

Homework 1

Ali Bemani
Oscar Bautista Gonzalez
Martin Hellkvist
Syed Aqeel Raza

Ali.Bemani@hig.se
Oscar.Bautista.Gonzalez@hig.se
Martin.Hellkvist@angstrom.uu.se
s.aqeelraza@gmail.com

Solutions are provided in formatted form for problems 1.1, 1.3 and 1.4. Due to time limitations we decided to not format the solutions for 1.2 and 1.5. These solutions are provided within this document, as photocopies of handwritten notes.

Problem 1.1

A differentiable function f is μ -strongly convex iff $\forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}, \mu > 0$

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2. \quad (1)$$

We here provide solutions for the questions of Problem 1.1 in the order they appear in the instructions.

(i) Prove that (1) is equivalent to a minimum positive curvature $\nabla^2 f(\mathbf{x}) \succcurlyeq \mu \mathbf{I}, \forall \mathbf{x} \in \mathcal{X}$.

Solution: We first assume $\mathcal{X} \subset \mathbb{R}^n$. For a twice differentiable function f we have from Taylor's theorem that

$$f(\mathbf{x}_2) = f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{1}{2} (\mathbf{x}_2 - \mathbf{x}_1)^T \nabla^2 f(\mathbf{x}_3) (\mathbf{x}_2 - \mathbf{x}_1), \quad (2)$$

for some \mathbf{x}_3 on the interval $[\mathbf{x}_1, \mathbf{x}_2]$.

Combining (1) with (2) we obtain

$$\begin{aligned} f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{1}{2} (\mathbf{x}_2 - \mathbf{x}_1)^T \nabla^2 f(\mathbf{x}_3) (\mathbf{x}_2 - \mathbf{x}_1) &\geq \\ f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2, \end{aligned} \quad (3)$$

which is equivalent to

$$(\mathbf{x}_2 - \mathbf{x}_1)^T (\nabla^2 f(\mathbf{x}_3) - \mu \mathbf{I}) (\mathbf{x}_2 - \mathbf{x}_1) \geq 0, \quad (4)$$

or

$$\nabla^2 f(\mathbf{x}) \succcurlyeq \mu \mathbf{I}, \quad \forall \mathbf{x} \in \mathcal{X}. \quad (5)$$

We have now proved that $\nabla^2 f(\mathbf{x}) \succcurlyeq \mu \mathbf{I} \Leftrightarrow (1)$.

(ii) Prove that (1) is equivalent to $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \geq \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$

Solution: From (1) we have

$$f(\mathbf{x}_2) - f(\mathbf{x}_1) - \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \geq \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2, \quad (6)$$

$$f(\mathbf{x}_1) - f(\mathbf{x}_2) - \nabla f(\mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2) \geq \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2, \quad (7)$$

If we sum (6) and (7), we obtain

$$\mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \leq (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1), \quad (8)$$

which concludes the proof.

(iii) Prove that (1) implies

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

Solution: (1) provides that $f(\mathbf{x}_2)$ is convex w.r.t. any fixed $\mathbf{x}_1 = \mathbf{x}$. We take the gradient of the righthand side of (1) w.r.t. \mathbf{x} and set it to zero to find the minimizing \mathbf{x}^* (minimizing the righthand side) as

$$\nabla f(\mathbf{x}) + \frac{\mu}{2} (2\mathbf{x}^* - 2\mathbf{x}) = 0, \quad (9)$$

$$\Rightarrow \mathbf{x}^* = \mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x}). \quad (10)$$

