

PART 1

(1)

- Motiv<sup>o</sup> opt<sup>o</sup> convexe  
en ML  $x_i \rightarrow y_i$

Regress<sup>o</sup> linéaire

$$\min_{\beta} \sum_i |y_i - \langle x_i, \beta \rangle| = f(\beta)$$

Classif. proba  $\frac{e^{\langle x, \beta \rangle}}{1 + e^{\langle x, \beta \rangle}}$

$$\min_{\beta} \sum_i l(\langle x_i, \beta \rangle, y_i)$$

logit ~~logistic~~ hinge

Arbres appli:  
an data science

$\Delta n \rightarrow \beta$  ds la suite!!

Aggr de contraintes, pénalités  
non lisses.

- F<sup>o</sup> convexe, Ensemble convexe, F<sup>o</sup> étendue  $f: \mathbb{R}^n \rightarrow \bar{\mathbb{R}}, \mathbb{R}$   
(indication)

ex:  $\max \|Ax - b\|^2$   $\|x\|_1$ ,  
 $\{x \mid f(x) \leq c\}$   
axe.

$$f \circ A$$

$$\min_{x \in C} \|x - y\|_1 + \epsilon$$

axe aux

$$\begin{cases} f \\ g \end{cases} \max(f, g)$$

écart

- Differentials & Gradients

$$m=1$$

$$\frac{f(u+\epsilon) - f(u)}{\epsilon} \rightarrow f'(u)$$

$$\boxed{m \geq 1} f(u+\epsilon h) = f(u) + \epsilon \langle \nabla f(u), h \rangle + o(\epsilon)$$

$$\nabla f: \mathbb{R}^n \rightarrow \mathbb{R}^n \quad (\text{vect. fct})$$

ex:  $f(u) = \frac{1}{2} \|Ax - b\|^2$   $f(u+\epsilon h) = \frac{1}{2} \|Ax - b + A\epsilon h\|^2 = f(u) + \frac{1}{2} \langle Ax - b, A\epsilon h \rangle + \epsilon^2 \|Ah\|^2$

ex: multiclass  $\logit(B_i - B_k)$ .

$$\rightarrow = \langle \epsilon, A^T(Au - b) \rangle$$

$$= \nabla P(u).$$

prediction loss  
prediction:

$$L(u, y) = \log \sum_i \exp(u_i) - y_i \log u_i$$

$$(u_i - y_i) = \langle \beta_i, x_i \rangle$$

$$\nabla L(u, y) = \left( \frac{e^{u_i}}{\sum_i e^{u_i}} - y_i \right)_i$$

observation  
( $y_i - y_{ik}$ ) proba vector

(2)

Chain rule:  $\nabla(f \circ A) = \cancel{A^T} \nabla f \circ A$

$$f(A(\cdot + \epsilon)) = f(Ax) + (\cancel{\nabla f(Ax)}, A\epsilon) + o(\epsilon).$$

$$\nabla f \circ g(x) = [\cancel{\nabla g(x)}]^T (\cancel{\nabla f(g(x))}) + R^n$$

$$f: R^P \rightarrow R$$

$$\underbrace{L}_{\in R^{n \times P}}$$

$$\epsilon \in R^P$$

$$g: R^n \rightarrow R^P$$

$$\epsilon \in R^{P \times n}$$

Differential:  $f(x + \epsilon) = f(x) + \epsilon \nabla f(x)[\epsilon] + o(\epsilon)$ .

$$g: R^m \rightarrow R^P$$

$(R^{n \times P}$  linear)

