

# OUTILS D'OPTIMISATION POUR LES SCIENCES DES DONNÉES ET DE LA DÉCISION

20 octobre 2023

Aujourd'hui : Accélération (cours + TD avec notebook  
d'illustration)

NB: Polycopié mis à jour avec le contenu des séances 4-6

# OPTIMISATION ET ACCÉLÉRATION

# CONVEXE

## Cadre

minimiser  $f(w)$   
 $w \in \mathbb{R}^d$

$f: \mathbb{R}^d \rightarrow \mathbb{R}$

et **convexe**

$\nabla f$  existe en  
tout  $w$

$\nabla f$  est  
lipschitzien

$\nabla f$  est  
 $L$ -lipschitz

$$\|\nabla f(v) - \nabla f(w)\| \leq L\|v - w\|$$

$f$  convexe :

- Tout minimum local est global
- $\|\nabla f(\bar{w})\| = 0 \Leftrightarrow \bar{w}$  est

un minimum global

$f$  convexe



$f$  non  
convexe

Objectif: Montrer de  
**convergence**  
du

que l'algorithme  
descente de gradient  
vers une solution  
problème

Rappel: Itération

de la descente de  
gradient

$$w_{k+1} = w_k$$

$$- \alpha_k \nabla f(w_k)$$

$$\alpha_k > 0$$

L> Analyse classique  
Garanties

en optimisation ..  
asymptotiques

Ex:

$$\| \nabla f(w_k) \| \rightarrow 0 \quad k \rightarrow \infty$$

⇒ Typique de l'optimisation  
non convexe

Optimisation  
convexe

$$f(w_k)$$

$$\rightarrow \min_{w \in \mathbb{R}^d} f(w)$$

↑

Optimisation  
fortement  
convexe

$$w_k$$

$$\rightarrow w^* \in \arg\min_w f(w)$$

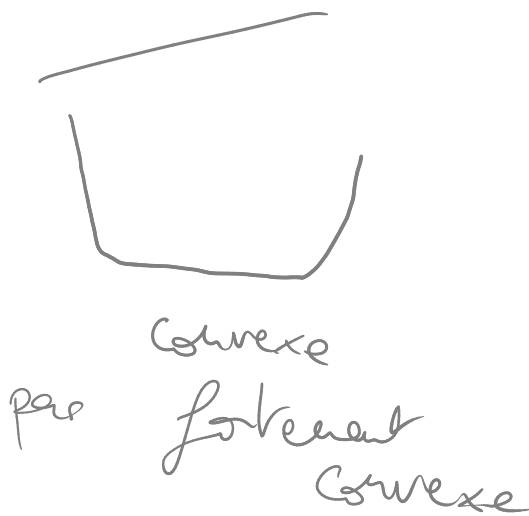
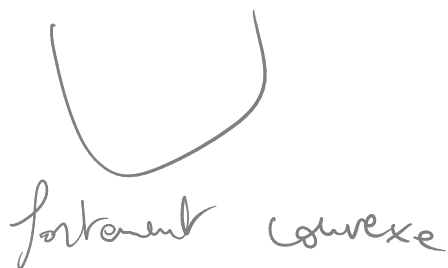
Rappel:

$f$  est

$\mu$ -fortement convexe

⇒  $f$  possède  
global

un unique minimum



L> En pratique, on  
fini (ex: nombre

dispose d'un temps  
max d'itérations)

et on recherche  
dites non-asymptotiques

• Complexité ;  
combien de temps  
pour atteindre la

Ex) Quel est  
de descente de  
pour obtenir

⇒ But: obtenir  
qui est une fonction

• Vitesse de  
donné un budget,  
espérer atteindre ?

