

Fundamentals of Machine Learning Over Networks

**Group 6
HW 1**

4-1

HW1 - Group 6

if $f(y) \succeq f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} \|y-x\|^2 \Rightarrow \nabla^2 f(x) \succeq M I, \forall x \in X$

if $g(t) \triangleq f((1-t)x + ty) \quad t \in [0,1] \rightarrow g'(t) = \nabla f((1-t)x + ty)^T (y-x)$
 $g''(t) = (y-x)^T \nabla^2 f((1-t)x + ty) (y-x)$

according to Taylor series:

$$g(t) = g(0) + g'(0)t + \frac{g''(0)}{2}t^2 + \frac{g'''(t)}{6}t^3$$

$$f((1-t)x + ty) = f(x) + \nabla f(x)^T (y-x)t + \frac{(y-x)^T \nabla^2 f(x) (y-x)}{2}t^2 + o(t^2) \quad (I)$$

from the \succeq condition:

$$f(z) \succeq f(x) + \nabla f(x)^T (z-x) + \frac{M}{2} \|z-x\|^2$$

$$\hookrightarrow f(z) \succeq f(x) + \nabla f(x)^T (y-x)t + t^2 \cdot \frac{M}{2} \|y-x\|^2 \quad (II)$$

$$I, II \Rightarrow (y-x)^T \nabla^2 f(x) (y-x) \frac{t^2}{2} \geq t^2 \cdot \frac{M}{2} \|y-x\|^2$$

$$\Rightarrow (y-x)^T \nabla^2 f(x) (y-x) - M (y-x)^T I (y-x) \geq 0$$

$A \succeq 0 \Leftrightarrow x^T A x \geq 0$ ↑ symmetric

$$\Rightarrow (y-x)^T (\nabla^2 f(x) - M I) (y-x) \geq 0 \Rightarrow \nabla^2 f(x) \succeq M I$$

4-2

$$(*) \quad g(y) \geq g(x) + \nabla g(x)^T (y-x) \quad g(x) = f(x) - \frac{\mu}{2} \|x\|^2$$

$$f(y) - \frac{\mu}{2} \|y\|^2 \geq f(x) - \frac{\mu}{2} \|x\|^2 + (\nabla f(x) - \mu x)^T (y-x)$$

$$f(y) \geq \frac{\mu}{2} (\|x\|^2 + \|y\|^2) + \nabla f(x)^T (y-x) - \mu x^T y + \mu \frac{x^T x}{\|x\|^2} + f(x)$$

$$f(y) \geq \frac{\mu}{2} \|x\|^2 + \frac{\mu}{2} \|y\|^2 + \nabla f(x)^T (y-x) - \mu x^T y + f(x)$$

$$f(y) \geq \frac{\mu}{2} \|x-y\|^2 + \nabla f(x)^T (y-x) + f(x)$$

f is strong convex $\Rightarrow (*)$ is valid.

in convex we have

$$(\nabla g(y) - \nabla g(x))^T (y-x) \geq 0$$

$$\rightarrow (\nabla f(y) - \mu y - \nabla f(x) + \mu x)^T (y-x) \geq 0$$

$$(\nabla f(y) - \nabla f(x))^T (y-x) \geq \mu y^T (y-x) - \mu x^T (y-x)$$

$$\quad \quad \quad \geq \mu \|y\|^2 - \mu y^T x - \mu x^T y + \mu \|x\|^2$$

$$(\nabla f(y) - \nabla f(x))^T (y-x) \geq \mu \|y-x\|^2$$

4-a

$$f(x) - f^* \leq \frac{1}{2M} \|\nabla f(x)\|_2^2, \forall x$$

we have

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|_2^2$$

gradient (f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|_2^2) = 0

$$\nabla f(x) + \mu(y-x) = 0$$

$$y = -\frac{\nabla f(x)}{\mu} + x$$

$$\rightarrow f^* \geq f(x) + \nabla f(x)^T (x - \frac{1}{\mu} \nabla f(x) + x) + \frac{\mu}{2} \|x - \frac{1}{\mu} \nabla f(x) + x\|_2^2$$

$$f(x) - f^* \leq \frac{1}{\mu} \|\nabla f(x)\|^2 - \frac{\mu}{2} \|\frac{1}{\mu} \nabla f(x)\|_2^2$$

