

# MLoN Homework1

Group 5

## 1 Problem 1

A differentiable function  $f$  is  $\mu$ -strongly convex iff  $\forall x_1, x_2 \in \chi, \mu > 0$

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

### 1.1

Prove Equation 1 is equivalent to a minimum positive curvature

$$\nabla^2 f(x) \geq \mu I_d, \forall x \in \chi \quad (2)$$

Ⓐ Equation 2  $\Rightarrow$  Equation 1:

Taylor expansion:

$$\begin{aligned} f(x_2) &= f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{2}(x_2 - x_1)^T \nabla^2 f(x_1)(x_2 - x_1) + R_2(X) \\ &= f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{2}(x_2 - x_1)^T \nabla^2 f(x_3)(x_2 - x_1) \\ &\geq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{\mu}{2} \|x_2 - x_1\|_2^2 \end{aligned}$$

Ⓑ Equation 1  $\Rightarrow$  Equation 2:

$$\begin{aligned} f(x_2) &= f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{2}(x_2 - x_1)^T \nabla^2 f(x_1)(x_2 - x_1) + R_2(X) \\ &= f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{2}(x_2 - x_1)^T \nabla^2 f(x_3)(x_2 - x_1) \\ &\geq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{\mu}{2} \|x_2 - x_1\|_2^2 \\ &\Rightarrow \\ &\quad (x_2 - x_1)^T \nabla^2 f(x_3)(x_2 - x_1) \geq \mu \|x_2 - x_1\|_2^2 \\ &\Rightarrow \\ &\quad \nabla^2 f(x) \geq \mu I_d, \forall x \in \chi \end{aligned}$$

Ⓐ & Ⓑ  $\Rightarrow$  Equation 1  $\equiv$  Equation 2

## 1.2

Prove Equation 1 is equivalent to

$$(\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) \geq \mu \|x_2 - x_1\|_2^2 \quad (3)$$

Ⓐ Equation 1  $\Rightarrow$  Equation 3:

$$\begin{aligned} f(x_2) &\geq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{\mu}{2} \|x_2 - x_1\|_2^2 \\ f(x_1) &\geq f(x_2) + \nabla f(x_2)^T (x_1 - x_2) + \frac{\mu}{2} \|x_1 - x_2\|_2^2 \\ \Rightarrow \\ (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) &\geq \mu \|x_2 - x_1\|_2^2 \end{aligned}$$

Ⓑ Equation 3  $\Rightarrow$  Equation 1:

$$\begin{aligned} (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) &\geq \mu \|x_2 - x_1\|_2^2 \\ \Rightarrow \\ (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) &\geq (x_2 - x_1)^T \mu (x_2 - x_1) \\ \Rightarrow \\ [(\nabla f(x_2) - \mu x_2) - (\nabla f(x_1) - \mu x_1)]^T (x_2 - x_1) &\geq 0 \\ \Rightarrow \\ (\nabla g(x_2) - \nabla g(x_1))^T (x_2 - x_1) &\geq 0, \\ g(x) &= f(x) - \frac{\mu}{2} \|x\|^2 \\ \Rightarrow \\ g(x) &\text{ is convex (the monotone gradient condition for convexity).} \\ g(x) &\text{ is convex iff. } (g(x_2) - g(x_1))^T (x_2 - x_1) \geq 0. \\ \Rightarrow \\ g(x_2) &\geq g(x_1) + \nabla g(x_1)^T (x_2 - x_1) \\ \Rightarrow \\ f(x_2) - \frac{\mu}{2} \|x_2\|^2 &\geq f(x_1) - \frac{\mu}{2} \|x_1\|^2 + (\nabla f(x_1) - \mu x_1)^T (x_2 - x_1) \\ \Rightarrow \\ f(x_2) &\geq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{\mu}{2} \|x_2\|^2 + \frac{\mu}{2} \|x_1\|^2 - \mu x_1^T x_2 \\ \Rightarrow \\ f(x_2) &\geq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{\mu}{2} \|x_2 - x_1\|_2^2 \end{aligned}$$

