

TDS OPTIMISATION – 2A (MMIS / IF)

Exercice 1 – Intersion : min-min vs. min-max.

a) Soit $f : X \times Y \rightarrow \mathbb{R}$; montrer que l'on peut choisir l'ordre de minimisation :

$$\inf_{(x,y) \in X \times Y} f(x,y) = \inf_{x \in X} \left(\inf_{y \in Y} f(x,y) \right) = \inf_{y \in Y} \left(\inf_{x \in X} f(x,y) \right).$$

b) En déduire un résultat sur la « régularisation de Moreau » : soit $g : X \rightarrow \mathbb{R}$ et $\gamma > 0$; montrer que

$$\inf_{x \in X} g(x) = \inf_{x \in X} G(x) \quad \text{où} \quad G(x) = \inf_{y \in X} g(y) + \frac{1}{2\gamma} \|x - y\|^2.$$

c) Attention : en général on ne peut pas « intervertir » min et max

$$\inf_{x \in X} \left(\sup_{y \in Y} f(x,y) \right) \neq \sup_{y \in Y} \left(\inf_{x \in X} f(x,y) \right),$$

comme le montre l'exemple dans $\mathbb{R} : X = [-1, 1], Y = \{-1, 1\}, f(x,y) = xy$. [Remarquez qu'on a tout de même toujours une inégalité]

Exercice 2 – Problème séparable.

a) Soient $f : X \rightarrow \mathbb{R}$ et $g : Y \rightarrow \mathbb{R}$; montrer que l'on peut « découpler » la minimisation de $f + g$:

$$\inf_{(x,y) \in X \times Y} f(x) + g(y) = \left(\inf_{x \in X} f(x) \right) + \left(\inf_{y \in Y} g(y) \right).$$

Montrer aussi que si le minimum est atteint pour f par $\bar{x} \in X$ et pour g par $\bar{y} \in Y$, alors (\bar{x}, \bar{y}) atteint le minimum de $f + g$ sur $X \times Y$.

b) Soient $c, \ell, u \in \mathbb{R}^n$; résoudre explicitement la minimisation de $c^\top x$ sous la contrainte $\ell \leq x \leq u$. Faire un dessin pour $n = 2$ pour visualiser le problème et la solution.

c) Même question quand la contrainte est la boule euclidienne. Quelle est la différence fondamentale ?

Exercice 3 – Conditions d'optimalité. Soient $f : \mathbb{R}^n \rightarrow \mathbb{R}$ une fonction différentiable, et $\bar{x} \in \mathbb{R}^n$.

a) Pour toute direction $u \in \mathbb{R}^n$, on définit l'application $q(t) := f(\bar{x} + t u)$ pour $t \in \mathbb{R}$. Calculer $q'(t)$.

b) Supposons que f soit deux fois différentiable. Calculer $q''(t)$.

On suppose que f admet un minimum local en \bar{x} , c'est-à-dire

$$\text{pour tout } x \text{ dans un voisinage de } \bar{x}, \quad f(x) \geq f(\bar{x}).$$

c) En utilisant le développement de Taylor-Young de la fonction q au premier ordre en 0, montrer que $\nabla f(\bar{x}) = 0$.

d) En utilisant le développement au second ordre, montrer que $\nabla^2 f(\bar{x})$ est « semidéfinie positive » (ce qui est aussi appelé « positive », c'est-à-dire que pour tout $u \in \mathbb{R}^n$, on a $u^\top \nabla^2 f(\bar{x}) u \geq 0$).

e) Pour le cas de la dimension $n = 2$, donner les conditions sur les dérivées partielles équivalentes aux deux propriétés des questions précédentes.

Exercice 4 – Sélecteur de Dantzig. On considère un modèle de régression $y = A\theta + \xi$ où le bruit est gaussien $\xi \sim \mathcal{N}(0, \sigma I_m)$. Les observations sont $A \in \mathbb{R}^{m \times n}$ et $y \in \mathbb{R}^m$; $\theta \in \mathbb{R}^n$ est le paramètre inconnu qu'on souhaite estimer. Dans le cas sur-paramétré (i.e. quand la taille n de θ est grande par rapport à m celle de y), le « sélecteur de Dantzig » consiste à résoudre le problème d'optimisation

$$\min_{\theta \in \mathbb{R}^n} \|\theta\|_1, \quad \text{sous contrainte } \|A^\top (A\theta - y)\|_\infty \leq \kappa \sigma$$

où $\kappa > 0$ est un hyper-paramètre.

