

OUTILS D' LES SDD

OPTIMISATION

28 octobre 2024

Aujourd'hui:

CM: Optimisation

non convexe / Projet

TD/TP: Optim

convexe

OPTIMISATION NON NÉGATIVE

CONVEXITÉ ET COURBURE

① Cadre de travail

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

↳ Difficultés en optimisation
non convexe (par comparaison avec
le cas convexe)

a) Distinction entre minima locaux
et globaux

b) Un point $\bar{w} \in \mathbb{R}^d$ pour lequel
 $\|\nabla f(\bar{w})\| = 0$ n'est pas forcément
un minimum

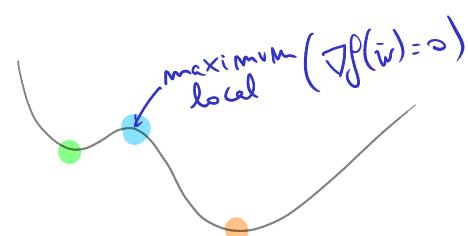
c) Si on applique la descente
on ne peut garantir qu'une
le gradient est nul !

↳ Malgré ces difficultés
(et notamment en sciences des
descents de gradient converge vers
des minima locaux voire globaux

↳ En s'intéressant à des
on peut expliquer ce phénomène
algorithmes encore plus efficaces
ces non convexes.

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

Ex)



Une fonction non convexe

de gradient à un problème non convexe,
convergence vers un point en lequel

théoriques, en pratique on observe souvent
données), on observe souvent que la
des minima locaux voire globaux
sous-classes de problèmes non convexe,
pratique, et même construire des
que la descente de gradient dans le

② Paysages d'optimisation

On s'intéresse à des fonctions

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

f de classe C^2

on peut définir $\nabla f(w) \in \mathbb{R}^d$ (gradient)
 et $\nabla^2 f(w) \in \mathbb{R}^{d \times d}$
 (matrice hessienne, qui est
 symétrique)

Pour de telles fonctions, on

$$\left[w^* \in \mathbb{R}^d \text{ minimum local de } f \right]$$

ce qui est équivalent à

$$\left[\bar{w} \in \mathbb{R}^d \text{ n'est pas un minimum local} \right]$$

sait que

$$\Rightarrow \left[\nabla f(\bar{w}) = 0 \text{ et } \nabla^2 f(\bar{w}) \succeq 0 \right]$$

\uparrow
 $v^T \nabla^2 f(\bar{w}) v \geq 0$
 $\forall v \in \mathbb{R}^d$

$$\left\langle \left. \left[\nabla f(\bar{w}) \neq 0 \text{ ou } \exists v \in \mathbb{R}^d, v^T \nabla^2 f(\bar{w}) v < 0 \right] \right\langle \right.$$

\uparrow
 v est une direction
 de courbure négative pour f en
 \bar{w}

Def. Paysage d'optimisation de f

Le paysage de f est la fonction des valeurs de ∇f

$\forall w \in \mathbb{R}^d$, on parle de

- Point non critique
- Point critique d'ordre 1
- Point selle
- Point critique d'ordre 2

classification des points de \mathbb{R}^d en
 et $\nabla^2 f$ en ces points.

si $\nabla f(w) \neq 0$

si $\nabla f(w) = 0$

si $\nabla f(w) = 0$ et $\nabla^2 f(w) \neq 0$

si $\nabla f(w) = 0$ et $\nabla^2 f(w) \succeq 0$

NB: $\{ \text{points critiques d'ordre } 1 \} =$

$\{ \text{points selle} \} \cup \{ \text{points critiques d'ordre } 2 \}$

$$\nabla^2 f(w) \succeq 0 \Leftrightarrow \forall v \in \mathbb{R}^d,$$

$$\nabla^2 f(w) \not\succeq 0 \Leftrightarrow \exists v \in \mathbb{R}^d,$$

$$v^T \nabla^2 f(w) v \geq 0$$

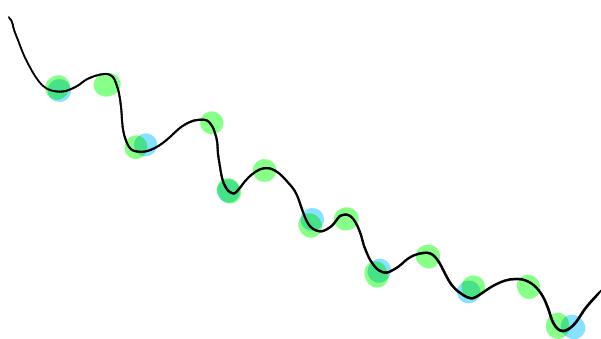
$$v^T \nabla^2 f(w) v < 0$$

$$\nabla^2 f(w) = \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \not\succeq 0$$

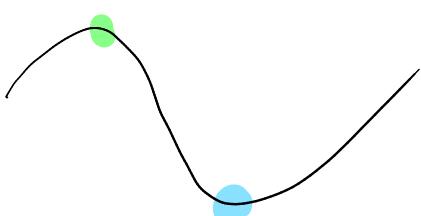
$$\begin{bmatrix} 1 & 0 \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \geq 0$$

$$\begin{bmatrix} 0 & 1 \end{bmatrix}^T \begin{bmatrix} 1 & 0 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} < 0$$

Mauvais paysage d'optimisation



Bon paysage d'optimisation



- Beaucoup de points critiques d'ordre 1 et d'ordre 2
- Beaucoup de minima locaux non globaux

- Peu de points critiques
- Tout minimum local est global

→ En sciences des données, il rencontrés ont en général un bon plus faciles à résoudre qu'un mauvais paysage.