$$f(x) - f^* \leq \frac{1}{2\mu} \|\nabla f(x)\|_2^2$$

4-b

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|_2^2 \xrightarrow{?} \|y-x\|_2 \leq \frac{1}{\mu} \|\nabla f(y) - \nabla f(x)\|_2$$

according to proof 4-2 in [P2] for strong convexity we have:

$$(\nabla f(y) - \nabla f(x))^T (y-x) \geq \mu \|y-x\|_2^2$$

according to holder's inequality $\|fg\|_1 \leq \|f\|_p \|g\|_q$ $\frac{1}{p} + \frac{1}{q} = 1$

$$\Rightarrow \|\nabla f(y) - \nabla f(x)\|_2 \|y-x\|_2 \geq (\nabla f(y) - \nabla f(x))^T (y-x) \geq \mu \|y-x\|_2^2$$

$$\Rightarrow \frac{1}{\mu} \|\nabla f(y) - \nabla f(x)\|_2 \geq \|y-x\|_2$$

[P3]

4-c

Prove $(\nabla f(y) - \nabla f(x))^T (y - x) \leq \frac{1}{\mu} \|\nabla f(y) - \nabla f(x)\|_2^2, \forall x, y$

$$\phi_\mu(z) = f(z) - \nabla f(x)^T z$$

first we prove that $\phi_\mu(z)$ is strong convex.

we use 4-2 in [P2] for strong convexity

$$\begin{aligned} & (\nabla \phi_\mu(z_2) - \nabla \phi_\mu(z_1))^T (z_2 - z_1) \geq \frac{1}{\mu} \|\nabla \phi_\mu(z_2) - \nabla \phi_\mu(z_1)\|_2^2 \quad (*) \\ & \left[(\nabla f(z_2) - \nabla f(x))^T - \nabla f(z_1)^T + \nabla f(x)^T \right]^T (z_2 - z_1) \\ & = (\nabla f(z_2) - \nabla f(z_1))^T (z_2 - z_1) \geq \mu \|z_2 - z_1\|_2^2 \end{aligned}$$

Since f is strong convex \Rightarrow inequality is established.

Since $\phi_\mu(z)$ is strong convex \rightarrow PL-inequality (4-a, [P3]) with $z^* = x$ can be applied.

$$\phi_\mu(y) - \phi_\mu(x) \leq \frac{1}{2\mu} \|\nabla \phi_\mu(y)\|_2^2$$

$$\begin{aligned} & \rightarrow (f(y) - \nabla f(x)^T y - f(x) + \nabla f(x)^T x) \leq \frac{1}{2\mu} \|\nabla f(y) - \nabla f(x)\|_2^2 \\ & \rightarrow f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{1}{2\mu} \|\nabla f(y) - \nabla f(x)\|_2^2 \end{aligned}$$

now interchanging x and y

$$f(x) \leq f(y) + \nabla f(y)^T (x - y) + \frac{1}{2\mu} \|\nabla f(x) - \nabla f(y)\|_2^2$$

Summing 2, 3 results

$$(\nabla f(y) - \nabla f(x))^T (y - x) \leq \frac{1}{\mu} \|\nabla f(y) - \nabla f(x)\|_2^2 \quad \checkmark$$

4-d

$$\begin{aligned} f \text{ is convex} & \rightarrow f(y) \geq f(x) + \nabla f(x)^T (y - x) \\ r \text{ is strong convex} & \rightarrow r(y) \geq r(x) + \nabla r(x)^T (y - x) + \frac{\mu}{2} \|y - x\|_2^2 \\ & = f(y) + r(y) \geq f(x) + r(x) + (\nabla f(x)^T + \nabla r(x)^T)(y - x) + \frac{\mu}{2} \|y - x\|_2^2 \\ & \rightarrow \text{Definition of convex} \end{aligned}$$

$$2.a) f(x_2) \leq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{L}{2} \|x_2 - x_1\|_2^2, \quad \forall x_1, x_2$$

we define g as $\Rightarrow g(t) = f(x_1 + t(x_2 - x_1)), \quad t \in \mathbb{R}$

