

Optimisation Convexe : DM3

① Interior Point Method

$$\bullet \phi_t(x) = t \left(\frac{1}{2} x^T \varphi x + p^T x \right) - \sum_{i=1}^m \log(b - a_i^T x)$$

où a_i est la $i^{\text{ème}}$ ligne de A .

$$\nabla \log(b - a_i^T x) = \frac{\nabla(b - a_i^T x)}{b - a_i^T x} = \frac{-a_i}{b - a_i^T x}$$

$$\text{dnc} \quad \boxed{\nabla \phi_t(x) = t(\varphi x + p) + \sum_{i=1}^m \frac{a_i}{b - a_i^T x}}$$

$$\nabla^T \left(\frac{1}{b - a_i^T x} \right) = - \frac{\nabla^T(b - a_i^T x)}{(b - a_i^T x)^2} = \frac{a_i^T}{(b - a_i^T x)^2}$$

$$\text{dnc} \quad \boxed{\nabla^2 \phi_t(x) = t\varphi + \sum_{i=1}^m \frac{a_i a_i^T}{(b - a_i^T x)^2}}$$

③ Support Vector Machine Problem

1) Primal: On cherche $\begin{pmatrix} w \\ z \end{pmatrix} \in \text{Int } \mathcal{D}$ tel que $y_i(w^T x_i) > 1 - z_i$ et $z_i > 0$, $i=1, \dots, m$

$$\text{Sci } \mathcal{D} = \mathbb{R}^d \times \mathbb{R}^m \quad (w \in \mathbb{R}^d, z \in \mathbb{R}^m) \quad \text{dnc } \text{Int } \mathcal{D} = \mathbb{R}^d \times \mathbb{R}^m$$

dnc l'ensemble des points strictement faisibles est: $\left\{ \begin{pmatrix} w \\ z \end{pmatrix} \in \mathbb{R}^d \times \mathbb{R}^m \mid z_i > 1 - y_i(w^T x_i) \text{ et } z_i > 0, \right.$
 $\left. i=1, \dots, m \right\}$

On peut dnc prendre par exemple: $w = 0_{\mathbb{R}^d}$ et $z = 2 \times \frac{1}{m} \left(\begin{smallmatrix} 1 \\ \vdots \\ 1 \end{smallmatrix} \right)_{m \times 1}$

car alors $z_i > 0$ et $z_i > 1 - y_i(0_{\mathbb{R}^d}^T x_i) = 1$.

Dual: On cherche $\lambda \in \text{Int } \mathcal{D}$ tel que $0 < \lambda < \frac{1}{2m}$

Sci $\mathcal{D} = \mathbb{R}^m$ dnc $\text{Int } \mathcal{D} = \mathbb{R}^m$. L'ensemble de points strictement faisibles est dnc:

$\left\{ \lambda \in \mathbb{R}^m \mid 0 < \lambda < \frac{1}{2m} \right\}$. On peut dnc prendre par exemple:

$$\lambda = \frac{1}{22m}$$

2) Reformulons le dual et le primal pour les écrire sous la forme $\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T Q x + p^T x$ st $Ax \leq b$:

Primal: $\frac{1}{2} \|w\|_2^2 + \frac{1}{2m} \sum_{i=1}^m z_i = \frac{1}{2} w^T w + \frac{1}{2m} 1_m^T z$

dnc on pose $Q = \begin{pmatrix} I_d & 0 \\ 0 & 0 \end{pmatrix}_{\substack{d \\ d \quad m}}$ où $I_d = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}_d$

$Q \in M_{m+d}$ et est symétrique positive

$$p = \begin{pmatrix} 0 \\ 0 \\ \frac{1}{2m} \mathbf{1}_m \end{pmatrix}^T d \quad \text{On a donc } \frac{1}{2} \|w\|_2^2 + \frac{1}{2m} \sum_{i=1}^m z_i = \frac{1}{2} (w, z) Q \begin{pmatrix} w \\ z \end{pmatrix} + p^T \begin{pmatrix} w \\ z \end{pmatrix}$$