→ Exemples: reconstruction de réseaux de neurones

a été montré que les problèmes paysage d'optimisation, et sont donc problème d'optimisation avec un

phase, apprentissage de dictionnaire, linéaires, ACP (analyse en composantes principales)

Notre exemple pour le projet:

Données: $M \in \mathbb{R}^{d \times d}$ matrice de

- \Rightarrow On dispose d'un ensemble de coefficients de
- \Rightarrow Cadre typique

BUT: Trouver la meilleure

Approche naïve: minimiser $W \in \mathbb{R}^{d \times d}$

⊕ Problème concrète

⊕ Solution est M si $S = \{-1, -d\} \times \{-1, -d\}$

⊖ Pas d'information concernant
qui ne sont pas observés

Approche moderne: On suppose que

\Rightarrow Conceptuellement:
sans avoir à connaître

\Rightarrow Mathématiquement:

Problème d'optimisation associé

Paramètre: minimiser $U \in \mathbb{R}^{d \times r}$
Approximation de U^*

$$\frac{1}{2} \sum_{(i,j) \in S} ([UU^T]_{ij} - M_{ij})^2$$

Factorisation de matrice de rang faible

données que l'on observe partiellement
sous-ensemble $S \subseteq \{-1, -d\} \times \{-1, -d\}$
 M

de systèmes de recommandation (ex: Netflix)

approximation possible de M

$$\frac{1}{2} \sum_{(i,j) \in S} (W_{ij} - M_{ij})^2$$

↑
On cherche la matrice W qui
colle le mieux possible aux
observations de M

$\rightarrow \{-1, -d\} \times \{-1, -d\}$ (si on connaît tous les
coefficients)

les coefficients de la matrice

M possède une structure de rang faible

On peut approcher M de manière précise
beaucoup de ses coefficients

On suppose qu'il existe $U^* \in \mathbb{R}^{d \times r}$
avec $r \ll d$ telle que $M = U^*(U^*)^T$

$$([UU^T]_{ij} - M_{ij})^2$$

↓
Approximation de M

④ Exploite la structure de M

⑤ Points de variétés

dr contre d^2 pour
l'approche naïve

→ On peut montrer que ce que les points critiques d'ordre des minima globaux.

→ Par ailleurs, si on échantillonne un minimum global.

Bilan :

. caractériser le paysage difficulté d'un problème difficulté à trouver ses

. De nombreux problèmes propriétés \ minima

C'est le cas pour la

③ Descente de gradient

Q. étant donné un problème que peut-on dire de la

Descente de gradient : w_{k+1}

④ Non convexe

⑤ Présence de points selle

problème (rappel de 2016-2017)
1 soit soit des points selle, soit

⑥ $O(\log(d))$ valeur, alors M est

suffisamment de valeur, mais moins que la taille de M !

d'optimisation permet de juger de la non convexe, et notamment de la minima globaux

en siéan des données possèdent la globaux } = { points critiques d'ordre 2 } factorisation de matrice)

et points critiques d'ordre 2

non convexe avec un bon paysage, descente de gradient sur ce problème ?

$$= w_k - \alpha_k \nabla f(w_k) \quad \text{avec } \alpha_k > 0$$

partant de $w_0 \in \mathbb{R}^d$

Si w_0 est un point stable ce n'est pas un minimum et w_0 !
 $w_1 = w_0, w_2 = w_1,$

→ Très facile de constater que le gradient ne converge pas
 → Mais en pratique, ces

Théorème (2015-2019) : La descente de gradient converge vers un point critique à l'ordre 2

$(\nabla f(w_0) = 0, \nabla^2 f(w_0) \neq 0)$, alors la descente de gradient est bloquée en ...

des exemples pour lesquels la descente vers un point critique d'ordre 2 exemples ne semblent pas fréquents

descente de gradient converge vers 2 pour presque tout $w_0 \in \mathbb{R}^d$.

④ Le projet

→ En individuel ou

→ Point de départ:
 . Notebook sur
 . Surfer avec des

→ But: Étudier un

Si $\nabla f(w_k) = 0$

en binôme

l'optimisation non convexe
question

algorithme de descente de gradient
+ courbure négative

et $\nabla^2 f(w_k) \neq 0$, alors

$$w_{k+1} = w_k + \alpha_k v_k$$

$$\text{avec } v_k^\top \nabla^2 f(w_k) v_k < 0$$

→ Converge plus vite

vers des points critiques d'ordre 1 que la descente de gradient

→ Converge même

en partant d'un point nulle !

↳ le projet:

- Questions
- Application
- Cadre

↳ Deadline ≈ fin janvier

Théoriques (démonstration, complexité)
à la factorisation de matrice
stochastique

2025