

Introduction aux problèmes d’optimisation différentiable

Guillaume Lécué¹

Dans cette section, on commence par présenter quelques problèmes d’optimisation (avec ou sans contrainte) issus de problèmes réels rencontrés en statistiques, traitement du signal, machine learning, etc.. Ces problèmes permettent d’inscrire le cours d’optimisation différentiable comme un outil pour les autres cours de l’ENSAE.

On présente ensuite une représentation géométrique des problèmes d’optimisation différentiable qui permettent de soulever les points clefs qui permettront de résoudre ces problèmes d’un point de vue mathématiques. En particulier, on insistera sur l’aspect “local” d’un problème d’optimisation qui demande une description locale d’une fonction et d’un ensemble (la contrainte). Ceci justifie l’adjectif “différentiable” dans le titre du cours car sous l’hypothèse de différentiation, décrire localement une fonction et l’ensemble des contraintes est plus facile.

Les deux points qui consistent à décrire localement une fonction et une contrainte dressent naturellement le plan de la première partie du cours qui concerne la résolution mathématique des problèmes d’optimisation. Dans la deuxième partie du cours, on intéressera à la résolution algorithmique de problèmes d’optimisation.

1 Quelques problèmes d’optimisation différentiable

L’optimisation est une branche des mathématiques cherchant à résoudre des problèmes de la forme

$$\min (f(x) : x \in K) \tag{1.1}$$

où $f : E \rightarrow \mathbb{R}$ et $K \subset E$. Dans ce cours, on prendra $E = \mathbb{R}^d$ ou, plus généralement, un ouvert U de \mathbb{R}^d . Les questions naturellement associées à (1.1) sont par exemple :

- un minimum existe-t’il ?
- le minimum est-il unique ?
- comment caractériser le minimum par des conditions nécessaires et suffisantes ou seulement nécessaires ?
- est-il possible de calculer ou d’approcher une solution au problème (1.1) ?

Les trois premières questions sont des questions mathématiques. La dernière est une question algorithmique qui soulève elle-même des questions d’ordre mathématiques comme sur la vitesse de convergence d’un algorithme donné.

L’objectif de ce cours est de donner des éléments de réponse à ces quatre questions dans le cas où f est différentiable et K est une contrainte définie par des fonctions différentiables. Le rôle joué par cette hypothèse de différentiation est de permettre de décrire “simplement” une fonction f et une contrainte K localement, c’est-à-dire en un point.

Remarque 1.1 Dans la formalisation du problème (1.1), il y a a priori un petit abus de notation car au moment de poser le problème, on ne sait pas a priori s’il a une solution ; on devrait alors parler du problème $\inf (f(x) : x \in K)$ plutôt que de poser le problème directement avec un min.

1. CREST, ENSAE. Bureau 3029, 5 avenue Henry Le Chatelier. 91 120 Palaiseau. Email : guillaume.lecue@ensae.fr.

Définition 1.2 La fonction f qu'on cherche à optimiser dans (1.1) est appelée **fonction objectif** et le sous-ensemble K est appelé la **contrainte**. Une **solution au problème** (1.1) est un élément $x^* \in K$ tel que $f(x^*) = \min(f(x) : x \in K)$. L'ensemble des solutions au problème (1.1) est noté $\operatorname{argmin}_{x \in K} f(x)$:

$$\operatorname{argmin}_{x \in K} f(x) = \{x^* \in K : f(x^*) = \min(f(x) : x \in K)\}.$$

(C'est un ensemble.) Dire que (1.1) **admet une solution** signifie que $\operatorname{argmin}_{x \in K} f(x)$ est non vide. De même, le **problème** (1.1) **admet une unique solution** signifie que $\operatorname{argmin}_{x \in K} f(x)$ contient un unique élément.

Idéalement, on voudrait pouvoir trouver une solution x^* au problème (1.1) et ensuite, si nécessaire, en déduire la valeur $f(x^*)$. Mais pour certains problèmes, on ne sait qu'estimer la valeur minimale de f sur K . Par ailleurs, d'un point de vue purement pratique, on ne sera pas, en général, en mesure de trouver exactement une solution x^* au problème mais on se contentera d'une approximation de x^* , par exemple, par des algorithmes itératifs.

Remarque 1.3 On ne s'intéresse qu'au problème de minimisation vu qu'en remplaçant f par $-f$ dans un problème de maximisation on retrouve un problème de minimisation :

$$\max(f(x) : x \in K) = -\min(-f(x) : x \in K)$$

et les ensembles de solutions sont les même :

$$\operatorname{argmax}_{x \in K} f(x) = \operatorname{argmin}_{x \in K} (-f(x)).$$

On présente maintenant quelques problèmes concrets qui s'écrivent sous la forme d'un problème d'optimisation.

1.1 Transport optimal de marchandises

Énoncé du problème : Des magasins (numérotés de 1 à n) ont une demande d_1, \dots, d_n pour un produit. On dispose de plusieurs entrepôts (numérotés de 1 à m) ayant chacun pour stock s_1, \dots, s_m de ce produit. Sachant que le coût unitaire de transport d'un produit de l'entrepôt j vers le magasin i est donné par c_{ij} , quelle stratégie minimise le coût de transport total pour satisfaire la demande de tous les magasins ?

Modélisation mathématique : Pour tout $1 \leq i \leq n, 1 \leq j \leq m$ on note x_{ij} le nombre de produits transportés de l'entrepôt j vers le magasin i . Étant donné une matrice de transport $x = (x_{ij}) \in \mathbb{R}^{n \times m}$, le coût total de transport est donné par $\sum_{i=1, \dots, n} \sum_{j=1, \dots, m} x_{ij} c_{ij}$. On définit ainsi la fonction objectif par

$$x \in \mathbb{R}^{n \times m} \rightarrow f(x) = \sum_{i,j} x_{ij} c_{ij}. \quad (1.2)$$

L'ensemble K de contraintes est un sous-ensemble de $\mathbb{R}^{n \times m}$ de toutes les matrices $x = (x_{ij})$ vérifiant des contraintes physiques et d'approvisionnement :

- chaque magasin $i = 1, \dots, n$ veut répondre à sa demande : $\sum_{j=1}^m x_{ij} \geq d_i$
- chaque entrepôt $j = 1, \dots, m$ ne peut pas fournir plus que son stock : $\sum_{i=1}^n x_{ij} \leq s_j$
- il ne peut y avoir qu'un nombre positif de produits déplacés : $x_{ij} \geq 0$. (On pourrait aussi préciser la contrainte avec $x_{ij} \in \mathbb{N}$, mais on verra par la suite pourquoi on ne le fait pas).

