

Problem 1. A differentiable function f is μ -strongly convex iff $\forall x_1, x_2 \in \mathcal{X}, \mu > 0$

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

Then proof that:

• (1) is equivalent to a minimum positive curvature $\nabla^{2} f(\mathbf{x}) \succeq \mu \mathbf{I}_{d}, \quad \forall \mathbf{x} \in \mathcal{X};$

Proof. According to (1), we have

$$f\left(oldsymbol{x}_{2}
ight)\geq f\left(oldsymbol{x}_{1}
ight)+
abla f\left(oldsymbol{x}_{1}
ight)^{T}\left(oldsymbol{x}_{2}-oldsymbol{x}_{1}
ight)+rac{\mu}{2}\left\|oldsymbol{x}_{1}
ight\|_{2}^{2}+rac{\mu}{2}\left\|oldsymbol{x}_{2}
ight\|_{2}^{2}-\muoldsymbol{x}_{1}^{T}oldsymbol{x}_{2},$$

or

$$f(x_2) - \frac{\mu}{2} ||x_2||_2^2 \ge f(x_1) - \frac{\mu}{2} ||x_1||_2^2 + (\nabla f(x_1) - \mu x_1)^T (x_2 - x_1).$$

Lets define $g(\mathbf{x}) \triangleq f(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2$. We have $\nabla g(\mathbf{x}) = \nabla f(\mathbf{x}) - \mu \mathbf{x}$ and $\nabla^2 g(\mathbf{x}) = \nabla^2 f(\mathbf{x}) - \mu \mathbf{I}_d$, therefore

$$g(\mathbf{x}_2) \ge g(\mathbf{x}_1) + \nabla g(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1),$$
 (2)

which is the first-order condition for convexity of $g(\mathbf{x})$. Hence, $\nabla^2 g(\mathbf{x}) \succeq \mathbf{0}$, or equivalently, $\nabla^2 f(\mathbf{x}) - \mu \mathbf{I}_d \succeq \mathbf{0}$. This implies $\nabla^2 f(\mathbf{x}) \succeq \mu \mathbf{I}_d$.

• (1) is equivalent to $(\nabla f(\boldsymbol{x}_2) - \nabla f(\boldsymbol{x}_1))^T (\boldsymbol{x}_2 - \boldsymbol{x}_1) \ge \mu \|\boldsymbol{x}_2 - \boldsymbol{x}_1\|_2^2$

Proof. We know that $g(\mathbf{x}) \triangleq f(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2$ is a convex function. It follows from the monotone gradient condition for convexity that $g(\mathbf{x})$ is convex if and only if $(\nabla g(\mathbf{x}_2) - \nabla g(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \geq 0$. Therefore

$$\left(\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)-\mu\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\geq0,$$

or

$$\left(\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\geq\mu\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2}^{2}.$$

• (1) implies

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

Proof. By taking minimization with respect to x_2 from both sides of (1), we have

$$f^* \ge f(x_1) + \nabla f(x_1)^T (x_2^* - x_1) + \frac{\mu}{2} ||x_2^* - x_1||_2^2,$$
 (3)

where $\nabla f(\boldsymbol{x}_1) + \mu(\boldsymbol{x}_2^* - \boldsymbol{x}_1) = 0$, or $\boldsymbol{x}_2^* = -\frac{1}{\mu} \nabla f(\boldsymbol{x}_1) + \boldsymbol{x}_1$. By substituting \boldsymbol{x}_2^* in (3), we have

$$f^* \ge f\left(oldsymbol{x}_1
ight) +
abla f\left(oldsymbol{x}_1
ight)^T \left(-rac{1}{\mu}
abla f\left(oldsymbol{x}_1
ight) + oldsymbol{x}_1 - oldsymbol{x}_1
ight) + rac{\mu}{2} \left\|-rac{1}{\mu}
abla f\left(oldsymbol{x}_1
ight) + oldsymbol{x}_1 - oldsymbol{x}_1
ight\|_2^2$$

$$= f\left(oldsymbol{x}_1
ight) - rac{1}{2\mu} \left\|
abla f\left(oldsymbol{x}_1
ight)
ight\|_2^2,$$



or

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

(b) $\|x_2 - x_1\|_2 \le \frac{1}{\mu} \|\nabla f(x_2) - \nabla f(x_1)\|_2$, $\forall x_1, x_2$;

Proof. Using Cauchy-Schwartz inequality on the equivalent condition $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \ge \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$, we have

$$\left\|\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2} \geq \left(\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right) \geq \mu\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2}^{2}.$$
 (4)