⚠  $f: R^n \rightarrow R$

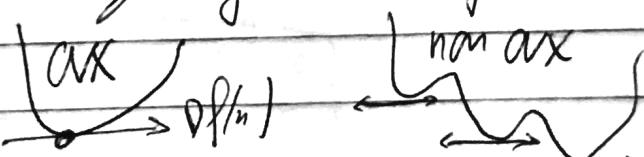
$$\nabla f(x) \in R^{1 \times m}$$

$$\nabla f(x) \in R^{n \times 1}$$

$$(\nabla f(x) = \nabla f(x)^T)$$

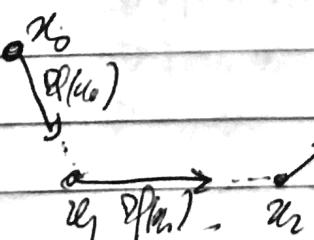
Fundamental thm  $x^* \in \text{arg min } f \Leftrightarrow \nabla f(x^*) = 0$

for x



Slope descent:

$$x_{k+1} = x_k - \tau \nabla f(x_k)$$



$$\nabla f(x^*) = 0$$

$$f(x - \tau \nabla f(x)) = f(x) - \tau \langle \nabla f(x), \nabla f(x) \rangle + O(\tau).$$

$\| \nabla f(x) \|^2 > 0$  if  $x$  not min  
 $\Rightarrow \exists \tau$  small enough such that  $f(x_{k+1}) < f(x_k)$

choice of  $\tau$   $\xrightarrow{\text{fixed}} \text{line search}$

Thm:  $f(x_{\text{min}}), \tau_k \in [t_{\text{min}}, t_{\text{max}}]$   
 $t_{\text{min}} / t_{\text{max}} \quad x_k \rightarrow x^*$

TP

Hessian:  $f(x) + \epsilon h = f(x) + \epsilon \nabla f(x) + \frac{\epsilon^2}{2} \nabla^2 f(x) h^2 + o(\epsilon^2)$

 $m=1$ 

$$\nabla f(x) + \epsilon \nabla^2 f(x) h = \nabla f(x) + \epsilon \langle \nabla^2 f(x) h, h \rangle$$

$$\nabla^2 f(x)^\top = \nabla^2 f(x) \in \mathbb{R}^{n \times n}, \text{ symmetric}$$

$$\nabla^2 f(x) = \left[ \frac{\partial^2 f}{\partial x_i \partial x_j} \right]$$

$$f(x+h) = f(x) + \epsilon \langle \nabla f(x), h \rangle + \frac{\epsilon^2}{2} \langle \nabla^2 f(x) h, h \rangle + o(\epsilon^2)$$

$\downarrow$  2x diff.

Then:  $f \text{ arx} \Leftrightarrow \nabla^2 f(x) \succ 0$ .

Rappd:  $A \succ 0 \Leftrightarrow \forall h, \langle Ah, h \rangle \geq 0 \Leftrightarrow \text{Spectrum}(A) \subset \mathbb{R}^+$

- Strong convexity  $\Rightarrow$  smoothness.
- Conditioning
- Convergence speed

Newton

BFGS

SGD

Aleksandrov

TP: Logistic damping

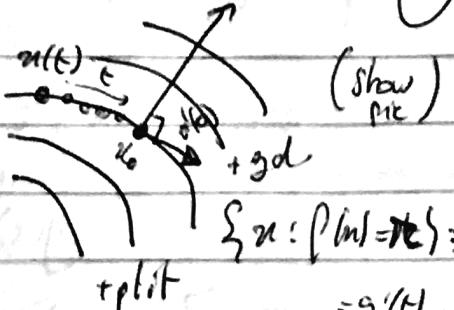
Next • Non smooth optim TP Projected Gradient  
TP Forward Backward  
TP Lasso TP  
TP Duality

• SGD - TP

• Automatic Differentiation - TP

(4)