On définit alors l'ensemble de contrainte par

$$K = \left\{ x = (x_{ij}) \in \mathbb{R}^{n \times m} : \forall j = 1, \dots, m, i = 1, \dots, n, \sum_{j=1}^m x_{ij} \geq d_i; \sum_{i=1}^n x_{ij} \leq s_j; x_{ij} \geq 0 \right\}. \quad (1.3)$$

Résoudre le problème d'allocation optimal de produit à partir de m entrepôts vers n magasins connaissant le coût unitaire de transport de chaque entrepôt vers chaque magasin ainsi que la demande de chaque magasin et le stock de chaque entrepôt se résout donc en trouvant une solution au problème d'optimisation de la forme (1.1) où f est donnée en (1.2) et K est donnée en (1.3).

Par ailleurs, ici, la fonction objectif et la contrainte sont définies par des fonctions linéaires. C'est un cas particulièrement favorable pour la résolution numérique du problème. On parle alors de **linear programming**.

Définition 1.4 On dit que (1.1) est un **problème d'optimisation convexe** quand la contrainte K est un ensemble convexe et quand la fonction objectif est convexe. On dit que (1.1) est un **problème d'optimisation linéaire** quand la fonction objectif est une fonction affine et quand la contrainte K est l'intersection d'hyperplans et de demi-plans.

Remarque 1.5 On aurait pu ajouter la contrainte $x_{ij} \in \mathbb{N}$ dans K si les produits à transporter sont insécables (ce qui n'est pas le cas si on transporte des liquides comme de l'essence par exemple). Le problème d'optimisation associé est alors beaucoup plus compliqué à résoudre que celui n'imposant que la contrainte $x_{ij} \geq 0$. En effet, la convexité de la contrainte est un critère important de faisabilité du problème en pratique. On aura donc tendance à "convexifier" les contraintes non-convexe : ici on passe de \mathbb{N} à \mathbb{R}_+ par convexification de \mathbb{N} . Le problème résultant de cette convexification est un linear programming, donc beaucoup plus facile à résoudre. On donne une solution à coordonnées entière en prenant la partie entière des coordonnées d'une solution au problème de linear programming.

1.2 Construction d'un filtre anti-spams

Énoncé du problème : On dispose d'un ensemble d'emails qu'on a nous même étiquetés comme *SPAM* ou *EMAIL*. C'est notre base de données d'apprentissage, qu'on note $(X_i, Y_i)_{i=1}^n$, où X_i est le i -ième email et $Y_i \in \{SPAM, EMAIL\}$ est le *label* ou étiquette associé à X_i . On souhaite construire un filtre anti-spam à partir de cette base de données. C'est-à-dire, étant donné un nouvel email, on souhaite être capable de dire automatiquement si c'est une *SPAM* ou un *EMAIL*.

Featuring : La première étape pour résoudre un problème d'apprentissage est de transformer les données d'entrées, les X_i , en un vecteur ; c'est ce qu'on appelle le featuring. Cette transformation doit être pertinente vis-à-vis du problème posé. On peut par exemple, pour chaque email construire des features comme 1) taille de l'email 2) nombre d'image 3) nombre de fautes d'orthographe 4) présence ou non du nom/prénom du destinataire 5) heure d'envoi de l'email 6) nombre de destinataires 7) présence de certains mots etc.. Disons, qu'on a construit p features. On se retrouve alors après le featuring avec n couples $(\tilde{X}_i, Y_i)_{i=1}^n$ où \tilde{X}_i est le vecteur des features de X_i associé à son label Y_i . Aussi pour simplifier, on va prendre $Y_i = 1$ quand c'est un *SPAM* et $Y_i = -1$ quand c'est un *EMAIL*.

Algorithme (ou modèle) : A partir de la base de données "featurisées" $(\tilde{X}_i, Y_i)_{i=1}^n$, on va séparer l'espace \mathbb{R}^p en deux parties disjointes $\mathbb{R}^p = \mathcal{S} \sqcup \mathcal{E}$: une partie telle que si on reçoit un email X tel que $\tilde{X} \in \mathcal{S}$ alors on décidera que X est un *SPAM* et une deuxième partie telle

que si $\tilde{X} \in \mathcal{E}$ alors on décidera que c'est un *EMAIL*. Une première approche est de chercher un séparateur linéaire.

Dans ce cas, on va chercher un vecteur $w \in \mathbb{R}^p$ et $b \in \mathbb{R}$ tels que l'hyperplan séparateur $\{x \in \mathbb{R}^p : \langle x, w \rangle + b = 0\}$ fasse le moins d'erreur possible sur la base de données $(\tilde{X}_i, Y_i)_{i=1}^n$. Il y a une erreur de commise sur la i -ième donnée quand $Y_i \neq \text{sgn}(\langle w, \tilde{X}_i \rangle + b)$. Le nombre total d'erreurs de l'hyper-plan séparateur associé au paramètre $(w, b) \in \mathbb{R}^p \times \mathbb{R}$ sur la base de données est donc donné par $\sum_{i=1}^n I(Y_i \neq \text{sgn}(\langle w, \tilde{X}_i \rangle + b))$. Pour minimiser ce nombre d'erreur total, on va chercher à résoudre le problème d'optimisation

$$\min \left(\sum_{i=1}^n I(Y_i \neq \text{sgn}(\langle w, \tilde{X}_i \rangle + b)) : w \in \mathbb{R}^p, b \in \mathbb{R} \right). \quad (1.4)$$

On retrouve bien un problème de la forme (1.1) où la contrainte K est tout l'espace $\mathbb{R}^p \times \mathbb{R}$ et la fonction objectif est

$$f : \begin{cases} \mathbb{R}^p \times \mathbb{R} & \rightarrow \mathbb{R} \\ (w, b) & \rightarrow \sum_{i=1}^n I(Y_i \neq \text{sgn}(\langle w, \tilde{X}_i \rangle + b)). \end{cases} \quad (1.5)$$

Comme la contrainte est tout l'espace des paramètres, on dit que le problème est **sans contrainte**.

On l'a déjà évoqué un peu plus haut, un point clef de la résolution pratique d'un problème d'optimisation est la convexité (de f et K). Ici, la fonction objectif (1.5) n'est pas convexe et le problème d'optimisation associé est en général très difficile à résoudre algorithmiquement. Une stratégie classique est alors de faire une **relaxation convexe** : on cherche une fonction convexe proche de la fonction objectif. Pour ce faire, on voit que $I(Y_i \neq \text{sgn}(\langle w, \tilde{X}_i \rangle + b)) = I(Y_i \text{sgn}(\langle w, \tilde{X}_i \rangle + b) \leq 0)$. On va alors remplacer la fonction $t \in \mathbb{R} \rightarrow I(t \leq 0)$ (qui est non convexe) par une fonction convexe qui fait sens pour le problème (càd qui charge plus les nombres négatif – car se sont des erreurs – que les nombres positifs) ; par exemple, on peut prendre $\phi : t \in \mathbb{R} \rightarrow \max(1 - t, 0)$. On obtient alors une nouvelle fonction objectif donnée par

$$\tilde{f} : \begin{cases} \mathbb{R}^p \times \mathbb{R} & \rightarrow \mathbb{R} \\ (w, b) & \rightarrow \sum_{i=1}^n \phi(Y_i \text{sgn}(\langle w, \tilde{X}_i \rangle + b)). \end{cases} \quad (1.6)$$

Le problème d'optimisation associé $\min(\tilde{f}(w, b) : (w, b) \in \mathbb{R}^p \times \mathbb{R})$ est convexe et peut donc se résoudre numériquement par des algorithmes classiques développés en *optimisation convexe*. Néanmoins, il n'est pas différentiable car ϕ n'est pas dérivable en 1. On peut néanmoins, choisir une autre relaxation convexe comme $\phi(t) = \exp(-t)$ et avoir une problème d'optimisation associé qui est convexe et différentiable.

