



LUND
UNIVERSITY

Department of
AUTOMATIC CONTROL

Exam in Optimization for Learning

2020-10-26

Points and grading

All answers must include a clear motivation. Answers should be given in English. The total number of points is 25. The maximum number of points is specified for each subproblem. Preliminary grading scales:

Grade 3: 12 points
4: 17 points
5: 22 points

Accepted aid

All material from the course.

Results

Solutions will be posted on the course webpage, and results will be registered in LADOK. Date and location for display of corrected exams will be posted on the course webpage.

1. Determine whether or not the functions below are convex.

- a. $f_1 : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f_1(x) = \frac{1}{g(x)}$$

for each $x \in \mathbb{R}^n$, where $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is concave and satisfies $g(x) > 0$ for each $x \in \mathbb{R}^n$. (1 p)

- b. $f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f(x) = \sqrt{x^T L^T L x}$$

for each $x \in \mathbb{R}^n$, where $L \in \mathbb{R}^{m \times n}$. (1 p)

- c. $f_3 : \mathbb{R}^2 \rightarrow \mathbb{R}$ such that

$$f_3(x) = \begin{cases} \sqrt{x_1 x_2} & \text{if } x = (x_1, x_2) \in \mathbb{R}_{++}^2, \\ \infty & \text{if } x \in \mathbb{R}^2 \setminus \mathbb{R}_{++}^2. \end{cases}$$

(1 p)

- d. $f_4 : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f_4(x) = \sum_{i=1}^r x_{(i)} = x_{(1)} + \dots + x_{(r)}$$

for each $x \in \mathbb{R}^n$, where r is an integer such that $1 \leq r \leq n$ and $x_{(i)}$ denote the i th largest component of x , i.e.,

$$x_{(1)} \geq \dots \geq x_{(n)}.$$

(1 p)

Remark: In each subproblem, you are allowed to assume that norms are convex. This remark also holds for all other problems in this exam.

Solution

- a. Convex. Define the convex function $h : \mathbb{R} \rightarrow \mathbb{R} \cup \{\infty\}$ such that

$$h(u) = \frac{1}{u} + \iota_{\mathbb{R}_{++}}(u)$$

for each $u \in \mathbb{R}$, where $\mathbb{R}_{++} = \{u \in \mathbb{R} : u > 0\}$ (e.g., show that h satisfies the second-order condition for convexity on $\text{dom } h$). The function f_1 can be written as

$$f_1 = h \circ g.$$

Since h is convex and nonincreasing, and g is concave, the composition rule gives that f_1 is convex.

b. Convex. Note that

$$\begin{aligned} f_2(x) &= \sqrt{(Lx)^T Lx} \\ &= \|Lx\|_2 \\ &= (\|\cdot\|_2 \circ L)(x) \end{aligned}$$

for each $x \in \mathbb{R}^n$. Thus, f_2 can be written as a composition between the convex function $\|\cdot\|_2$ and the affine mapping given by L . The composition rule gives that f_2 is convex.

c. Not convex. Let $x = (1, 1)$ and $y = (4, 1)$. Then $f_3(x) = 1$ and $f_3(y) = 2$. However, considering the convex combination

$$\frac{1}{2}x + \frac{1}{2}y = (2.5, 1)$$

gives

$$f_3\left(\frac{1}{2}x + \frac{1}{2}y\right) = \sqrt{2.5} > 1.5 = \frac{1}{2}f_3(x) + \frac{1}{2}f_3(y),$$

which violates the definition of convexity.

d. Convex. Note that the function f_4 can be written as

$$f_4(x) = \max \{x_{i_1} + \dots + x_{i_r} : \forall i_1, \dots, i_r \in \mathbb{N}, 1 \leq i_1 < \dots < i_r \leq n\}$$

for each $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. Note that each function in the maximum can be written as $a_{i_1, \dots, i_r} : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$a_{i_1, \dots, i_r}(x) = x_{i_1} + \dots + x_{i_r}$$

for each $x = (x_1, \dots, x_n) \in \mathbb{R}^n$. Note that each of the a_{i_1, \dots, i_r} are affine, and therefore also convex. We conclude that f_4 is given by a point-wise supremum of convex functions, and therefore itself convex.

