

Introduction to Continuous optimization

Assessment
(6th January 2021)

Exercise I

We denote $\mathbb{R}_{\text{sym}}^{n \times n}$ the space of dimension $n(n+1)/2$ of symmetric $n \times n$ matrices. We consider the scalar product $X : Y = \sum_{i,j} X_{i,j} Y_{i,j} = \text{Tr}(XY)$ (or $\text{Tr}(X^T Y)$ but it is the same here since X, Y are symmetric).

Let $\mathcal{S}_+ \subset \mathbb{R}_{\text{sym}}^{n \times n}$ be the set of $n \times n$ symmetric, positive semidefinite matrices: $X = X^T$, $(X\xi) \cdot \xi \geq 0$ for any $\xi \in \mathbb{R}^n$. Let \mathcal{S}_{++} be the interior of \mathcal{S}_+ , that is, the set of positive definite matrices: $(X\xi) \cdot \xi > 0$ for all $\xi \neq 0$.

We let, for $X \in \mathbb{R}_{\text{sym}}^{n \times n}$:

$$h(X) := \begin{cases} -\ln \det X & \text{if } X \in \mathcal{S}_{++}, \\ +\infty & \text{else.} \end{cases}$$

1. Let $X \in \mathcal{S}_{++}$, H a symmetric matrix. Show that for $t \in \mathbb{R}$ with $|t|$ small enough, $X + tH \in \mathcal{S}_{++}$.

For ξ a vector, one has $((X+tH)\xi) \cdot \xi = (X\xi) \cdot \xi + t(H\xi) \cdot \xi \geq (\lambda_1(X) - t\|H\|)|\xi|^2 > 0$ if $t < \lambda_1(X)/\|H\|$, where $\lambda_1(X) > 0$ is the smallest eigenvalue of X .

2. Using $X + tH = X(I + tX^{-1}H)$, show that

$$\nabla h(X) = -X^{-1}.$$

We recall that $\det(I + A) = 1 + \text{Tr} A + o(\|A\|)$.

We have:

$$\begin{aligned} h(X + tH) &= -\ln \det(X(I + tX^{-1}H)) \\ &= -\ln(\det(X) \det(I + tX^{-1}H)) = h(X) - \ln \det(I + tX^{-1}H) \\ &= h(X) - \ln(1 + t\text{Tr}(X^{-1}H) + o(t)) = h(X) - t\text{Tr}(X^{-1}H) + o(t) = h(X) - tX^{-1} : H + o(t) \end{aligned}$$

which shows the claim.

3. One now wants to compute the conjugate $h^*(Y) = \sup_X X : Y - h(X)$.

Let $Y \in \mathbb{R}_{\text{sym}}^{n \times n}$ and assume e is an eigenvector of Y with eigenvalue $\lambda \in \mathbb{R}$ (and $|e| = 1$).

Considering first X of the form $te \otimes e + \varepsilon I$ (where for $e \in \mathbb{R}^n \setminus \{0\}$ with $|e| = 1$, $e \otimes e$ is the matrix $e_i e_j$ which has eigenvector e with eigenvalue 1), $\varepsilon > 0$, $t \rightarrow +\infty$, show that $h^*(Y) = +\infty$ if $\lambda \geq 0$.

Deduce that $\text{dom } h^* \subset \{Y : -Y \in \mathcal{S}_{++}\}$.

One has, for $X = te \otimes e + \varepsilon I$, $\det X = \varepsilon^{n-1}(t + \varepsilon)$ and $X : Y - h(X) = \varepsilon \text{Tr} X + t\lambda + (n-1) \ln \varepsilon + \ln(\varepsilon + t)$ which goes to $+\infty$ as $t \rightarrow \infty$ if $\lambda \geq 0$. Hence, $h^*(Y) < +\infty$ only if all eigenvalues of Y are strictly negative, that is $-Y \in \mathcal{S}_{++}$.

4. Now, assuming $-Y > 0$ we admit (even if it is quite easy to show) that $\sup_X X : Y - h(X)$ is reached at some positive matrix X .