HW1-B-a

so we have $\Rightarrow f(x_2) - f(x_1) = g(1) - g(0) = \int_0^1 g'(t) dt =$

$$\int_0^1 (x_2 - x_1) \nabla f(x_1 + t(x_2 - x_1))^T dt = \underbrace{(x_2 - x_1) \nabla f(x_1)^T}_{\text{add both } \oplus \text{ \& } \ominus \text{ part, } \ominus \text{ is under integral}} +$$

$$(x_2 - x_1) \left[\int_0^1 [\nabla f(x_1 + t(x_2 - x_1)) - \nabla f(x_1)]^T dt \right] \leq \text{quach-schwarz inequality}$$

$$(x_2 - x_1) \nabla f(x_1)^T + \|x_2 - x_1\|_2 \left\| \int_0^1 [\nabla f(x_1 + t(x_2 - x_1)) - \nabla f(x_1)] dt \right\|_2$$

$$\leq (x_2 - x_1) \nabla f(x_1)^T + \|x_2 - x_1\|_2 \int_0^1 \underbrace{\|\nabla f(x_1 + t(x_2 - x_1)) - \nabla f(x_1)\|_2}_{L\text{-smooth} \Rightarrow \mathbb{L} \|x_1 + t(x_2 - x_1) - x_1\|_2} dt$$

$$\leq (x_2 - x_1) \nabla f(x_1)^T + \|x_2 - x_1\|_2 \int_0^1 L t \|x_2 - x_1\|_2 dt$$

$$\leq (x_2 - x_1) \nabla f(x_1)^T + \frac{L}{2} \|x_2 - x_1\|_2^2 \Rightarrow$$

$$f(x_2) - f(x_1) \leq (x_2 - x_1) \nabla f(x_1)^T + \frac{L}{2} \|x_2 - x_1\|_2^2$$

HW1(b)-(b):

For any $z \in \mathbb{R}^n$, we have

$$f \text{ is convex} \Rightarrow f(z) \geq f(x) + \nabla f(x)^T (z-x)$$

$$f \text{ is } L\text{-smooth} \Rightarrow f(z) \leq f(x) + \nabla f(x)^T (z-x) + \frac{L}{2} \|z-x\|_2^2 \quad \leftarrow \text{HW1(b)-(a)}$$

$$f(x_2) - f(x_1) = f(x_2) - f(z) + f(z) - f(x_1) \geq -\nabla f(x_2)^T (z-x_2) - \frac{L}{2} \|z-x_2\|_2^2 + \nabla f(x_1)^T (z-x_1)$$

$$= -\frac{L}{2} \|z-x_2\|_2^2 + (\nabla f(x_1) - \nabla f(x_2))^T (z-x_2) + \nabla f(x_1)^T (x_2-x_1)$$

$$= -\left\| \sqrt{\frac{L}{2}} (z-x_2) - \frac{1}{\sqrt{2L}} (\nabla f(x_1) - \nabla f(x_2)) \right\|_2^2 + \frac{1}{2L} \|\nabla f(x_1) - \nabla f(x_2)\|_2^2 + \nabla f(x_1)^T (x_2-x_1)$$

$$\stackrel{\downarrow}{=} \nabla f(x_1)^T (x_2-x_1) + \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \Rightarrow \checkmark$$

$$\text{for } z = x_2 + \frac{1}{L} (\nabla f(x_1) - \nabla f(x_2))$$

2.c) $(\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) \geq \frac{1}{L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2, \forall x_1, x_2$

we use smoothness for both $f(x_1)$ & $f(x_2) \Rightarrow$

$$f(x_2) - f(x_1) \geq \nabla f(x_1)^T (x_2 - x_1) + \frac{1}{2L} \|\nabla f(x_1) - \nabla f(x_2)\|_2^2$$

$$f(x_1) - f(x_2) \geq \nabla f(x_2)^T (x_1 - x_2) + \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2$$

we add both
sides \Rightarrow

$$\longrightarrow (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) \geq \frac{1}{L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2$$

HW1-Bc

HW1(c) -

(Group 6)

Consider:

$$\min \frac{1}{N} \sum_{i \in [N]} f_i(x_i)$$

$$\text{s.t. } Ax = b$$

$$\text{for } A \in \mathbb{R}^{p \times N} \text{ and } x = [x_1, \dots, x_N]^T.$$

a) Assume: strong convexity and smoothness; $N = 1000$

To solve (a) first we form a Lagrangian (dual) function:

$$L(x, \lambda) = \frac{1}{N} \sum_{i \in [N]} f_i(x_i) + \lambda(b - Ax); \text{ where } \lambda \text{ is Lagrangian multiplier.}$$

