

Éléments d'optimisation.

Machine Learning:

Trouver des paramètres $\hat{\theta}$ qui minimisent un critère empirique $\hat{\theta} \in \operatorname{argmin}_{\theta} \hat{m}(\theta)$.

⚠ Il est parfois possible de résoudre ce problème avec $\nabla_{\theta} \hat{m}(\theta) = 0$, mais que se passe-t-il lorsque :

- (i) On rajoute des contraintes,
- (ii) L'équation a des solutions complexes ?

I. Dualité de Lagrange

Considérons le problème (appelé problème primal).

$$\begin{aligned} & \min_{\theta} f(\theta) \quad \theta \in \mathbb{H} \\ \text{s.t.} \quad & g_i(\theta) \leq 0 \quad i=1, \dots, m \\ & h_j(\theta) = 0 \quad j=1, \dots, p \end{aligned} \tag{P}$$

Hyp: $(P) < +\infty$

1) Fonction Lagrangienne, Problème Dual, Dualité Faible

②

$\forall \theta \in \mathbb{H}, (\mu_1, \dots, \mu_m) \in \mathbb{R}^m, (\lambda_1, \dots, \lambda_p) \in \mathbb{R}^p,$

on définit

$$L(\theta, \mu, \lambda) = f(\theta) + \sum_{i=1}^m \mu_i g_i(\theta) + \sum_{j=1}^p \lambda_j h_j(\theta) \quad (\text{Lagrangien}).$$

Mémoire Si θ satisfait les contraintes de (P) et $\mu \geq 0$,

$$L(\theta, \mu, \lambda) \leq f(\theta).$$

Définition : $p(\theta) = \sup_{\mu \geq 0, \lambda} L(\theta, \mu, \lambda) = \begin{cases} p_f(\theta) & \text{si } \theta \text{ satisfait } (P) \\ +\infty & \text{sinon.} \end{cases}$

Proposition : $(P) = \inf_{\theta \in \mathbb{H}} p(\theta) = \inf_{\theta \in \mathbb{H}} \sup_{\mu \geq 0, \lambda} L(\theta, \mu, \lambda)$

 Le problème dual est obtenu par inversion du inf et du sup.

Définition : ~~$d(\mu, \lambda) = \inf_{\theta \in \mathbb{H}} L(\theta, \mu, \lambda)$~~

$$(D) = \sup_{\mu \geq 0, \lambda} d(\mu, \lambda)$$

③

Théorème (Dualité faible): $(D) \leq (P)$

démonstration:

Si θ satisfait les contraints de (P) , et $p \geq 0$,

$$L(\theta, p, \lambda) \leq P(\theta)$$

donc $d(p, \lambda) \leq (P)$

donc $\sup_{p \geq 0, \lambda} d(p, \lambda) \leq (P)$

□

Exercice: Trouve le dual de $\inf_{x \in \mathcal{X}} c^T x$
 $Ax \leq b$

Solution: $L(x, p) = c^T x + p^T (Ax - b) = -b^T p + (A^T p + c)^T x$

donc $d(p) = \inf_x L(x, p) = -b^T p + \inf_x (A^T p + c)^T x$

$$= \begin{cases} -b^T p & \text{si } A^T p + c = 0 \\ -\infty & \text{sinon.} \end{cases}$$

donc $(D) = \sup_{\substack{p \geq 0 \\ A^T p + c = 0}} -b^T p$

2) Lagrangien et points de selle

Définition: On dit que $(\theta^*, \mu^*, \lambda^*)$ est un point de selle du Lagrangien si

~~satisfait les contraintes~~ ≥ 0

de $(P) \in \mathbb{H}$

$$L(\theta^*, \mu^*, \lambda^*) \leq L(\theta^*, \mu^*, \lambda^*) \leq L(\theta, \mu^*, \lambda^*) \quad \forall \lambda, \forall \mu > 0$$

et $\forall \theta$ satisfaisant
les contraintes de (P) .

