underset{\underset{\text{population}}{\|}}{f(w)} = \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)]$

$\longrightarrow$ If we want to apply GD Framework,

then $w_{k+1} = w_k - \alpha_k \cdot \nabla \hat{f}_n(w_k)$

$= w_k - \alpha_k \cdot \dfrac{1}{n} \underbrace{\sum_{i=1}^{n} \nabla F(w_k, z_i)}_{n - \text{summation}}$

$\Rightarrow$ ☆ 1 update $\longleftrightarrow$ n calculation (gradient)

$\Downarrow$

inefficient when n is big (large dataset)

$\triangle$

## 2. SGD $\longrightarrow$ Stochastic Gradient Descent

① original problem:

$$w^* = \underset{w}{argmin} \ f(w) := \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)]$$

$$\longrightarrow \nabla f(w) = \nabla \mathbb{E}_{z \sim \mathcal{D}} [F(w, z)]$$

(under strong regularity
of $F(\cdot, \cdot)$ )
$$= \mathbb{E}_{z \sim \mathcal{D}} [\nabla_w F(w, z)]$$
$\hookrightarrow$ not always correct !

② ERM Surrogate :

$$\hat{w} = \underset{w}{argmin} \ \hat{f}_n(w) := \frac{1}{N} \sum_{i=1}^{N} F(w, z_i)$$

$$\longrightarrow \nabla \hat{f}_n(w) = \frac{1}{N} \sum_{i=1}^{N} \nabla_w F(w, z_i)$$

$\Rightarrow$ O(n) computational & storage cost per update

③ SGD $\longrightarrow$ using $\boxed{\begin{array}{c} \nabla_w F(w, z_I) \\ I \sim uniform[1,....N] \end{array}}$ $\approx \boxed{\nabla_w \hat{f}_n(w)}$

if we fix $z_1,...., z_n$.
then this is a constant

$\Rightarrow$ only use one data point per update ( O(1) )

consider the relationship:

a) $\nabla_w \hat{f}_n(w) \longleftrightarrow \nabla_w f(w)$

R.V. with respect to $\{z_1, \ldots, z_N\}$

$$\Rightarrow \mathbb{E}_z[\nabla_w \hat{f}_n(w)] = \mathbb{E}_{z_i}[\nabla_w F(w, z_i)]$$

☆

$$\text{(under regularity)} = \nabla_w \mathbb{E}_{z_i}[F(w, z_i)]$$

$$= \nabla_w f(w)$$

unbiased estimator

b) $\nabla_w \hat{f}_{SGD}(w) := \boxed{\nabla_w F(w, z_I)} \longleftrightarrow \nabla_w \hat{f}_n(w)$

$\longrightarrow$ stochastic gradient

R.V. with respect to $I \sim \text{uniform}([n])$

$$\Rightarrow \mathbb{E}_I[\nabla_w \hat{f}_{SGD}(w)] = \mathbb{E}_I[\nabla_w F(w, z_I)]$$

$$= \sum_{i=1}^{n} \frac{1}{n} \cdot \nabla_w F(w, z_i)$$

$$= \hat{f}_n(w)$$

$\Rightarrow$ Stochastic gradient descent (SGD) update:

$$w_{k+1} = w_k - \alpha_k \cdot \nabla_w F(w_k, z_{I_k})$$

where $I_k \sim \text{uniform}(1, 2, \ldots, N) \longleftrightarrow$ stochastic!

☆

for each update, we only require $O(1)$ computation !!!

Previously: $w_{k+1} = w_k - \alpha_k \cdot \frac{1}{N} \sum_{i=1}^{N} \nabla_w F(w_k, z_i)$

# 3. Convergence Result of SGD

Recap: under convexity regularity of $f(\cdot)$, GD method can achieve convergence to minimizer with sufficiently small step length.