Ⓐ & Ⓑ  $\Rightarrow$  Equation 1  $\equiv$  Equation 3

### 1.3 a

Implies Polyak-Łojasiewicz (PL) Inequality

$$f(x) - f^* \leq \frac{1}{2\mu} \|\nabla f(x)\|_2^2, \forall x \quad (4)$$

Taking minimization respect to  $x_2$  on both sides of Equation 1:

$$\begin{aligned} f(x_2) &\geq f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{\mu}{2} \|x_2 - x_1\|_2^2 \\ \Rightarrow \\ f(x^*) &\geq f(x_1) + \nabla f(x_1)^T (x^* - x_1) + \frac{\mu}{2} \|x^* - x_1\|_2^2 \\ (f'(x^*) = 0 &\Rightarrow x^* = x_1 - \frac{1}{\mu} \nabla f(x_1)) \\ \Rightarrow \\ f(x^*) &\geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|_2^2 \\ \Rightarrow \\ &\text{Equation 4} \end{aligned}$$

### 1.4 b

Implies

$$\|x_2 - x_1\|_2 \leq \frac{1}{\mu} \|\nabla f(x_2) - \nabla f(x_1)\|_2, \forall x_1, \forall x_2 \quad (5)$$

Applying Cauchy-Schwarz inequality  $|\langle u, v \rangle| \leq \|u\| \|v\|$  on the equivalent condition Equation 3.

$$\begin{aligned} (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) &\geq \mu \|x_2 - x_1\|_2^2 \\ \Rightarrow \\ \|\nabla f(x_2) - \nabla f(x_1)\|_2 \|x_2 - x_1\|_2 &\geq \mu \|x_2 - x_1\|_2^2 \\ \Rightarrow \\ &\text{Equation 5} \end{aligned}$$

### 1.5 c

Implies

$$(\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1) \leq \frac{1}{\mu} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2, \forall x_1, \forall x_2 \quad (6)$$

Multiplying Equation 5 and  $\|\nabla f(x_2) - \nabla f(x_1)\|_2$   
 $\Rightarrow$   
 $\|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \geq \mu \|x_2 - x_1\|_2 \|\nabla f(x_2) - \nabla f(x_1)\|_2$   
Applying Cauchy–Schwarz inequality  
 $\Rightarrow$   
 $\|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \geq \mu (\nabla f(x_2) - \nabla f(x_1))^T (x_2 - x_1)$   
 $\Rightarrow$   
Equation 6

## 1.6 d

Implies  $f(x) + r(x)$  is strongly convex for any convex  $f$  and strongly convex  $r$ .

$f$  is convex:

$$f(y) \geq f(x) + \nabla f(x)^T (x_2 - x_1) \quad (7)$$

$r$  is strongly convex:

$$r(y) \geq r(x) + \nabla r(x)^T (x_2 - x_1) + \frac{\mu}{2} \|y - x\|^2 \quad (8)$$

$$\begin{aligned} g(x) &= f(x) + r(x) \\ &\geq f(x) + \nabla f(x)^T (x_2 - x_1) + r(x) + \nabla r(x)^T (x_2 - x_1) + \frac{\mu}{2} \|y - x\|^2 \\ &= g(x) + \nabla g(x)^T (x_2 - x_1) + \frac{\mu}{2} \|y - x\|^2 \\ \Rightarrow \\ g(x) &= f(x) + r(x) \text{ is strongly convex.} \end{aligned}$$

## 2 Problem 2

A function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -smooth iff it is differentiable and its gradient is  $L$ -Lipschitz-continuous (usually w.r.t. norm-2):

$$\forall x_1, \forall x_2 \in \mathbb{R}^d, \|\nabla f(x_2) - \nabla f(x_1)\|_2 \leq L \|x_2 - x_1\|_2 \quad (9)$$

### 2.1 a

Implies

$$f(x_2) \leq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{L}{2} \|x_2 - x_1\|_2^2 \quad (10)$$

$f$  is convex and  $L$ -smooth.