Show that $X = -Y^{-1}$. Deduce the expression of h^* . Deduce also that h is convex.

At the maximum X one has $\nabla_X(X : Y - h(X)) = Y - (-X^{-1}) = 0$ so that $Y = (-X^{-1}) \Leftrightarrow X = -Y^{-1}$. Then,

$$X : Y - h(X) = \text{Tr}(-Y^{-1}Y) + \ln \det(-Y^{-1}) = -n - \ln \det(-Y).$$

In particular, the function

$$Y \mapsto \begin{cases} -n - \ln \det(-Y) & \text{if } -Y \in \mathcal{S}_{++} \\ +\infty & \text{else.} \end{cases}$$

is convex, and so is $h(X) = n + h^*(-X)$.

5. We consider the problem $\min_{X \in \mathcal{S}_+} C : X$ and the Bregman distance

$$D_h(X, Y) = h(X) - h(Y) - \nabla h(Y) : (X - Y)$$

induced by h , defined for $X, Y \in \mathcal{S}_{++}$. Write the expression of an iteration of non-linear gradient descent for the problem, with step $\tau > 0$, relative to the Bregman distance D_h . Why can we always assume that C is symmetric? What assumption is needed on C in order for the problem to have a solution (and the algorithm to be well defined for all k)?

X^{k+1} is obtained as (if it exists)

$$X^{k+1} = \arg \min_X \frac{1}{\tau} D_h(X, X^k) + C : X$$

and satisfies: $-(X^{k+1})^{-1} = -(X^k)^{-1} - \tau C$, that is

$$X^{k+1} = ((X^k)^{-1} + \tau C)^{-1}$$

If C is not symmetric then $C : X = C^T : X^T = C^T : X = (C + C^T) : X/2$ so one can replace C with its symmetric part. If C has a negative eigenvalue, as in the analysis of the previous question, the minimum problem has no solution (the value is $-\infty$) as soon as $(X^0)^{-1} + \tau k C$ has a negative eigenvalue. If $C \geq 0$, the iterates are $X^k = ((X^0)^{-1} + \tau k C)^{-1}$

Exercice II - prox

Compute the proximity operator (for some parameter $\tau > 0$):

$$\text{prox}_{\tau g}(x) = \arg \min_z g(z) + \frac{1}{2\tau} |z - x|^2$$

for the convex functions:

1. $g_1(x) = -\ln x$ for $x > 0$, $+\infty$ else;

The equation is $z - x - \tau/z = 0$, that is $z^2 - zx - \tau = 0$, that is $z = (x + \sqrt{x^2 + 4\tau})/2$.

2. $g_2(x) = \sum_{i=1}^n \frac{1}{3}|x_i|^3$, ($x \in \mathbb{R}^n$);

The problem is:

$$\min_z \sum_{i=1}^n \frac{1}{3}|z_i|^3 + \frac{1}{2\tau}|z_i - x_i|^2$$

and can be minimized independently for each i : the minimizer satisfies

$$\tau \text{sign}(z_i)z_i^2 + z_i - x_i = 0, \quad i = 1, \dots, n.$$

Clearly the sign of z_i is the same as the sign of x_i (as $\tau \text{sign}(z_i)z_i^2 + z_i = z_i(\tau|z_i| + 1)$ has the same sign as z_i). Solving the equation one obtains:

$$z_i = \text{sign}(x_i) \frac{\sqrt{1 + 4\tau|x_i|} - 1}{2\tau}.$$

3. $g_3(x) = \sum_{i=1}^n \frac{2}{3}|x_i|^{3/2}$, ($x \in \mathbb{R}^n$);

$g_3 = g_2^*$ and one has the Moreau identity:

$$\text{prox}_{\tau g_3}(x) = x - \tau \text{prox}_{\frac{1}{\tau} g_3^*}\left(\frac{x}{\tau}\right).$$

Hence,

$$\begin{aligned} z_i &= x_i - \tau \text{sign}(x_i) \frac{\sqrt{1 + (4/\tau)|x_i/\tau|} - 1}{2/\tau} = \text{sign}(x_i) \left(|x_i| + \frac{\tau^2}{2} - \tau \frac{\sqrt{\tau^2 + 4|x_i|}}{2} \right) \\ &= \text{sign}(x_i) \frac{4|x_i| + \tau^2 + \tau^2 - 2\tau\sqrt{\tau^2 + 4|x_i|}}{4} = \text{sign}(x_i) \left(\frac{\sqrt{\tau^2 + 4|x_i|} - \tau}{2} \right)^2 \end{aligned}$$

The last expression is the one which is obtained directly if one solves the minimization problem (without using Moreau's identity).