1.3 Séparateur linéaire à forte marge

Énoncé du problème : A l'image du problème précédent, on dispose d'un ensemble de points étiquetés $\{(X_i, Y_i) : i = 1, \dots, n\}$ où $X_i \in \mathbb{R}^p$ et $Y_i \in \{-1, 1\}$. Cette fois, on suppose que les données sont linéairement séparables, càd qu'il existe un hyper-plan de \mathbb{R}^p permettant de séparer les données étiquetées $+1$ des données étiquetées -1 . Un tel hyper-plan est appelé hyper-plan séparateur. En général, il en existe plusieurs et on va en choisir un tel que les deux nuages de points $\{X_i : Y_i = 1\}$ et $\{X_i : Y_i = -1\}$ soient les plus éloignés possible de l'hyper-plan séparateur comme dans la Figure 1.

Procédure : Un hyper-plan peut être paramétré par un vecteur normal $w \in \mathbb{R}^p$ et un intercept $b \in \mathbb{R}$ ainsi $\Delta(w, b) = \{x \in \mathbb{R}^p : \langle w, x \rangle + b = 0\}$. Dans un premier temps, il faut qu'on détermine

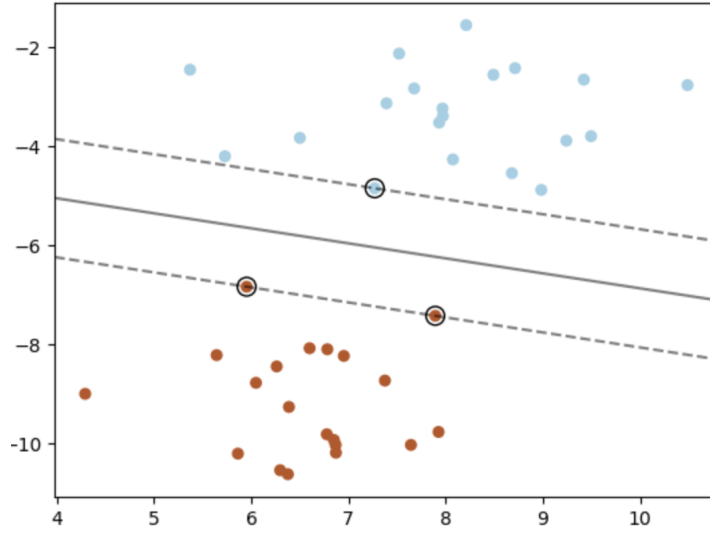


FIGURE 1 – Séparateur linéaire à forte marge.

la distance entre un point $x \in \mathbb{R}^p$ et un hyper-plan $\Delta(w, b)$. On rappelle d'abord, la définition de distance d'un point à un ensemble : si $x \in \mathbb{R}^p$ et $E \subset \mathbb{R}^p$ alors on définit la distance (euclidienne) de x à E par

$$d(x, E) = \inf_{y \in E} \|x - y\|_2.$$

Quand E est un hyper-plan de la forme $\Delta(w, b)$, on peut expliciter cette distance en fonction de x, w et b . Pour cela, on va montrer le lemme suivant.

Lemme 1.6 Soit $w, x \in \mathbb{R}^p$ et $b \in \mathbb{R}$. La distance entre x et $\Delta(w, b)$ est donnée par

$$d(x, \Delta(w, b)) = \frac{|\langle x, w \rangle + b|}{\|w\|_2}.$$

Preuve. On note par $p(x)$ la projection de x sur $\Delta(w, b)$. C'est par définition l'unique point tel que $p(x) \in \arg\min_{y \in \Delta(w, b)} \|y - x\|_2$. On a donc $p(x) \in \Delta(w, b)$ et pour tout $y \in \Delta(w, b)$, $\|p(x) - x\|_2^2 \leq \|y - x\|_2^2$. En particulier,

$$\|p(x) - x\|_2^2 \leq \|y - x\|_2^2 = \|y - p(x)\|_2^2 + 2\langle y - p(x), p(x) - x \rangle + \|p(x) - x\|_2^2$$

donc $2\langle y - p(x), x - p(x) \rangle \leq \|y - p(x)\|_2^2$. Ceci étant vrai pour tout $y \in \Delta(w, b)$ et $\{y - p(x) : y \in \Delta(w, b)\} = \text{vect}(w)^\top$, on a pour tout $z \in \text{vect}(w)^\top$, $2\langle z, x - p(x) \rangle \leq \|z\|_2^2$. Or étant donné $z \in \text{vect}(w)^\top$, on a pour tout $\lambda \in \mathbb{R}$, $\lambda z \in \Delta(w, b)$, donc $2\lambda\langle z, x - p(x) \rangle \leq |\lambda|^2 \|z\|_2^2$. Ceci étant vrai pour tout $\lambda \in \mathbb{R}$, on a $\langle z, x - p(x) \rangle = 0$ pour tout $z \in \text{vect}(w)^\top$. Autrement dit $x - p(x) \in \text{vect}(w)$.

On note $\mu \in \mathbb{R}$ tel que $x - p(x) = \mu w$. Comme $p(x) \in \Delta(w, b)$, on a $\langle p(x), w \rangle + b = 0$ et donc $\langle x + \mu w, w \rangle + b = 0$. On obtient alors, $-\mu \|w\|_2^2 = \langle x, w \rangle + b$ et donc

$$p(x) = x - \mu w = x + \frac{\langle x, w \rangle + b}{\|w\|_2^2} w.$$

Par ailleurs,

$$d(x, \Delta(w, b)) = \|x - p(x)\|_2 = |\mu| \|w\|_2 = \frac{|\langle x, w \rangle + b|}{\|w\|_2}.$$

■

Notre objectif est de trouver un hyper-plan séparateur ayant la plus forte marge. On écrit d'abord les conditions que doivent vérifier un hyper-plan séparateur : $\Delta(w, b)$ sépare les deux nuages de points $\{X_i : Y_i = 1\}$ et $\{X_i : Y_i = -1\}$ quand

$$Y_i(\langle w, X_i \rangle + b) \geq 0, \forall i = 1, \dots, n.$$

La **marge** est par définition la plus petite distance entre les nuages de points $\{X_i : Y_i = 1\}$ et $\{X_i : Y_i = -1\}$ et l'hyper-plan séparateur $\Delta(w, b)$. Pour un $w \in \mathbb{R}^p$ et $b \in \mathbb{R}$ donnés, elle vaut

$$\min_{i=1, \dots, n} \frac{|\langle X_i, w \rangle + b|}{\|w\|_2}. \quad (1.7)$$