2. Determine if the following sets are convex or not:

a. $S_1 = \{x \in \mathbb{R}^n : x_1 + \dots + x_n = 1\}$. (1 p)

b. $S_2 = \{x \in \mathbb{R}^n : \|x - a\|_2 \leq \|x - b\|_2\}$, where $a, b \in \mathbb{R}^n$ and $a \neq b$. (1 p)

c. $S_3 = \left\{x \in \mathbb{R}^3 : 2x_1 \geq \sqrt{x_2^2 + x_3^2}\right\}$. (1 p)

d. $S_4 = \left\{x \in \mathbb{R}^2 : 2 \leq e^{x_1^2 + x_2^2} \leq 4\right\}$. (1 p)

e. $S_5 = \{x \in \mathbb{R}^n : x^T y \leq 1, \forall y \in C\}$, where $C \subseteq \mathbb{R}^n$. (1 p)

Solution

- a.** Convex. The set

$$S_1 = \{x \in \mathbb{R}^n : \mathbf{1}^T x = 1\}$$

defines a hyperplane in \mathbb{R}^n , which we know is convex.

- b.** Convex. Since norms are nonnegative, we have that

$$\begin{aligned} \|x - a\|_2 &\leq \|x - b\|_2 \\ \Leftrightarrow \\ \|x - a\|_2^2 &\leq \|x - b\|_2^2 \\ \Leftrightarrow \\ (x - a)^T(x - a) &\leq (x - b)^T(x - b) \\ \Leftrightarrow \\ 2(b - a)^T x &\leq \|b\|_2^2 - \|a\|_2^2, \end{aligned}$$

which defines a halfspace in $x \in \mathbb{R}^n$. Thus, the set

$$S_2 = \{x \in \mathbb{R}^n : 2(b - a)^T x \leq \|b\|_2^2 - \|a\|_2^2\}$$

is convex.

- c.** Convex. The set S_3 can be written as the zero:th sublevel set of the function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$ such that

$$f(x) = \|Lx\|_2 - 2x_1$$

for each $x = (x_1, x_2, x_3) \in \mathbb{R}^3$, where

$$L = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

I.e.,

$$S_3 = \{x \in \mathbb{R}^3 : f(x) \leq 0\}.$$

The first term $\|Lx\|_2$ of f is convex by **1.b.**. The second term $-2x_1$ is affine and therefore convex. There, f is convex since it is a sum of convex functions. However, the set S_3 is then a sublevel set of a convex function, and therefore a convex set.

- d.** Not convex. The set S_4 can be written as

$$S_4 = \{x \in \mathbb{R}^2 : \log 2 \leq x_1^2 + x_2^2 \leq \log 4\}.$$

The set S_4 is nonempty since, e.g., $(0, \sqrt{\log 4}) \in S_4$. Consider any $x \in S_4$. Then $-x \in S_4$. However, the convex combination

$$\frac{1}{2}x + \frac{1}{2}(-x) = 0,$$

is not in S_4 . Thus, the set S_4 is not convex.

- e. Convex. Note that the set S_5 can be written as

$$S_5 = \bigcap_{y \in C} \left\{ x \in \mathbb{R}^n : x^T y \leq 1 \right\},$$

i.e., an intersection of halfspaces, which we know is convex, since halfspaces are convex and arbitrary intersections of convex sets are convex.

3. Consider the convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f(x) = \left(a^T x - b \right)^2$$

for each $x \in \mathbb{R}^n$, where $a \in \mathbb{R}^n \setminus \{0\}$, $b \in \mathbb{R}$ and $n \geq 2$.