Écrivent pour les contraintes: on pose $A = \begin{pmatrix} -\text{diag}(y)X & -I_m \\ 0 & -I_m \end{pmatrix} \begin{matrix} d \\ m \\ m \end{matrix}$, $A \in M_{2m, m+d}$

où $\text{diag}(y) = \begin{pmatrix} y_1 & & 0 \\ & \ddots & \\ 0 & & y_m \end{pmatrix}$ donc $\text{diag}(y)X = \begin{pmatrix} x_1^T y_1 \\ \vdots \\ x_n^T y_m \end{pmatrix}$

et $I_m = \begin{pmatrix} 1 & & 0 \\ & \ddots & \\ 0 & & 1 \end{pmatrix}_m$

$b = \begin{pmatrix} -\frac{1}{2m} \mathbf{1}_m \\ 0_{2m} \end{pmatrix}$ et les contraintes s'écrivent alors: $A \begin{pmatrix} w \\ z \end{pmatrix} \leq b$

On a donc $\min_{(w,z)} \frac{1}{2} \|w\|_2^2 + \frac{1}{2m} \sum_{i=1}^m z_i$ st $y_i(w^T x_i) \geq 1 - z_i, z_i \geq 0 \quad \forall i=1, \dots, m$

qui s'écrit: $\min_{(w,z)} \frac{1}{2} (w, z) Q \begin{pmatrix} w \\ z \end{pmatrix} + p^T \begin{pmatrix} w \\ z \end{pmatrix}$ st $A \begin{pmatrix} w \\ z \end{pmatrix} \leq b$

Dual: On prend $-\left(\frac{1}{2} \left\| \sum_{i=1}^m \lambda_i y_i x_i \right\|_2^2\right) + \frac{1}{m} \mathbf{1}^T \lambda$ pour avoir le problème

de minimisation équivalent

$$\frac{1}{2} \left\| \sum_{i=1}^m \lambda_i y_i x_i \right\|_2^2 = \frac{1}{2} \left(\sum_{i=1}^m \sum_{j=1}^m \lambda_i \lambda_j y_i y_j x_i^T x_j \right)$$

On pose $B = \text{diag}(y)X$ et alors on pose $Q = BB^T$, $Q \in M_m$ et est symétrique positive

$p = \mathbf{1}_m$

On a donc $\frac{1}{2} \left\| \sum_{i=1}^m \lambda_i y_i x_i \right\|_2^2 - \frac{1}{n} \mathbf{1}^T \lambda = \frac{1}{2} \lambda^T Q \lambda + p^T \lambda$

Puis pour les contraintes : on pose $A = \begin{pmatrix} I_m \\ -I_m \end{pmatrix}_m \quad A \in \mathbb{M}_{2m, m}$

et $b = \begin{pmatrix} \frac{1}{2m} \mathbf{1}_m \\ 0_{\mathbb{R}^m} \end{pmatrix}_m$ et les contraintes s'écrivent alors : $A \lambda \leq b$

On a donc $\min_{\lambda} \frac{1}{2} \left\| \sum_{i=1}^m \lambda_i y_i x_i \right\|_2^2 - \frac{1}{n} \mathbf{1}^T \lambda$ st $0 \leq \lambda \leq \frac{1}{2m}$ qui est

équivalent au dual qui s'écrit : $\min_{\lambda} \frac{1}{2} \lambda^T Q \lambda + p^T \lambda$ st $A \lambda \leq b$.

Implémentation / Commentaires:

Tous les graphiques sont en annexe. On utilise toujours une tolérance de 10^{-3} .

Test des méthodes sur une fonction très simple: (fichiers test-damped.py et test-LS.py)

On va commencer par étudier une exemple très simple avant la question 3).

On étudie le problème: $\min_x \frac{1}{2}x^2 + x$ st $x \geq 0$, $x \in \mathbb{R}$

On a donc $q=1$, $p=1$, $A=-1$ et $b=0$

On cherche son dual: le Lagrangien s'écrit, $L(x, \lambda) = \frac{1}{2}x^2 + x - \lambda x$ pour $\lambda \geq 0$
et $x \in \mathbb{R}$

La fonction dual est: $g(\lambda) = \inf_{x \in \mathbb{R}} \left(\frac{1}{2}x^2 + x(1-\lambda) \right)$, $\lambda \geq 0$ (L convexe en x)

$$\frac{\partial L(x, \lambda)}{\partial x} = 0 \Leftrightarrow x = \lambda - 1.$$

$$\text{Donc } g(\lambda) = \frac{1}{2}(\lambda-1)^2 - (\lambda-1)^2 = -\frac{1}{2}(\lambda-1)^2 = -\frac{1}{2}\lambda^2 + \lambda - \frac{1}{2}$$

et ainsi, le problème dual est: $\max_{\lambda \in \mathbb{R}} -\frac{1}{2}\lambda^2 + \lambda - \frac{1}{2}$ st $\lambda \geq 0$.