4. $g_4(x) = \frac{1}{2} \sum_i x_i^2$ if $x_i \geq 0$, $i = 1, \dots, n$, and $+\infty$ else, defined for $x \in \mathbb{R}^n$ (and with domain $\text{dom } g_4 = [0, +\infty)^n$).

One solves:

$$\min_{z_i \geq 0} \sum_i \frac{1}{2}z_i^2 + \frac{1}{2\tau}|z_i - x_i|^2$$

which is solved independently for each i . The solution is $z_i = 0$ if $x_i < 0$, otherwise, $(1 + \tau)z_i - x_i = 0$, that is, $z_i = x_i/(1 + \tau)$. Hence, $z_i = x_i^+/(1 + \tau)$.

Exercise III - rate for the proximal point algorithm

We consider M a maximal-monotone operator, defined in a Hilbert space \mathcal{X} . Given $x^0 \in \mathcal{X}$, we let for $k \geq 0$:

$$x^{k+1} = (I + M)^{-1}x^k$$

that is, the iterations of the proximal-point algorithm.

1. Let x^* be a zero, that is, a point such that $Mx^* \ni 0$ (we assume the set $M^{-1}(0)$ is not empty). Show that $x^* = (I + M)^{-1}(x^*)$ and that

$$|x^{k+1} - x^*|^2 + |x^k - x^{k+1}|^2 \leq |x^k - x^*|^2.$$

One has $x^* + 0 = x^*$, hence $x^* + Mx^* \ni x^*$, that is $x^* = (I + M)^{-1}(x^*)$. Denoting $J_M = (I + M)^{-1}$ we know that J_M is “firmly non expansive”:

$$|J_M x - J_M x'|^2 + |(I - J_M)x - (I - J_M)x'|^2 \leq |x - x'|^2.$$

With $x = x^k$ and $x' = x^*$ this gives the desired inequality.

2. Show that $|x^{k+1} - x^k|$ is a decreasing function of $k \geq 0$.

This is even easier: if $k \geq 1$, $|x^{k+1} - x^k| = |J_M x^k - J_M x^{k-1}| \leq |x^k - x^{k-1}|$ since J_M is one-Lipschitz.

3. Deduce that

$$|x^{k+1} - x^k| \leq \frac{|x^0 - x^*|}{\sqrt{k+1}}.$$

We sum the inequality of the first question:

$$|x^{k+1} - x^*|^2 + \sum_{l=0}^k |x^l - x^{l+1}|^2 \leq |x^0 - x^*|^2$$

then we use the second question to observe that $\sum_{l=0}^k |x^l - x^{l+1}|^2 \geq (k+1)|x^k - x^{k+1}|^2$.

4. Let x^{k_i} be a (weakly) converging subsequence, to some point \bar{x} . Show that for any $x' \in \mathcal{X}$ and $y' \in Mx'$,

$$\langle x' - \bar{x}, y' \rangle \geq 0.$$

Deduce that $0 \in M\bar{x}$.

Since M is monotone and $x^k - x^{k+1} \in Mx^{k+1}$, $\langle x' - x^{k_i}, y' - (x^{k_i} - x^{k_i+1}) \rangle \geq 0$ and in the limit, thanks to the previous estimate, we obtain the inequality (we have a product (weak convergence) \times (strong convergence)).

Since M is maximal-monotone, it means that $0 \in M\bar{x}$ (otherwise one could extend the graph). (Of course, using Opial’s lemma, one can then show that $x^k \rightarrow \bar{x}$, weakly.)