C'est la fonction objectif de notre problème d'optimisation qu'on cherche à maximiser parmi tous les hyper-plans séparateurs. On aboutit alors au problème d'optimisation suivant :

$$\max_{w \in \mathbb{R}^p, b \in \mathbb{R}} \left(\min_{i=1, \dots, n} \frac{|\langle X_i, w \rangle + b|}{\|w\|_2} : Y_i(\langle w, X_i \rangle + b) \geq 0, \forall i = 1, \dots, n \right). \quad (1.8)$$

C'est bien un problème de la forme (1.1) mais la fonction objectif n'est pas convexe. Le problème (1.8) est donc difficile à résoudre numériquement. On va alors le réécrire de manière plus simple. En fait, on va montrer qu'il existe un problème d'optimisation convexe qui permet de résoudre (1.8). Dans ce cas, on dira que les deux problèmes d'optimisation sont équivalents lorsque les solutions de l'un permettent d'obtenir les solutions de l'autre.

On considère le problème d'optimisation

$$\min_{w \in \mathbb{R}^p, b \in \mathbb{R}} (\|w\|_2 : Y_i(\langle w, X_i \rangle + b) \geq 1, \forall i = 1, \dots, n). \quad (1.9)$$

On montre que les deux problèmes (1.8) et (1.9) sont équivalents dans le lemme suivant.

Lemme 1.7 *Les deux problèmes (1.8) et (1.9) sont équivalents :*

- i) Si (\hat{w}, \hat{b}) est solution de (1.8) alors $(\hat{w}/\alpha, \hat{b}/\alpha)$ pour $\alpha = \min_{i=1, \dots, n} |\langle X_i, \hat{w} \rangle + \hat{b}|$ est solution de (1.9).
- ii) Si (\hat{w}, \hat{b}) est solution de (1.9) alors pour tout $\alpha > 0$, $(\hat{w}/\alpha, \hat{b}/\alpha)$ est solution de (1.8).

Preuve. Preuve de i). On note $\alpha = \min_{i=1, \dots, n} |\langle X_i, \hat{w} \rangle + \hat{b}|$. On a alors $(\tilde{w}, \tilde{b}) = (\hat{w}/\alpha, \hat{b}/\alpha)$ est aussi solution de (1.8) vu que si (w, b) est séparateur alors $(w/\alpha, b/\alpha)$ l'est aussi et que la fonction objectif est homogène. Or

$$\min_{i=1, \dots, n} \frac{|\langle X_i, \tilde{w} \rangle + \tilde{b}|}{\|\tilde{w}\|_2} = \frac{1}{\|\tilde{w}\|_2}$$

et $Y_i(\langle \tilde{w}, X_i \rangle + \tilde{b}) \geq 1$ (car $Y_i(\langle \tilde{w}, X_i \rangle + \tilde{b}) > 0$ et $\min_{i=1, \dots, n} |\langle \tilde{w}, X_i \rangle + \tilde{b}| = 1$) donc $(\tilde{w}, \tilde{b}) \in \operatorname{argmin}_{w, b} (1/\|w\|_2 : Y_i(\langle w, X_i \rangle + b) \geq 1)$. Par ailleurs,

$$\operatorname{argmax}_{w \in \mathbb{R}^p, b \in \mathbb{R}} \left(\frac{1}{\|w\|_2} : Y_i(\langle w, X_i \rangle + b) \geq 1 \right) = \operatorname{argmin}_{w \in \mathbb{R}^p, b \in \mathbb{R}} (\|w\|_2 : Y_i(\langle w, X_i \rangle + b) \geq 1). \quad (1.10)$$

Donc (\tilde{w}, \tilde{b}) est bien solution de (1.9).

Preuve de ii). Si (\hat{w}, \hat{b}) est solution de (1.9) alors d'après (1.10), (\hat{w}, \hat{b}) est aussi solution de

$$\max_{w \in \mathbb{R}^p, b \in \mathbb{R}} \left(\frac{1}{\|w\|_2} : Y_i(\langle w, X_i \rangle + b) \geq 1 \right).$$

De plus pour tout $\alpha > 0$, $(\tilde{w}, \tilde{b}) = (\hat{w}/\alpha, \hat{b}/\alpha)$ est solution de

$$\max_{w \in \mathbb{R}^p, b \in \mathbb{R}} \left(\frac{\alpha}{\|w\|_2} : Y_i(\langle w, X_i \rangle + b) \geq \alpha \right).$$

En particulier, on a $\min_{i=1, \dots, n} |\langle \tilde{w}, X_i \rangle + \tilde{b}| = \alpha$ (sinon, on pourrait augmenter la fonction objectif dans le problème d'optimisation précédent). De plus, $Y_i(\langle w, X_i \rangle + b) \geq \min_{i=1, \dots, n} |\langle w, X_i \rangle + b|$ pour tout $i = 1, \dots, n$ si et seulement si $Y_i(\langle w, X_i \rangle + b) \geq 0$ pour tout $i = 1, \dots, n$ car $|Y_i| = 1$. Donc (\tilde{w}, \tilde{b}) est bien solution de (1.8). ■

D'après le Lemme 1.7, il suffit de résoudre (1.9) pour résoudre (1.8). Or (1.9) est un problème d'optimisation dont la fonction objectif est convexe et les contraintes sont affines. C'est donc un problèmes qu'on pourra résoudre assez facilement en pratique. En fait, on peut même voir que (1.9) entre dans la classe des problèmes quadratiques pour lesquels il existe des algorithmes efficaces.

Définition 1.8 *Un problème d'optimisation est dit **quadratique (quadratic programming)** quand il peut s'écrire sous la forme :*

$$\min_{x \in \mathbb{R}^p} \left(\frac{1}{2} x^\top Q x + \langle q, x \rangle : A x \leq b \right)$$

où $Q \in \mathcal{S}_d$ est une matrice symétrique de $\mathbb{R}^{d \times d}$, $q \in \mathbb{R}^d$, $A \in \mathbb{R}^{p \times d}$ et $b \in \mathbb{R}^p$. Par convention, l'inégalité " $Ax \leq b$ " signifie que pour tout $j = 1, \dots, p$, $(Ax)_j \leq b_j$.