- a. Prove or disprove that f is strongly convex. (1 p)
- b. Find the conjugate function f^* . (2 p)
- c. Let $\gamma > 0$. Find the proximal operator $\text{prox}_{\gamma f}$. (1 p)

Solution

- a. Note that

$$\nabla f(x) = 2a a^T x - 2ab \quad \text{and} \quad \nabla^2 f(x) = 2a a^T$$

for each $x \in \mathbb{R}^n$. Since $a \in \mathbb{R}^n$ and $n \geq 2$, we know that there exists a vector $z \in \mathbb{R}^n \setminus \{0\}$ such that $a^T z = 0$. Let $x \in \mathbb{R}^n$ and note that

$$z^T \nabla^2 f(x) z = 2 \left(a^T z \right)^2 = 0.$$

This shows that the second-order condition for strong convexity fails for f . Hence, f is not strong convexity.

- b. Let $s \in \mathbb{R}^n$. Note that we can write s as

$$s = \frac{a^T s}{\|a\|_2^2} a + \left(s - \frac{a^T s}{\|a\|_2^2} a \right)$$

where the first term is parallel to a and the second term is orthogonal to a , i.e., $a^T \left(s - \frac{a^T s}{\|a\|_2^2} a \right) = 0$. Then

$$\begin{aligned} f^*(s) &= \sup_{x \in \mathbb{R}^n} \left(s^T x - f(x) \right) \\ &= \sup_{x \in \mathbb{R}^n} \left(\frac{a^T s}{\|a\|_2^2} a^T x + \left(s - \frac{a^T s}{\|a\|_2^2} a \right)^T x - \left(a^T x - b \right)^2 \right). \end{aligned} \quad (1)$$

First, suppose that

$$s - \frac{a^T s}{\|a\|_2^2} a \neq 0.$$

Then, by picking $x = t\left(s - \frac{a^T s}{\|a\|_2^2} a\right)$ for $t \in \mathbb{R}$ in (1), we get that

$$\begin{aligned} f^*(s) &\geq t \frac{a^T s}{\|a\|_2^2} \underbrace{a^T \left(s - \frac{a^T s}{\|a\|_2^2} a\right)}_{=0} + t \left\| s - \frac{a^T s}{\|a\|_2^2} a \right\|_2^2 - \left(\underbrace{t a^T \left(s - \frac{a^T s}{\|a\|_2^2} a\right)}_{=0} - b \right)^2 \\ &= t \underbrace{\left\| s - \frac{a^T s}{\|a\|_2^2} a \right\|_2^2}_{>0} - b^2 \rightarrow \infty \quad \text{as } t \rightarrow \infty. \end{aligned}$$

Thus, $f^*(s) = \infty$ in this case.

Second, suppose that

$$s - \frac{a^T s}{\|a\|_2^2} a = 0.$$

Then (1) becomes

$$f^*(s) = \sup_{x \in \mathbb{R}^n} \left(\frac{a^T s}{\|a\|_2^2} a^T x - (a^T x - b)^2 \right). \quad (2)$$

By the orthogonal decomposition theorem, any $x \in \mathbb{R}^n$ can uniquely be decomposed into

$$x = \alpha a + c,$$

for some $\alpha \in \mathbb{R}$ and $c \in \mathbb{R}^n$ such that $a^T c = 0$. Using this observation, (2) can be written as

$$\begin{aligned} f^*(s) &= \sup_{\alpha \in \mathbb{R}} \left(\alpha a^T s - (\alpha \|a\|_2^2 - b)^2 \right) \\ &= \sup_{\alpha \in \mathbb{R}} \left(-\|a\|_2^4 \alpha^2 + (2b \|a\|_2^2 + a^T s) \alpha - b^2 \right) \\ &= -\inf_{\alpha \in \mathbb{R}} \left(\|a\|_2^4 \alpha^2 - (2b \|a\|_2^2 + a^T s) \alpha + b^2 \right). \end{aligned}$$