5. Let $T : \mathcal{X} \rightarrow \mathcal{X}$ be a 1-Lipschitz operator and, for $\theta \in (0, 1)$, let $T_\theta = (1 - \theta)I + \theta T$. Let x^* be a fixed point of T (and therefore also of T_θ for any θ). We now consider the algorithm given by

$$x^{k+1} = T_\theta x^k.$$

Use the parallelogram identity to show that:

$$|x^{k+1} - x^*|^2 \leq |x^k - x^*|^2 - \theta(1 - \theta)|Tx^k - x^k|^2$$

One has

$$\begin{aligned} |x^{k+1} - x^*|^2 &= |(1 - \theta)(x^k - x^*) + \theta(Tx^k - x^*)|^2 \\ &= (1 - \theta)|x^k - x^*|^2 + \theta|Tx^k - x^*|^2 - \theta(1 - \theta)|Tx^k - x^k|^2 \end{aligned}$$

and one uses $|Tx^k - x^*| = |Tx^k - T^*| \leq |x^k - x^*|$ to conclude.

6. As before, deduce that:

$$|Tx^k - x^k| \leq \frac{|x^0 - x^*|}{\sqrt{\theta(1 - \theta)}\sqrt{k + 1}}.$$

(Remark: in this framework, one can show [Baillon-Bruck 1996] that a similar estimate holds in any metric space, but it is much harder).

As in question 2., one has $|x^{k+1} - x^k| \leq |x^k - x^{k-1}|$ for $k \geq 1$, using that T_θ is 1-Lipschitz. We deduce $|Tx^k - x^k| \leq |Tx^{k-1} - x^{k-1}|$. Thus,

$$\theta(1 - \theta)(k + 1)|Tx^k - x^k| \leq \theta(1 - \theta) \sum_{l=0}^k |Tx^l - x^l| + |x^{k+1} - x^*|^2 \leq |x^0 - x^*|^2$$

7. Application: show that the over-relaxed proximal point algorithm:

$$\begin{aligned} x^{k+\frac{1}{2}} &= (I + M)^{-1}x^k \\ x^{k+1} &= x^k + \lambda(x^{k+\frac{1}{2}} - x^k) \end{aligned}$$

for $1 < \lambda < 2$ is a converging method.

We know that $(I + M)^{-1} = I/2 + R/2$ for a 1-Lipschitz map R . Then,

$$x^{k+1} = x^k + \frac{\lambda}{2}(Rx^k - x^k) = (1 - \frac{\lambda}{2})x^k + \frac{\lambda}{2}Rx^k = R_{\frac{\lambda}{2}}x^k$$

is the iteration of an averaged operator, and we can use the previous results to show that $x^{k+1} - x^k \rightarrow 0$. Then, we can conclude as in question 4.

Exercise IV - Yosida approximation

Let A be a maximal monotone operator in a Hilbert space, and defined the Yosida approximation, for $\lambda > 0$, as

$$A_\lambda x = \frac{x - J_{\lambda A}x}{\lambda}$$

where $J_{\lambda A} = (I + \lambda A)^{-1}$ is the resolvent.

1. Show that A_λ is a monotone operator.

This is because $J_{\lambda A}$ is 1-Lipschitz. Then, for any x, y ,

$$\langle A_\lambda x - A_\lambda y, x - y \rangle = \frac{1}{\lambda}(|x - y|^2 - \langle J_{\lambda A}x - J_{\lambda A}y, x - y \rangle) \geq 0$$

using that $\langle J_{\lambda A}x - J_{\lambda A}y, x - y \rangle \leq |J_{\lambda A}x - J_{\lambda A}y||x - y| \leq |x - y|^2$.

2. Show that $A_\lambda x = J_{A^{-1}/\lambda}(x/\lambda)$. Deduce that it is $(1/\lambda)$ -Lipschitz. Bonus: show that it is λ -co-coercive.