Dividing both sides by $\|\boldsymbol{x}_2 - \boldsymbol{x}_1\|_2$ (assuming $\boldsymbol{x}_2 \neq \boldsymbol{x}_1$) gives

$$\|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \ge \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2.$$
 (5)

 $\left(\mathrm{c}\right)\ \left(\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\leq\frac{1}{\mu}\left\|\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2},\quad\forall\boldsymbol{x}_{1},\boldsymbol{x}_{2};$

Proof. From (4) and (5) we have

$$\frac{\left(\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)}{\left\|\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}}\leq\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2}\leq\frac{1}{\mu}\left\|\nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}.$$

Hence,
$$\left(\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\right) \leq \frac{1}{\mu}\left\|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2}$$
.

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

Proof. Assuming differentiability of f(x) and r(x), we have by definition of convexity of f and strong convexity of r

$$f(\boldsymbol{x}_2) \ge f(\boldsymbol{x}_1) + \nabla f(\boldsymbol{x}_1)^T (\boldsymbol{x}_2 - \boldsymbol{x}_1),$$
 (6)

$$r(x_2) \ge r(x_1) + \nabla r(x_1)^T (x_2 - x_1) + \frac{\mu}{2} ||x_2 - x_1||_2^2.$$
 (7)

By summing (6) and (7), we have

$$f\left(\boldsymbol{x}_{2}\right)+r\left(\boldsymbol{x}_{2}\right)\geq f\left(\boldsymbol{x}_{1}\right)+r\left(\boldsymbol{x}_{1}\right)+\nabla(r\left(\boldsymbol{x}_{1}\right)+f\left(\boldsymbol{x}_{1}\right))^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)+\frac{\mu}{2}\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2}^{2},$$

which is the definition of strong convexity for $f(\mathbf{x}) + r(\mathbf{x})$.



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

$$\forall \boldsymbol{x}_{1}, \boldsymbol{x}_{2} \in \mathbb{R}^{d}, \|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\|_{2} \leq L \|\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\|_{2}. \tag{8}$$

For all x_1, x_2 , prove that (8) implies

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

Proof. To begin, let's define the helper function $g(\mathbf{x}) = \frac{L}{2}\mathbf{x}^T\mathbf{x} - f(\mathbf{x})$ and prove that it is convex. $g(\mathbf{x})$ is convex iff the first-order condition holds [1]

$$g(\mathbf{x}_2) \ge g(\mathbf{x}_1) + \nabla g(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1).$$
 (9)

Of course, the variables in (9) can be switched to create

$$g(\mathbf{x}_1) \ge g(\mathbf{x}_2) + \nabla g(\mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2). \tag{10}$$

Taking $g(\mathbf{x}_1)$ from (10) and inserting it into (9) gives

$$g(x_2) \ge (g(x_2) + \nabla g(x_2)^T (x_1 - x_2)) + \nabla g(x_1)^T (x_2 - x_1),$$

and after rearranging that becomes

$$(\nabla g(\mathbf{x}_2) - \nabla g(\mathbf{x}_1))^T(\mathbf{x}_2 - \mathbf{x}_1) \ge 0.$$
(11)

Inserting the definition of the helper function into (11) gives

$$(L\boldsymbol{x}_2 - \nabla f(\boldsymbol{x}_2) - (L\boldsymbol{x}_1 - \nabla f(\boldsymbol{x}_1))^T(\boldsymbol{x}_2 - \boldsymbol{x}_1) \ge 0.$$

This can be rearranged to resemble (8):

$$L \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \ge (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1).$$
 (12)

Multiplying both sides of (8) with $\|x_2 - x_1\|_2$ and then using Cauchy-Schwarz inequality on the L.H.S. gives (12). Thus, (8) implies (12), which in turn proves that the first-order conditions for g(x) holds.