Analysis:

$\rightarrow f(\omega)$ in optimization framework

$$\underbrace{\nabla F(\omega, z_I)}_{\rightarrow \text{ R.V. with respect to } I \sim \text{unif}(1,...,N)} = \nabla \boxed{\hat{f}_n(\omega)} + \zeta \qquad \mathbb{E}_I[\zeta] = 0$$

Assume: $\mathbb{E}_I[\zeta^2] \leq \sigma^2$

Thm:

$f$ is $c$-strongly convex, $\nabla f$ is $L$-Lipschitz.

then with $\boxed{\text{fixed step length } \alpha} \leq C/L^2$, we have, $Q$-linear

$\boxed{\text{Linear Convergence}}$

$$\longrightarrow \mathbb{E}\|\omega_n - \omega^*\|_2^2 \leq (1 - c\alpha)^n \, \mathbb{E}\|\omega_0 - \omega^*\|_2^2 + \boxed{\frac{\sigma^2 \alpha}{c}}$$

$\boxed{\omega^* = \arg\min_\omega f(\omega)}$

$\boxed{\text{Trade off}}$

noise

intuitively, this term comes from all the randomness of sampling for each update

Note: $f(\omega) = \dfrac{1}{N} \sum\limits_{i=1}^{N} F(\omega, z_i)$

$\boxed{\text{Pf Sketch:}}$ ① $\omega_n - \omega^* = \omega_{n-1} - \omega^* - \alpha \nabla_\omega F(\omega_{n-1}, z_{I_{n-1}})$

$$= \omega_{n-1} - \omega^* - \alpha \nabla_\omega f(\omega_{n-1})$$
$$- \alpha \zeta_{n-1}$$

② $\mathbb{E}_{I_{n-1}} \|\omega_n - \omega^*\|_2^2$ (conditioned on $\omega_{n-1}$)

$$= \boxed{\|\omega_{n-1} - \omega^* - \alpha \nabla_\omega f(\omega_{n-1})\|_2^2} \rightarrow \text{previous result}$$

$$+ \boxed{\mathbb{E}_{I_{n-1}}[\zeta_{n-1}] \cdot *} \longrightarrow 0$$

$$+ \boxed{\alpha^2 \, \mathbb{E}_{I_{n-1}}[\zeta_{n-1}^2]} \leq \alpha^2 6^2$$

$$① + ② \Rightarrow \mathbb{E}_{I_{n-1}} \|\omega_n - \omega^*\|_2^2 \leq (1 - \alpha c) \|\omega_{n-1} - \omega^*\|_2^2$$

$$+ \alpha^2 6^2$$

$$\Rightarrow \mathbb{E}_{I_{n-2}, I_{n-1}} \left[ \|\omega_n - \omega^*\|_2^2 \mid \omega_{n-2} \right]$$

$$\leq (1 - \alpha c)^2 \|\omega_{n-2} - \omega^*\|_2^2 + \left[ (1 - \alpha c) + 1 \right] \alpha^2 6^2$$

$$\Rightarrow \cdots \quad \mathbb{E}_{I_0, \cdots I_{n-1}} \left[ \|\omega_n - \omega^*\|_2^2 \right]$$

$$= \mathbb{E}_{I_0} \left[ \mathbb{E}_{I_1, \cdots, I_{n-1}} \left[ \|\omega_n - \omega^*\|_2^2 \mid \omega_0 \right] \right]$$

$$\leq (1 - \alpha c)^n \|\omega_0 - \omega^*\|_2^2 + \left[ 1 + \cdots + (1 - \alpha c)^{n-1} \right] \alpha^2 6^2$$

$$= (1 - \alpha c)^n \|\omega_0 - \omega^*\|_2^2 + \frac{1 - (1 - \alpha c)^n}{\alpha c} \cdot \alpha^2 6^2$$

$$= (1 - \alpha c)^n \|\omega_0 - \omega^*\|_2^2 + \frac{\alpha}{c} \cdot 6^2$$

$$\#$$

Rmk: There is <mark>no guarantee</mark> that $\omega_n \longrightarrow \omega^*$ as $n \to +\infty$

since we only have $\mathbb{E}_I \|\omega_n - \omega^*\|_2^2 \leq \frac{6^2}{c} \cdot \alpha$ as $n \to +\infty$

# 4. Convergence Result of SGD with different step length

↳ not fixed step length

Notation

$$\begin{cases} S_{1,n} = \sum_{k=1}^{n} \alpha_k \\ S_{2,n} = \sum_{k=1}^{n} \alpha_k^2 \end{cases}$$

**Theorem:**

→ if $f$ is $c$-strongly convex & $\nabla f$ is $L$-Lip.

then with step length $\alpha_k \leq c/L^2$, we have:

→ no need to be constant

$$\Rightarrow \min_{k \in [n]} \mathbb{E}_I \| w_k - w^* \|_2^2 \leq \frac{\mathbb{E}_I \| w_0 - w^* \|_2^2 + \sigma^2 S_{2,n}}{c S_{1,n}}$$

ISSUE: convergence might be slow!

