

Optimisation-RO: Optimisation convexe

Clément Royer

Certificat Chef de Projet IA - Université Paris Dauphine-PSL

15 novembre 2021



- 1 Convexité
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe

- 1 Convexité
 - Ensembles et fonctions convexes
 - Convexité forte
 - Analyse convexe
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe

- 1 Convexité
 - Ensembles et fonctions convexes
 - Convexité forte
 - Analyse convexe
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe

Ensemble convexe

Un ensemble $\mathcal{C} \in \mathbb{R}^d$ est dit **convexe** si

$$\forall(\mathbf{u}, \mathbf{v}) \in \mathcal{C}^2, \forall t \in [0, 1], \quad t\mathbf{u} + (1 - t)\mathbf{v} \in \mathcal{C}.$$

Ensemble convexe

Un ensemble $\mathcal{C} \in \mathbb{R}^d$ est dit **convexe** si

$$\forall (\mathbf{u}, \mathbf{v}) \in \mathcal{C}^2, \forall t \in [0, 1], \quad t\mathbf{u} + (1 - t)\mathbf{v} \in \mathcal{C}.$$

Exemples :

- \mathbb{R}^d ;
- Droite : $\{t\mathbf{w} | t \in \mathbb{R}\}$ pour tout $\mathbf{w} \in \mathbb{R}^d$;
- Boule : $\left\{ \mathbf{w} \in \mathbb{R}^d \mid \|\mathbf{w}\|^2 = \sum_{i=1}^d [\mathbf{w}]_i^2 \leq 1 \right\}$.

Définition

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est dite **convexe** si

$$\forall(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1 - t)\mathbf{v}) \leq t f(\mathbf{u}) + (1 - t) f(\mathbf{v}).$$

Définition

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est dite **convexe** si

$$\forall (\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2, \forall t \in [0, 1], \quad f(t\mathbf{u} + (1-t)\mathbf{v}) \leq t f(\mathbf{u}) + (1-t) f(\mathbf{v}).$$

Exemples :

- Fonction linéaire : $f(\mathbf{w}) = \mathbf{a}^T \mathbf{w} + b$;
- Norme au carré : $f(\mathbf{w}) = \|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$.

Convexité et gradient

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^1 est convexe si et seulement si

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}).$$

Convexité et gradient

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^1 est convexe si et seulement si

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}).$$

L'autre inégalité clé en optimisation.

Convexité et gradient

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^1 est convexe si et seulement si

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T(\mathbf{v} - \mathbf{u}).$$

L'autre inégalité clé en optimisation.

Convexité et matrice hessienne

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^2 est dite convexe si et seulement si $\nabla^2 f(\mathbf{w}) \succeq 0$ pour tout vecteur $\mathbf{w} \in \mathbb{R}^d$, .

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}), f \text{ convex.}$$

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}), f \text{ convex.}$$

Theorem

Every local minimum of f is a global minimum.

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}), f \text{ convex.}$$

Theorem

Every local minimum of f is a global minimum.

Corollary

If f is continuously differentiable, every point \mathbf{w}^* such that $\|\nabla f(\mathbf{w}^*)\| = 0$ is a global minimum of f .

- 1 Convexité
 - Ensembles et fonctions convexes
 - Convexité forte
 - Analyse convexe
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe

Définition

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est μ -fortement convexe si pour tous $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ et $t \in [0, 1]$, on a

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|^2.$$

$\mathbf{w} \mapsto \frac{\mu}{2}\|\mathbf{w}\|^2$ est μ -fortement convexe.

Définition

Une fonction $f : \mathbb{R}^d \rightarrow \mathbb{R}$ est μ -fortement convexe si pour tous $(\mathbf{u}, \mathbf{v}) \in (\mathbb{R}^d)^2$ et $t \in [0, 1]$, on a

$$f(t\mathbf{u} + (1-t)\mathbf{v}) \leq tf(\mathbf{u}) + (1-t)f(\mathbf{v}) - \frac{\mu}{2}t(1-t)\|\mathbf{v} - \mathbf{u}\|^2.$$

$\mathbf{w} \mapsto \frac{\mu}{2}\|\mathbf{w}\|^2$ est μ -fortement convexe.

Théorème

- Une fonction fortement convexe a au plus un minimum global.
- Une fonction continue fortement convexe a un unique minimum global.