Using the now established convexity of $g(\mathbf{x}) = \frac{L}{2} \mathbf{x}^T \mathbf{x} - f(\mathbf{x})$, we can insert it into the first-order

condition (9) to show the final result: $\frac{L}{2}\boldsymbol{x}_{2}^{T}\boldsymbol{x}_{2} - f\left(\boldsymbol{x}_{2}\right) \geq \frac{L}{2}\boldsymbol{x}_{1}^{T}\boldsymbol{x}_{1} - f\left(\boldsymbol{x}_{1}\right) + (L\boldsymbol{x}_{1} - \nabla f\left(\boldsymbol{x}_{1}\right))^{T}(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}).$

Rearranging terms and noting that $\frac{L}{2}(\boldsymbol{x}_{2}^{T}\boldsymbol{x}_{2}+\boldsymbol{x}_{1}^{T}\boldsymbol{x}_{1})-L\boldsymbol{x}_{1}^{T}\boldsymbol{x}_{2}=\frac{L}{2}\left\|\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right\|_{2}^{2}\text{ yields }f\left(\boldsymbol{x}_{2}\right)\leq f\left(\boldsymbol{x}_{1}\right)+L\left(\boldsymbol{x}_{1}^{T}\boldsymbol{x}_{2}\right)+L\left(\boldsymbol{x}_{1}^{T}$

$$abla f\left(oldsymbol{x}_{1}
ight)^{T}\left(oldsymbol{x}_{2}-oldsymbol{x}_{1}
ight)+rac{L}{2}\left\|oldsymbol{x}_{2}-oldsymbol{x}_{1}
ight\|_{2}^{2}.$$



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

Proof. Let us consider f is convex. Therefore, the function $\Phi_{x_1}(z) \triangleq f(z) - \nabla f(x_1)^T z$, will also be convex and its optimum occurs at $z^* = x_1$. Moreover, we know $\nabla \Phi_{x_1}(z) = \nabla f(z) - \nabla f(x_1)$. Hence,

$$\nabla \Phi_{\boldsymbol{x}_{1}}(\boldsymbol{z}) - \nabla \Phi_{\boldsymbol{x}_{1}}(\boldsymbol{x}_{2}) = \nabla f(\boldsymbol{z}) - \nabla f(\boldsymbol{x}_{2}). \tag{13}$$

By substituting (13) in (8), we have

$$\|\nabla \Phi_{x_1}(z) - \nabla \Phi_{x_1}(x_2)\|_2 \le L \|z - x_2\|_2.$$
 (14)

According to the part (a), (14) is equivalent to

$$\Phi_{oldsymbol{x}_1}\left(oldsymbol{z}
ight) \leq \Phi_{oldsymbol{x}_1}\left(oldsymbol{x}_2
ight) +
abla \Phi_{oldsymbol{x}_1}\left(oldsymbol{x}_2
ight)^T \left(oldsymbol{z} - oldsymbol{x}_2
ight) + rac{L}{2} \left\|oldsymbol{z} - oldsymbol{x}_2
ight\|_2^2.$$

Taking minimization with respect to z on both sides, yields,

$$\begin{split} f\left(\boldsymbol{x}_{1}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)^{T} \boldsymbol{x}_{1} &\leq f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)^{T} \boldsymbol{x}_{2} + \nabla \Phi_{\boldsymbol{x}_{1}} \left(\boldsymbol{x}_{2}\right)^{T} \left(\boldsymbol{z}^{*} - \boldsymbol{x}_{2}\right) + \frac{L}{2} \left\|\boldsymbol{z}^{*} - \boldsymbol{x}_{2}\right\|_{2}^{2} \\ &= f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)^{T} \boldsymbol{x}_{2} + \nabla \Phi_{\boldsymbol{x}_{1}} \left(\boldsymbol{x}_{2}\right)^{T} \left(-\frac{1}{L} \nabla \Phi_{\boldsymbol{x}_{1}} \left(\boldsymbol{x}_{2}\right)\right) + \frac{L}{2} \left\|-\frac{1}{L} \nabla \Phi_{\boldsymbol{x}_{1}} \left(\boldsymbol{x}_{2}\right)\right\|_{2}^{2} \\ &= f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)^{T} \boldsymbol{x}_{2} - \frac{1}{2L} \left\|\nabla \Phi_{\boldsymbol{x}_{1}} \left(\boldsymbol{x}_{2}\right)\right\|_{2}^{2} \\ &= f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)^{T} \boldsymbol{x}_{2} - \frac{1}{2L} \left\|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2} \end{split}$$

Re-arranging gives

$$f(x_2) \ge f(x_1) + \nabla f(x_1)^T (x_2 - x_1) + \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2.$$
 (15)

(c) $\left(\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\right) \geq \frac{1}{L}\left\|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2}$

Proof. Using the result from b), we know that (8) implies (15). Using (15) as a starting point, we begin by rearranging the terms to get

$$\nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \le f(\mathbf{x}_2) - f(\mathbf{x}_1) - \frac{1}{2L} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2.$$
 (16)