Théorème: $(\theta^*, \mu^*, \lambda^*)$ est un point de selle ssi

$$\left\{ \begin{array}{l} \bullet L(\theta^*, \mu^*, \lambda^*) = \inf_{\theta \in \mathbb{H}} L(\theta, \mu^*, \lambda^*) \\ \bullet \forall i, g_i(\theta^*) \leq 0, \forall j, h_j(\theta^*) = 0 \\ \bullet \forall i, \mu^* g_i(\theta^*) = 0 \end{array} \right.$$

démonstration:

$$\Rightarrow \text{Si } g_i(\theta^*) \geq 0, \mu_i^* \rightarrow +\infty \Rightarrow \infty \leq L(\theta^*, \mu^*, \lambda^*).$$

$$\Rightarrow \forall i, g_i(\theta^*) \leq 0$$

Si $h_j(\theta^*) \neq 0$, on obtient une absurdité similaire

$$\Rightarrow \forall j, h_j(\theta^*) = 0$$

La troisième condition est immédiate.

(\Leftarrow) Immédiat.

□

Théorème: Si $(\theta^*, p^*, \lambda^*)$ est un point de selle, alors

- θ^* est une soln. de (P)
- (p^*, λ^*) est une soln. de (D)
- $(D) = (P)$.

démonstr.:

Si $(\theta^*, p^*, \lambda^*)$ est un point de selle, alors le théorème précédent assure que θ^* satisfait les contraintes de (P) . De plus, on peut écrire

$$f(\theta^*) = f(\theta^*) + \underbrace{\sum_{i=1}^m p_i^* g_i(\theta^*)}_{=0} + \sum_{j=1}^l \lambda_j^* \underbrace{h_j(\theta^*)}_{=0}$$

$$= L(\theta^*, p^*, \lambda^*)$$

$$= \inf_{\theta \in \mathbb{M}} L(\theta, p^*, \lambda^*)$$

$$= d(p^*, \lambda^*)$$

d'où le résultat par dualité-fiblé.

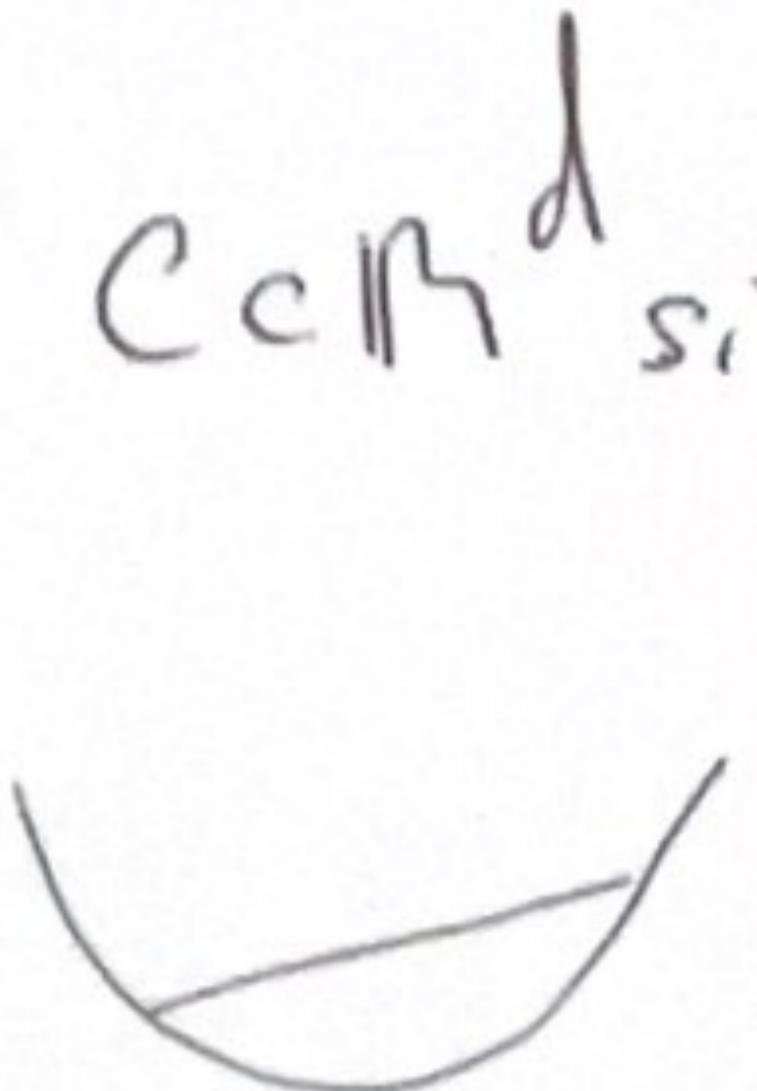
□

3) La convexe, Condition de Slater

Definition: Une fonction $f: \mathbb{R}^d \rightarrow \mathbb{R}$ est dite convexe sur $C \subset \mathbb{R}^d$ si