Gradient et convexité forte

Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^1 . Alors,

$$\forall \mathbf{u}, \mathbf{v} \in \mathbb{R}^d, \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \nabla f(\mathbf{u})^T (\mathbf{v} - \mathbf{u}) + \frac{\mu}{2} \|\mathbf{v} - \mathbf{u}\|^2.$$

Hessienne et convexité forte

Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ de classe \mathcal{C}^2 . Alors

$$f \text{ est } \mu\text{-fortement convexe} \iff \nabla^2 f(\mathbf{w}) \succeq \mu \mathbf{I} \quad \forall \mathbf{w} \in \mathbb{R}^d.$$

Minimisation d'une quadratique convexe

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad \mathbf{A} \succeq 0.$$

- $\nabla^2 f(\mathbf{w}) = \mathbf{A}$;
- Fortement convexe si $\mathbf{A} \succ 0$ avec $\mu = \lambda_{\min}(\mathbf{A})$.

Minimisation d'une quadratique convexe

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) := \frac{1}{2} \mathbf{w}^T \mathbf{A} \mathbf{w} + \mathbf{b}^T \mathbf{w}, \quad \mathbf{A} \succeq 0.$$

- $\nabla^2 f(\mathbf{w}) = \mathbf{A}$;
- Fortement convexe si $\mathbf{A} \succ 0$ avec $\mu = \lambda_{\min}(\mathbf{A})$.

Projection sur un ensemble fermé convexe

$$\underset{\mathbf{w} \in \mathcal{X}}{\text{minimiser}} \frac{1}{2} \|\mathbf{w} - \mathbf{a}\|^2, \quad \mathcal{X} \text{ fermé convexe.}$$

- Generalise le cas $\mathcal{X} = \mathbb{R}^d$;
- L'objectif est 1-fortement convexe \Rightarrow il existe une unique solution.

- 1 Convexité
 - Ensembles et fonctions convexes
 - Convexité forte
 - Analyse convexe
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe

Définition

Une fonction est dite **non lisse** si elle n'est pas dérivable partout.

NB: Non lisse \neq Discontinue.

Exemples de fonctions non lisses

- $w \mapsto |w|$ de \mathbb{R} dans \mathbb{R} ;
- $w \mapsto \|w\|_1 = \sum_{i=1}^d |w_i|$ de \mathbb{R}^d dans \mathbb{R} ;
- ReLU: $w \mapsto \max\{w, 0\}$ de \mathbb{R}^d dans \mathbb{R} .

- Un algorithme d'optimisation lisse est basé sur des dérivées;
- Non lisse \Leftrightarrow Pas de dérivée en certains points;
- On utilise des notions plus générales de dérivée.

Alternatives

- Si possible, reformuler en un problème lisse :

Ex) minimiser $w \in \mathbb{R}$ $|w|$ se ré-écrit

$$\text{minimiser}_{w, t^+, t^- \in \mathbb{R}} t^+ - t^- \quad \text{s. c.} \quad w = t^+ - t^-, t^+ \geq 0, t^- \geq 0.$$

- Si la *fonction* est lipschitzienne, elle possède un gradient en presque tous les points (**mais** souvent pas en les minima).

Ex) Fonction d'activation des réseaux de neurones

$$\text{ReLU}(\mathbf{w}) = [\max\{w_i, 0\}].$$

Définition

Soit $f : \mathbb{R}^d \rightarrow \mathbb{R}$ une fonction convexe. Un vecteur $\mathbf{g} \in \mathbb{R}^d$ est un **sous-gradient** de f en $\mathbf{w} \in \mathbb{R}^d$ si

$$\forall \mathbf{z} \in \mathbb{R}^d, \quad f(\mathbf{z}) \geq f(\mathbf{w}) + \mathbf{g}^T(\mathbf{z} - \mathbf{w}).$$

L'ensemble des sous-gradients de f en \mathbf{w} s'appelle le *sous-différentiel* de f en \mathbf{w} : on le note $\partial f(\mathbf{w})$.

- Si f dérivable en \mathbf{w} , $\partial f(\mathbf{w}) = \{\nabla f(\mathbf{w})\}$;
- $0 \in \partial f(\mathbf{w}) \Leftrightarrow \mathbf{w}$ minimum de f !

Exemple : Soit $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(w) = |w|$.

$$\partial f(w) = \begin{cases} -1 & \text{si } w < 0 \\ 1 & \text{si } w > 0 \\ [-1, 1] & \text{si } w = 0. \end{cases}$$

Itération pour minimiser $\mathbf{w} f(\mathbf{w})$, f convexe nonsmooth

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \mathbf{g}_k, \quad \mathbf{g}_k \in \partial f(\mathbf{w}_k).$$

- Dépend du choix du sous-gradient;
- Choix de α_k plus technique (f peut croître dans la direction d'un sous-gradient !).

À retenir

- Certaines méthodes dites “de gradient” se basent en fait sur des sous-gradients;
- Ceux-ci sont bien compris pour des problèmes simples, et dans le cas convexe.