- a) Notons $e(\theta) = 1/2 \|A\theta - y\|_2^2$ l'erreur quadratique du modèle. Observer que $\nabla e(\theta) = A^\top(A\theta - y)$.
- b) En introduisant des variables supplémentaires, reformuler ce problème comme un problème linéaire.
- c) Construire les vecteurs et matrices (c, G, h) pour écrire ce problème linéaire sous forme canonique (pour pouvoir ensuite le résoudre par un solveur disponible)

$$\min_x c^\top x \quad \text{sous contrainte } Gx \leq h$$

Exercice 5 – Lemme de descente. Soit une fonction $f : \mathbb{R}^n \rightarrow \mathbb{R}$ différentiable telle que son gradient $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ est L -Lipchitz, c'est-à-dire, que pour tous $x, y \in \mathbb{R}^n$,

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|.$$

- a) Se rappeler le « théorème fondamental de l'analyse » qui donne ici, pour tous $x, y \in \mathbb{R}^n$,

$$f(x) - f(y) = \int_0^1 (x - y)^\top \nabla f(y + t(x - y)) dt$$

- b) En déduire que pour tous $x, y \in \mathbb{R}^n$

$$f(x) \leq f(y) + (x - y)^\top \nabla f(y) + \frac{L}{2} \|x - y\|^2$$

- c) Donner une fonction f pour laquelle l'égalité est atteinte dans l'inégalité ci-dessus (et une autre pour laquelle elle ne l'est pas).
- d) Appliquer l'inégalité dans le cas $y = x - \gamma \nabla f(x)$ avec $\gamma < \frac{2}{L}$.

Exercice 6 – Ensembles convexes de matrices. Soit \mathcal{S}_n^{++} l'ensemble des matrices définies positives

$$\mathcal{S}_n^{++} = \{X \in \mathcal{S}_n : w^\top X w > 0, \text{ pour tout } w \in \mathbb{R}^n, w \neq 0\},$$

et \mathcal{S}_n^+ l'ensemble des matrices semi-définies positives

$$\mathcal{S}_n^+ = \{X \in \mathcal{S}_n : w^\top X w \geq 0, \text{ pour tout } w \in \mathbb{R}^n\}.$$

- a) Montrer que \mathcal{S}_n^{++} et \mathcal{S}_n^+ sont convexes. Montrer que \mathcal{S}_n^+ est un cône. Qu'en est-il de \mathcal{S}_n^{++} ?
- b) Montrer que \mathcal{S}_n^+ est fermé. Quelle est l'adhérence de \mathcal{S}_n^{++} ?

Exercice 7 – Fonctions convexes. Donner le domaine où les fonctions suivantes sont convexes

- a) $f(x, y) = x + 2y + y^2$ définie sur $\text{dom } f = \mathbb{R}^2$.
- b) $f(x, y) = x + 2y + y^2/x$ définie sur $\text{dom } f = \{(x, y) \in \mathbb{R}^2 : x \neq 0\}$.
- c) Montrer aussi que $f(x) = -|x| + \alpha x^2$ n'est jamais convexe (quelque soit $\alpha > 0$)

Exercice 8 – Fonctions supports. Soit C un sous-ensemble de \mathbb{R}^n ; on définit la fonction-support de C de la manière suivante :

$$\sigma_C(x) := \sup_{y \in C} x^\top y \quad \text{pour } x \in \mathbb{R}^n.$$

- a) Montrer que $\sigma_C : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{+\infty\}$ est convexe.
- b) Calculer la fonction-support pour des sous-ensembles de \mathbb{R}^n suivants :
 — C la boule euclidienne de rayon 1 ;
 — $C = (\mathbb{R}^+)^n$ l'orthant positif ;
 — $C = [a, b]$ le segment joignant deux points a et b dans \mathbb{R}^n .