From (1) we have

$$f^* \geq f(\mathbf{x}^*) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 \quad (11)$$

$$\begin{aligned} &= f(\mathbf{x}) + \nabla f(\mathbf{x})^T \left(\mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x}) - \mathbf{x} \right) \\ &\quad + \frac{\mu}{2} \left\| \mathbf{x} - \frac{1}{\mu} \nabla f(\mathbf{x}) - \mathbf{x} \right\|_2^2 \end{aligned} \quad (12)$$

$$= f(\mathbf{x}) - \frac{1}{\mu} \nabla f(\mathbf{x})^T \nabla f(\mathbf{x}) + \frac{\mu}{2} \left\| \frac{1}{\mu} \nabla f(\mathbf{x}) \right\|_2^2 \quad (13)$$

$$= f(\mathbf{x}) - \frac{1}{\mu} \|\nabla f(\mathbf{x})\|_2^2 + \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2 \quad (14)$$

$$= f(\mathbf{x}) - \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2, \quad (15)$$

for any $\mathbf{x} \in \mathcal{X}$. Rearranging the terms concludes the proof:

$$f(\mathbf{x}) - f^* \leq \frac{1}{2\mu} \|\nabla f(\mathbf{x})\|_2^2. \quad (16)$$

$$(b) \quad \|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2, \quad \forall \mathbf{x} \in \mathcal{X}.$$

Solution: From (ii), we have $(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \geq \mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2$. Together with Cauchy-Schwarz inequality $\mathbf{v}^T \mathbf{w} \leq \|\mathbf{v}\|_2 \|\mathbf{w}\|_2$ we obtain

$$\mu \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \leq \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \|\mathbf{x}_2 - \mathbf{x}_1\|_2. \quad (17)$$

Dividing by $\|\mathbf{x}_2 - \mathbf{x}_1\|_2$ and assuming $\mathbf{x}_2 \neq \mathbf{x}_1$ so not to have $\|\mathbf{x}_2 - \mathbf{x}_1\|_2 = 0$ concludes the proof:

$$\|\mathbf{x}_2 - \mathbf{x}_1\|_2 \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2. \quad (18)$$

$$(c) (\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2, \forall \mathbf{x}_1, \mathbf{x}_2.$$

Solution: Cauchy-Schwarz together with the result in (iii) (b) gives

$$(\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1))^T (\mathbf{x}_2 - \mathbf{x}_1) \leq \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \|\mathbf{x}_2 - \mathbf{x}_1\|_2 \quad (19)$$

$$\leq \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2 \quad (20)$$

$$= \frac{1}{\mu} \|\nabla f(\mathbf{x}_2) - \nabla f(\mathbf{x}_1)\|_2^2 \quad (21)$$

$$(d) f(\mathbf{x}) + r(\mathbf{x}) \text{ is strongly convex for any convex } f \text{ and strongly convex } r.$$

Solution: Assume r is μ -strongly convex. By the first-order condition on convexity we have for f that

$$f(\mathbf{x}_2) \geq f(\mathbf{x}_1) + \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1). \quad (22)$$

By (1) we have for $r(\mathbf{x})$ that

$$r(\mathbf{x}_2) \geq r(\mathbf{x}_1) + \nabla r(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2. \quad (23)$$

Let $g(\mathbf{x})$ denote $g(\mathbf{x}) = f(\mathbf{x}) + r(\mathbf{x})$. It follows that $\nabla g(\mathbf{x}) = \nabla f(\mathbf{x}) + \nabla r(\mathbf{x})$. By the addition of (22) and (23), we have for $g(\mathbf{x})$ that

$$g(\mathbf{x}_2) = f(\mathbf{x}_2) + r(\mathbf{x}_2) \quad (24)$$

$$\begin{aligned} &\geq f(\mathbf{x}_1) + r(\mathbf{x}_1) \\ &+ \nabla f(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \nabla r(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) \end{aligned} \quad (25)$$

$$\begin{aligned} &+ \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2 \\ &= g(\mathbf{x}_1) + \nabla g(\mathbf{x}_1)^T (\mathbf{x}_2 - \mathbf{x}_1) + \frac{\mu}{2} \|\mathbf{x}_2 - \mathbf{x}_1\|_2^2. \end{aligned} \quad (26)$$

Which proves that $f(\mathbf{x}) + r(\mathbf{x})$ is strongly convex if f is convex and r is strongly convex.

problem 1.2

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):

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

For all x_1, x_2 , prove that

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

Solution:

Let define $g(t) \triangleq f(x_1 + t(x_2 - x_1))$

we know that:

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

It then follows that:

$$f(x_2) - f(x_1) - \nabla f(x_1)^T (x_2 - x_1) = \int_0^1 \nabla f(x_1 + t(x_2 - x_1))^T (x_2 - x_1) dt - \nabla f(x_1)^T (x_2 - x_1)$$

$$\Rightarrow f(x_2) - f(x_1) - \nabla f(x_1)^T (x_2 - x_1) = \int_0^1 (\nabla f(x_1 + t(x_2 - x_1)) - \nabla f(x_1))^T (x_2 - x_1) dt$$

From the Cauchy-Schartz inequality we can get:

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

f has L -Lipschitz continuous gradient

we then have:

$$\begin{aligned}\Rightarrow f(x_2) - f(x_1) - \nabla f(x_1)^T (x_2 - x_1) &\leq \int_0^1 t L \|x_2 - x_1\|_2 \cdot \|x_2 - x_1\|_2 dt \\ &\leq \int_0^1 t L \|x_2 - x_1\|_2^2 dt \\ &\leq L \|x_2 - x_1\|_2^2 \int_0^1 t dt \\ &\leq \frac{L}{2} \|x_2 - x_1\|_2^2\end{aligned}$$

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

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

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

Let define $g_{x_1}(x_2) \triangleq f(x_2) - f(x_1) - \nabla f^T(x_1)(x_2 - x_1)$

Since f is convex therefore: $f(x_2) \geq f(x_1) + \nabla f^T(x_1)(x_2 - x_1)$
 $\Rightarrow f(x_2) - f(x_1) - \nabla f^T(x_1)(x_2 - x_1) \geq 0$
 $\Rightarrow g_{x_1}(x_2) \geq 0$

In particular $g_{x_1}(x_1) = 0 \Rightarrow g_{x_1}(x_1) = \min_x g_{x_1}(x_2)$

and $\nabla g_{x_1}(x_1) = -\nabla f(x_1) + \nabla f(x_1) = 0$

from the optimality of x_1 , it then follows that

$$\begin{aligned} g_{x_1}(x_1) &\leq \min_{x_2} g_{x_1}(x_2 - \eta \nabla g_{x_1}(x_2)) \\ (*) \quad &= \min_{x_2} f(x_2 - \eta \nabla g_{x_1}(x_2)) - f(x_1) - \nabla f^T(x_1)(x_2 - \eta \nabla g_{x_1}(x_2) - x_1) \end{aligned}$$

By definition of L -smooth we have:

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

In the follows from (*) we have:

$$\begin{aligned} g_{x_1}(x_1) &\leq \min_{x_2} f(x_2) + \nabla f^T(x_2)(-\eta \nabla g_{x_1}(x_2)) + \frac{\lambda}{2} \|\eta \nabla g_{x_1}(x_2)\|_2^2 \\ &\quad - f(x_1) - \nabla f^T(x_1)(x_2 - x_1 - \eta \nabla g_{x_1}(x_2)) \end{aligned}$$

$$g_{x_1}^{(x_1)} \leq \min_{\eta} g_{x_1}^{(x_2)} + \frac{1}{2} \|\eta \nabla g_{x_1}^{(x_2)}\|_2^2 - \eta \nabla g_{x_1}^{(x_2)}^\top (\nabla f(x_2) - \nabla f(x_1))$$

$$g_{x_1}^{(x_1)} \leq \min_{\eta} g_{x_1}^{(x_2)} + \frac{1}{2} \eta^2 \|\nabla g_{x_1}^{(x_2)}\|_2^2 - \eta \|\nabla g_{x_1}^{(x_2)}\|_2^2$$

