## Mathématiques pour l'intelligence artificielle UE MIA, M2 IAAA, AMU, 2022-2023

Thomas Schatz Jeudi 3 novembre 2022

9 séances de 3h

Partie 1 (DM1) : algèbre linéaire et probabilités

1. Notions de bases sur les preuves (+ Algèbre linéaire?)

2. Algèbre linéaire (+ Probabilités?)

3. Probabilités

Partie 2 (DM2): statistique et optimisation

4. Statistiques

5. Optimisation

Partie 3 (DM3):

6. Optimisation sous contraintes

7. Optimisation stochastique



'Backtracking' avec la règle d'Armijo le long de l'arc de projection

Input: 
$$x_0 \in \mathbf{R}^n, f: U \subset \mathbf{R}^n \to \mathbf{R}$$
  
 $U$  convexe, for

Iteration: 
$$x_{k+1} := p_k(\beta^{m_k}s)$$

$$p_k(r) = [x_k - r\nabla f(x_k)]_U \in$$

de classe  $C^1, s > 0, 0 < \beta < 1, 0 < \sigma < 1$ . ermé, non-vide

et  $m_k$  plus petit entier m tel que  $f(x_k) - f(x_{k+1}) \ge \sigma \nabla f(x_k)^T (x_k - x_{k+1})$ 

'Backtracking' avec la règle d'Armijo le long de l'arc de projection



'Backtracking' avec la règle d'Armijo le long de l'arc de projection

Input: 
$$x_0 \in \mathbf{R}^n, f: U \subset \mathbf{R}^n \to \mathbf{R}$$
  
 $U$  convexe, for

Iteration: 
$$x_{k+1} := p_k(\beta^{m_k}s)$$

$$p_k(r) = [x_k - r\nabla f(x_k)]_U \in$$

de classe  $C^1, s > 0, 0 < \beta < 1, 0 < \sigma < 1$ . ermé, non-vide

et  $m_k$  plus petit entier m tel que  $f(x_k) - f(x_{k+1}) \ge \sigma \nabla f(x_k)^T (x_k - x_{k+1})$ 

minimize f(x)



## Dualité

#### subject to $x \in X$ , $g_j(x) \leq 0$ , $j = 1, \ldots, r$ , $f: \mathbf{R}^n \to \mathbf{R}$ $g_j: \mathbf{R}^n \to \mathbf{R}$ $X \subset \mathbf{R}^n$

 $f^* = \inf_{\substack{x \in X \\ g_j(x) \le 0, j=1,\dots,r}} f(x).$ 

Il existe  $x^*$ , tel que  $f(x^*) = f^*$ 

Т j=1

## Dualité

# $L(x,\mu) = f(x) + \sum \mu_j g_j(x) = f(x) + \mu' g(x).$



9

**Fonction duale** 

**Problème dual** maximize  $q(\mu)$ subject to  $\mu \ge 0, \ \mu \in D$ ,

> **Proposition 5.1.2:** The domain D of the dual function q is convex and q is concave over D.

Proposition 5.1.3: (Weak Duality Theorem) We have

## Dualité

- $q(\mu) = \inf_{x \in X} L(x, \mu).$ 
  - $D = \left\{ \mu \mid q(\mu) > -\infty \right\}.$

 $q^* \leq f^*$ .

#### Dualité forte et qualification des contraintes

Example :

dual:  $q^* = f^*$ 

## Dualité

Si f est convexe et les contraintes  $g_j$  sont linéaires, alors il n'y a pas d'écart

## Dualité

**Proposition 5.1.5: (Optimality Conditions)**  $(x^*, \mu^*)$  is an optimal solution-Lagrange multiplier pair if and only if

 $x^* \in X, \quad g(x^*) \leq 0, \quad (\text{Primal Feasibility}), \quad (5.4)$  $\mu^* \geq 0, \quad (\text{Dual Feasibility}), \quad (5.5)$  $x^* = \arg\min_{x \in X} L(x, \mu^*), \quad (\text{Lagrangian Optimality}), \quad (5.6)$  $\mu_j^* g_j(x^*) = 0, \quad j = 1, \dots, r, \quad (\text{Complementary Slackness}). \quad (5.7)$ 

9 séances de 3h

Partie 1 (DM1) : algèbre linéaire et probabilités

1. Notions de bases sur les preuves (+ Algèbre linéaire?)

2. Algèbre linéaire (+ Probabilités?)

3. Probabilités

Partie 2 (DM2): statistique et optimisation

4. Statistiques

5. Optimisation

Partie 3 (DM3):

6. Optimisation sous contraintes

7.Optimisation stochastique

9 séances de 3h

Partie 1 (DM1) : algèbre linéaire et probabilités

1. Notions de bases sur les preuves (+ Algèbre linéaire?)

2. Algèbre linéaire (+ Probabilités?)

3. Probabilités

Partie 2 (DM2): statistique et optimisation

4. Statistiques

5. Optimisation

Partie 3 (DM3):