Since the Problem is now unconstrained, we can use the descent methods.

$$\nabla L = \frac{1}{N} \begin{bmatrix} \frac{\partial f_1}{\partial x_1} - \lambda a_1 \\ \frac{\partial f_2}{\partial x_2} - \lambda a_2 \\ \vdots \\ \frac{\partial f_n}{\partial x_n} - \lambda a_n \end{bmatrix}$$

still

a) For $N = 1000$; Since the Problem dimension is ^{still} small ($N = 1000$), we can use the Gradient descent method.
(GD)

b) Due to the difficulty of finding Proper co-ordinate for high-dimensional Problems GD is not a good choice for $N = 10^3$ any more. Here, for $N \geq 10^3$, we can use Newton method since the Hessian matrix is diagonal and finding the inverse is easy.

$$\nabla^2 L = \begin{bmatrix} \frac{\partial^2 f_1}{\partial x_1^2} & & & \\ & \frac{\partial^2 f_2}{\partial x_2^2} & & \\ & & \ddots & \\ & & & \frac{\partial^2 f_n}{\partial x_n^2} \end{bmatrix}$$

c) for the case $p_2 1, N_2 1$? \rightarrow Yes we can use Newton method because the Hessian ^{matrix} is diagonal and the inverse is easy to find.

for $1 < p < N, N=1 \rightarrow$ still Newton method is acceptable since:

$$\nabla^2 L = \frac{1}{N} \begin{bmatrix} \frac{\partial^2 f_1}{\partial x_1^2} & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & \frac{\partial^2 f_N}{\partial x_N^2} & 0 \end{bmatrix} - \begin{bmatrix} \lambda_1 a_{11} + \dots + \lambda_p a_{p1} \\ \vdots \\ \lambda_{N-p+1} a_{N-p+1, N} + \dots + \lambda_p a_{pN} \end{bmatrix}$$

again $\nabla^2 L$ is diagonal.

d)

d-a) $L(x, \lambda) = \frac{1}{N} \sum f_i(x_i) + r(x) + \lambda(b - Ax)$

$$\nabla L = \frac{1}{N} \begin{bmatrix} \frac{\partial f_1}{\partial x_1} + \frac{\partial r}{\partial x_1} - \lambda a_1 \\ \frac{\partial f_2}{\partial x_2} + \frac{\partial r}{\partial x_2} - \lambda a_2 \\ \vdots \\ \frac{\partial f_N}{\partial x_N} + \frac{\partial r}{\partial x_N} - \lambda a_N \end{bmatrix} \rightarrow \text{GD can solve the Problem.}$$

d-b) $N_2 1, p_2 1$

$$\nabla^2 L_2 = \begin{bmatrix} \frac{\partial^2 f_1}{\partial x_1^2} & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & \dots & \frac{\partial^2 f_N}{\partial x_N^2} & 0 \end{bmatrix} + \begin{bmatrix} \frac{\partial^2 r}{\partial x_1^2} & \frac{\partial^2 r}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 r}{\partial x_1 \partial x_N} \\ \frac{\partial^2 r}{\partial x_1 \partial x_2} & \frac{\partial^2 r}{\partial x_2^2} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 r}{\partial x_1 \partial x_N} & \dots & \dots & \frac{\partial^2 r}{\partial x_N^2} \end{bmatrix} \quad (2)$$

d-b) cont.

Because $\nabla^2 L$ is a Hermitian and Positive semi definite, we can use

Cholesky decomposition to decompose $\nabla^2 L$ into a lower Rank matrix and

its conjugate transpose. Hence, finding the inverse is easy and Newton works well.

d-c) still it is possible to use Newton. Because, we can use Cholesky decomposition and finding the inverse of lower rank matrix is easy.