Exercice 9 – Meilleure approximation polynômiale. Considérons N points dans le plan : $(x_i, y_i) \in \mathbb{R}^2$ pour $i = 1, \dots, N$. On cherche un polynôme de degré au plus n qui approche le nuage de points au mieux au sens des moindres carrés; en d'autres termes, on cherche $p(x) = a_0 + a_1x + \dots + a_nx^n$ pour minimiser la quantité $\sum_{i=1}^N |p(x_i) - y_i|^2$.

a) Montrer que ce problème s'écrit

$$\min_{z \in \mathbb{R}^d} \|Vz - b\|_2^2,$$

avec une variable z dont la dimension d est à préciser. Donner explicitement V et b .

- b) Écrire les conditions d'optimalité de ce problème. Préciser si elles sont nécessaires et/ou suffisantes.
- c) On suppose que les colonnes de V sont indépendantes. Montrer qu'il existe une unique solution au problème ; la donner explicitement.
- d) Montrer que si $N \geq n+1$ et si les x_i sont tous différents, alors les colonnes de V sont indépendantes. Que se passe-t-il si $N < n+1$?

Exercice 10 – Vendeur de journaux : modélisation, convexité, et résolution explicite. On considère un problème classique en logistique : le problème du vendeur de journaux. Dans sa forme la plus basique, ce problème se présente comme suit. Tous les matins, un vendeur de journaux achète un stock de quotidiens qu'il va revendre à ses clients tout au long de la journée. Le vendeur achète x journaux au prix unitaire c et qu'il vend ensuite au prix unitaire $p > c$.

- a) Supposons que la demande pour la journée à venir est connue ; quelle est la meilleure décision x pour le vendeur ?
- b) La difficulté est ainsi que la demande est inconnue au moment de prise de décision. En introduisant la demande comme une variable aléatoire positive ξ , écrire le problème du vendeur comme la maximisation sur x de son espérance de gain (ou la minimisation de son espérance de perte).
- c) Reformuler le problème (avec un changement de signe si besoin) comme la minimisation de la fonction $f(x) = cx - p\mathbb{E}[\min\{\xi, x\}]$. Montrer que f est convexe.

Pour la fin de l'exercice, on suppose que la décision est une variable continue $x \in \mathbb{R}^+$ et que l'incertitude suit une loi continue, dont on connaît la fonction de répartition H (et la densité h) :

$$H(u) = \mathbb{P}(\xi \leq u) = \int_0^u h(t) dt.$$

- d) Montrer que f est dérivable avec $f'(x) = c - p(1 - H(x))$.
- e) Résoudre le problème. [Hint : commencer par vérifier que $f'(0) < 0$]

Exercice 11 – Projection par optimisation. Dans cet exercice, on retrouve la propriété fondamentale des espaces de Hilbert à l'aide de l'optimisation, dans le cas particulier de la dimension finie. Soit un ensemble $C \subset \mathbb{R}^n$ convexe et fermé et un point $a \in \mathbb{R}^n$; intéressons-nous au problème :

$$(P) \quad \begin{cases} \min & \frac{1}{2}\|x - a\|^2 \\ x \in & C. \end{cases}$$

- a) Montrer qu'il existe une unique solution à (P).
- b) Montrer que la solution \bar{x} est caractérisée par la propriété

$$\bar{x} \in C \quad \text{et} \quad (a - \bar{x})^\top (x - \bar{x}) \leq 0, \quad \text{pour tout } x \in C.$$

- c) Énoncer les résultats des deux questions précédentes en termes géométriques. Donner une illustration sur un dessin.
- d) Question bonus : montrer que si $a \notin C$, on peut "séparer" C de a , c'est-à-dire qu'il existe $b \in \mathbb{R}^n$ et $\beta \in \mathbb{R}$ tels que

$$b^\top x < \beta < b^\top a, \quad \text{pour tout } x \in C.$$

- e) Question bonus, application de la question précédente : montrer que l'enveloppe convexe fermée d'un ensemble $A \subset \mathbb{R}^n$ est exactement l'intersection de tous les demi-espaces contenant A .