- ❶ La sphère $\{\mathbf{w} \mid \|\mathbf{w}\|^2 = 1\}$ est-elle un ensemble convexe ?
- ❷ Quelle est la différence entre une fonction convexe et une fonction fortement convexe du point de vue des minima ?
- ❸ Quel objet mathématique remplace le gradient en optimisation non lisse ?

- 1 Convexité
- 2 Programmation convexe
 - Classes principales de problèmes
 - Dualité
 - Résolution de programmes convexes
- 3 Algorithmes pour l'optimisation convexe

- 1 Convexité
- 2 Programmation convexe
 - Classes principales de problèmes
 - Dualité
 - Résolution de programmes convexes
- 3 Algorithmes pour l'optimisation convexe

Forme de base

$$\begin{cases} \text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & f(\mathbf{w}) \\ \text{s. c.} & g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, m \\ & h_i(\mathbf{x}) = 0, \quad i = 1, \dots, \ell, \end{cases}$$

avec $f : \mathbb{R}^d \rightarrow \mathbb{R}$ convexe, $g_i : \mathbb{R}^d \rightarrow \mathbb{R}$ convexe pour tout i et $h_i : \mathbb{R}^d \rightarrow \mathbb{R}$ affine pour tout i .

- Tout programme convexe peut être mis sous cette forme;
- On considère ici que les fonctions ne prennent que des valeurs finies (mais théorie plus générale).

Exemple : Programmes linéaires (LP)

$$\begin{array}{ll}\text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & \mathbf{c}^T \mathbf{w} \\ \text{s. c.} & \mathbf{A} \mathbf{w} = \mathbf{b} \\ & \mathbf{w} \geq 0,\end{array}$$

avec $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$.

- Problème canonique en optimisation;
- Possible de résoudre des problèmes avec des millions/milliards de variables.

Exemple : Programmes quadratiques (QP)

$$\begin{array}{ll}\text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & \frac{1}{2} \mathbf{w}^T \mathbf{H} \mathbf{w} + \mathbf{c}^T \mathbf{w} \\ \text{s. c.} & \mathbf{A} \mathbf{w} = \mathbf{b} \\ & \mathbf{w} \geq 0,\end{array}$$

avec $\mathbf{H} \in \mathbb{R}^{d \times d}$, $\mathbf{c} \in \mathbb{R}^d$, $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$.

- Extension du cas sans contraintes avec contraintes linéaires;
- Peut avoir une solution même si $\mathbf{H} \succeq 0$!

Exemples : Programmes coniques d'ordre deux (SOCP)

$$\begin{array}{ll} \text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & \mathbf{c}^T \mathbf{w} \\ \text{s. c.} & \|\mathbf{A}_i \mathbf{w} + \mathbf{b}_i\| \leq \mathbf{c}_i^T \mathbf{w} + d_i, \quad i = 1, \dots, n. \end{array}$$

avec pour tout $i = 1, \dots, m$, $\mathbf{A}_i \in \mathbb{R}^{n_i \times d}$, $\mathbf{b}_i \in \mathbb{R}^{n_i}$, $\mathbf{c}_i \in \mathbb{R}^d$ et $d_i \in \mathbb{R}$.

- Généralisent LP et QP;
- Peuvent être résolus efficacement.

$$\begin{array}{ll}\text{minimiser}_{\mathbf{W} \in \mathbb{R}^{d \times d}} & \text{trace}(\mathbf{C}^T \mathbf{W}) \\ \text{s. c.} & \text{trace}(\mathbf{A}_i^T \mathbf{X}) = b_i, \quad i = 1, \dots, n \\ & \mathbf{W} = \mathbf{W}^T \succcurlyeq 0,\end{array}$$

avec $\mathbf{C}, \mathbf{A}_1, \dots, \mathbf{A}_n \in \mathbb{R}^{d \times d}$ et $\mathbf{b} = [b_i] \in \mathbb{R}^n$.

- Généralisent LP et QP;
- Possible de les résoudre en temps polynomial, mais les calculs algébriques peuvent être coûteux;
- Problèmes souvent de grandes tailles (relaxations de problèmes continus ou combinatoires).

- 1 Convexité
- 2 Programmation convexe
 - Classes principales de problèmes
 - Dualité
 - Résolution de programmes convexes
- 3 Algorithmes pour l'optimisation convexe

Principe

- Définir une fonction qui combine l'objectif et les contraintes;
- Maximiser cette fonction conduit à un nouveau problème d'optimisation dit dual.