$$\begin{aligned} & \|\nabla f(x_2) - \nabla f(x_1)\|_2 \leq L \|x_2 - x_1\|_2 \\ \Rightarrow & \\ & \|\nabla f(x_2) - \nabla f(x_1)\|_2 \cdot \|x_2 - x_1\|_2 \leq L \|x_2 - x_1\|_2^2 \\ & \text{(Applying Cauchy-Schwartz)} \\ \Rightarrow & \\ & (\nabla f(x_2) - \nabla f(x_1))^T(x_2 - x_1) \leq L \|x_2 - x_1\|_2^2 \\ \Rightarrow & \\ & (\nabla f(x_2) - \nabla f(x_1))^T(x_2 - x_1) \leq L \|x_2 - x_1\|_2^2 = L(x_2 - x_1)^T(x_2 - x_1) \\ \Rightarrow & \\ & (Lx_2 - \nabla f(x_2) - Lx_1 + \nabla f(x_1))^T(x_2 - x_1) \geq 0 \\ \Rightarrow & \\ & (\nabla g(x_2) - \nabla g(x_1))^T(x_2 - x_1) \geq 0 \text{ (define } g(x) = \frac{L}{2}x^Tx - f(x)) \\ \Rightarrow & \\ & g(x) \text{ is convex (the monotone gradient condition).} \\ \Rightarrow & \\ & g(x_2) \geq g(x_1) + \nabla g(x_1)^T(x_2 - x_1) \\ \Rightarrow & \\ & \frac{L}{2}x_2^Tx_2 - f(x_2) \geq \frac{L}{2}x_1^Tx_1 - f(x_1) + (Lx_1 - \nabla f(x_1))^T(x_2 - x_1) \\ \Rightarrow & \\ & f(x_2) \leq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{L}{2}x_2^Tx_2 + \frac{L}{2}x_1^Tx_1 - Lx_1^Tx_2 \\ \Rightarrow & \\ & f(x_2) \leq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{L}{2} \|x_2 - x_1\|_2^2 \end{aligned}$$

## 2.2 b

Implies

$$f(x_2) \geq f(x_1) + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \quad (11)$$

$f$  is convex and  $L$ -smooth.

$$f(x_2) - f(x_1) = f(x_2) - f(x_3) + f(x_3) - f(x_1)$$

$$(\text{define } x_3 = x_2 - \frac{1}{L}(\nabla f(x_1) - \nabla f(x_2)))$$

$\Rightarrow$

$$f(x_2) - f(x_3) \geq -\nabla f(x_2)^T(x_3 - x_2) - \frac{L}{2} \|x_2 - x_3\|_2^2 \quad (\text{quadratic upper bound})$$

$$f(x_3) - f(x_1) \geq \nabla f(x_1)^T(x_3 - x_1) \quad (\text{convex})$$

$\Rightarrow$

$$f(x_2) - f(x_1) \geq -\nabla f(x_2)^T(x_3 - x_2) - \frac{L}{2} \|x_2 - x_3\|_2^2 + \nabla f(x_1)^T(x_3 - x_1)$$

$\Rightarrow$

$$\begin{aligned} f(x_2) - f(x_1) &\geq -\frac{1}{L} \nabla f(x_2)^T(\nabla f(x_1) - \nabla f(x_2)) - \frac{1}{2L} \|\nabla f(x_1) - \nabla f(x_2)\|_2^2 \\ &\quad + \nabla f(x_1)^T(x_2 - x_1) + \frac{1}{L} \nabla f(x_1)^T(\nabla f(x_1) - \nabla f(x_2)) \quad (\text{substitute } z) \end{aligned}$$

$\Rightarrow$

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

## 2.3 c

Implies

$$(\nabla f(x_2) - \nabla f(x_1))^T(x_2 - x_1) \geq \frac{1}{L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \quad (12)$$

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

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

$\Rightarrow$

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

### 3 Problem 3

Define, discuss the benefits, and give examples for the different convergence rates of a sequence of updates  $x_k$ : (a) Sublinear, (b) Linear, (c) SuperLinear, (d) Quadratic.

#### 3.1 Sublinear

The sequence of updates  $x_k$  sublinearly converges to  $L$  if satisfies the following equation:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - L|}{|x_k - L|} = 1$$

In this case,  $x_k$  converges slowly. For example, the  $x_k$  below converges to 0 sublinearly:

$$x_k = \frac{1}{k+1}$$

Because we have:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - 0|}{|x_k - 0|} = \lim_{k \rightarrow \infty} \frac{\left| \frac{1}{k+2} \right|}{\left| \frac{1}{k+1} \right|} = \lim_{k \rightarrow \infty} \frac{\left| \frac{\partial}{\partial k} \left( \frac{1}{k+2} \right) \right|}{\left| \frac{\partial}{\partial k} \left( \frac{1}{k+1} \right) \right|} = \lim_{k \rightarrow \infty} \frac{|(k+2)^2|}{|(k+1)^2|} = 1$$