Exercice 12 – Approche par modèle : exemple du gradient. Soit $f: \mathbb{R}^n \rightarrow \mathbb{R}$ une fonction L -fortement différentiable (*i.e.* différentiable avec gradient L -Lipshitz) et \mathcal{X} un convexe fermé de \mathbb{R}^n .

a) Soit $x \in \mathbb{R}^n$; montrer que l'opérateur suivant est bien défini :

$$T(x) = \operatorname{argmin}_{y \in \mathcal{X}} \{ \langle \nabla f(x), y - x \rangle + L/2 \|x - y\|^2 \}.$$

- b) Dans le cas sans contrainte ($\mathcal{X} = \mathbb{R}^n$), retrouver que l'itération $x_{k+1} = T(x_k)$ correspond à l'algorithme du gradient (à pas constant $1/L$).
- c) Dans le cas général, expliciter l'itération $x_{k+1} = T(x_k)$ en écrivant les conditions d'optimalité du problème. Quel algorithme reconnaissez-vous ?
- d) En vous aidant de la question b de l'exercice 5, montrer qu'il s'agit bien d'un algorithme de descente.

Exercice 13 – Stabiliser Newton. Nous souhaitons résoudre une équation $F(x) = 0$ avec $F: \mathbb{R}^n \rightarrow \mathbb{R}^n$ est une fonction de classe C^1 . On suppose que le jacobien $JF(x)$ est toujours inversible, et on propose d'utiliser la « méthode de Newton pour résoudre les équations » pour le problème $F(x) = 0$. Une itération de cette méthode est :

$$x_{k+1} = x_k - [JF(x_k)]^{-1} F(x_k). \quad (\text{N})$$

On considère maintenant la fonction

$$\varphi: \begin{cases} \mathbb{R}^n \longrightarrow \mathbb{R} \\ x \longmapsto \|F(x)\|_2^2. \end{cases}$$

- a) Calculer le gradient $\nabla \varphi(x)$.
- b) Montrer que $d_k = -[JF(x_k)]^{-1} F(x_k)$ est une direction de descente de φ en x_k (si $F(x_k) \neq 0$).
- c) Comment pourrait-on stabiliser l'itération (N) ?

Exercice 14 – Exo du cours avec une contrainte en plus. Soient les deux vecteurs de \mathbb{R}^n , $e = [1, \dots, 1]^\top$ et $c = [1, 0, \dots, 0]^\top$. Résoudre le problème

$$\begin{cases} \max & c^\top x \\ & e^\top x = 0 \\ & \|x\|^2 \leq 1. \end{cases}$$

Visualiser, sur un dessin, le problème et la solution en dimension $n = 2$.

Exercice 15 – Contraintes actives. Soit l'ensemble $C = \{(x, y) \in \mathbb{R}^2 : x + y \leq 1, x \geq 0, y \geq 0\}$.

- a) Dessiner C . En considérant les contraintes actives, exhiber 7 zones dans C .
- b) Écrire les conditions d'optimalité de KKT de la minimisation de la fonction $f(x, y) = \exp(x - y) - x - y$ sur C . Sont-elles nécessaires et/ou suffisantes ?
- c) Trouver le minimum global de cette fonction sur C .

Exercice 16 – Problèmes sur-paramétrés. Les réseaux de neurones profonds introduisent des problèmes sur-paramétrés (avec plus de variables que de contraintes) pour lesquels les méthodes d'optimisation introduisent un biais implicite. Illustrons ce phénomène dans un cas très simple.

Soit le système linéaire $Ax = b$ avec une matrice A de taille $m \times n$ ($m \leq n$) et un vecteur $b \in \operatorname{Im} A$.

a) Quel est l'ensemble des solutions du problème sur-paramétré $Ax = b$?

Une approche pour résoudre ce problème est de le formuler comme un problème d'optimisation

$$(P) \begin{cases} \min & \frac{1}{2} \|Ax - b\|^2 \\ & x \in \mathbb{R}^n. \end{cases}$$

b) Écrire l'algorithme de gradient pour résoudre (P). Montrer qu'en prenant par exemple $x_0 = 0$ on a $x_k \in \operatorname{Im} A^\top$ pour tout $k \geq 0$.

- c) Quel pas de descente choisir a priori pour garantir la convergence? Supposons dans la suite avoir convergence de l'algorithme vers \bar{x} . Montrer que \bar{x} est solution de

$$(P) \begin{cases} \min & \frac{1}{2}\|x\|^2 \\ & Ax = b, \end{cases}$$

c'est-à-dire que \bar{x} est la plus petite solution du système sur-paramétré. [hint : Ecrire les conditions d'optimalité de KKT ; discuter nécessaires/suffisantes...].