Définition

Soit le problème d'optimisation

$$\begin{cases} \text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & f(\mathbf{w}) \\ \text{s. c.} & g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, m \\ & h_i(\mathbf{w}) = 0, \quad i = 1, \dots, \ell, \end{cases}$$

Le **lagrangien** du problème est défini par

$$\begin{aligned} \forall (\mathbf{w}, \mathbf{v}, \mathbf{s}) \in \mathbb{R}^d \times \mathbb{R}^\ell \times (\mathbb{R}_+)^m, \\ \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s}) &:= f(\mathbf{w}) + \mathbf{v}^T \mathbf{h}(\mathbf{w}) + \mathbf{s}^T \mathbf{g}(\mathbf{w}) \\ &= f(\mathbf{w}) + \sum_{i=1}^{\ell} v_i h_i(\mathbf{w}) + \sum_{i=1}^m s_i g_i(\mathbf{w}). \end{aligned}$$

Définition

Soit le problème d'optimisation

$$\begin{cases} \text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} & f(\mathbf{w}) \\ \text{s. c.} & g_i(\mathbf{w}) \leq 0, \quad i = 1, \dots, m \\ & h_i(\mathbf{w}) = 0, \quad i = 1, \dots, \ell, \end{cases}$$

Le **lagrangien** du problème est défini par

$$\begin{aligned} \forall (\mathbf{w}, \mathbf{v}, \mathbf{s}) \in \mathbb{R}^d \times \mathbb{R}^\ell \times (\mathbb{R}_+)^m, \\ \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s}) &:= f(\mathbf{w}) + \mathbf{v}^T \mathbf{h}(\mathbf{w}) + \mathbf{s}^T \mathbf{g}(\mathbf{w}) \\ &= f(\mathbf{w}) + \sum_{i=1}^{\ell} v_i h_i(\mathbf{w}) + \sum_{i=1}^m s_i g_i(\mathbf{w}). \end{aligned}$$

- Agrégation de l'objectif et des contraintes;
- Cette fonction peut prendre des valeurs infinies !

Définition: maximiser $d(\mathbf{v}, \mathbf{s})$, $d(\mathbf{v}, \mathbf{s}) := \inf_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s})$
 $\mathbf{v} \in \mathbb{R}^\ell, \mathbf{s} \in (\mathbb{R}_+)^m$

- Problème **concave** (opposé convexe) !
- Le problème originel est dit primal et s'écrit

$$\text{minimiser}_{\mathbf{w} \in \mathbb{R}^d} \sup_{\mathbf{v} \in \mathbb{R}^\ell, \mathbf{s} \in (\mathbb{R}_+)^m} \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s})$$

Définition: maximiser $d(\mathbf{v}, \mathbf{s})$, $d(\mathbf{v}, \mathbf{s}) := \inf_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s})$
 $\mathbf{v} \in \mathbb{R}^\ell, \mathbf{s} \in (\mathbb{R}_+)^m$

- Problème **concave** (opposé convexe) !
- Le problème originel est dit primal et s'écrit

$$\text{minimiser} \quad \sup_{\mathbf{v} \in \mathbb{R}^\ell, \mathbf{s} \in (\mathbb{R}_+)^m} \mathcal{L}(\mathbf{w}, \mathbf{v}, \mathbf{s})$$

À quoi sert le dual?

- Si on a dualité forte, les deux problèmes ont même valeur optimale;
- Un triplet $(\mathbf{w}^*, \mathbf{v}^*, \mathbf{s}^*)$ solution primale-duale vérifie alors :

$$\mathcal{L}(\mathbf{w}^*, \mathbf{v}, \mathbf{s}) \leq \mathcal{L}(\mathbf{w}^*, \mathbf{v}^*, \mathbf{s}^*) \leq \mathcal{L}(\mathbf{w}, \mathbf{v}^*, \mathbf{s}^*).$$

- Pour résoudre le problème, on cherche alors un point selle du lagrangien qui vérifie les contraintes.

Conditions de Karush-Kuhn-Tucker (KKT)

Problème: minimiser $_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ s. c. $\mathbf{g}(\mathbf{w}) \leq 0_{\mathbb{R}^m}, \mathbf{h}(\mathbf{w}) = 0_{\mathbb{R}^\ell}$.

Énoncé

On suppose qu'il y a dualité forte et que $f, \mathbf{g}, \mathbf{h}$ sont de classe \mathcal{C}^1 . Alors, si $\mathbf{w}^* \in \mathbb{R}^d$ est une solution du problème, il existe $\mathbf{v}^* \in \mathbb{R}^\ell$ et $\mathbf{s}^* \in \mathbb{R}^m$ tels que

$$\begin{aligned}\nabla f(\mathbf{w}^*) + \mathbf{J}_h(\mathbf{w}^*)\mathbf{v}^* + \mathbf{J}_g(\mathbf{w}^*)\mathbf{s}^* &= 0 \\ \mathbf{h}(\mathbf{w}^*) &= 0_{\mathbb{R}^\ell} \\ \mathbf{g}(\mathbf{w}^*) &\leq 0_{\mathbb{R}^m} \\ g_i(\mathbf{w}^*)s_i^* &= 0 \quad \forall i = 1, \dots, m.\end{aligned}$$

Conditions de Karush-Kuhn-Tucker (KKT)

Problème: minimiser $\mathbf{w} \in \mathbb{R}^d$ $f(\mathbf{w})$ s. c. $\mathbf{g}(\mathbf{w}) \leq 0_{\mathbb{R}^m}$, $\mathbf{h}(\mathbf{w}) = 0_{\mathbb{R}^\ell}$.