Then, we switch x_2 with x_1 in (16) to get:

$$\nabla f(x_2)^T (x_1 - x_2) \le f(x_1) - f(x_2) - \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2.$$
 (17)

The sum of (16) and (17) is

$$\left(\nabla f\left(\boldsymbol{x}_{1}\right)-\nabla f\left(\boldsymbol{x}_{2}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\leq-\frac{1}{L}\left\Vert \nabla f\left(\boldsymbol{x}_{2}\right)-\nabla f\left(\boldsymbol{x}_{1}\right)\right\Vert _{2}^{2}.$$

Multiplying both sides by -1 gives the final result

$$\left(\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\right) \geq \frac{1}{L}\left\|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2}.$$
(18)

To summarize, b) shows that (8) implies (15), and c) shows that (15) implies (18), thus (8) implies (18). \Box

Assume convexity if needed.



Problem 3. Define, discuss the benefits, and give examples for the different convergence rates of a sequence of updates $\{x_k\}$:

First, we define a limit which will be useful to express the convergence rate of a sequence of updates $\{x_k\}$. Let's assume that the sequence will eventually converge, so that

$$\lim_{k\to\infty} \boldsymbol{x}_k = \boldsymbol{x}^*$$

holds. Then, the following expression can be used to represent how quickly the sequence converges:

$$\lim_{k \to \infty} \frac{|\boldsymbol{x}_{k+1} - \boldsymbol{x}^*|}{|\boldsymbol{x}_k - \boldsymbol{x}^*|^p} = m.$$
(19)

The values of p and m determines the convergence rate of the sequence.

(a) Sublinear

Sublinear is the slowest form of convergence, and is defined by a series which converge according to (19) with m=1 and p=1. For simplicity, the example will be a series with scalar $\{x_k\}$. The following sequence has sublinear convergence rate:

$$a_k = \frac{1}{k}. (20)$$

This sequence converges to zero, we can verify the sublinear convergence by inserting it into (19) with p = 1:

$$\lim_{k \to \infty} \frac{1/(k+1) - 0}{1/k - 0} = \lim_{k \to \infty} \frac{k}{k+1} = 1.$$

(b) Linear

Linear convergence is faster than sublinear, and is defined by a series which converge according to (19) with 0 < m < 1 and p = 1. Gradient descent over an L-smooth and μ -strongly convex function has a linear convergence rate with

$$m = 1 - \frac{2}{1 + L/\mu}.$$

Since $\mu < L$ by definition, we see that 0 < m < 1 which makes the convergence linear.

(c) Superlinear

Superlinear convergence is faster than linear, and is defined by a series which converge according to (19) with m = 0 and p = 1. To find an example of a sequence with superlinear convergence, we need a series where the ratio of two adjacent updates goes to zero as k approaches infinity. An example of such a series is:

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

This sequence converges to one, we can verify the superlinear convergence by inserting it into (19) with p = 1:

$$\lim_{k \to \infty} \frac{1 + 1/(k+1)^{k+1} - 1}{1 + 1/k^k - 1} = \lim_{k \to \infty} \frac{k^k}{(k+1)^{k+1}} = 0.$$

(d) Quadratic

Quadratic convergence the fastest of the four forms of convergence covered here, and is defined by a series which converge according to (19) with m > 0 and p = 2. Given that x_k is sufficiently close to the minimum, the Newton method experiences quadratic convergence. However, in practice, another algorithm has to be ran for a couple of steps before the Newton method can be applied to ensure that the current x_k is sufficiently close to the minimum.



Problem 4. Consider

minimize
$$\frac{1}{N} \sum_{i \in [N]} f_i(x_i)$$
 subject to
$$A\mathbf{x} = \mathbf{b},$$

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

(a) Assume strong-convexity and smoothness on f. How would you solve this problem when N = 1000?

Proof. This is a convex optimization problem with equality constraint. There are a few ways to handle it:

(a) Eliminating equality constraints. This can be done by finding a matrix F and vector \hat{x} that parameterize the feasible set:

$$\{\boldsymbol{x}|A\boldsymbol{x}=\boldsymbol{b}\} = \{F\boldsymbol{z} + \hat{\boldsymbol{x}}|\boldsymbol{z} \in \mathbb{R}^{n-p}\}$$
(22)

Then, the problem can be transformed into

$$\min \hat{f}(z) = f(Fz + \hat{x}). \tag{23}$$

For the above problem, we can use gradient descent to solve it. Since f is strong-convex and smooth, we could achieve linear convergence.