In case of P_2 we can have a closed expression of inverse matrix

and in case of $1 < P \leq N$, still the matrix has a good structure for finding the inverse.

$$(w) \Rightarrow 4) (\nabla f(x) - \nabla f(y))^T (x-y) \geq \frac{\mu L}{\mu+L} \|x-y\|_2^2 + \frac{1}{\mu+L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

we define $g(x)$ as $g(x) = f(x) - \frac{\mu}{2} \|x\|^2 \rightarrow g$ is convex, $\nabla g(x) = \nabla f(x) - \mu x$

$g(x)$ is also L -smooth \rightarrow with parameter $L-\mu$
 e.g. $\Rightarrow g(x)$ is $(L-\mu)$ -smooth \Rightarrow [proof is ~~also~~ in the last page]

$$(\nabla g(x) - \nabla g(y))^T (x-y) \geq \frac{1}{L-\mu} \|\nabla g(x) - \nabla g(y)\|_2^2 \xrightarrow{\text{replace}}$$

$$(\nabla f(x) - \mu x - (\nabla f(y) - \mu y))^T (x-y) \geq \frac{1}{L-\mu} \|\nabla f(x) - \nabla f(y) - \mu(x-y)\|_2^2$$

we use the result of HW1, prob. 2 \rightarrow

$$\Rightarrow \left((\nabla f(x) - \nabla f(y)) - \mu(x-y) \right)^T \left((\nabla f(x) - \nabla f(y)) - \mu(x-y) \right) =$$

$$\|\nabla f(x) - \nabla f(y)\|_2^2 + \mu^2 \|x-y\|_2^2 - \mu(x-y)^T (\nabla f(x) - \nabla f(y))$$

$$- \mu(x-y)(\nabla f(x) - \nabla f(y))^T = \|\nabla f(x) - \nabla f(y)\|_2^2 + \mu^2 \|x-y\|_2^2 - 2\mu(x-y)(\nabla f(x) - \nabla f(y))^T$$

$$\Rightarrow (\nabla f(x) - \nabla f(y))^T (x-y) - \mu \|x-y\|_2^2 \geq \frac{1}{L-\mu} \left[\|\nabla f(x) - \nabla f(y)\|_2^2 + \mu^2 \|x-y\|_2^2 - 2\mu(x-y)(\nabla f(x) - \nabla f(y))^T \right]$$

$$2\mu(x-y)(\nabla f(x) - \nabla f(y))^T \xrightarrow[\text{simple}]{\text{make it}} (\nabla f(x) - \nabla f(y))^T (x-y) \left[1 + \frac{2\mu}{L-\mu} \right] \geq$$

$$\|x-y\|_2^2 \left(\mu + \frac{\mu^2}{L-\mu} \right) + \frac{1}{L-\mu} \|\nabla f(x) - \nabla f(y)\|_2^2 \left(\frac{L+\mu}{L-\mu} \right) (\nabla f(x) - \nabla f(y))^T (x-y) \geq$$

$$\frac{\mu L}{L-\mu} \|x-y\|_2^2 + \frac{1}{L-\mu} \|\nabla f(x) - \nabla f(y)\|_2^2 \Rightarrow (\nabla f(x) - \nabla f(y))^T (x-y) \geq \checkmark$$

$$\frac{\mu L}{L+\mu} \|x-y\|_2^2 + \frac{1}{L+\mu} \|\nabla f(x) - \nabla f(y)\|_2^2$$

proof of smoothness of $g(x) \Rightarrow$

$$g(x) = f(x) - \frac{\mu}{2} \|x\|_2^2 \Rightarrow \nabla g(x) = \nabla f(x) - \mu x \Rightarrow$$

$$(\nabla g(y) - \nabla g(x)) = (\nabla f(y) - \mu y - \nabla f(x) + \mu x) \Rightarrow$$

$$[\nabla f(y) - \mu y - (\nabla f(x) - \mu x)]^T [\nabla f(y) - \mu y - (\nabla f(x) - \mu x)] =$$

$$\|\nabla f(y) - \nabla f(x)\|_2^2 - 2\mu (y-x)^T (\nabla f(y) - \nabla f(x)) + \mu^2 \|y-x\|_2^2 \leq$$

$$L^2 \|y-x\|_2^2 + \mu^2 \|y-x\|_2^2 - \underbrace{2\mu (y-x)^T [\nabla f(y) - \nabla f(x)]}_{\leq L^2 \|y-x\|_2^2}$$

$$\Rightarrow \|y-x\|_2^2 (L^2 + \mu^2 - 2L^2) \leq \|y-x\|_2^2 (L^2 + \mu^2 - 2\mu L) = \underbrace{(L-\mu)^2}_{\downarrow} \|y-x\|_2^2$$

because $L \geq \mu \Rightarrow -2L^2 \leq -2\mu L$
we have

g is $L-\mu$ -Smooth