Énoncé

On suppose qu'il y a dualité forte et que $f, \mathbf{g}, \mathbf{h}$ sont de classe \mathcal{C}^1 . Alors, si $\mathbf{w}^* \in \mathbb{R}^d$ est une solution du problème, il existe $\mathbf{v}^* \in \mathbb{R}^\ell$ et $\mathbf{s}^* \in \mathbb{R}^m$ tels que

$$\begin{aligned}\nabla f(\mathbf{w}^*) + \mathbf{J}_h(\mathbf{w}^*)\mathbf{v}^* + \mathbf{J}_g(\mathbf{w}^*)\mathbf{s}^* &= 0 \\ \mathbf{h}(\mathbf{w}^*) &= 0_{\mathbb{R}^\ell} \\ \mathbf{g}(\mathbf{w}^*) &\leq 0_{\mathbb{R}^m} \\ g_i(\mathbf{w}^*)s_i^* &= 0 \quad \forall i = 1, \dots, m.\end{aligned}$$

- Les trois premières conditions correspondent à annuler le gradient du lagrangien par rapport à \mathbf{w}^* , \mathbf{v}^* et \mathbf{s}^* .
- \mathbf{v}^* et \mathbf{s}^* sont des variables duales aussi appelées multiplicateurs de Lagrange.

- 1 Convexité
- 2 Programmation convexe
 - Classes principales de problèmes
 - Dualité
 - Résolution de programmes convexes
- 3 Algorithmes pour l'optimisation convexe

Pour les classes majeures

- LP, QP, SDP, SOCP;
- D'autres classes non traitées ici.

Dans la suite

- On illustre les notions dans le cadre des programmes linéaires;
- On donne l'idée derrière les méthodes dites de points intérieurs.

Problème de base (dit primal)

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \mathbf{c}^T \mathbf{w} \quad \text{s. c.} \quad \mathbf{A} \mathbf{w} = \mathbf{b}, \mathbf{w} \geq 0,$$

$$\mathbf{A} \in \mathbb{R}^{m \times d}, \mathbf{b} \in \mathbb{R}^m.$$

Problème dual (aussi LP !)

$$\underset{\mathbf{v} \in \mathbb{R}^m, \mathbf{s} \in \mathbb{R}^d}{\text{maximiser}} \mathbf{b}^T \mathbf{v} \quad \text{s. c.} \quad \mathbf{A}^T \mathbf{v} + \mathbf{s} = \mathbf{c}, \mathbf{s} \geq 0.$$

Conditions de KKT

Si $\mathbf{w}^* \in \mathbb{R}^d$ est une solution, il existe $\mathbf{v}^* \in \mathbb{R}^m$ et $\mathbf{s}^* \in \mathbb{R}^d$ tels que

$$\begin{aligned} \mathbf{A} \mathbf{w}^* &= \mathbf{b}, & \mathbf{w}^* &\geq 0, \\ \mathbf{A}^T \mathbf{v}^* + \mathbf{s}^* &= \mathbf{c}, & \mathbf{s}^* &\geq 0, \\ w_i s_i &= 0 \quad \forall i = 1, \dots, d. \end{aligned}$$

Théorème

Pour tout programme linéaire,

- i) Soit ce problème et son dual ont chacun une solution, et l'ensemble vérifie les conditions de KKT;
- ii) Soit un des deux problèmes n'a pas de solution, et l'ensemble réalisable de l'autre est vide (problème irréalisable);
- iii) Soit les deux ensembles réalisables sont vides.

Théorème

Pour tout programme linéaire,

- i) Soit ce problème et son dual ont chacun une solution, et l'ensemble vérifie les conditions de KKT;
- ii) Soit un des deux problèmes n'a pas de solution, et l'ensemble réalisable de l'autre est vide (problème irréalisable);
- iii) Soit les deux ensembles réalisables sont vides.

Conséquences

- Trouver une solution revient à trouver un triplet vérifiant les contraintes de chaque problème !
- C'est ce que font les approches de points intérieurs, de manière primale-duale ou duale.