(b) Solving equality constrained problems via the dual. The dual problem is given by

$$g(\mathbf{v}) = -\mathbf{b}^T \mathbf{v} - f^* \left(-A^T \mathbf{v} \right) \tag{24}$$

(c) We could solve the problem directly using Newton's method with equality constraints. In this case, the Newton step Δx_{nt} is characterized by

$$\begin{bmatrix} \nabla^2 f(\boldsymbol{x}) & A^T \\ A & 0 \end{bmatrix} \begin{bmatrix} \Delta x_{nt} \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(\boldsymbol{x}) \\ 0 \end{bmatrix}$$
 (25)

(d) We could also solve it using ADMM. This is because the variables in the objective function is separate. Therefore, we could update one variable at one time, while fixing the others.

(b) What if $N = 10^9$?

Proof. In this case, Newton's method and ADMM can still be applied.

(c) Can we use Newton's method for $N = 10^9$? Try efficient method for computing $\nabla^2 f(x_k)$ for p = 1 and b = 1 (probability simplex constraint). Extend it to $1 \le p \ll N$.

Proof. Since x_i is separate in the objective function, the Hessian matrix $\nabla^2 f(x)$ is a diagonal matrix, i.e.,

$$\nabla^{2} f(\boldsymbol{x}) = \frac{1}{N} \begin{bmatrix} \frac{\partial^{2} f_{1}}{\partial x_{1}^{2}} & \cdots & 0 \\ 0 & \frac{\partial^{2} f_{i}}{\partial x_{i}^{2}} & 0 \\ 0 & \cdots & \frac{\partial^{2} f_{N}}{\partial x_{i}^{2}} \end{bmatrix}$$
(26)

Therefore, we can still use Newton's method.

(d) Now, add twice differentiable r(x) to the objective and solve (a)–(c).

Proof. Here it is unclear whether $r(\mathbf{x})$ is convex or not! If $r(\mathbf{x})$ is convex, the similar procedure can be performed after adding $r(\mathbf{x})$.



Problem 5. In the convergence proof of GD with constant step size and strongly convex objective function (see slides), prove the coercivity of the gradient:

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

Proof. Let us consider the function $g(\boldsymbol{x}) = f(\boldsymbol{x}) - \frac{\mu}{2} \|\boldsymbol{x}\|_2^2$. As already noted in Problem 1 Equation (2), since $f(\boldsymbol{x})$ is μ -strongly convex, the function $g(\boldsymbol{x})$ is convex, i.e.

$$g(\mathbf{x}_2) \ge g(\mathbf{x}_1) + \nabla g(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1),$$
 (27)

Also, we observe that since f(x) is L-smooth, we have that g(x) is also smooth with smoothness parameter $L - \mu$:

$$\|\nabla g(\mathbf{x}_{2}) - \nabla g(\mathbf{x}_{1})\|_{2} = \|\nabla f(\mathbf{x}_{2}) - \nabla f(\mathbf{x}_{1}) - \mu(\mathbf{x}_{2} - \mathbf{x}_{1})\|_{2}$$

$$\leq \|\nabla f(\mathbf{x}_{2}) - \nabla f(\mathbf{x}_{1})\|_{2} - \mu \|\mathbf{x}_{2} - \mathbf{x}_{1}\|_{2}$$

$$\leq (L - \mu) \|\mathbf{x}_{2} - \mathbf{x}_{1}\|_{2}$$
(28)

$$\left(\nabla g\left(\boldsymbol{x}_{2}\right)-\nabla g\left(\boldsymbol{x}_{1}\right)\right)^{T}\left(\boldsymbol{x}_{2}-\boldsymbol{x}_{1}\right)\geq\frac{1}{L-\mu}\left\|\nabla g\left(\boldsymbol{x}_{2}\right)-\nabla g\left(\boldsymbol{x}_{1}\right)\right\|_{2}^{2}$$
(29)

By expressing it in function of f(x), we have

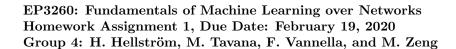
$$\left(\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right) - \mu(\boldsymbol{x}_{2} - \boldsymbol{x}_{1})\right)^{T} \left(\boldsymbol{x}_{2} - \boldsymbol{x}_{1}\right) \geq \frac{1}{L - \mu} \left\|\nabla f\left(\boldsymbol{x}_{2}\right) - \nabla f\left(\boldsymbol{x}_{1}\right)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

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

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

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

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





References

[1] Boyd, Stephen, Stephen P. Boyd, and Lieven Vandenberghe. Convex optimization. Cambridge university press, 2004.