6. Optimisation sous contraintes

7.Optimisation stochastique

#### **Optimisation stochastique**

**Contexte :** minimisation du risque empirique pour une fonction de coût "séparable par point de donnée"

 $R_n(u)$ 

#### **Descente de gradient stochastique**

 $w_1 \in \mathbb{R}^d$  given

 $w_{k+1} \leftarrow w$ 

 $i_k$  is chosen randomly from  $\{1, \ldots, n\}$  and  $\alpha_k$  is a positive stepsize

$$w) = \frac{1}{n} \sum_{i=1}^{n} f_i(w)$$

$$v_k - \alpha_k \nabla f_{i_k}(w_k)$$

#### **Optimisation stochastique**

#### **Exemple de garantie de convergence**

(cf. Bottou, Curtis et Nocedal (2018) Optimisation Methods for Large-Scale Machine Learning)

Si



Lipschitz continuous with Lipschitz constant L > 0, i.e.,

 $\|\nabla F(w) - \nabla F(\overline{w})\|_2 \leq L \|w - \overline{w}\|_2$  for all  $\{w, \overline{w}\} \subset \mathbb{R}^d$ .

plus des conditions de régularité pas très contraignantes

Alors

$$\hat{\alpha}_k^2 < \infty$$

Assumption 4.1 (Lipschitz-continuous objective gradients). The objective function F:  $\mathbb{R}^d \to \mathbb{R}$  is continuously differentiable and the gradient function of F, namely,  $\nabla F : \mathbb{R}^d \to \mathbb{R}^d$ , is

$$\liminf_{k \to \infty} \mathbb{E}[\|\nabla F(w_k)\|_2^2] = 0$$

9 séances de 3h

Partie 1 (DM1) : algèbre linéaire et probabilités

1. Notions de bases sur les preuves (+ Algèbre linéaire?)

2. Algèbre linéaire (+ Probabilités?)

3. Probabilités

Partie 2 (DM2): statistique et optimisation

4. Statistiques

5. Optimisation

Partie 3 (DM3):

6. Optimisation sous contraintes

7.Optimisation stochastique

9 séances de 3h

Partie 1 (DM1) : algèbre linéaire et probabilités

1. Notions de bases sur les preuves (+ Algèbre linéaire?)

2. Algèbre linéaire (+ Probabilités?)

3. Probabilités

Partie 2 (DM2): statistique et optimisation

4. Statistiques

5. Optimisation

Partie 3 (DM3):

6. Optimisation sous contraintes

7.Optimisation stochastique

## Théorie de l'apprentissage : TLDR

High-dimensional statistics : non asymptotic

- Bias-variance decomposition of risk (same approach as in classical case)
- Explicit computation impossible for typical ML problems
- Upper bounding :
  - bias : Rademacher (or Gaussian) averages
  - variance : concentration inequalities

High-dimensional statistics : asymptotic

• Different asymptotic regimes, e.g.  $n \rightarrow +\infty$  et  $n/p \rightarrow constante$ 

Deep learning theory?

- Good generalisation properties of minimum-norm interpolations?

• Over-parametrisation + SGD -> computationally tractable minimum-norm interpolation?

## **Application 2: Normes de** matrices

$$\|A\| = \sup_{\mathbf{x}\neq\mathbf{0}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|A\|$$
$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2\right)^{1/2}.$$

A matrix norm is a function  $\|\cdot\|: \mathbb{R}^{m \times n} \to \mathbb{R}$  that has the following properties: •  $||A|| \ge 0$  for any  $A \in \mathbb{R}^{m \times n}$ , and ||A|| = 0 if and only if A = 0•  $\|\alpha A\| = |\alpha| \|A\|$  for any  $m \times n$  matrix A and scalar  $\alpha$ •  $||A + B|| \le ||A|| + ||B||$  for any  $m \times n$  matrices A and B

> $A\mathbf{x}$  $||A||_2 = \sigma_1$

$$||A||_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$$

## **Application 2: Normes de** matrices

$$\|A\| = \sup_{\mathbf{x}\neq\mathbf{0}} \frac{\|A\mathbf{x}\|}{\|\mathbf{x}\|} = \max_{\|\mathbf{x}\|=1} \|A\|$$
$$\|A\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2\right)^{1/2}.$$

A matrix norm is a function  $\|\cdot\|: \mathbb{R}^{m \times n} \to \mathbb{R}$  that has the following properties: •  $||A|| \ge 0$  for any  $A \in \mathbb{R}^{m \times n}$ , and ||A|| = 0 if and only if A = 0•  $\|\alpha A\| = |\alpha| \|A\|$  for any  $m \times n$  matrix A and scalar  $\alpha$ •  $||A + B|| \le ||A|| + ||B||$  for any  $m \times n$  matrices A and B

> $A\mathbf{x}$  $||A||_2 = \sigma_1$

$$||A||_F = \sqrt{\sigma_1^2 + \dots + \sigma_r^2}$$