Remarque 1.9 *Ce qu'on a fait ici sur l'espace \mathbb{R}^p s'étend facilement à n'importe quel espace de Hilbert, et en particulier, au espace de Hilbert de dimension infinie. Dans ce cas, il existe toujours un hyperplan séparant un nombre fini (ici n) de points. Etant donné des données d'entrées $X_i, i = 1, \dots, n$ labélisées à valeurs dans un espace quelconque \mathcal{X} (ça peut être par exemple l'ensemble de tous les mots sur le dictionnaire $\{A, T, C, G\}$ pour les problèmes de génétique), on peut les plonger dans un espace de Hilbert de dimension infinie $X_i \in \mathcal{X} \rightarrow \tilde{X}_i \in H$ et chercher les séparateurs linéaires à forte marge dans H pour les données $(\tilde{X}_i, Y_i)_{i=1}^n$. Il existe une multitude de plongement de ce type et ces techniques sont appelées **les méthodes à noyaux** en machine learning car le plongement $\mathcal{X} \rightarrow H$ se fait grâce à un noyau $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ par $x \in \mathcal{X} \rightarrow K(x, \cdot) \in H$.*

1.4 Estimateur des moindres carrés

Énoncé et modélisation du problème : On souhaite prédire la consommation de gaz en fonction de la température extérieure. On dispose d'une historique de consommation et de température $(T_i, C_i)_{i=1}^n$ où T_i est la température du jour i et C_i est la température du jour i . On veut pouvoir inférer la consommation C étant donnée une température T .

Procédure statistique : L'approche classique pour résoudre ce problème est de chercher une droite dans \mathbb{R}^2 qui minimise la somme des résidus de la droite au nuage de points $\{(T_i, C_i) : i = 1, \dots, n\}$. On fait donc ici l'hypothèse que la consommation de gaz peut s'expliquer de manière affine en fonction de la température (c'est une hypothèse forte). Une droite de \mathbb{R}^2 est

habituellement paramétrée par un coefficient directeur $w \in \mathbb{R}$ et une ordonnée à l'origine $b \in \mathbb{R}$. Étant donné une droite $T \rightarrow wT + b$, le résidu à la i -ième donnée est $(C_i - (wT_i + b))^2$ (on pourrait aussi prendre une erreur en $|C_i - (wT_i + b)|$; ici on prends le carré car il y a un modèle probabiliste derrière qui le justifie, mais qui n'est pas présente ici par soucis de concision). En effet, si on avait prédit la consommation du jour i par $wT_i + b$ alors on aurait commis une erreur $(C_i - (wT_i + b))^2$. La somme totale des erreurs nous donne la fonction objectif à minimiser :

$$f : \begin{cases} \mathbb{R}^p \times \mathbb{R} & \rightarrow \mathbb{R} \\ (w, b) & \rightarrow \sum_{i=1}^n (Y_i - (wT_i + b))^2. \end{cases} \quad (1.11)$$

et le problème d'optimisation associé est $\min(f(w, b) : (w, b) \in \mathbb{R}^2 \times \mathbb{R})$, qui est donc un problème d'optimisation sans contrainte, convexe et différentiable.

Critique du modèle et amélioration : Dans l'approche précédente, on fait l'hypothèse que la consommation s'explique bien par une fonction affine de la température. C'est une hypothèse forte qui n'est en fait pas vérifiée en pratique. Une manière d'affiner l'inférence est d'augmenter le nombre de features. Dans le modèle précédent, il n'y a qu'une seule feature, la température; on peut en ajouter d'autres par exemple, les températures et consommations des jours précédents, la force du vent, des features calendaires (jour de la semaine, jour férié, jour de vacance, etc.). Disons qu'on aboutit à un nombre p de features. Pour chaque jour i , on a donc un vecteur $X_i \in \mathbb{R}^p$ de features qui, on espère, devrait permettre de bien prédire la consommation C_i . Si on fait l'hypothèse que C_i peut être prédit de manière affine à partir de X_i , on va alors chercher un vecteur $w \in \mathbb{R}^p$ et un **intercept** $b \in \mathbb{R}$ tel que la somme des erreurs soit petite

$$\sum_{i=1}^n (C_i - (\langle X_i, w \rangle + b))^2.$$

On est alors amené à résoudre la problème d'optimisation suivant

$$(\hat{w}, \hat{b}) \in \underset{(w, b) \in \mathbb{R}^p \times \mathbb{R}}{\operatorname{argmin}} \sum_{i=1}^n (C_i - (\langle X_i, w \rangle + b))^2$$

qui est aussi sans contrainte, convexe et différentiable. L'estimateur associé de la consommation est la fonction $X \in \mathbb{R}^p \rightarrow \langle \hat{w}, X \rangle + \hat{b}$; on l'appelle **l'estimateur des moindres carrés**

1.5 Estimation d'un signal constant par bloc

Énoncé et modélisation du problème : On fait des observations bruitées d'un signal constant par bloc et on cherche à reconstruire ce signal. On observe $y_i = f(i/n) + g_i$ pour $i = 1, \dots, n$ où g_1, \dots, g_n sont des Gaussiennes i.i.d. centrées de variance σ^2 et $f : [0, 1] \rightarrow \mathbb{R}$ est une fonction constante par morceau. On souhaite reconstruire le vecteur $(f(i/n))_{i=1, \dots, n}$.

Procédure statistique : Pour résoudre ce genre de problème, on a souvent recours à des procédures qui cherchent à faire un meilleur compromis entre un terme d'adéquation aux données (souvent donné par la vraisemblance du modèle) et l'a priori qu'on a sur le signal à estimer; ici, le signal est constant par bloc. Le terme d'adéquation aux données est ici le critère des moindres carrés : $x \in \mathbb{R}^n \rightarrow (1/n) \sum_{i=1}^n (x_i - y_i)^2$ (comme vu précédemment). Pour induire, une structure "constant par bloc", on introduit ce qui est appelé la "norme en variation totale" : $x \in \mathbb{R}^n \rightarrow \|x\|_{TV} = \sum_{i=1}^{n-1} |x_{i+1} - x_i|$. L'idée de cette "norme" est que si $\|x\|_{TV}$ est petit alors les coordonnées successives de x seront proches et donc elle induit une structure "constant par bloc"

dans un vecteur – c’est ce qu’on cherche à faire ici. On va alors chercher un signal qui est à la fois proche des données et qui a une petite norme TV, par exemple en cherchant

$$\hat{x}_n \in \operatorname{argmin}_{x \in \mathbb{R}^n} \left(\frac{1}{n} \sum_{i=1}^n (x_i - y_i)^2 + \lambda \|x\|_{TV} \right) \quad (1.12)$$

où λ est un paramètre appelé paramètre de régularisation qui quantifie le trade-off entre le terme d’adéquation aux données et le terme induisant la structure “constant par bloc” de la norme TV.

On retrouve bien un problème de la forme (1.1) où la fonction objectif est $x \in \mathbb{R}^n \rightarrow (1/n) \sum_{i=1}^n (x_i - y_i)^2 + \lambda \|x\|_{TV}$ et la contrainte $K = \mathbb{R}^n$ – c’est donc un problème d’optimisation sans contrainte. C’est un problème d’optimisation convexe mais pas différentiable car $x \rightarrow \|x\|_{TV}$ n’est pas différentiable.