We use Moreau's identity:

$$x = (I + \lambda A)^{-1}x + \lambda(I + \frac{1}{\lambda}A^{-1})^{-1}(\frac{x}{\lambda}) = J_{\lambda A}x + \lambda J_{A^{-1}/\lambda}(\frac{x}{\lambda}).$$

We conclude using that J_\bullet is 1-Lipschitz.

3. Let $x \in \text{dom } A$ (that is, $Ax \neq \emptyset$). Show that

$$\lim_{\lambda \rightarrow 0} A_\lambda x = A_0 x := \arg \min_{p \in Ax} |p|.$$

Hint: first, show that if $p_\lambda = A_\lambda x$ then $p_\lambda \in A(x - \lambda p_\lambda)$. Using the monotonicity of A , deduce that for any $p \in Ax$, $|p_\lambda|^2 \leq \langle p_\lambda, p \rangle$, hence that $|p_\lambda| \leq |p|$. Conclude by using that A is maximal.

One has $p_\lambda = (x - J_{\lambda A}x)/\lambda$, hence $(I + \lambda A)(x - \lambda p_\lambda) \ni x$, that is, $p_\lambda \in A(x - \lambda p_\lambda)$. Since A is monotone, for any y and $q \in Ay$,

$$\langle q - p_\lambda, y - x + \lambda p_\lambda \rangle \geq 0. \quad (*)$$

In particular for $y = x, q = p \in Ax$,

$$\langle p, p_\lambda \rangle \geq |p_\lambda|^2 \Rightarrow |p_\lambda| \leq |p|.$$

Hence in the limit, if $p_{\lambda_k} \rightarrow \bar{p}$, we find from $(*)$ that

$$\langle q - \bar{p}, y - x \rangle \geq 0$$

and $|\bar{p}| \leq |p|$ for any $p \in Ax$. Since A is maximal, we deduce that $\bar{p} \in Ax$, so that it is the (unique) element of minimal norm, and the whole sequence p_λ converges to \bar{p} .

4. Contraction semigroup: since A_λ is Lipschitz, by the Cauchy-Lipschitz theorem, one can solve for all $x \in \mathcal{X}$:

$$\begin{cases} \dot{X}^\lambda(t, x) = -A_\lambda X^\lambda(t, x) & t > 0, \\ X^\lambda(0, x) = x \end{cases}$$

and the solution, which is at least C^1 in time, satisfies:

$$X^\lambda(t, x) = x - \int_0^t A_\lambda X^\lambda(s, x) ds = X^\lambda(t', x) - \int_0^{t-t'} A_\lambda X^\lambda(s, X^\lambda(t', x)) ds$$

for any $t' < t$. In particular, $X^\lambda(t, x) = X^\lambda(t - t', X^\lambda(t', x))$.

Show that for any $x, y \in \mathcal{X}$, $t \mapsto |X^\lambda(t, x) - X^\lambda(t, y)|^2$ is non-increasing. Deduce that $|X^\lambda(t, x) - X^\lambda(t, y)| \leq |x - y|$ for all $t \geq 0$.

One simply observes that because A_λ is monotone:

$$\frac{1}{2} \frac{\partial}{\partial t} |X^\lambda(t, x) - X^\lambda(t, y)|^2 = -\langle X^\lambda(t, x) - X^\lambda(t, y), A_\lambda X^\lambda(t, x) - A_\lambda X^\lambda(t, y) \rangle \leq 0.$$

Then, we deduce

$$|X^\lambda(t, x) - X^\lambda(t, y)| \leq |X^\lambda(0, x) - X^\lambda(0, y)| = |x - y|.$$

5. Show that $t \mapsto |A_\lambda X^\lambda(t, x)|$ is nonincreasing. If $x \in \text{dom } A$, show that $|A_\lambda X(t, x)| \leq |A_0 x|$ for all $\lambda > 0$ and $t \geq 0$.

Hint: use that

$$X^\lambda(t+h, x) - X^\lambda(t, x) = X^\lambda(t-t', X^\lambda(t'+h, x)) - X^\lambda(t-t', X^\lambda(t', x))$$

for any $h > 0$, $t, t' < t$, and the previous question.