- Rappel sur le gradient + descente
- Corrèct° exo séance précédente (4)
- Stepper descent dir



$$x(t) \in L_n \Leftrightarrow \underbrace{f(x(t))}_{\equiv g(t)} = n$$

$$g(t+\epsilon) = f(x(t) + \epsilon x'(t)) = f(x(t)) + \epsilon \langle \nabla f(x(t)), x'(t) \rangle + o(\epsilon)$$

done  $g'(t) = \langle \nabla f(x(t)), x'(t) \rangle$

mais  $\hat{c} g = ct \Rightarrow g'(t) = 0 \Rightarrow g'(0) = \langle \nabla f(x_0), \underbrace{x'(0)}_{\text{tangente}} \rangle = 0$

$$\underset{\|x_0-x_\epsilon\|}{\operatorname{argmin}} f(x_0 + \epsilon d) = f(x_0) + \epsilon \langle \nabla f(x_0), d \rangle + o(\epsilon).$$

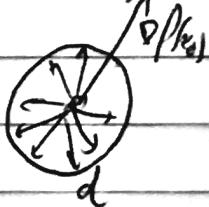
$$x_\epsilon = x_0 + \epsilon d$$

$$\|d\| \leq 1$$

need to solve in the limit

$$\bullet \underset{\epsilon B_\epsilon(x_0)}{x}$$

$$\operatorname{argmin} \langle \nabla f(x_0), d \rangle \rightsquigarrow d = -\frac{\nabla f(x_0)}{\|\nabla f(x_0)\|}$$



→ Rappels hémisphère  $\nabla f$

$C^2$  smooth function:

$$\|\nabla^2 f\| \leq L$$

$$\|\nabla f\| = \max_{\|x\| \leq 1} \frac{|\nabla f|}{\|x\|}$$

i.e.  $|Ax_i - Ay_j| \leq \|A\|_F \|y\|_2 \|x\|_2$   
↳ Lipschitz constant

contrac<sup>o</sup>-ratio

If M symmetric

= max eigenvalues(M).

$$Ai = \lambda_i e_i$$

↳ BON.

Rapp:  $f(x) = \frac{1}{2} \|Ax-y\|^2 \quad \nabla f(x) = A^T(Ax-y) \quad \nabla^2 f(x) = A^T A$

$$L = \|A^T A\| \leq \|A\|^2$$

ex:  $f(x) = \frac{1}{2} (x_1^2 + y_2^2) \quad \nabla f(x) = \begin{pmatrix} x_1 \\ y_2 \end{pmatrix} \quad \nabla^2 f(x) = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$   
if  $y > 1 \rightarrow \|\nabla^2 f(x)\| = y$

(5)

Thm de convergence :  $u_{k+1} = u_k - \tau \nabla f(u_k)$

si  $0 < \tau < \frac{2}{L}$  alors  $u_k \rightarrow u^*$  e.s. et  $f(u_k) - f(u^*) \leq C/k$  (difficile).

• Strong convexity :  $\nabla^2 f(x) \geq \mu I_d \Rightarrow \nabla^2 f(x) - \mu I_d \geq 0$   
 $\Rightarrow$  strict CVX ie  $\langle \nabla^2 f(x), h \rangle \geq \|h\|^2/\mu$ .  
 $\Rightarrow$  unique du Min ie  $\text{eigen}(\nabla^2 f(x)) \geq \mu$ .

$$\text{ex} : f(x) = \frac{1}{2} \|Ax - y\|^2 \quad \nabla^2 f(x) = A^* A \quad \mu = \min \text{eigen}(A^* A)$$

$\mu > 0 \Rightarrow A$  injective ie  $\ker A = \{0\}$

ie e.s. et unique solution  
 $x^* = (A^* A)^{-1} A^* y$ . "least" square

Thm de convergence Linien : si  $0 < \mu \leq L < +\infty$

$$|u_k - u^*| \leq \underbrace{\rho}_{\text{uniqueness}} \cdot \underbrace{\rho^k}_{\text{con.}} \cdot \|u_0 - u_k\|$$

Dans le cas quadratique

$$\text{si } \tau = \frac{2}{L+\mu} \text{ alors } \rho = \frac{L-\mu}{L+\mu} = 1 - \frac{2\epsilon}{1+\epsilon} \quad \epsilon = \frac{\mu}{L} < 1$$