Principe

- Partir d'un triplet $(\mathbf{w}, \mathbf{v}, \mathbf{s})$ réalisable pour le problème et son dual avec $\mathbf{w} > 0$ et $\mathbf{s} > 0$.
- Appliquer la méthode de Newton au système de KKT

$$\begin{cases} \mathbf{A}\mathbf{w}^* &= \mathbf{b} \\ \mathbf{A}^T \mathbf{v}^* + \mathbf{s}^* &= \mathbf{c} \\ w_i s_i &= 0 \quad \forall i = 1, \dots, d. \end{cases}$$

- Obtenir un nouveau point "intérieur" avec une valeur $\mathbf{w}^T \mathbf{s}$ réduite.

Principe

- Partir d'un triplet $(\mathbf{w}, \mathbf{v}, \mathbf{s})$ réalisable pour le problème et son dual avec $\mathbf{w} > 0$ et $\mathbf{s} > 0$.
- Appliquer la méthode de Newton au système de KKT

$$\begin{cases} \mathbf{A}\mathbf{w}^* &= \mathbf{b} \\ \mathbf{A}^T \mathbf{v}^* + \mathbf{s}^* &= \mathbf{c} \\ w_i s_i &= 0 \quad \forall i = 1, \dots, d. \end{cases}$$

- Obtenir un nouveau point "intérieur" avec une valeur $\mathbf{w}^T \mathbf{s}$ réduite.

Ça marche !

- Vitesses de convergence en théorie;
- Très bonnes performances en pratique;
- Implémentations très sophistiquées.

- 1 Citer deux classes de problèmes d'optimisation convexe.
- 2 Quel objet mathématique concatène l'objectif et les contraintes d'un problème d'optimisation ? Que représente une solution du problème par rapport à cet objet ?
- 3 Quelle technique s'applique aussi bien à la résolution de programmes linéaires qu'à la résolution de programmes quadratiques ?

- 1 Convexité
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe
 - Descente de gradient et optimisation convexe
 - Accélération
 - Régularisation

- 1 Convexité
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe
 - Descente de gradient et optimisation convexe
 - Accélération
 - Régularisation

$$\underset{\mathbf{x} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{x}), \quad f \in \mathcal{C}_L^{1,1}.$$

Descente de gradient

- Itération: $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$, terminer si $\nabla f(\mathbf{w}_k) = 0$.
- Choix typique en théorie : $\alpha_k = \frac{1}{L}$.

Avec la convexité

Hypothèse : $f^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$ est atteint.

- Garanties relativement à un minimum **global**;
- On peut montrer que $f(\mathbf{w}_k) \rightarrow f^*$;
- On peut aussi montrer une convergence vers l'argmin

Cas non convexe

- Critère : $\|\nabla f(\mathbf{w}_k)\|$;
- Idée : être proche d'un point stationnaire.

Cas convexe/fortement convexe

- : $f(\mathbf{w}_k) - f^* \leq \epsilon$, avec $f^* = \min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w})$;
- Idée : être proche de la valeur à l'optimum.
- Valeur liée à $\|\mathbf{w}_k - \mathbf{w}^*\|$ dans le cas fortement convexe.

Théorème

Si $f \in \mathcal{C}_L^{1,1}$ est convexe et $\alpha_k = \frac{1}{L}$, la descente de gradient calcule \mathbf{w}_k tel que $f(\mathbf{w}_k) - f^* \leq \epsilon$ en au plus

- $\mathcal{O}(\epsilon^{-1})$ itérations;
- $\mathcal{O}\left(\frac{L}{\mu} \ln(\epsilon^{-1})\right)$ itérations si f est μ -fortement convexe.

- Cas non convexe : $\mathcal{O}(\epsilon^{-2})$ pour garantir $\|\nabla f(\mathbf{w}_k)\| \leq \epsilon$.
- On dit que la descente de gradient possède une meilleure complexité dans le cas convexe/fortement convexe.

Théorème

Si $f \in \mathcal{C}_L^{1,1}$ est convexe et $\alpha_k = \frac{1}{L}$, pour tout $K \in \mathbb{N}$, on a :

$$f(\mathbf{w}_K) - f^* \leq \frac{L \max_{\mathbf{w} \in \arg\min_{\mathbf{v}} f(\mathbf{v})} \|\mathbf{w}_0 - \mathbf{w}\|}{2} \frac{1}{K}$$

pour f convexe, et

$$f(\mathbf{w}_K) - f^* \leq \left(1 - \frac{\mu}{L}\right)^K (f(\mathbf{w}_0) - f^*).$$

pour f μ -fortement convexe.

- Cas non convexe : $\mathcal{O}(1/\sqrt{K})$ pour $\min_{0 \leq k \leq K-1} \|\nabla f(\mathbf{w}_k)\|$.
- On dit que la descente de gradient converge plus rapidement dans le cas fortement convexe que dans le cas convexe.