- d) Refaire le raisonnement précédent avec SGD, l'algorithme de gradient stochastique ou incrémental, pour arriver à la même conclusion.

Exercice 17 – Norme duale par KKT. Soient $p > 1$ et $q = p/(p-1)$ tels que $1/p + 1/q = 1$. Pour un vecteur non-nul à coefficients positifs $y \in (\mathbb{R}_+)^n$ fixé, considérons le problème d'optimisation

$$(P) \begin{cases} \max & x^\top y \\ & x \geq 0, \sum_{i=1}^n x_i^p \leq 1. \end{cases}$$

- a) Montrer que les hypothèses sont réunies pour pouvoir appliquer le théorème de KKT. Les conditions d'optimalité seront-elles nécessaires et/ou suffisantes? Pour l'optimalité globale ou locale?
- b) Écrire les conditions d'optimalité pour (P).
- c) Considérons les différentes possibilités pour les contraintes actives. Montrer que les multiplicateurs associés aux contraintes $x_i \geq 0$ sont nuls et que le multiplicateur associé à la contrainte $\sum_{i=1}^n x_i^p \leq 1$ est non-nul, si bien que les conditions se réduisent aux $n+1$ équations suivantes en $(x, \mu) \in (\mathbb{R}_+)^{n+1}$

$$\begin{cases} -y_i + p\mu x_i^{p-1} = 0 & \text{pour tout } i = 1, \dots, n \\ \sum_{i=1}^n x_i^p = 1. \end{cases}$$

- d) Donnez la solution optimale et la valeur optimale du problème (P).
[Indication : pas besoin d'explicitier μ , travailler avec $c = p\mu > 0$ suffit.]
- e) En déduire l'inégalité de Hölder : pour tous $u, v \in \mathbb{R}^n$,

$$\sum_{i=1}^n |u_i||v_i| \leq \|u\|_p \|v\|_q.$$

Exercice 18 – Entropie et dualité. Nous souhaitons estimer un vecteur inconnu $\bar{x} \in \mathbb{R}^n$ à coefficients positifs, dont nous ne connaissons que certaines de ses « réalisations », c'est-à-dire nous ne connaissons que $a_i^\top \bar{x} = b_i$, pour des $a_i \in \mathbb{R}^n$ et $b_i \in \mathbb{R}$ connus ($i = 1, \dots, m$), avec $m < n$. On note A la matrice dont les lignes sont les a_i^\top , et $b \in \mathbb{R}^m$ le vecteur des b_i .

Parmi tous les vecteurs vérifiant ces conditions, on décide de préférer un vecteur « entropique » défini comme une solution du problème (où \log est le logarithme népérien)

$$(P) \begin{cases} \min & \sum_{k=1}^n x_k \log(x_k) \\ & Ax = b, \\ & x \geq 0. \end{cases}$$

- a) Considérons la fonction $\varphi_\alpha : \mathbb{R}_+ \rightarrow \mathbb{R}$ définie pour une constante $\alpha \in \mathbb{R}$ donnée par

$$\varphi_\alpha(t) = \begin{cases} t \log(t) + \alpha t & \text{si } t > 0 \\ 0 & \text{si } t = 0. \end{cases}$$

Montrer que φ_α est convexe et coercive sur \mathbb{R}_+^* .

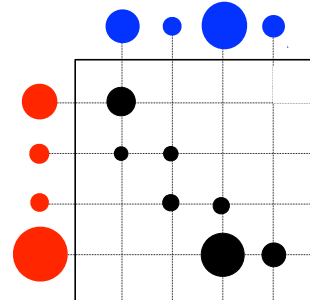
- b) Noter que φ_α est continue sur \mathbb{R}_+ tout entier ; en déduire que φ_α est convexe sur \mathbb{R}_+ .
- c) Donner le t qui atteint le minimum de φ_α sur \mathbb{R}_+ , ainsi que la valeur de ce minimum.

Revenons à présent au problème (P), que l'on va résoudre par dualité. On suggère de dualiser uniquement la contrainte couplante $Ax = b$.