On a $x_0=1$ et $\lambda_0=1$ des points strictement faisables pour le primal et le dual respectivement.

$$\eta = f_0(x_0) - g(\lambda_0) = \frac{1}{2} + 1 - 0 = \frac{3}{2} \quad \text{le "duality gap" de départ.}$$

$$\text{On prend donc } t_0 = \frac{m \leftarrow \text{nb de contraintes}}{2} = \frac{1}{2} = \frac{2}{3} \quad \text{comme } t \text{ de départ de la}$$

méthode log-barrière. On prend aussi x_0 au initialisation.

On voit on utilise les KKT conditions que:

$$\begin{cases} -x \leq 0 \\ \lambda \geq 0 \\ -\lambda x = 0 \\ x = 1-1 \end{cases} \Leftrightarrow \begin{cases} x \geq 0 \\ \lambda \geq 0 \\ \lambda(1-1) = 0 \\ x = 1-1 \end{cases} \Leftrightarrow \begin{cases} x \geq 0 \\ \lambda \geq 0 \\ \lambda = 0 \text{ ou } \lambda = 1 \\ x = 1-1 \end{cases}$$

on a le x optimal $x^* = 0$. On va lancer la méthode log-barrière on utilise

"damped Newton" puis une autre fois on utilise "LS Newton" ^{déjà} pour vérifier que

deux méthodes fonctionnent dans ce cas très simple. (On utilise $\alpha = 0,01$ et $\beta = 0,5$ pour "LS Newton")

On trace également les courbes "duality gap versus iterations" pour $p = 2, 15, 50$ et 100 .

On observe déjà que le x optimal est très proche de 0 pour chaque p et pour les deux méthodes (utilisant "damped Newton" et "LS Newton").

On peut observer ensuite sur les graphes que pour les deux méthodes, quand p est petit, il y a plus d'"outer iterations" et moins d'"inner iterations" et inversement quand

p est grand. On observe également que pour $p = 2$, l'algorithme se termine dans les 25-30

itérations dans les deux méthodes alors qu'il se termine en une quinzaine d'itérations

(pour les deux méthodes) pour $p = 15, 50$ et 200 .

On peut donc déjà être rassuré sur l'implémentation des deux méthodes qui donnent des

résultats cohérents dans cet exemple très simple (et notamment sur le fait que

le nombre d'itérations ne varie pas beaucoup pour p au delà de 10¹ ^{du moins jusqu'à 100} et qu'il est plus grand pour p petit).

3) On va maintenant tester plus sérieusement nos algorithmes en effectuant la classification du dataset Iris (pour les classes Iris-versicolre et Iris-virginica).

Tout d'abord, on centre nos données en leur soustrayant la moyenne des données. On affecte ensuite à la classe Iris-versicolre le label -1 et le label 1 à la classe Iris-virginica. On sélectionne ensuite aléatoirement les données (80% du dataset) avec lesquelles on va chercher l'hyperplan (passant par l'origine) séparant le "meilleur" ces données selon notre critère de risque. Le reste des données sera utilisé pour tester la performance du classifieur.

Pour initialiser la méthode log-barrière on va utiliser les points strictement faisables qu'on avait trouvés précédemment (question 1) SVM problem) afin de calculer la dualité gap initial g et on prend alors $t_0 = \frac{m}{2}$ ^{nombre de contraintes}. On prend également le point strictement faisable du primal en initialisation ainsi que $p = 10$ (ce qu'on a vu en cours et vérifié ici dans l'exemple très simple que le choix de p n'était pas critique pour p dans un assez large intervalle (environ 3 à 100) donc prendre p vers 10-20 est satisfaisant).

On prend aussi dans le cas où LS-Newton est utilisée $\alpha = 0,01$ et $\beta = 0,5$.

On obtient à la fin de la méthode, une solution w qu'on utilise pour classer les points des 20% du dataset restant avec la règle: $w^T x_i > 0 \rightarrow \text{label } 1, w^T x_i \leq 0 \rightarrow \text{label } -1$.