Let $g : \mathbb{R} \rightarrow \mathbb{R}$ denote the objective function of the minimization problem above, i.e.,

$$g(\alpha) = \|a\|_2^4 \alpha^2 - (2b \|a\|_2^2 + a^T s) \alpha + b^2$$

for each $\alpha \in \mathbb{R}$. Note that

$$\nabla^2 g(\alpha) = 2 \|a\|_2^4 > 0$$

for each $\alpha \in \mathbb{R}$. The second-order condition for (strong) convexity gives that g is (strongly) convex. Fermat's rule gives that $\alpha^* \in \mathbb{R}$ is a minimizer to the

minimization problem above if and only if

$$\begin{aligned}
0 &= \nabla g(\alpha^*) \\
&\Leftrightarrow \\
0 &= 2 \|a\|_2^4 \alpha^* - 2b \|a\|_2^2 - a^T s \\
&\Leftrightarrow \\
\alpha^* &= \frac{2b \|a\|_2^2 + a^T s}{2 \|a\|_2^4},
\end{aligned}$$

since g is differentiable and convex. Therefore,

$$\begin{aligned}
f^*(s) &= -g(\alpha^*) \\
&= -\|a\|_2^4 \left(\frac{2b \|a\|_2^2 + a^T s}{2 \|a\|_2^4} \right)^2 + (2b \|a\|_2^2 + a^T s) \left(\frac{2b \|a\|_2^2 + a^T s}{2 \|a\|_2^4} \right) - b^2 \\
&= -\frac{(2b \|a\|_2^2 + a^T s)^2}{4 \|a\|_2^4} + \frac{(2b \|a\|_2^2 + a^T s)^2}{2 \|a\|_2^4} - b^2 \\
&= \frac{(2b \|a\|_2^2 + a^T s)^2}{4 \|a\|_2^4} - b^2 \\
&= \frac{4b^2 \|a\|_2^4 + 4b \|a\|_2^2 a^T s + (a^T s)^2}{4 \|a\|_2^4} - b^2 \\
&= b \frac{a^T s}{\|a\|_2^2} + \frac{1}{4} \left(\frac{a^T s}{\|a\|_2^2} \right)^2.
\end{aligned}$$

To summarize, we have that

$$f^*(s) = \begin{cases} b \frac{a^T s}{\|a\|_2^2} + \frac{1}{4} \left(\frac{a^T s}{\|a\|_2^2} \right)^2 & \text{if } s = \frac{a^T s}{\|a\|_2^2} a, \\ \infty & \text{if } s \neq \frac{a^T s}{\|a\|_2^2} a. \end{cases}$$

c. Let $z \in \mathbb{R}^n$. Then

$$\text{prox}_{\gamma f}(z) = \underset{x \in \mathbb{R}^n}{\text{argmin}} \left((a^T x - b)^2 + \frac{1}{2\gamma} \|x - z\|^2 \right).$$

Note that the objective in the minimization problem above is differentiable and convex. Fermat's rule gives that $x^* \in \mathbb{R}^n$ is a minimizer of the optimization problem above if and only if

$$\begin{aligned}
0 &= 2a(a^T x^* - b) + \frac{1}{\gamma}(x^* - z) \\
&\Leftrightarrow \\
0 &= 2\gamma a a^T x^* - 2\gamma a b + x^* - z \\
&\Leftrightarrow \\
(I + 2\gamma a a^T) x^* &= 2\gamma a b + z.
\end{aligned}$$

Note that $(I + 2\gamma aa^T)$ is invertible—its smallest eigenvalue is greater or equal to 1. Therefore,

$$x^* = (I + 2\gamma aa^T)^{-1} (2\gamma ab + z)$$

and

$$\text{prox}_{\gamma f}(z) = (I + 2\gamma aa^T)^{-1} (2\gamma ab + z).$$

4. Let $C \subseteq \mathbb{R}^n$ be a nonempty, closed and convex set. Its support function $\sigma_C : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is defined as

$$\sigma_C(y) = \sup_{x \in C} y^T x$$

for each $y \in \mathbb{R}^n$.