- 1 Convexité
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe
 - Descente de gradient et optimisation convexe
 - Accélération
 - Régularisation

Motivation

- En optimisation non convexe, $\mathcal{O}(1/\sqrt{K})$ est la meilleure vitesse de convergence pour une méthode type gradient;
- Dans le cas convexe, c'est $\mathcal{O}(1/K^2)$, mieux que la descente de gradient en $\mathcal{O}(1/K)$.

Comment obtenir cette meilleure vitesse ?

- Stratégies de **gradient accéléré**, basées sur l'idée de **momentum**;
- **Principe** : Réutiliser l'information de l'itération précédente.

Algorithme($\mathbf{w}_0, \mathbf{w}_{-1} = \mathbf{w}_0$)

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1})) + \beta_k(\mathbf{w}_k - \mathbf{w}_{k-1}).$$

- Un appel de gradient par itération;
- Terme de momentum : $\mathbf{w}_k - \mathbf{w}_{k-1}$ (pas précédent).

Version à deux suites($\mathbf{w}_0, \mathbf{z}_0 = \mathbf{w}_0$)

$$\begin{cases} \mathbf{w}_{k+1} &= \mathbf{z}_k - \alpha_k \nabla f(\mathbf{z}_k) \\ \mathbf{z}_{k+1} &= \mathbf{w}_{k+1} + \beta_{k+1}(\mathbf{w}_{k+1} - \mathbf{w}_k). \end{cases}$$

Longueur de pas α_k

- $\alpha_k = \frac{1}{L}$;
- Autres : décroissante, recherche linéaire, etc.

Momentum β_k

- f μ -fortement convexe : $\beta_k = \frac{\sqrt{L}-\sqrt{\mu}}{\sqrt{L}+\sqrt{\mu}}$;
- f convexe : Utiliser deux suites

$$t_{k+1} = \frac{1}{2}(1 + \sqrt{1 + 4t_k^2}), t_0 = 0, \quad \beta_k = \frac{t_k - 1}{t_{k+1}}.$$

Fonctions convexes

- Descente de gradient : $f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K}\right)$;
- Gradient accéléré : $f(\mathbf{w}_K) - f^* \leq \mathcal{O}\left(\frac{1}{K^2}\right)$.

Fonctions μ -fortement convexes

- Descente de gradient :

$$f(\mathbf{w}_K) - f^* \leq \left(1 - \frac{\mu}{L}\right)^K (f(\mathbf{w}_0) - f^*).$$

- Gradient accéléré

$$f(\mathbf{w}_K) - f^* \leq C \left(1 - \sqrt{\frac{\mu}{L}}\right)^K (f(\mathbf{w}_0) - f^*).$$

Méthode de la boule lestée (*Heavy ball*, Polyak, 1964)

- $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k) + \beta(\mathbf{w}_k - \mathbf{w}_{k+1});$
- Optimale sur des quadratiques fortement convexes, mais ne converge pas toujours pour f fortement convexe!
- “Précurseur” du gradient accéléré.

Méthode du gradient conjugué (Hestenes and Stiefel, 1952)

- $\mathbf{w}_{k+1} = \mathbf{w}_k + \alpha_k p_k, p_k = -\nabla f(x_k) + \beta_k p_{k-1};$
- Développée pour les quadratiques fortement convexes, optimale **sans connaître L ou μ !**
- D'autres versions pour les fonctions fortement convexes, convexes et non convexes efficaces en pratique.

- 1 Convexité
- 2 Programmation convexe
- 3 Algorithmes pour l'optimisation convexe
 - Descente de gradient et optimisation convexe
 - Accélération
 - Régularisation

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} \underbrace{f(\mathbf{w})}_{\text{attache données}} + \underbrace{\lambda \Omega(\mathbf{w})}_{\text{régularisation}} .$$

où $\lambda > 0$ est un paramètre de régularisation.

Exemple : Régularisation *ridge* ou *écrêtée*

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \frac{\lambda}{2} \|\mathbf{w}\|^2 .$$

Interprétations :

- Revient à ajouter une contrainte sur $\|\mathbf{w}\|^2 = \sum_{i=1}^d w_i^2$;
- Favorise les vecteurs \mathbf{w} avec des composantes uniformément faibles en amplitude;
- Réduit la variance des solutions par rapport aux données;
- Problème fortement convexe pour λ suffisamment grand.

Problème régularisé et exemples (2)

Régularisation ℓ_0

- Objectif : Trouver $\mathbf{w} \in \mathbb{R}^d$ qui colle aux données avec le plus de coefficients nuls possible;
- Problème idéal : $\min_{\mathbf{w}} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_0$, avec $\|\mathbf{v}\|_0 := |\{i | [\mathbf{v}]_i \neq 0\}|$.
Mais la fonction $\|\cdots\|_0$ est non lisse, discontinue et introduit de la combinatoire.