#### 3.2 Linear

The sequence of updates  $x_k$  linearly converges to  $L$  if satisfies the following equation:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - L|}{|x_k - L|} = \mu, \mu \in (0, 1)$$

In this case,  $x_k$  converges faster than the case with sublinear convergence. For example, the  $x_k$  below converges to 0 linearly:

$$x_k = \left( \frac{1}{2} \right)^k$$

Because we have:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - 0|}{|x_k - 0|} = \lim_{k \rightarrow \infty} \frac{\left| \left( \frac{1}{2} \right)^{k+1} \right|}{\left| \left( \frac{1}{2} \right)^k \right|} = \lim_{k \rightarrow \infty} \frac{1}{2} = \frac{1}{2}$$

### 3.3 Superlinear

The sequence of updates  $x_k$  sublinearly converges to  $L$  if satisfies the following equation:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - L|}{|x_k - L|} = 0$$

In this case,  $x_k$  converges faster than the case with linear convergence. For example, the  $x_k$  below converges to 0 superlinearly:

$$x_k = \left(\frac{1}{2}\right)^{2^k}$$

Because we have:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - 0|}{|x_k - 0|} = \lim_{k \rightarrow \infty} \frac{\left|\left(\frac{1}{2}\right)^{2^{k+1}}\right|}{\left|\left(\frac{1}{2}\right)^{2^k}\right|} = \lim_{k \rightarrow \infty} \left|\left(\frac{1}{2}\right)^{2^k}\right| = 0$$

### 3.4 Quadratic

Quadratic convergence is one type of superlinear convergence. The sequence of updates  $x_k$  converges to  $L$  with the order 2. Specifically, it should satisfies the following equation where  $M$  is a positive constant value.

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - L|}{|x_k - L|^2} < M$$

The example we show in the superlinear above also converge to 0 quadratically. Because we have:

$$\lim_{k \rightarrow \infty} \frac{|x_{k+1} - 0|}{|x_k - 0|^2} = \lim_{k \rightarrow \infty} \frac{\left|\left(\frac{1}{2}\right)^{2^{k+1}}\right|}{\left|\left(\frac{1}{2}\right)^{2^k + 2^k}\right|} = \lim_{k \rightarrow \infty} \frac{\left|\left(\frac{1}{2}\right)^{2^{k+1}}\right|}{\left|\left(\frac{1}{2}\right)^{2^{k+1}}\right|} = 1$$

So we can choose any constant value  $M$  larger than 1 to satisfy the definition.

### 3.5 Discussion

In machine learning, in order to find the optimal solution, we are happy when the convergence rate is faster. In other words, we prefer linear or superlinear convergence rate rather than sublinear convergence rate. From the machine learning aspect, we want to minimize the loss function  $f(\omega)$  according to the specific problem where  $\omega$  is the model parameter. We denote the optimal solution as  $\omega^*$ . A variation of definition of convergence rate is:

$$\lim_{k \rightarrow \infty} \frac{f(\omega^{k+1}) - f(\omega^*)}{f(\omega^k) - f(\omega^*)} = \rho$$



$$\begin{aligned}\rho = 1 &\rightarrow \textit{sublinear} \\ \rho \in (0, 1) &\rightarrow \textit{linear} \\ \rho = 0 &\rightarrow \textit{superlinear}\end{aligned}$$

For some nice designed loss functions  $f$ , gradient descent algorithm can have a linear convergence rate. A good example is least square problem. When we have a convex function  $f$ , we could introduce an L2 regularization term to make  $f$  to be strongly-convex which has Polyak-Lojasiewicz (PL) inequality property. It implies that we could improve gradient methods from sublinear rate to a linear rate. At the same time, when  $f$  is fixed, we may try other optimization algorithms like Newton's method which achieves superlinear convergence rate (however requires expensive computation when the dimension is large).