- a. Show that the support function σ_C is convex, independent of the convexity of the set C . (1 p)
- b. Show that $\sigma_C^* = \iota_C$. (1 p)
- c. Find an expression for $\text{prox}_{\gamma \sigma_C}$, where $\gamma > 0$, that involves Π_C , i.e., the (Euclidean) projection onto the set C . (1 p)

Solution

- a. *Alternative 1:* Note that

$$\sigma_C(y) = \sup_{x \in \mathbb{R}^n} (y^T x - \iota_C(x)) = \iota_C^*(y)$$

for each $y \in \mathbb{R}^n$. I.e., σ_C is equal to the conjugate function to the indicator function of C . Thus, σ_C is a convex function since conjugate functions are always convex.

Alternative 2: Note that σ_C is a points-wise supremum of convex functions (linear to be precise) and therefore itself a convex function.

- b. Note that $\sigma_C^* = \iota_C^{**} = \iota_C$ since ι_C is a proper, closed and convex function.
- c. Let $z \in \mathbb{R}^n$. Moreau decomposition gives that

$$\begin{aligned} \text{prox}_{\gamma \sigma_C}(z) &= z - \gamma \text{prox}_{\gamma^{-1} \sigma_C^*}(\gamma^{-1} z) \\ &= z - \gamma \text{prox}_{\gamma^{-1} \iota_C}(\gamma^{-1} z) \\ &= z - \gamma \text{prox}_{\iota_C}(\gamma^{-1} z) \\ &= z - \gamma \Pi_C(\gamma^{-1} z). \end{aligned}$$

5. Consider the problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \frac{1}{2} \|Ax - b\|_2^2$$

where $A \in \mathbb{R}^{n \times n}$ satisfies

$$A = \mathbf{diag}(a_1, \dots, a_n), \quad a_i \neq 0, \quad \forall i \in \{1, \dots, n\},$$

and $b = (b_1, \dots, b_n) \in \mathbb{R}^n$. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

$$f(x) = \frac{1}{2} \|Ax - b\|_2^2$$

for each $x \in \mathbb{R}^n$.

- a. Give a closed-form expression of the solution. (0.5 p)
- b. Show that $\beta = \max_{i \in \{1, \dots, n\}} a_i^2$ is a smoothness constant for f . (1 p)
- c. Show that $\beta_i = a_i^2$ is a coordinate-wise smoothness constants for f , for each coordinate $i = 1, \dots, n$. (1 p)
- d. Consider the gradient method with step-size $1/\beta$, where β is the smoothness constant in **b.** Suppose you are given the iterate $x^k \in \mathbb{R}^n$, where $k \in \mathbb{N}_0$ is the iteration number. For each coordinate $i = 1, \dots, n$, provide the update formula for x_i^k . Utilize that $A = \mathbf{diag}(a_1, \dots, a_n)$. (1 p)
- e. Let $b_i = 0$ for each $i = 1, \dots, n$ and provide an exact linear convergence rate for each of the coordinates for the gradient method in **d.** This means, find the $\rho_i \in [0, 1)$ such that

$$\|x_i^{k+1}\|_2 = \rho_i \|x_i^k\|_2,$$

for each coordinate $i = 1, \dots, n$. (Each coordinate will converge linearly to $x_i^* = 0$ in this case.) (1 p)