2 Représentation graphique / visualisation géométrique

But : Le but de cette section est de décrire une manière de visualiser une fonction $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ qui permet de comprendre le type d’équation que va vérifier une solution au problème (1.1). **L’idée principale est de décrire une fonction à partir de ces lignes de niveaux** à la manière d’une carte IGN, voir Figure 2 .

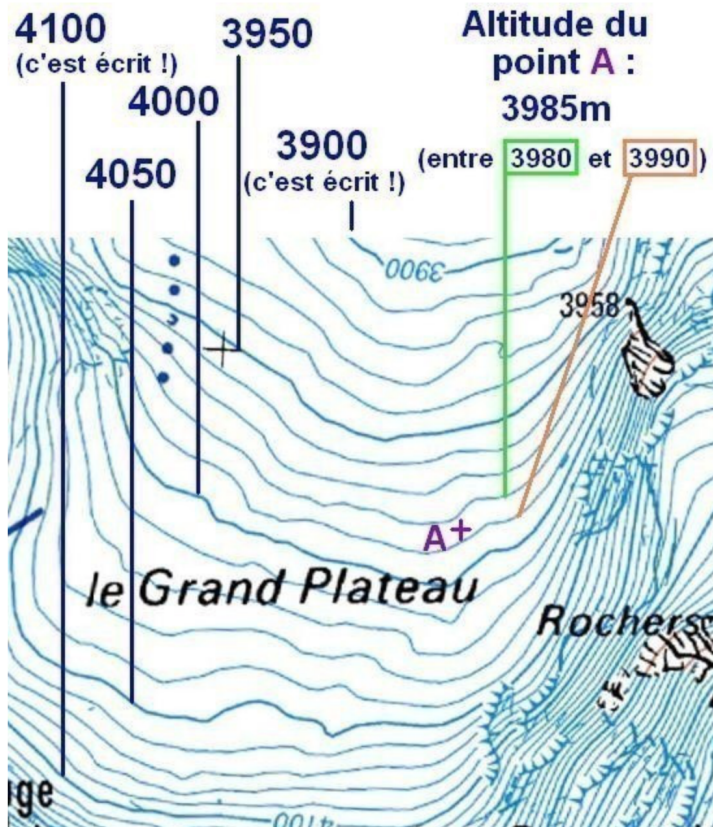


FIGURE 2 – Représentation d’un relief par les lignes de niveaux de la fonction altitude.

Définition 2.1 Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$. Soit $\alpha \in \mathbb{R}$. La ligne de niveau α de f est

$$\mathcal{L}_f(\alpha) = \{x \in \mathbb{R}^p : f(x) = \alpha\}.$$

Par exemple, si $f(x, y) = x^2 + y^2$ alors $\mathcal{L}_f(\alpha)$ est le cercle Euclidien centré en 0 et de rayon $\sqrt{\alpha}$ si $\alpha \geq 0$ et sinon, $\mathcal{L}_f(\alpha)$ est vide.

L’approche classique pour la représentation d’une fonction $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ est de tracer une surface dans \mathbb{R}^3 où \mathbb{R}^3 est décomposé en un plan “xy” pour la variable $x \in \mathbb{R}^2$ et une droite “z” où sont représentées les valeurs prises par f (voir Figure 3). Cette manière de visualiser une fonction $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ ne permet pas de visualiser facilement les idées clefs derrière la résolution des problèmes du type (1.1). On va plutôt utiliser une visualisation de $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ par lignes de niveaux (voir Figure 4). Pour cela, on donne quelques lignes sous `python` permettant de générer les deux types de visualisation (on encourage le lecteur à les reproduire et à en reproduire d’autres).

Exemple : visualisations de la différence de deux densités Gaussiennes.

On génère un réseau de points du plan ‘xy’ et les valeurs ‘z’ associées prise par la fonction “différence de deux Gaussiennes” par les lignes de code suivantes :

```
delta = 0.025
x = np.arange(-3.0, 3.0, delta)
y = np.arange(-2.0, 2.0, delta)
X, Y = np.meshgrid(x, y)
Z1 = mlab.bivariate_normal(X, Y, 1.0, 1.0, 0.0, 0.0)
Z2 = mlab.bivariate_normal(X, Y, 1.5, 0.5, 1, 1)
Z = 10.0 * (Z2 - Z1)
```

En première approche, on peut faire une visualisation 3D de cette surface comme dans la Figure 3. Ici, on représente tout l’espace \mathbb{R}^3 . On obtient ce type de représentation sous `python` avec les lignes de code suivantes :

```
fig = plt.figure()
ax = fig.add_subplot(111, projection='3d')
ax.plot_surface(X, Y, Z, color='b')
```

L’approche qu’on va utiliser est celle basée sur la visualisation dans \mathbb{R}^2 des lignes de niveau de f , comme dans la Figure 4. Ici, on travail dans \mathbb{R}^2 , l’espace de la variable $x \in \mathbb{R}^2$. On obtient ce type de représentation sous `python` avec les lignes de code suivantes :

```
plt.figure()
CS = plt.contour(X, Y, Z)
plt.clabel(CS, inline=1, fontsize=10)
```

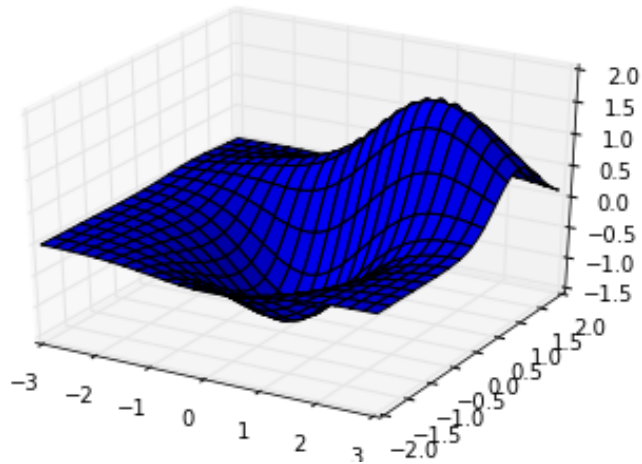


FIGURE 3 – Représentation en 3D d’une fonction $f : \mathbb{R}^2 \rightarrow \mathbb{R}$.