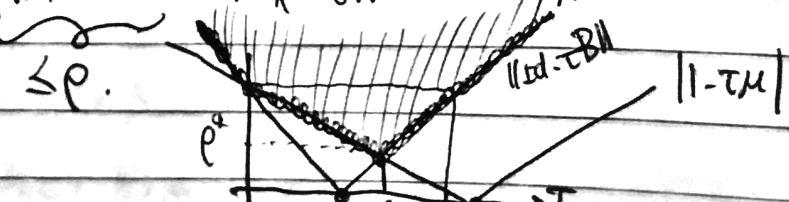
Conditionnement

en effet :  $f(x) = \frac{1}{2} \langle Bx, x \rangle - b^* \cdot x \rightarrow Bx^* = b$ . (ratio valeur propre).  
 $\nabla f(x) = Bx - b$ .

$$u_{k+1} = u_k - \tau (Bx_k - b)$$

$$(u_{k+1} - x^*) = (u_k - x^*) - \tau B(u_k - x^*) = (Id - \tau B)(u_k - x^*)$$

$$\Rightarrow \|u_k - x^*\| \leq \underbrace{\|Id - \tau B\|}_{\leq \rho} \times \|u_k - x^*\|$$



Newton :  $u_{k+1} = u_k - \nabla^2 f(u_k)^{-1} \nabla f(u_k)$  nif Quad. 1 iter!  $\frac{1}{\mu} \frac{2}{L+\mu} \frac{1}{\mu} \leq \frac{2}{L} \rho^k$

⚠  $\beta \leftrightarrow w$

⑥

TP classif logistic: input  $(x_i, y_i)_{i=1}^n$

Parametr:  $\beta \in \mathbb{R}^p$   $\epsilon \mathbb{R}^{p+1}, +1$ .

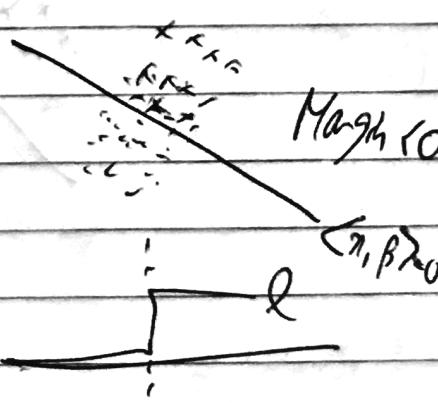
Frontière de décision  $\langle x, \beta \rangle = 0$

Décision ok if  $\text{sign}(\langle x_i, \beta \rangle \cdot y_i) > 0$

0-1 loss: on paye 1 par mauvais choix

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \text{sign}(\langle x_i, \beta \rangle \cdot y_i)$$

$\ell(y_i \langle x_i, \beta \rangle)$



$\Rightarrow$  non convexe

Ridge:  ~~$\ell(s) = \frac{1}{2} s^2$~~   $\Rightarrow \ell(s) = \text{Proba}(y_i f(s) = +1) \Rightarrow$  smooth decision

Affine decision  $\langle x, \beta \rangle = b \Leftrightarrow \langle \begin{bmatrix} x \\ 1 \end{bmatrix}, \begin{bmatrix} \beta \\ b \end{bmatrix} \rangle = 0$

$\Rightarrow$  feature augmenting

$$\min_{\beta} f(\beta) = \frac{1}{n} \sum_i \ell(y_i \langle x_i, \beta \rangle)$$

$\boxed{\text{exo}}$   $\nabla f(\beta) = \frac{1}{n} \sum_i y_i \cdot \ell'(y_i \langle x_i, \beta \rangle) x_i$