- d) Mettre le problème sous la forme du cours. Donner l'expression du lagrangien, et donner la définition de la fonction duale et du problème dual. On notera λ la variable duale, θ la fonction duale.
- e) Observer que la maximisation du lagrangien (à λ fixé) se découple en n problèmes. Donner l'unique $x_\lambda \in \mathbb{R}^n$ qui maximise le lagrangien (à λ fixé).
- f) Montrer que θ est différentiable, et donner l'expression de $\theta(\lambda)$ et de $\nabla\theta(\lambda)$ en fonction de x_λ .
- g) Supposons qu'il existe une solution duale qu'on note λ^* . Comment pourrait-on la calculer ?
- h) Montrer que x_{λ^*} est réalisable dans le primal. En déduire que x_{λ^*} est une solution primale.
- i) Expliquer finalement comment répondre à la question initiale de cet exercice.

Exercice 19 – Transport optimal et distance de Wasserstein. Soient deux vecteurs positifs $a \in \mathbb{R}_+^n$ et $b \in \mathbb{R}_+^m$ tels que $\sum_{i=1}^n a_i = 1$ $\sum_{j=1}^m b_j = 1$ (représentant ainsi des densités de probabilité discrètes). On souhaite faire un transport optimal de a vers b : on veut trouver une matrice $P = (P_{ij}) \in \mathbb{R}_+^{n \times m}$ qui représente comment se répartit chaque a_i vers les b_j étant donnés des coûts associés $C_{ij} \geq 0$. Ce problème se formule comme

$$W(a, b) = \begin{cases} \min & \sum_{i=1}^n \sum_{j=1}^m C_{ij} P_{ij} \\ & \sum_{j=1}^m P_{ij} = a_i, \quad \text{pour tout } i = 1, \dots, n \\ & \sum_{i=1}^n P_{ij} = b_j, \quad \text{pour tout } j = 1, \dots, m \\ & P_{ij} \geq 0 \quad \text{pour tout } i = 1, \dots, n \text{ et } j = 1, \dots, m \end{cases}$$



Sur la figure, la distribution discrète $a \in \mathbb{R}^4$ est en rouge et $b \in \mathbb{R}^4$ en bleu. On a $n = m = 4$ avec $a_4 \geq a_1 \geq a_2 \geq a_3$ et $b_3 \geq b_1 \geq b_4 \geq b_2$. Les points noirs représentent les coefficients non-nuls de P .

- a) On considère la dualisation de toutes les contraintes sur les lignes ($a_i - \sum_{j=1}^m P_{ij} = 0$ pour tout i) et sur les colonnes ($b_j - \sum_{i=1}^n P_{ij} = 0$ pour tout j). Mettre le problème sous la forme du cours, introduire le lagrangien associé, et définir la fonction duale. On notera les variables duales $\lambda^a = (\lambda_i^a)_{i=1, \dots, n} \in \mathbb{R}^n$ et $\lambda^b = (\lambda_j^b)_{j=1, \dots, m} \in \mathbb{R}^m$.
- b) Montrer qu'il n'y a pas de saut dual. En déduire que

$$W(a, b) = \begin{cases} \max & a^\top \lambda^a + b^\top \lambda^b \\ & \lambda_i^a + \lambda_j^b \leq C_{ij}, \quad \text{pour tout } i = 1, \dots, n \text{ et } j = 1, \dots, m \end{cases}$$

Remarque culturelle : $W(a, b)$ définit une distance ; on l'appelle distance de Wasserstein entre les deux distributions a et b .

- c) Ecrire les conditions de KKT de ce problème. [ne pas essayer de les résoudre !]

Exercice 20 – Transport optimal régularisé par entropie. Reprenons le problème de transport optimal vu en cours, auquel nous allons ajouter la régularisation entropique :

$$H(P) = \sum_{i=1}^n \sum_{j=1}^m P_{ij} (\log(P_{ij}) - 1) \quad (\text{où } \log \text{ est le logarithme népérien}).$$