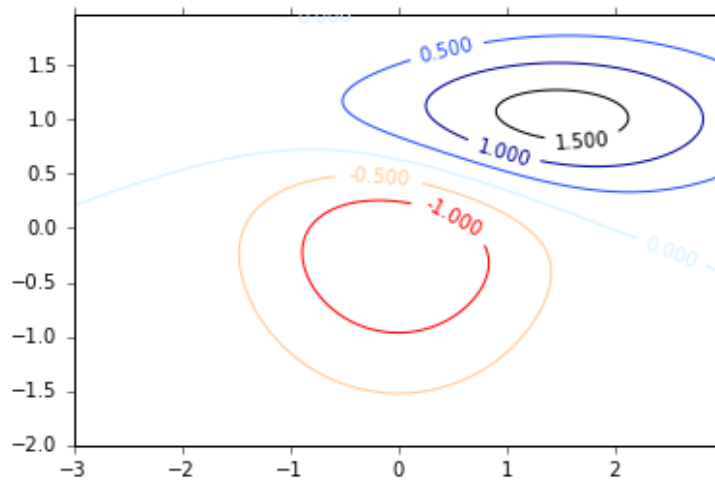


FIGURE 4 – Représentation par lignes de niveaux d’une fonction $f : \mathbb{R}^2 \rightarrow \mathbb{R}$.

3 Les problèmes d’optimisation sont des problèmes locaux

But : Dans cette section, on souhaite faire passer l’idée que les problèmes d’optimisation du type (1.1) sont des problèmes “locaux”. C’est-à-dire que pour les résoudre, on a besoin de décrire localement la fonction f et la contrainte K . On commence d’abord par dire ce qu’on entend par “décrire localement une fonction et un sous-ensemble de \mathbb{R}^d ”. Cette section est informelle et a pour but de faire passer les principales idées du cours et de son organisation qui en découle. On formalisera ces concepts dans les 5 ou 6 cours suivants.

Décrire localement une fonction signifie qu’on regarde à quoi ressemble une fonction autour d’un point donné. C’est ici que l’hypothèse de “différentiation” (inscrite dans le titre du cours) aide beaucoup car quand une fonction est différentiable en un point x^* alors la fonction f ressemble localement à une fonction affine de vecteur normal donné par le **gradient de f en x^*** , noté $\nabla f(x^*)$, (on reviendra sur cette notion fondamentale tout au long du cours) et par un **intercept**.

En équation, quand f est différentiable en x^* , f est décrite localement en x^* par la fonction affine

$$F_{x^*} : x \in \mathbb{R}^d \rightarrow \langle \nabla f(x^*), x - x^* \rangle + f(x^*). \quad (3.1)$$

Cela signifie que la fonction affine ci-dessus est une approximation d'ordre 1 de f en x^* . En, particulier, on représentera la ligne de niveau $\alpha = f(x^*)$ de cette fonction affine. On voit que

$$\mathcal{L}_{F_{x^*}}(f(x^*)) = \{x : \langle \nabla f(x^*), x - x^* \rangle = 0\} = x^* + \text{vect}(\nabla f(x^*))^\top, \quad (3.2)$$

C'est donc l'hyperplan qui passe par x^* et qui est orthogonal à $\nabla f(x^*)$. On verra par la suite que $\nabla f(x^*)$ est aussi un vecteur normal de la ligne de niveau $f(x^*)$ de f en x^* (ce qui est assez intuitif car $\mathcal{L}_{F_{x^*}}(f(x^*))$ est aussi une approximation du premier ordre de $\mathcal{L}_f(f(x^*))$ vu que F_{x^*} est une approximation au premier ordre de f en x^*).

La deuxième étape est maintenant de décrire localement une sous-ensemble de \mathbb{R}^d . Décrire localement un ensemble $K \subset \mathbb{R}^d$ est assez compliqué même sous l'hypothèse de différentiation. Cette hypothèse signifie que la contrainte K s'écrit sous la forme

$$K = \left\{ x \in U : \begin{array}{l} g_1(x) = \dots = g_r(x) = 0 \\ h_1(x) \leq 0, \dots, h_l(x) \leq 0 \end{array} \right\} \quad (3.3)$$

où $g_1, \dots, g_r : U \rightarrow \mathbb{R}$ et $h_1, \dots, h_l : U \rightarrow \mathbb{R}$ sont des classes \mathcal{C}^1 (et U est un ouvert de \mathbb{R}^d). Pour décrire localement en un point x^* un ensemble K , on aura recours à la notion de **cône tangent à K en x^*** , noté $T_K(x^*)$, qui est une manière de généraliser la notion d'hyper-plan tangent d'une surface différentiable. Il se trouve que cette notion de cône tangent est elle-même difficile à manipuler pour le problème (1.1) qu'on souhaite résoudre, on aura alors recours à une hypothèse (qui a été beaucoup étudiée en optimisation), c'est **l'hypothèse de qualification**. Cette dernière hypothèse dit littéralement que le cône tangent à K (pour un K de la forme (3.3)) en x^* se décrit facilement par

$$T_K(x^*) = \left\{ v \in \mathbb{R}^d : \begin{array}{l} \langle \nabla g_1(x^*), v \rangle = \dots = \langle \nabla g_r(x^*), v \rangle = 0 \\ \langle \nabla h_1(x^*), v \rangle \leq 0, \dots, \langle \nabla h_l(x^*), v \rangle \leq 0 \end{array} \right\} \quad (3.4)$$

Une fois qu'on dispose des outils décrivant localement la fonction f et la contrainte K , on peut se représenter visuellement des conditions nécessaires satisfaites par le ou les points x^* de l'espace \mathbb{R}^d (dans le cas $d = 2$) qui seront solution du problème (1.1) (voir Figure 3) :

- on part du minimum global de f . S'il est dans K alors ce minimum est solution de (1.1), sinon on fait croître les lignes de niveau de f , jusqu'à celle qui est la première à rencontrer la contrainte K . C'est en un tel point x^* que sera atteint le ou la solution au problème (1.1).
- En même temps qu'on fait croître les lignes de niveaux de f , on fait propager en son front, ces descriptions locales, en particulier, leur vecteur normal c'est à dire le gradient de f en chaque point dans le cas différentiable.
- Au point x^* de rencontre de la première ligne de niveau de f avec la contrainte K , on ne regarde plus que la description locale de f donnée par son gradient et celle de la contrainte donnée par $T_K(x^*)$. On voit alors que l'équation

$$-\nabla f(x^*) \in \left\{ z \in \mathbb{R}^d : \langle z, v \rangle \leq 0, \forall v \in T_K(x^*) \right\} := N_K(x^*)$$

est satisfaite. Cette équation est centrale dans l'étude des problème d'optimisation ; on l'appellera “ **condition nécessaire d'Euler/Péano/Kantorovitch** ”.

Dans la troisième étape, l'ensemble $N_K(x^*)$ est apparu. On l'appelle le **cône normal à K en x^*** . C'est une généralisation du vecteur normal à l'hyper-plan tangent à une fonction qui est le gradient $\nabla f(x^*)$ dans le cas différentiable. Sous l'hypothèse de qualification, ce cône normale s'écrit facilement :

$$N_K(x^*) = \left\{ \sum_{i=1}^r \beta_i \nabla g_i(x^*) + \sum_{j \in J(x^*)} \lambda_j \nabla h_j(x^*) : \beta_i \in \mathbb{R}, \lambda_j \geq 0 \right\} \quad (3.5)$$

et $J(x^*) = \{j \in \{1, \dots, l\} : h_j(x^*) = 0\}$. Dans ce cas, la condition $-\nabla f(x^*) \in N_K(x^*)$ est appelée la **condition de Karush-Kuhn-Tucker**, connue sous le nom de **condition KKT**.