The contraction semi-group property shows that

$$\begin{aligned} |X^\lambda(t+h, x) - X^\lambda(t, x)| &= |X^\lambda(t-t', X^\lambda(t'+h, x)) - X^\lambda(t-t', X^\lambda(t', x))| \\ &\leq |X^\lambda(t'+h, x) - X^\lambda(t', x)| \end{aligned}$$

and dividing by h and sending $h \rightarrow 0$ it follows

$$|A_\lambda X(t, x)| \leq |A_\lambda X(t', x)|.$$

Then from the inequality $|p_\lambda| \leq |p|$ of question **2.** we deduce $|A_\lambda X(t, x)| \leq |A_0 x|$.

6. Using question **2.**, show that for $\lambda, \mu > 0$ and for any $x \in \mathcal{X}$, $A_\lambda x = A_\mu(x + (\mu - \lambda)A_\lambda x)$. Deduce that:

$$\frac{\partial}{\partial t} |X^\lambda(t, x) - X^\mu(t, x)|^2 \leq (\mu - \lambda)(|A_\lambda X^\lambda(t, x)|^2 - |A_\mu X^\mu(t, x)|^2)$$

(or any similar estimate) and in particular that if $x \in \text{dom } A$, $|X^\lambda(t, x) - X^\mu(t, x)| \leq C|A_0 x| \sqrt{|\mu - \lambda|t}$ for some constant $C > 0$.

What can you conclude? (Without justifying everything, unless there is still time.)

If $z = A_\lambda x = J_{A^{-1}/\lambda}(x/\lambda)$ (question **2.**), then

$$\begin{aligned} z + \frac{1}{\lambda} A^{-1} z \ni \frac{x}{\lambda} &\Leftrightarrow \frac{\lambda}{\mu} z + \frac{1}{\mu} A^{-1} z \ni \frac{x}{\mu} \\ &\Leftrightarrow z + \frac{1}{\mu} A^{-1} z \ni \frac{x}{\mu} + (1 - \frac{\lambda}{\mu})z = \frac{x + (\mu - \lambda)z}{\mu} \\ &\Leftrightarrow z = J_{A^{-1}/\mu}(\frac{x + (\mu - \lambda)z}{\mu}) = A_\mu(x + (\mu - \lambda)A_\lambda x) \end{aligned}$$

As a consequence,

$$\begin{aligned} \frac{\partial}{\partial t} |X^\lambda(t, x) - X^\mu(t, x)|^2 &= -2 \langle X^\lambda - X^\mu, A_\mu(X^\lambda + (\mu - \lambda)A_\lambda X^\lambda) - A_\mu X^\mu \rangle \\ &= -2 \langle (X^\lambda + (\mu - \lambda)A_\lambda X^\lambda) - X^\mu, A_\mu(X^\lambda + (\mu - \lambda)A_\lambda X^\lambda) - A_\mu X^\mu \rangle \\ &\quad + 2(\mu - \lambda) \langle A_\lambda X^\lambda, A_\lambda X^\lambda - A_\mu X^\mu \rangle \\ &\leq 2(\mu - \lambda) \langle A_\lambda X^\lambda, A_\lambda X^\lambda - A_\mu X^\mu \rangle. \end{aligned}$$

Symmetrically,

$$\frac{\partial}{\partial t} |X^\lambda(t, x) - X^\mu(t, x)|^2 \leq 2(\lambda - \mu) \langle A_\mu X^\mu, A_\mu X^\mu - A_\lambda X^\lambda \rangle$$

and averaging the two estimates we get the answer. Hence the time derivative is bounded by $|\lambda - \mu||A_0 x|^2$ and the estimate follows with $C = 1$ (integrating from 0 to t).

It follows that as $\lambda \rightarrow 0$, $X^\lambda(x, t)$ is a Cauchy sequence in $C^0([0, T]; \mathcal{X})$ (for any $T > 0$), which converges uniformly to some continuous path $X(t, x)$. This path is a solution of $\partial_t X + AX \ni 0$: precisely one can show that satisfies for all $t \geq 0$:

$$X(t, x) = x - \int_0^t A_0 X(s, x) ds.$$