Nous considérons donc le problème, avec $\varepsilon > 0$,

$$(P) \quad \min_{P \in \mathcal{U}(a, b)} \sum_{i=1}^n \sum_{j=1}^m C_{ij} P_{ij} + \varepsilon H(P)$$

où $\mathcal{U}(a, b)$ est l'ensemble des plans de transports

$$\mathcal{U}(a, b) = \left\{ P \in \mathbb{R}^{n \times m} \text{ tel que } \begin{cases} \sum_{j=1}^m P_{ij} = a_i & \text{pour tout } i = 1, \dots, n, \\ \sum_{i=1}^n P_{ij} = b_j & \text{pour tout } j = 1, \dots, m \\ P_{ij} \geq 0 & \text{pour tout } i, j = 1, \dots, n \text{ et } j = 1, \dots, m \end{cases} \right\}$$

entre deux vecteurs $a \in \mathbb{R}_+^n$ et $b \in \mathbb{R}_+^m$ tels que $\sum_{i=1}^n a_i = 1$ $\sum_{j=1}^m b_j = 1$.

a) Soit la fonction $\varphi: \mathbb{R}_+ \rightarrow \mathbb{R}$ définie et continue sur \mathbb{R}_+

$$\varphi(t) = \begin{cases} t \log(t) & \text{si } t > 0 \\ 0 & \text{si } t = 0. \end{cases}$$

Montrer que φ est strictement convexe sur \mathbb{R}_+ .

b) Montrer que φ est en fait strictement convexe sur \mathbb{R}_+ tout entier. [Hint : on peut observer que, pour $0 < \alpha < 1$ et $t > 0$, on a $\log(\alpha t) < \log(t)$.]

c) En déduire que la fonction $H: \mathbb{R}_+^{n \times m} \rightarrow \mathbb{R}$ est continue et strictement convexe.

d) Montrer qu'il existe une unique solution à (P). Notons-la P_ε .

e) En introduisant la matrice $K = (K_{ij}) \in \mathbb{R}^{n \times m}$ définie par $K_{ij} = \exp(-C_{ij}/\varepsilon)$ pour tout (i, j) , réarranger l'objectif pour montrer que

$$P_\varepsilon = \operatorname{argmin}_{P \in \mathcal{U}(a,b)} \sum_{i=1}^n \sum_{j=1}^m P_{ij} \log(P_{ij}/K_{ij}).$$

Remarque culturelle : ceci signifie que P_ε s'interprète comme la projection, au sens de la divergence de Kullback-Leiber, de K (appelée noyaux de Gibbs) sur $\mathcal{U}(a, b)$. Indice pour la question : utiliser que la somme des P_{ij} est constante.

Reprenons le problème (P) initial et étudions une approche duale. Nous allons dualiser toutes les contraintes d'égalité de $\mathcal{U}(a, b)$ (et pas les contraintes de positivité). La fonction $\varphi_\alpha: \mathbb{R}_+ \rightarrow \mathbb{R}$ définie pour $\alpha \in \mathbb{R}$ par

$$\varphi_\alpha(t) = \begin{cases} \varepsilon t \log(t) + \alpha t & \text{si } t > 0 \\ 0 & \text{si } t = 0 \end{cases}$$

apparaîtra dans les développements.

f) Reformuler (P) sous la forme du cours (en changeant le signe pour avoir un max). Définir le lagrangien et la fonction duale θ . On notera les variables duales $\lambda^a \in \mathbb{R}^n$ et $\lambda^b \in \mathbb{R}^m$.

g) Montrer que

$$\theta(\lambda^a, \lambda^b) = -a^\top \lambda^a - b^\top \lambda^b - \sum_{i=1}^n \sum_{j=1}^m \min_{P_{ij} \geq 0} \varphi_{\alpha_{ij}}(P_{ij})$$

pour des $\alpha_{ij} \in \mathbb{R}$ que l'on explicitera.

h) Calculer le minimum sur \mathbb{R}_+ de la fonction φ_α .

i) En déduire que

$$\theta(\lambda^a, \lambda^b) = -a^\top \lambda^a - b^\top \lambda^b + \varepsilon \sum_{i=1}^n \sum_{j=1}^m \exp((-C_{ij} + \lambda_i^a + \lambda_j^b)/\varepsilon)$$

Comparer avec le dual du problème non-régularisé ($\varepsilon = 0$) vu en cours. Interpréter l'impact de la régularisation sur le dual.

j) En déduire que θ est différentiable et donner les expressions de $\frac{\partial}{\partial \lambda_i^a} \theta(\lambda^a, \lambda^b)$ pour tout i , ainsi que $\frac{\partial}{\partial \lambda_j^b} \theta(\lambda^a, \lambda^b)$ pour tout j .