- f. Now drop the assumption that $b_i = 0$ for each $i = 1, \dots, n$. Consider the coordinate gradient method (i.e., no proximal operator) with step-sizes $1/\beta_i$, where β_i are the coordinate smoothness constants in **c.** Provide an update formula for each coordinate $i = 1, \dots, n$. Utilize that $A = \mathbf{diag}(a_1, \dots, a_n)$. Show that $x_i^{k+1} = x_i^*$ with x_i^* from **a.**, independent on $x^k \in \mathbb{R}^n$. (1 p)

Solution

- a. The objective function f in the problem is convex and differentiable. Fermat's rule gives that $x^* \in \mathbb{R}^n$ is a minimizer of the problem if and only if

$$\begin{aligned} 0 &= \nabla f(x^*) \\ &\Leftrightarrow \\ 0 &= A^T (Ax^* - b) \\ &\Leftrightarrow \\ x^* &= (A^T A)^{-1} (A^T b) \\ &\Leftrightarrow \end{aligned}$$

$$x_i^* = \frac{b_i}{a_i}, \quad \text{for each } i = 1, \dots, n,$$

since $A^T A = \mathbf{diag}(a_1^2, \dots, a_n^2)$ is invertible as $a_i \neq 0$ for each $i = 1, \dots, n$.

b. Alternative 1: Note that

$$\begin{aligned}
\|\nabla f(x) - \nabla f(y)\|_2^2 &= \|A^T(Ax - b) - A^T(Ay - b)\|_2^2 \\
&= \|A^T A(x - y)\|_2^2 \\
&= \|\mathbf{diag}(a_1^2, \dots, a_n^2)(x - y)\|_2^2 \\
&= \sum_{i=1}^n (a_i^2(x_i - y_i))^2 \\
&= \sum_{i=1}^n a_i^4(x_i - y_i)^2 \\
&\leq \left(\max_{i \in \{1, \dots, n\}} a_i^4 \right) \sum_{j=1}^n (x_j - y_j)^2 \\
&= \left(\max_{i \in \{1, \dots, n\}} a_i^4 \right) \|x - y\|_2^2.
\end{aligned}$$

for each $x, y \in \mathbb{R}^n$. Taking square root gives the result, i.e.,

$$\begin{aligned}
\|\nabla f(x) - \nabla f(y)\|_2 &\leq \sqrt{\max_{i \in \{1, \dots, n\}} a_i^4} \|x - y\|_2 \\
&= \left(\max_{i \in \{1, \dots, n\}} a_i^2 \right) \|x - y\|_2
\end{aligned}$$

for each $x, y \in \mathbb{R}^n$, which is the definition of $\max_{i \in \{1, \dots, n\}} a_i^2$ -smoothness.

Alternative 2: Since f is convex and twice differentiable, $\beta \geq 0$ is a smoothness constant if and only if

$$\nabla^2 f(x) \preceq \beta I$$

for each $x \in \mathbb{R}^n$. This is equivalent to that

$$0 \leq y^T(\beta I - \nabla^2 f(x))y$$

for each $x, y \in \mathbb{R}^n$. Note that

$$\begin{aligned}
\nabla^2 f(x) &= A^T A \\
&= \mathbf{diag}(a_1^2, \dots, a_n^2)
\end{aligned}$$

for each $x \in \mathbb{R}^n$. The condition on β reduces to

$$\begin{aligned}
0 &\leq y^T \mathbf{diag}(\beta - a_1^2, \dots, \beta - a_n^2) y \\
&= \sum_{i=1}^n y_i^2 (\beta - a_i^2)
\end{aligned}$$

for each $y \in \mathbb{R}^n$. This holds if and only if

$$\beta \geq \max_{i \in \{1, \dots, n\}} a_i^2$$

and the smallest such β is

$$\beta = \max_{i \in \{1, \dots, n\}} a_i^2.$$

- c. Fix a coordinate $i = 1, \dots, n$. Coordinate-wise smoothness with parameter $\beta_i \geq 0$ can be written as

$$|(\nabla f(x))_i - (\nabla f(y))_i| \leq \beta_i |x_i - y_i|$$

for each $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ and $y = (y_1, \dots, y_n) \in \mathbb{R}^n$. In our setting, we have