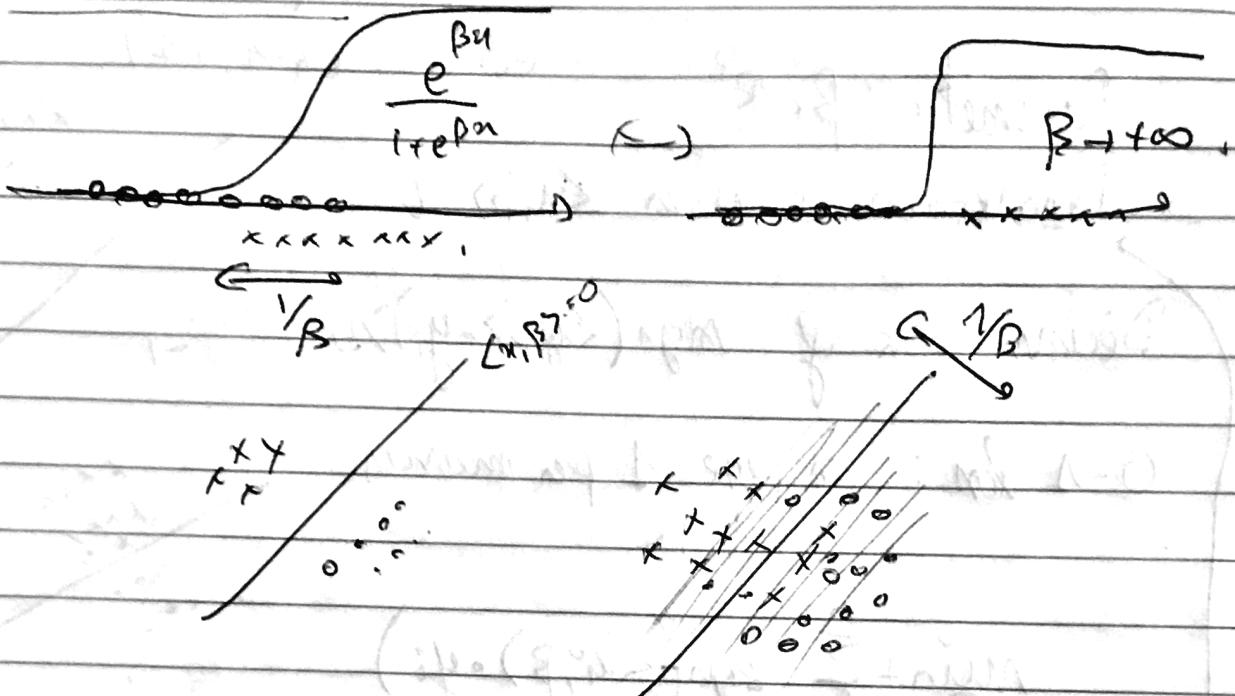
Preuve:  $f(\beta + \delta) = \frac{1}{n} \sum_i \ell(y_i \langle x_i, \beta \rangle + \epsilon y_i \langle x_i, \delta \rangle)$

Logistic  $\ell(s) = \frac{e^s}{1+e^s}$

$$\begin{aligned} &= \frac{1}{n} \sum_i \ell(y_i \langle x_i, \beta \rangle + \epsilon y_i \ell'(y_i \langle x_i, \beta \rangle) \langle x_i, \delta \rangle) + o(\epsilon) \\ &= f(\beta) + \epsilon \langle \delta, -\frac{1}{n} \sum_i x_i \times y_i \cdot \ell'(y_i \langle x_i, \beta \rangle) \rangle + o(\epsilon) \end{aligned}$$

6'

# logistic loss intuition



(7)

Autre notation

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n \times p}$$

$$f(\beta) = \mathcal{L}(X\beta, y)$$

$$\text{Prop: } \nabla f(\beta) = X^T \nabla \mathcal{L}(X\beta, y)$$

$$\text{Preuve: } \mathcal{L}(X(\beta + \varepsilon \delta), y) = \mathcal{L}(X\beta + \varepsilon X\delta, y).$$

$$= \mathcal{L}(X\beta, y) + \varepsilon \left\langle \nabla \mathcal{L}(X\beta, y), X\delta \right\rangle + o(\varepsilon)$$

$$\text{Id: } \mathcal{L}(u, y) = \frac{1}{n} \sum l(u_i, y_i).$$

$$\nabla \mathcal{L}(u, y) = -\frac{1}{n} (y_1 l'(u_1, y_1), y_2 l'(u_2, y_2), \dots, y_n l'(u_n, y_n))$$