Ex) Si on effectue  
descente de gradient,  
au pire cas de

⇒ Fonction

Remarque: On peut  
vitesse de  
complexité, et

donc des garanties

Étant donné  $\varepsilon > 0$ ,  
doit-on "attendre"  
précision  $\varepsilon$  ?

le nombre d'itérations  
gradient nécessaires

$f(w_k) - \min_{w \in \mathbb{R}^d} f(w) \leq \varepsilon$  ?

un nombre d'itérations  
de  $\varepsilon$

convergence. Étant  
quelle précision puis-je

$K$  itérations de  
quelle est la valeur

$f(w_k) - \min_{w \in \mathbb{R}^d} f(w)$  ?

de  $K$

toujours transformer une  
convergence en  
vice-versa

# ① Résultats pour gradient

Cas convexe :  $f$  (M L) atteint son

On note

la descente de convexe qui

minimum

$$f^* = \min_{w \in \mathbb{R}^d} f(w)$$

## Théorème (Vitesse descente de gradient)

Supposons que l'on de descente de gradient

$$\alpha_k = \frac{1}{L} (> 0). \text{ Alors}$$

$$f(w_k) - f^*$$

où  $C > 0$  est dépend de  $w_0$ ,  $L$

On dit que la a une vitesse de

N.B. : Le résultat reste choix de  $\alpha_k$  (ex:

de convergence de la son problèmes convexes,

effectue  $K$  itérations avec  $K \geq 1$  et

$$\leq \frac{C}{K}$$

une constante qui mais pas de  $K$

descente de gradient convergence a  $O\left(\frac{1}{K}\right)$

vrai pour d'autres recherche linéaire)

## Interprétation

- Garantie au pire de  $f(w_k) - f^*$  ("distance à la valeur optimale")
- Implique la  $0 \leq f(w_k) - f^* \leq O\left(\frac{1}{k}\right)$

- Garantie plus nombre d'itérations à la

$$f(w_k) - f^* \leq O\left(\frac{1}{k}\right)$$

Alors avec  $10k$  itérations,

- Résultat de complexité

Soit  $\varepsilon > 0$ . Alors

on a au plus

$$O\left(\frac{1}{\varepsilon}\right)$$

## Cas fortement convexe

NB: On a toujours  $\mu \leq L$ .

car son la valeur à la valeur optimale")

convergence asymptotique

$$\Rightarrow f(w_k) - f^* \rightarrow 0 \quad k \rightarrow \infty$$

$$f(w_k) \xrightarrow{k \rightarrow \infty} f^*$$

précise qui relie la précision

pour  $k$  fixé

$$f(w_{10k}) - f^* \leq O\left(\frac{1}{10k}\right)$$

équivalent :

$$\text{on a } f(w_k) - f^* \leq \varepsilon$$

itérations

$f \in C_{\mathcal{L}}^{1,1}$ ,  $\mu$ -fortement convexe avec  $\mu > 0$

On note  $\{w^*\} = \operatorname{argmin}_{w \in \mathbb{R}^d} f(w)$

$$\text{et } f^* = f(w^*)$$

Théorème (Vitese de  
de gradient sur les  
convexes)

Supposons qu'on  
de descente de  
Alors

$$f(w_k) - f^* \leq C \left(1 - \frac{\mu}{L}\right)^k$$

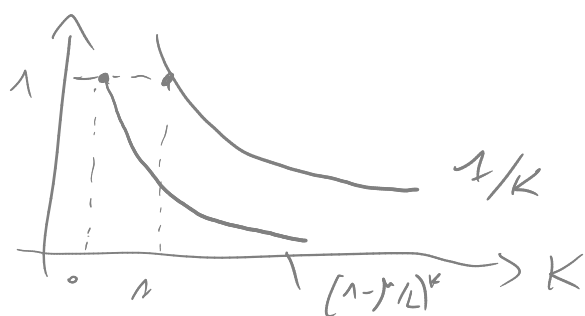
avec  $C = f(w_0) - f^* \geq 0$

On dit que la  
converge en vitese

(vitese exponentielle)

Interpretation

• Meilleure vitese  
dans le cas convexe!