Exemple : Trouver géométriquement la solution au problème

$$\min (x^2 + y^2 : (x - 5)^2 + (y - 5)^2 \leq 4) .$$

On représente d'abord la contrainte : c'est une boule Euclidienne de centre $(5, 5)$ et de rayon 2. Le minimum global de f sur \mathbb{R}^2 est atteint en $(0, 0)$ et la fonction objectif $f(x, y) = x^2 + y^2$ y prends la valeur 0. On fait alors croître les lignes de niveau de f à partir du niveau 0 jusqu'à rencontrer la contrainte. Le premier point de rencontre est en $x^* = (4, 4)$. Si on décrit la contrainte de la manière canonique (3.3) $K = \{(x, y) \in \mathbb{R}^2 : h_1(x, y) \leq 0\}$ où $h_1(x, y) = (x - 5)^2 + (y - 5)^2 - 4$, on a bien, en x^* ,

$$-\nabla f(x^*) \in \{\lambda \nabla h_1(x^*) : \lambda \geq 0\}$$

où $\{\lambda \nabla h_1(x^*) : \lambda \geq 0\}$ est le cône normal à K en x^* (c'est le cône normal au cône tangent à K en x^*). En effet, $\nabla f(x^*) = (2x, 2y) = (8, 8)$ et $\nabla h_1(x^*) = (2(x - 5), 2(y - 5)) = -(2, 2)$ pour $x^* = (x, y) = (4, 4)$. ■

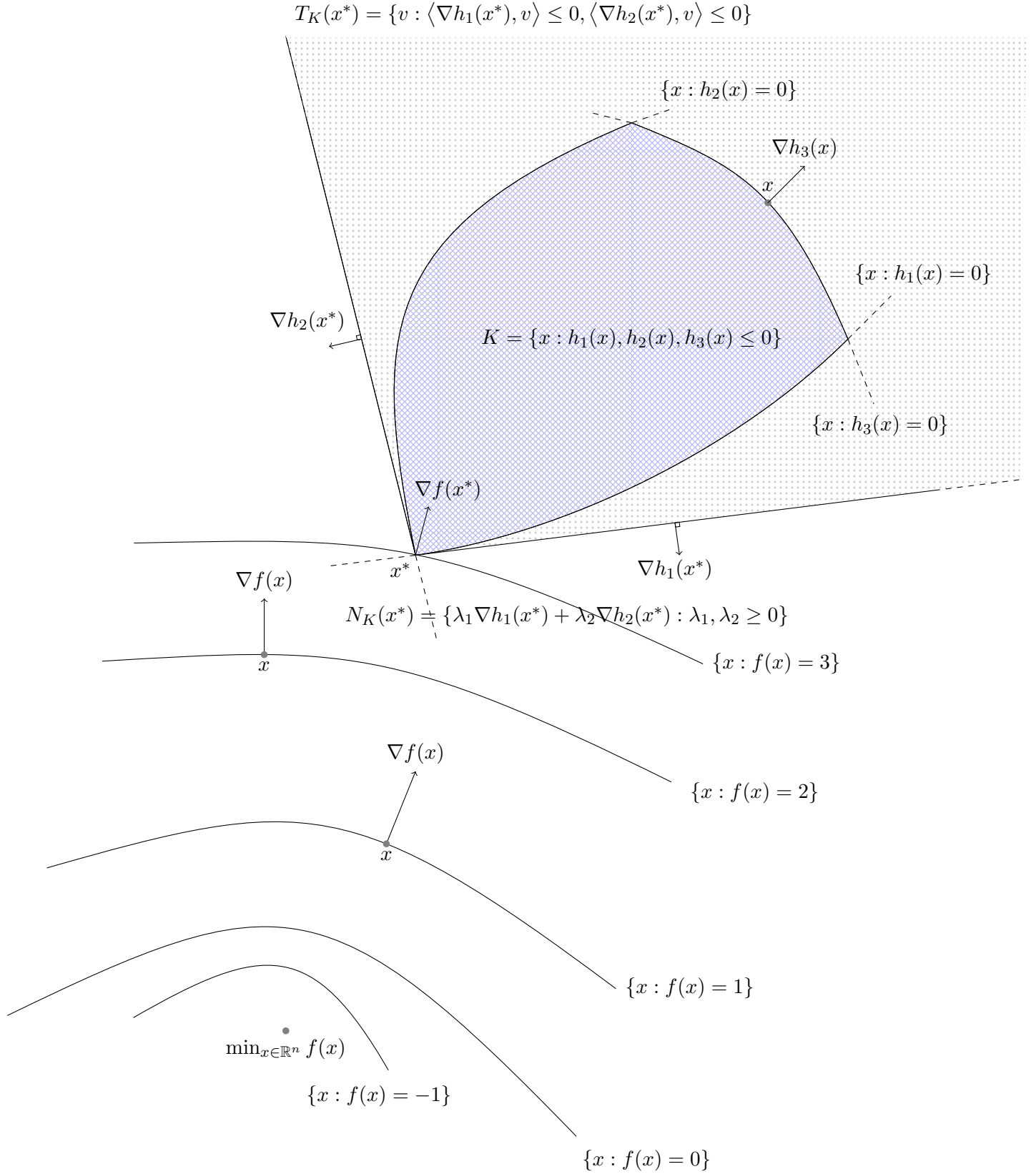


FIGURE 5 – Visualisation de la condition d'Euler/Péano/Kantorovitch et de la description locale de f par son gradient $\nabla f(x^*)$ en x^* et de K par son cône normale $N_K(x^*)$ en x^* (qui est aussi le cône normale de $T_K(x^*)$).

En résumé : Il existe une multitude de problèmes concrets qui s'énoncent sous forme d'un problème d'optimisation comme en (1.1). Pour résoudre, ce problème d'un point de vue mathématiques, on aura besoin d'outils tels que le gradient d'une fonction, le cône tangent et le cône normal à une contrainte K , la condition d'Euler/Péano/Kantorovitch et d'étudier l'hypothèse de qualification qui donne une description simple des cônes tangents et normaux à une contraintes. Dans ce cas, la condition d'Euler/Péano/Kantorovitch est appelée la condition KKT, qui est la condition la plus utilisée en optimisation. On verra aussi que la condition d'Euler/Péano/Kantorovitch qui est une condition sur les cônes normaux de la ligne de niveau $f(x^*)$ de f et de K en x^* peut aussi se retrouver par une autre dualité appelée dualité Lagrangienne. C'est le programme des 5 ou 6 prochains cours.