The minimum solution η to minimize this quadratic problem is:

$$\frac{1}{2} \eta^2 \|\nabla g_{x_1}^{(x_2)}\|_2^2 - \|\nabla g_{x_1}^{(x_2)}\|_2^2 = 0$$

$$\eta^* = \frac{1}{L}$$

$$\Rightarrow \text{minimum solution: } g_{x_1}^{(x_2)} - \frac{1}{2L} \|\nabla g_{x_1}^{(x_2)}\|_2^2$$

thus from our definition of $g_{x_1}^{(x_2)}$ it follows:

$$g_{x_1}^{(x_1)} \leq g_{x_1}^{(x_2)} - \frac{1}{2L} \|\nabla g_{x_1}^{(x_2)}\|_2^2$$

$$\circ \quad \leq f(x_2) - f(x_1) - \nabla f(x_1)^\top (x_2 - x_1) - \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2$$

$$\Rightarrow -f(x_2) \leq -f(x_1) - \nabla f(x_1)^\top (x_2 - x_1) - \frac{1}{2L} \|\nabla f(x_2) - \nabla f(x_1)\|_2^2$$

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

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

Let define two convex functions f_{x_1}, f_{x_2} with \mathbb{R}^n domain

$$\begin{cases} f_{x_1}(z) = f(z) - \nabla f(x_1)^T z \\ f_{x_2}(z) = f(z) - \nabla f(x_2)^T z \end{cases}$$

This two functions have L -Lipschitz continuous gradient
we know that if $f: \mathbb{R}^n \rightarrow \mathbb{R}$ and f has a minimizer
 x^* then from the inequality in problem 1.2 b we have:

$$\begin{aligned} f(z) &\geq f(x^*) + \nabla f(x^*)^T (z - x^*) + \frac{1}{2L} \|\nabla f(z) - \nabla f(x^*)\|_2^2 \\ \Rightarrow f(z) - f(x^*) &\geq \frac{1}{2L} \|\nabla f(z)\|_2^2 \end{aligned}$$

$$\Rightarrow z = x_1 \text{ minimize } f_{x_1}(z)$$

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

$$\text{similarly } z = x_2 \text{ minimize } f_{x_2}(z)$$

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

now we can combine 1, 2 inequality:

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

Problem 1.3

A key performance measure of any iterative algorithm is its rate of convergence, i.e., how many iterations is needed to obtain a certain level of accuracy.

We now define the rate of convergence, or convergence rate. Let the set $\{\mathbf{x}_k\}$ be a sequence of updates produced by an iterative algorithm. If the algorithm converges, it will reach \mathbf{x}^* as $k \rightarrow \infty$. Let the error $e_k = \|\mathbf{x}_k - \mathbf{x}^*\|$ denote how far the current value \mathbf{x}_k is from the optimal value \mathbf{x}^* .

Michelle Schatzman (2002), Numerical analysis: a mathematical introduction, Clarendon Press, Oxford. ISBN 0-19-850279-6. check KTH library!

To find the rate of convergence for an algorithm we investigate the following limit:

$$\lim_{k \rightarrow \infty} \frac{e_{k+1}}{e_k^p} = \lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^p}. \quad (27)$$

Different values of this limit defines different convergence rates:

- **Sublinear:** If for $p = 1$

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 1, \quad (28)$$

then the sequence $\{\mathbf{x}_k\}$ has sublinear convergence to \mathbf{x}^* .

- **Linear:** If for $p = 1$ and there exists some $\mu \in (0, 1)$ that fulfills,

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = \mu, \quad (29)$$

then the sequence $\{\mathbf{x}_k\}$ has linear convergence to \mathbf{x}^* .

- **Superlinear:** If for $p = 1$

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|} = 0, \quad (30)$$

then the sequence $\{\mathbf{x}_k\}$ has superlinear convergence to \mathbf{x}^* .

- **Quadratic:** If for $p = 2$

$$\lim_{k \rightarrow \infty} \frac{\|\mathbf{x}_{k+1} - \mathbf{x}^*\|}{\|\mathbf{x}_k - \mathbf{x}^*\|^2} < M > 0, \quad (31)$$

then the sequence $\{\mathbf{x}_k\}$ has quadratic convergence to \mathbf{x}^* .