On a fait cela pour les deux méthodes \log -barrière (avec damped Newton et LS-Newton) et pour Z allant de 0.1 à 10 par pas de 0.1. On trace alors la courbe représentant l'erreur de classification des 20% restant en fonction de Z . On répète un certain

nombre de fois (6) tout l'algorithme pour déterminer une tendance générale de l'erreur

en fonction de Z . (Remarque: taux d'erreur = $\frac{1}{m_1} \sum_{i=1}^{m_1} \frac{1}{|y_i + g(x_i)|}$ où m_1 : nb de données dans les 20% de dataset restant
 $g(x_i) = 2 \frac{1}{w^T x_i + 0.5} - 1$)

On observe sur les graphiques pour les deux méthodes une tendance générale qui semble

être que le taux d'erreur augmente rapidement avec Z entre 0.1 et 2 environ puis plafonne ensuite. Cependant, même si le taux semble augmenter rapidement ^{entre 0.1 et 2} il ne

semble pas dépasser généralement les 15%. De cette tendance générale sur nos 12 graphiques

on en retire donc que 'il semblerait judicieux de prendre un Z petit ce qu'on fera pour la question 4) en prenant $Z = 0.1$.

On termine cette question en générant aléatoirement un nouveau dataset 300.

Pour cela on génère des données centrées par des gaussiennes de \mathbb{R}^4 avec comme matrice

de variance-covariance, la matrice de variance-covariance empirique de nos données du dataset d'origine. On utilise ensuite l'hyperplan séparateur obtenu (avec $Z=0.1$) via les 80% de notre dataset d'origine afin de classer les données générées. Puis, on "décentre" nos nouvelles données en leur ajoutant la moyenne de nos données d'origine.

4) Dans cette question, toujours en utilisant le 80% dataset, on va tracer le "duality gap versus iterations" pour les problèmes primal et dual résolus par la méthode \log -boîte. On utilise donc pour être consistant avec la question précédente les 80% du dataset. On prend la même t_0 et les mêmes points initiaux strictement faisables que précédemment ainsi que $Z=0.1$ (et $\lambda=0.01$ et $\beta=0.5$ pour LS-Newton). On fait cela pour $p=2, 15, 50$ et 100 .

On obtient quatre graphiques, donc par chaque méthode \log -boîte (damped et LS-Newton) (un pour le primal et un pour le dual pour chaque méthode).

Ponctuel: Pour tracer ces courbes, on utilise le fait qu'on a un duality gap égal à $\frac{m}{T}$ à chaque outer iteration (et on le laisse constant pendant les inner iterations).

Dans tous les cas, on observe toujours le fait que quand p est petit on a

plus d'outer iterations et moins d'inner iterations et inversement quand p est grand.

Pour la variante damped Newton, on observe que pour les p "extrêmes" 2 et 100

l'algorithme se termine vers 175 iterations pour le primal et 145 iterations pour le dual

alors que pour $p = 15$ et 50 on a vers 140-150 iterations pour le primal et ^{vers} 110-120

iterations pour le dual. On a donc globalement une performance légèrement meilleure de

30 iterations environ dans la résolution du dual qui peut donc nous donner $d^* = p^*$

par dualité forte un peu plus rapidement.

Pour la variante LS Newton, on observe que les performances sont similaires pour le

primal et le dual, vers 60 iterations pour $p = 2$ et vers 40 iterations pour $p = 15, 50$ et 100.

On a donc une meilleure performance de cette variante par rapport à celle de damped Newton.

Cela peut paraître cohérent du fait que dans la variante LS on cherche sur la ligne

$x + t \Delta x$ un t approximativement optimal (dans le sens de la minimisation).

Remarque: On a adopté la "backtracking line search" en s'assurant d'abord que $x + t \Delta x$

est dans le domaine de la fonction à minimiser avant de chercher un t approximativement optimal.

Ainsi que la ^{variante} damped Newton permet directement d'être dans le domaine de la fonction

à minimiser en posant $t = \frac{1}{1 + l(x)}$ dans $x + t \Delta x$ qui agit
Newton déclinant

comme une pénalisation pour éviter le risque d'être hors domaine lorsque $\lambda(x)$ est trop grand. Cependant, le choix de $t = \frac{1}{1+\lambda(x)}$ n'est pas forcément optimal dans le cas de la minimisation sur la ligne $x+t\Delta x$. La variante LS semble donc une fois adoptée pour ne pas sortir du domaine plus efficace.

Pour finir, en termes de choix de p , pour la variante LS $p=2$ nous donne un surplus d'une vingtaine d'itérations par rapport à $p=15, 50$ et 100 . On peut donc alors penser que le choix d'un p entre 10 et 100 environ n'aura pas un impact critique sur la performance (et sera même qu'un p petit comme $p=2$).

Par ailleurs, pour la variante damped Newton c'est $p=15$ et 50 qui semblent être un peu plus performants. On pourrait alors penser dans ce cas p entre ces deux valeurs (15 et 50) environ.