$$\forall x_1, x_2 \in C, \forall \lambda \in [0,1]$$

$$f(\lambda x_1 + (1-\lambda)x_2) \leq \lambda f(x_1) + (1-\lambda)f(x_2).$$



Remarque: Pour que cette définition passe sens, il faut que
 $\forall x_1, x_2 \in C, \forall \lambda \in [0,1], \lambda x_1 + (1-\lambda)x_2 \in C$.

Un ensemble est dit convexe.

Definition: Le problème (P) est dit convexe si

- f est convexe
- $(U) = \mathbb{R}^d$
- $\forall i, g_i$ est convexe
- $\forall j, h_j$ est affine

Théorème (Condition de Slater)

Si (P) est convexe et qu'il existe un point strictement faisable ($\exists \theta \in U, \forall i, g_i(\theta) < 0$ et $\forall j, h_j(\theta) = 0$) alors $(D) = (P)$.

II. Conditions de KKT

💡 L'étude des points de selle du Lagrangien permet de déduire les conditions de KKT

1) Qualification des contraintes.

• Contraints linéaires (ou affins)

• (LICQ) [Linearly independent constraints qualification]

Un point θ^* satisfait les contraintes LICQ si

$$\{\nabla h_j, V_j\} \cup \{\nabla g_i \mid \nabla g_i(\theta^*) = 0\}$$

est linéairement indépendant.

• (Slate) Un point θ^* satisfait les contraintes de Slater si

les conditions du théorème du même nom sont vérifiées (memer que que cette contrainte dépend uniquement de (P)).

2) Conditions de KKT.

Si θ^* est un optimum lagrangien de (P) et si θ^* satisfait une des conditions du paragraphe précédent, alors $\exists \mu \geq 0$

(stationnarité) $\bullet \nabla f(\theta^*) + \sum_{i=1}^m \mu_i \nabla g_i(\theta^*) + \sum_j \lambda_j \nabla h_j(\theta^*) = 0$

(F. (P)) $\bullet V_i, g_i(\theta^*) \leq 0, V_i, h_j(\theta^*) = 0$

(F. (D)) $\bullet V_i, \mu_i \geq 0$

(complémentarité) $\bullet V_i, \mu_i g_i(\theta^*) = 0$

⑧

Exercice Étudier $\min_{x_1, x_2} x_2$

$$\text{tq} \quad (x_1 - 1)^2 + x_2^2 \leq 1$$

$$(x_1 + 1)^2 + x_2^2 \leq 1$$

3) Conditions Suffisantes

⚠ Souvent, les conditions de KKT ne sont que des conditions nécessaires. Cependant, elles sont suffisantes si (P) est convexe et satisfait les conditions hypothétiques de la condition de Slater.

III Algorithmes d'optimisation

Parfois, il n'est pas possible d'obtenir une solution exacte, on utilise alors des algorithmes d'approximation.

Déscente de gradient :

Pour résoudre $\min_{\theta} P(\theta)$, on utilise l'algorithme

$$\theta_{t+1} \leftarrow \theta_t - \alpha_{t+1} \nabla P(\theta_t) \quad (\text{Déscente de gradient})$$

(3)

En effet, comme (sous des hypothèses faibles)

$$f(\theta + \Delta\theta) = f(\theta) + \langle \nabla f(\theta), \Delta\theta \rangle + o(\|\Delta\theta\|),$$

L'algorithme suit l'occident la direction de décroissance maximale.

Exercice :

Décrire l'algorithme de descente de gradient pour le problème de SUM

$$\hat{\theta} \leftarrow \arg \min_{\theta} - \sum_{i=1}^n (f(y_i; \theta^\top x_i)) + \lambda \|\theta\|_2^2$$

Algorithme du gradient stochastique

Pour accélérer l'algorithme, il est possible de remplacer le gradient par un estimateur, plus facile à calculer.

Exercice :

Que devient l'algorithme précédent si à chaque étape, le gradient est calculé sur un sous-ensemble du jeu de données unique?