" =  $-\frac{1}{n} (y_1 \circ l'(u_1, y_1), y_2 \circ l'(u_2, y_2), \dots, y_n \circ l'(u_n, y_n))$

$\hookrightarrow$  note matlab/python

$$\text{Logistic } l' = \Theta$$

$$\text{Hessian: } l(s) = \log(1 + e^{-s}) \quad l'(s) = \frac{e^{-s}}{1 + e^{-s}} = \frac{1}{1 + e^{-s}}$$

$$l''(s) = \frac{e^{-s}}{(1 + e^{-s})^2}$$

$$f(\beta) = L(y \circ \Theta(X\beta)) \quad \nabla f(\beta) = X^T \cdot \text{diag}(y) \cdot \nabla L(y \circ X\beta)$$

$$\partial^2 f(\beta) = X^T \cdot \underbrace{\text{diag}(y)}_{I \times I} \cdot \partial^2 L(y \circ X\beta) \cdot \text{diag}(y) \cdot X.$$

$$\|\partial^2 f(\beta)\| \leq \|X\|^2 \cdot \underbrace{\|\partial^2 L(y \circ X\beta)\|}_{\leq 1/4} \cdot \|X\|$$

8

Multi-class Logistic : Histogram represent.

K classes, K vectors  $(\beta, \underline{\beta_k}) = \underline{\beta} \in \mathbb{R}^{P \times K}$

Data :  $(x_i, d_i)$

$\sum_k$  : histo

$\prod_i^k \prod_i^k$

$\prod_{i=1}^k \prod_{j=1}^k$  desk

LSE

Loss between histo:  $l(u, d) = \log \left[ \sum_k \exp(u_k) \right] - \sum_k u_k d_k$

$$\nabla l(u, d) = \left[ \frac{e^{u_k}}{\sum_k e^{u_k}} - d_k \right]_k$$

soft max.

Bmg:  $\underbrace{\text{LSE}(u_k) - c}_{\max(u_k)} = \text{LSE}((u_k)_b) - c$

Regression lineare:  $y_i = \langle x_i, \beta \rangle$

$$\min f(\beta) = \frac{1}{n} \sum_i |y_i - \langle x_i, \beta \rangle|^2 = \frac{1}{n} \|X\beta - y\|^2$$

$$\beta \text{ sol} \Leftrightarrow X^*(X\beta - y) = 0$$

$$\text{in } \ker(X) = \{0\}, \text{ sol} = \beta = (X^*X)^{-1}X^*y$$

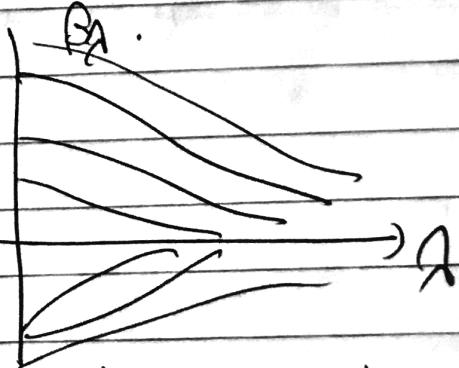
Regul<sup>0</sup>-Quad:  $\min_{\beta} \frac{1}{2} \|X\beta - y\|^2 + \lambda \|\beta\|^2$  (Ridge).

stetig convex

$$\beta_\lambda = (X^*X + \lambda \text{Id})^{-1}X^*y$$

→ LU / Cholesky

→ conjugate gradient



Cross val.:  $\min_{\lambda} \|X\beta_\lambda - \bar{y}\|$