A better approach: LASSO regularization

LASSO=Least Absolute Shrinkage and Selection Operator

$$\min_{\mathbf{w} \in \mathbb{R}^d} f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1, \quad \|\mathbf{w}\|_1 = \sum_{i=1}^d |w_i|.$$

- $\|\cdot\|_1$ convexe, continue, norme;
- Non lisse mais possède des sous-gradients.

Cadre : Optimisation composite

$$\underset{\mathbf{w} \in \mathbb{R}^d}{\text{minimiser}} f(\mathbf{w}) + \lambda \Omega(\mathbf{w}).$$

- $f \in \mathcal{C}^{1,1}$;
- Ω convexe.

Approche proximale

- Classique en optimisation : remplacer un problème par une suite de sous-problèmes plus simples;
- Ici on exploite la “douceur” de f et la structure de Ω pour construire des problèmes que l’on peut résoudre **efficacement**.

Itération

$$\mathbf{w}_{k+1} = \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^d} \left\{ f(\mathbf{w}_k) + \nabla f(\mathbf{w}_k)^T (\mathbf{w} - \mathbf{w}_k) + \frac{1}{2\alpha_k} \|\mathbf{w} - \mathbf{w}_k\|_2^2 + \lambda \Omega(\mathbf{w}) \right\}.$$

- Si $\Omega \equiv 0$, la solution est $\mathbf{w}_{k+1} = \mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$: **c'est l'itération de la descente de gradient !**
- En général, une itération coûte un calcul de gradient + **1 résolution de sous-problème.**

Propriétés

- Complexité/Vitesses de convergence;
- Règles de choix de longueur de pas;
- Variantes accélérées/avec sous-gradients.

Contexte

- Résoudre minimiser $\mathbf{w} \in \mathbb{R}^d$ $f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$;
- Problème classique en traitement du signal/de l'image;
- La méthode du gradient proximal a une **forme explicite**.

Contexte

- Résoudre minimiser $\mathbf{w} \in \mathbb{R}^d$ $f(\mathbf{w}) + \lambda \|\mathbf{w}\|_1$;
- Problème classique en traitement du signal/de l'image;
- La méthode du gradient proximal a une **forme explicite**.

Itération ISTA : Iterative Soft-Thresholding Algorithm

Définit \mathbf{w}_{k+1} par coordonnées: pour tout $i \in \{1, \dots, d\}$,

$$[\mathbf{w}_{k+1}]_i = \begin{cases} [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i + \alpha_k \lambda & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i < -\alpha_k \lambda \\ [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i - \alpha_k \lambda & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i > \alpha_k \lambda \\ 0 & \text{si } [\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)]_i \in [-\alpha_k \lambda, \alpha_k \lambda]. \end{cases}$$

Exemple de méthode proximale : ISTA (2)

Mise à jour dans ISTA

- Part du pas de gradient $\mathbf{w}_k - \alpha_k \nabla f(\mathbf{w}_k)$;
- Applique l'opérateur de *soft-thresholding* $s_{\alpha_k \lambda}(\bullet)$ à chaque coordonnée, avec

$$s_{\mu}(t) = \begin{cases} t + \mu & \text{si } t < -\mu \\ t - \mu & \text{si } t > \mu \\ 0 & \text{sinon.} \end{cases}$$

- Favorise les composantes nulles.

Variantes de ISTA

- Changement de longueur de pas;
- Ajout de momentum : FISTA (la plus utilisée en pratique).

- ➊ Peut-on appliquer l'algorithme de descente de gradient à un problème convexe ?
- ➋ Sur quoi repose l'idée de momentum ?
- ➌ Quel est l'intérêt d'une régularisation ℓ_1 ?

Convexité et optimisation

-
- Solveurs efficaces.

Algorithmes d'optimisation convexe

- Un même algorithme possède de meilleures propriétés sur un problème convexe !
- Méthodes accélérées optimales pour les problèmes convexes;
- Stratégies de régularisation : apportent de la structure, “convexifient” souvent le problème.

Ouvrages :

- S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- A. Beck, *First-order methods in optimization*, MOS-SIAM Series on Optimization, 2017.

Codes :

- CVX, CVXPY : logiciels académiques pour l'optimisation convexe (+IPOPT pour le cas non convexe).
- Gurobi, CPLEX, MOSEK : logiciels commerciaux pour les programmes linéaires (en variables continues et discrètes) et au-delà (programmes quadratiques, etc).

Notebook

- Les effets de l'accélération...
- ...et quelques solveurs convexes.

Demain : Méthodes stochastiques avec Pierre Ablin.

Notebook

- Les effets de l'accélération...
- ...et quelques solveurs convexes.

Demain : Méthodes stochastiques avec Pierre Ablin.

Merci beaucoup !