$$\begin{aligned} |(\nabla f(x))_i - (\nabla f(y))_i| &= |a_i(a_i x_i - b_i) - a_i(a_i y_i - b_i)| \\ &= |a_i a_i (x_i - y_i)| \\ &\leq a_i^2 |x_i - y_i|. \end{aligned}$$

for each $x = (x_1, \dots, x_n) \in \mathbb{R}^n$ and $y = (y_1, \dots, y_n) \in \mathbb{R}^n$. So we identify a_i^2 as a coordinate smoothness constant for coordinate i .

- d. The gradient method with step-size $\gamma = 1/\beta$ reads

$$\begin{aligned} x^{k+1} &= x^k - \frac{1}{\beta} \nabla f(x^k) \\ &= x^k - \frac{1}{\max_{j \in \{1, \dots, n\}} a_j^2} A^T (A x^k - b). \end{aligned}$$

Since $A = \mathbf{diag}(a_1, \dots, a_n)$, this reads as

$$x_i^{k+1} = x_i^k - \frac{1}{\max_{j \in \{1, \dots, n\}} a_j^2} a_i (a_i x_i^k - b_i)$$

for each coordinate $i = 1, \dots, n$.

- e. Fix a coordinate $i = 1, \dots, n$. Note that

$$\begin{aligned} \|x_i^{k+1}\|_2 &= \left\| x_i^k - \frac{1}{\max_{j \in \{1, \dots, n\}} a_j^2} a_i^2 x_i^k \right\|_2 \\ &= \underbrace{\left(1 - \frac{a_i^2}{\max_{j \in \{1, \dots, n\}} a_j^2} \right)}_{=\rho_i} \|x_i^k\|_2 \end{aligned}$$

where $\rho_i \in [0, 1)$.

- f. Fix a coordinate $i = 1, \dots, n$. The coordinate gradient method when updating coordinate i is

$$\begin{aligned} x_i^{k+1} &= x_i^k - \frac{1}{\beta_i} (\nabla f(x^k))_i \\ &= x_i^k - \frac{1}{a_i^2} a_i (a_i x_i^k - b_i) \\ &= \frac{b_i}{a_i} \\ &= x_i^*. \end{aligned}$$

6. Consider minimizing a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, with minimizer $x^* \in \mathbb{R}^n$, using a stochastic optimization algorithm, starting at some predetermined (deterministic) point $x_0 \in \mathbb{R}^n$. Analysis of the algorithm resulted in the following inequality

$$\mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \mid x_k \right] \leq \|x_k - x^*\|_2^2 - 2\gamma(f(x_k) - f(x^*)) + \gamma^2 G, \quad \forall k \in \mathbb{N}_0,$$

where G is a deterministic positive constant and γ is a deterministic fixed positive step-size of the algorithm. In particular, $(x_k)_{k \in \mathbb{N}_0}$ is a stochastic process.

- a. Apply an expectation to the above inequality to derive a Lyapunov inequality for the algorithm. (1 p)
- b. Use the obtained Lyapunov inequality to show that

$$\sum_{i=0}^k \mathbb{E}[f(x_i) - f(x^*)] \leq \frac{\|x_0 - x^*\|_2^2 + G(k+1)\gamma^2}{2\gamma}, \quad \forall k \in \mathbb{N}_0. \quad (3)$$

(1.5 p)

- c. The upper bound (3) goes to infity as $k \rightarrow \infty$ unless $G = 0$. Consider the step-size $\gamma = \theta/\sqrt{K+1}$, where $K \in \mathbb{N}_0$ is the total number of iterations we wish to run the algorithm and $\theta > 0$. Show that we get a $\mathcal{O}(1/\sqrt{K+1})$ convergence bound. In particular, show that

$$\min_{i \in \{0, \dots, K\}} \mathbb{E}[f(x_i) - f(x^*)] \leq \frac{\|x_0 - x^*\|_2^2 + G\theta^2}{2\theta\sqrt{K+1}}.$$