k) Montrer que l'unique solution optimisant le lagrangien, pour (λ^a, λ^b) fixés,

$$(P_{\lambda^a, \lambda^b})_{ij} = \exp((-C_{ij} + \lambda_i^a + \lambda_j^b)/\varepsilon) \quad \text{pour tout } (i, j).$$

Ré-écrire les dérivées partielles de θ en (λ^a, λ^b) en fonction de P_{λ^a, λ^b} .

l) Ecrire le problème dual. Que proposez-vous pour le résoudre numériquement ?

m) Supposons avoir les solutions duales $(\bar{\lambda}^a, \bar{\lambda}^b)$; montrer que $P_{\bar{\lambda}^a, \bar{\lambda}^b}$ est réalisable. En déduire qu'il n'y a pas de saut dual et que $P_\varepsilon = P_{\bar{\lambda}^a, \bar{\lambda}^b}$.

n) En déduire l'expression classique de P_ε , avec la matrice K de la question e :

$$P_\varepsilon = \text{diag}(\exp(\bar{\lambda}^a/\varepsilon))K \text{diag}(\exp(\bar{\lambda}^b/\varepsilon)).$$

Notation : pour un vecteur λ , on note $\text{diag}(\exp(\lambda))$ la matrice diagonale avec les coefficients $\exp(\lambda_i)$ sur la diagonale.

Exercice 21 – Identité de Moreau. On rappelle une propriété du cours : si $f: \mathbb{R} \rightarrow \mathbb{R}$ est une fonction convexe (vérifiant une hypothèse technique supplémentaire), alors elle est égale à sa « bi-conjuguée » $f = (f^*)^*$ et on peut « inverser » les sous-différentiels

$$u \in \partial f(x) \iff x \in \partial f^*(u). \quad (\text{F})$$

On considère dans cet exercice une fonction convexe f vérifiant cette propriété. Soit $x \in \mathbb{R}^n$ fixé.

- a) Écrire les conditions d'optimalité du problème définissant $\text{prox}_f(x)$. Sont-elles nécessaires et/ou suffisantes ? Mêmes questions pour $\text{prox}_{f^*}(x)$.
- b) Montrer que $x - \text{prox}_f(x) = \text{prox}_{f^*}(x)$. [Hint : utiliser (F)]. On vient d'établir l'identité de Moreau

$$\text{Id} = \text{prox}_f + \text{prox}_{f^*}.$$

- c) Que donne cette identité dans le cas où f est indicatrice d'un sous-espace vectoriel de \mathbb{R}^n ? (On pourra utiliser les notations : $V \subset \mathbb{R}^n$ et $f(x) = i_V(x)$)

Exercice 22 – Florilège. Les propositions suivantes, fausses, sont extraites de copies d'examen des années passées... Trouver des contre-exemples ou des arguments qui montrent que les énoncés suivants sont **faux**. Les modifier ensuite pour les rendre vrais. Et puis, ne plus faire ces erreurs...

- a) f est quadratique, donc convexe ;
- b) f est convexe, donc il existe un minimum ;
- c) f est continue sur un fermé et minorée, donc il existe un minimum ;
- d) f a un seul minimum local donc il est global ;
- e) f est affine, donc coercive ;
- f) f est convexe et C compact, donc il y a unicité de la solution.
- g) f différentiable et C convexe, alors f est convexe sur C ;
- h) f est positive donc admet un minimum ;
- i) l'image réciproque d'un compact par une application continue est compact ;
- j) $x \mapsto |x|$ est linéaire ;
- k) $\{x \in \mathbb{R}^n : \|Ax - b\| \leq 1\}$ est borné ;
- l) \mathbb{R}^n est compact ;
- m) la composée de deux fonctions convexes est convexe.
- n) le problème dual est convexe, donc il n'y a pas de saut dual ;
- o) la fonction $(x, y) \mapsto xy - 1$ est convexe ;
- p) on a $\|\nabla f(x)\| = 0$ donc x est un minimum local de f sur \mathbb{R}^n ;
- q) x est un minimum local de f sur C donc $\|\nabla f(x)\| = 0$;
- r) il existe une unique solution à un problème de moindres carrés $\min_x \|Ax - b\|^2$;
- s) la méthode de gradient à pas constant converge ;
- t) la méthode de Newton converge.