convergence de la descente  
fonctions  $\mu$ -fortement

effective  $K$  iterations  
gradient avec  $\alpha_k = \frac{1}{L}$

descente de gradient

$$O\left(\left(1 - \frac{\mu}{L}\right)^k\right)$$

$$\text{en } K \quad \left(1 - \frac{\mu}{L}\right)^k = e^{K \ln\left(1 - \frac{\mu}{L}\right)}$$

de convergence que

$$0 < \mu \leq L \Rightarrow 1 - \frac{\mu}{L} < 1$$

$$\left(1 - \frac{\mu}{L}\right)^k \rightarrow 0$$

plus vite que  
 $\frac{1}{K}$

• La vitesse de convergence dépend des propriétés de  $f$

• Implique la

• On peut aussi

$$\|w_k - w^*\|$$

$\Rightarrow$  Convergence à optimale et vers

(Plus fort que les convexe)

convergence dépend ( $L$  et  $\mu$ )

convergence asymptotique

montrer que

$$\leq O\left(\left(1 - \frac{\mu}{L}\right)^k\right)$$

la fois vers la valeur la solution optimale!

garanties dans le cas

## (2) Accélération

$\hookrightarrow$  Que ce soit dans fortement convexe, les de la descente de optimales  $\Rightarrow$  il

qui calculent un et obtiennent de

le cas convexe ou vitesses de convergence gradient ne sont pas existe des algorithmes gradient par itération meilleures vitesses

Descente de gradient

Convexe

$$O\left(\frac{1}{k}\right)$$

fortement convexe

$$O\left(\left(1 - \frac{\mu}{L}\right)^k\right)$$

Optimal

$$O\left(\frac{1}{k^2}\right)$$

$$O\left(\left(1 - \sqrt{\frac{\mu}{L}}\right)^k\right)$$



↳ les algorithmes  
convergence optimale  
l'idée du momentum

avec vitesse  $s$  de  
sont basés sur

Première proposition :

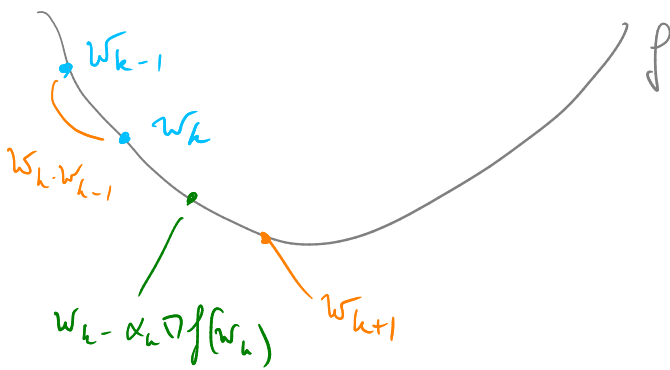
Heavy ball ("balle lourde"  
Polyak (~1964) "levée")

$\forall k \geq 1$

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k)$$

$$+ \beta_k (w_k - w_{k-1})$$

$\beta_k > 0$   
Paramètre de momentum  
Terme de momentum



• Polyak prouve la  
vitesse de convergence optimale  
pour  $f$  quadratique fortement  
convexe

• MAIS cet algorithme peut  
ne pas converger sur des  
fonctions générales

Seconde proposition :

Gradient accéléré  
(Yu. Nesterov 1983)

$$w_{k+1} = w_k - \alpha_k \nabla f(w_k +$$

$$\beta_k (w_k - w_{k-1})) + \beta_k (w_k - w_{k-1})$$

↳  $\neq$  avec Polyak :  
évalué après le pas

le gradient est  
de momentum

↳ Atteint les bornes

optimales de complexité

• pour  $f$  fortement  
 $(1 - \sqrt{\frac{\mu}{L}})^k$  avec  $c$

Convexe

$\beta_k = \beta > 0$   $\beta$  dépend  
de  $L$  et  $\mu$

• pour  $f$  convexe

$\frac{1}{k^2}$  avec  $c$

$\{\beta_k\}$  choisie indépendamment  
de  $f$  et de  $w_0$   
(et de  $\alpha_n$ )

NB: le gradient accéléré  
au cas proximal

peut se généraliser

(ex: ISTA  $\Rightarrow$  Fast ISTA)