Problem 1.4

Consider

$$\min_x f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x_i), \quad (32)$$

$$\text{s.t. } \mathbf{A}\mathbf{x} = \mathbf{b}, \quad (33)$$

for $\mathbf{A} \in \mathbb{R}^{p \times N}$, $\mathbf{x} \in \mathbb{R}^{N \times 1}$, $\mathbf{b} \in \mathbb{R}^{p \times 1}$.

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

Solution: We can solve this using constraint free gradient descent (GD) by eliminating the equality constraint. Assuming $p < N$ (otherwise, a solution to $\mathbf{A}\mathbf{x} = \mathbf{b}$ is unique), find \mathbf{F} which spans the nullspace of \mathbf{A} , i.e., $\mathbf{AF} = 0 \in \mathbb{R}^{p \times (n-1)}$. The matrix \mathbf{F} parametrizes the feasible set as

$$\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\} = \{\mathbf{Fz} + \hat{\mathbf{x}} : \mathbf{z} \in \mathbb{R}^{(n-p) \times 1}\}.$$

Choosing a particular solution to $\mathbf{Ax} = \mathbf{b}$, e.g., $\hat{\mathbf{x}} = \mathbf{A}^+ \mathbf{b}$, where \mathbf{A}^+ denotes the Moore-Penrose pseudoinverse, leads to the following minimization problem:

$$\min_z f(\mathbf{Fz} + \hat{\mathbf{x}}). \quad (34)$$

An optimal solution \mathbf{z}^* fulfills $\mathbf{x}^* = \mathbf{Fz}^* + \hat{\mathbf{x}}$, where \mathbf{x}^* clearly is a feasible solution:

$$\mathbf{Ax}^* = \mathbf{AFz}^* + \mathbf{A}\hat{\mathbf{x}} = 0 + \mathbf{A}\hat{\mathbf{x}} = \mathbf{b}.$$

The problem in (34) can be solved using GD, by computing the gradients with respect to \mathbf{z} which are $\nabla_{\mathbf{z}} f(\mathbf{Fz} + \hat{\mathbf{x}}) = \mathbf{F}^T \nabla f(\mathbf{Fz} + \hat{\mathbf{x}})$, where the gradient in the righthand side is with respect to the argument of f .

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

Solution: Assuming $p \ll N$, it can be more efficient to solve the Lagrange dual problem instead of (32):

$$\max_{\nu} -\mathbf{b}^T \nu - f^*(-\mathbf{A}^T \nu), \quad (35)$$

where $f^*(x)$ is the conjugate function to f . This is an optimization problem in p variables rather than N as in (32).

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

Solution: If we approximate f by its second order Taylor polynomial, the newton step to minimize this approximation is characterized by $\Delta \mathbf{x}$ [?, p. 526]:^w

$$\begin{bmatrix} \nabla^2 f(\mathbf{x})^2 & \mathbf{A}^T \\ \mathbf{A} & 0 \end{bmatrix} \begin{bmatrix} \Delta \mathbf{x} \\ w \end{bmatrix} = \begin{bmatrix} -\nabla f(\mathbf{x}) \\ 0 \end{bmatrix} \quad (36)$$

It will be computationally infeasible to compute the inverse of the matrix on the lefthand side. We could do stochastic GD on reduced problem in (34).

- (d) Now, add twice differentiable $r(\mathbf{x})$ to f and solve (a)–(c).

Solution: Assuming $r(\mathbf{x})$ is convex, our solutions in (a), (b) and (c) still apply.

problem 1.5

In the convergence proof of GD with constant step size and strongly convex objective function proof the coercivity of the gradient:

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

Solution:

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

from strong convexity of $f(x)$, we get $g(x)$ is convex based on convergence of GD with constant step size and μ -strongly convex and L -smooth, we can say that the function $f(x)$, is differentiable and its gradient is L -Lipschitz continuous,

so we get $g(x)$ is also L -Lipschitz continuous and smooth with parameter $(L-\mu)$

now we can apply inequality in problem 1.2C to $g(x)$: co-coercivity

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

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

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

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

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

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

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

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