(1 p)

Solution

- a. We start from the inequality

$$\mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \mid x_k \right] \leq \|x_k - x^*\|_2^2 - 2\gamma(f(x_k) - f(x^*)) + \gamma^2 G, \quad \forall k \in \mathbb{N}_0.$$

By monotonicity and linearity of expectation, we get that

$$\begin{aligned} \mathbb{E} \left[\mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \mid x_k \right] \right] &\leq \mathbb{E} \left[\|x_k - x^*\|_2^2 - 2\gamma(f(x_k) - f(x^*)) + \gamma^2 G \right] \\ &= \mathbb{E} \left[\|x_k - x^*\|_2^2 \right] - 2\gamma \mathbb{E} [f(x_k) - f(x^*)] + \gamma^2 G, \end{aligned}$$

holds for each $k \in \mathbb{N}_0$. The law of total expectation yields

$$\mathbb{E} \left[\|x_{k+1} - x^*\|_2^2 \right] \leq \mathbb{E} \left[\|x_k - x^*\|_2^2 \right] - 2\gamma \mathbb{E} [f(x_k) - f(x^*)] + \gamma^2 G, \quad \forall k \in \mathbb{N}_0.$$

This is the Lyapunov inequality we pick.

b. Recursively applying the Lyapunov inequality above gives

$$\begin{aligned}\mathbb{E} \left[\|x_{k+1} - x^\star\|_2^2 \right] &\leq \mathbb{E} \left[\|x_0 - x^\star\|_2^2 \right] - 2\gamma \sum_{i=0}^k \mathbb{E} [f(x_i) - f(x^\star)] + G\gamma^2(k+1) \\ &= \|x_0 - x^\star\|_2^2 - 2\gamma \sum_{i=0}^k \mathbb{E} [f(x_i) - f(x^\star)] + G\gamma^2(k+1), \quad \forall k \in \mathbb{N}_0,\end{aligned}$$

since $\|x_0 - x^\star\|_2^2$ is deterministic. Again, by monotonicity of expectation, we know that

$$0 \leq \mathbb{E} \left[\|x_{k+1} - x^\star\|_2^2 \right], \quad \forall k \in \mathbb{N}_0,$$

since

$$0 \leq \|x_{k+1} - x^\star\|_2^2, \quad \forall k \in \mathbb{N}_0.$$

We conclude that

$$0 \leq \|x_0 - x^\star\|_2^2 - 2\gamma \sum_{i=0}^k \mathbb{E} [f(x_i) - f(x^\star)] + G\gamma^2(k+1), \quad \forall k \in \mathbb{N}_0,$$

or by rearranging

$$\sum_{i=0}^k \mathbb{E} [f(x_i) - f(x^\star)] \leq \frac{\|x_0 - x^\star\|_2^2 + G\gamma^2(k+1)}{2\gamma}, \quad \forall k \in \mathbb{N}_0, \quad (4)$$

as desired.

c. We first note that

$$(K+1) \min_{i \in \{0, \dots, K\}} \mathbb{E} [f(x_i) - f(x^\star)] \leq \sum_{i=0}^K \mathbb{E} [f(x_i) - f(x^\star)].$$

Using this in the bound in **b.** with $\gamma = \theta/\sqrt{K+1}$ and $k = K$ give

$$\begin{aligned}\min_{i \in \{0, \dots, K\}} \mathbb{E} [f(x_i) - f(x^\star)] &\leq \frac{\|x_0 - x^\star\|_2^2 + G(K+1)\gamma^2}{2\gamma(K+1)} \\ &= \frac{\|x_0 - x^\star\|_2^2 + G\theta^2}{2\theta\sqrt{K+1}},\end{aligned}$$

as desired.