☆ **Rmk:** based on this thm, we would like to choose $\{\alpha_k\}_{k=1}^{n}$

such that $\begin{cases} S_{1,\infty} = \sum_{k=1}^{\infty} \alpha_k = \infty \\ S_{2,\infty} = \sum_{k=1}^{\infty} \alpha_k^2 < \infty \end{cases}$ to guarantee convergence
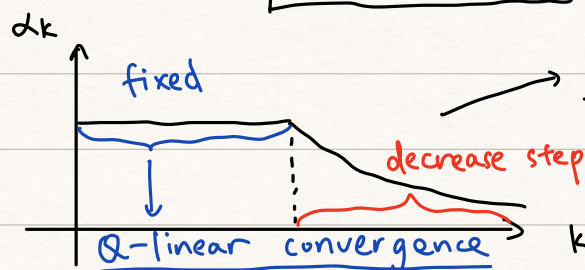
get rid of noise!

$$\lim_{n \to \infty} \min_{k \in [n]} \mathbb{E}_I \| w_k - w^* \|_2^2 = 0$$

→ choice: $\alpha_k = \alpha_0 \frac{1}{k^t}$ $\quad t \in (\frac{1}{2}, 1]$ → thm guarantee

$\quad \downarrow$

$$\alpha_k = \alpha_0 \cdot \frac{1}{k}$$

→ choice 2:

Fixed + decrease schedule

decrease step length ⇒ guarantee convergence



fixed

Q-linear convergence

$\alpha_k$ (vertical axis), $k$ (horizontal axis)

# 5. Mini-batch SGD

① $\quad \nabla_w F(w, z_I) = \nabla_w \hat{f}_n(w) + \zeta \qquad \underline{\mathbb{E}_I [\zeta] = 0}$

$$\boxed{\mathbb{E}_I [\zeta^2] \longleftrightarrow \sigma^2}$$

<u>previously</u>, we have the result:

$$\Downarrow$$

$$\underline{\mathbb{E}_I \| w_n - w^* \|_2^2 \leq (1 - \alpha \cdot c)^n \, \mathbb{E}_I \| w_0 - w^* \|_2^2 + \frac{\sigma^2}{c} \cdot \alpha}$$

$\longrightarrow$ to <u>improve</u> <u>performance</u> $(w_n \to w^*)$, we can <mark>decrease</mark> $\sigma^2$

$\longrightarrow$ <mark>idea:</mark> previously, we use $\nabla_w F(w, z_I) \approx \nabla_w \hat{f}_n(w)$

$\quad \Big\downarrow$ <u>Bagging</u> !!!

why not use $\quad \frac{1}{B} \sum\limits_{b=1}^{B} \nabla_w F(w, z_{I_b}) \approx \nabla_w \hat{f}_n(w)$

$$\downarrow$$

$\Big\{$ same expectation as $\nabla_w F(w, z_I)$

but with $\frac{1}{B}$ variance $\left(\frac{\sigma^2}{B}\right)$

② <u>mini-batch SGD</u> : $\boxed{1 \ll B \ll N}$ $\Big\{$ <mark>$B \ll N$ : computation</mark>

$\qquad\qquad \Big\downarrow \qquad\qquad\qquad\qquad\qquad$ <mark>$1 \ll B$ : reduce var.</mark>

$$w_{k+1} = w_k - \alpha_k \cdot \frac{1}{B} \sum\limits_{b=1}^{B} \nabla_w F(w, z_{I_k^b})$$

$$\boxed{I_k^1 \cdots\cdots I_k^B \sim \text{Uniform}(1, 2, \ldots, N)}$$

# 6. Momentum GD → utilize momentum information in previous step

## Framework:

Goal: $\hat{w} = \underset{w}{\arg\min} f(w)$

Update: $w_{k+1} = w_k - \alpha_k \cdot m_k$

where $m_k = \beta \cdot m_{k-1} + (1-\beta) \cdot \nabla_w f(w)$

$\boxed{\beta \in (0, 1)}$

$m_{k-1} \to \boxed{\text{previous state}}$

Rmk:

① common choice of $\beta \to 0.9$

② variants of momentum SGD → $\begin{cases} ADAM \\ AdaGrad \end{cases}$

③ converge faster since it can accumulate "speed" ✰

④ can help to escape bad local minima

⑤ Momentum SGD can be viewed as:

→ increase batch size (since it takes previous update into consideration) to reduce variance