## 4 Problem 4

Consider

$$\begin{aligned} & \text{minimize } \frac{1}{N} \sum_{i \in [N]} f_i(x_i) \\ & \text{s.t. } \mathbf{Ax} = \mathbf{b} \end{aligned} \quad (13)$$

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

### 4.1 a,b,c

We have strong-convexity and smoothness on  $f$ , assuming  $f$  is twice differentiable then we could calculate the Hessian matrix of  $f(\mathbf{x}) = \frac{1}{N} \sum_{i \in [N]} f_i(x_i)$ :

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} \nabla^2 f_1(x_1) & 0 & \dots & 0 \\ 0 & \nabla^2 f_2(x_2) & & \\ \vdots & & \ddots & \\ 0 & & & \nabla^2 f_N(x_N) \end{bmatrix} \quad (14)$$

since it is diagonal matrix, the inverse of Hessian matrix is fast no matter when  $N = 1000$  or  $N = 10^9$ . We could use Newton's method with equality constraints to solve it:

1. select a starting point  $\mathbf{x}$  that satisfies  $\mathbf{Ax} = \mathbf{b}$
2. Compute the Newton step  $\Delta \mathbf{x}_{nt}$  by solving:

$$\begin{bmatrix} \nabla^2 f(\mathbf{x}) & \mathbf{A}^T \\ \mathbf{A} & 0 \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x}_{nt} \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}) \\ 0 \end{bmatrix} \quad (15)$$

where  $w$  is the associated optimal dual variable for the quadratic problem. And compute  $\lambda^2(\mathbf{x})$  as:

$$\lambda^2(\mathbf{x}) = \Delta \mathbf{x}_{nt}^T \nabla^2 f(\mathbf{x}) \Delta \mathbf{x}_{nt} \quad (16)$$

3. Quit if  $\lambda^2(\mathbf{x})/2 \leq \epsilon$
4. Update  $\mathbf{x} = \mathbf{x} + \Delta \mathbf{x}_{nt}$

The computing Hessian matrix is efficient since it is diagonal.

If  $f$  is not twice differentiable, we could only calculate the gradient of  $f$ , thus we could use projected gradient descent to solve it.

### 4.2 d

At this point, the objective function  $f(\mathbf{x}) = \frac{1}{N} \sum_{i \in [N]} f_i(x_i) + r(\mathbf{x})$ . Since we do not know if  $\frac{\partial^2 r(\mathbf{x})}{\partial x_i \partial x_j}$  is zero, then the computing of Hessian matrix and its inverse for each iteration would become too expensive. So we could use quasi-Newton's method to update the Hessian by analyzing successive gradients.

## 5 Problem 5

### Problem

In the convergence proof of GD with constant step size and strongly convex objective function (see slides), prove the coercivity of the gradient:

$$\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq \frac{\mu L}{\mu + L} \|x - y\|_2^2 + \frac{1}{\mu + L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2 \quad (17)$$

### Proof

Extend function  $f$  to:

$$g(x) = f(x) + \frac{\mu}{2} \|x\|^2$$

According to  $\mu$ -convexity of  $f$ ,  $g$  is convex which is proven in Problem 1. Now the  $L$ -smoothness of  $f$  can be extended to  $L - \mu$  smoothness of  $g$  according to the definition of Lipschitz-smoothness. Then:

$$\begin{aligned} \langle \nabla g(x) - \nabla g(y), x - y \rangle &\geq \frac{1}{(L - \mu)} \|\nabla g(x) - \nabla g(y)\|_2^2 \\ &\Rightarrow \\ \langle \nabla f(x) - \nabla f(y) - \mu(x - y), x - y \rangle &\geq \frac{1}{(L - \mu)} \|\nabla f(x) - \nabla f(y) - \mu(x - y)\|_2^2 \\ &\Rightarrow (Denote(*)) = \langle \nabla f(x) - \nabla f(y), x - y \rangle \\ (*) - \mu \|x - y\|_2^2 &\geq \frac{1}{(L - \mu)} \|\nabla f(x) - \nabla f(y)\|_2^2 + \frac{\mu^2}{(L - \mu)} \|x - y\|_2^2 - 2 \frac{\mu}{L - \mu} (*) \\ &\Rightarrow \\ \frac{L + \mu}{L - \mu} (*) &\geq \frac{1}{(L - \mu)} \|\nabla f(x) - \nabla f(y)\|_2^2 + \frac{\mu L}{(L - \mu)} \|x - y\|_2^2 \\ &\Rightarrow \text{Strong coercivity} \end{aligned}$$