

# Introduction à l'apprentissage supervisé

Stéphane Canu

[asi.insa-rouen.fr/enseignants/~scanu](http://asi.insa-rouen.fr/enseignants/~scanu)

[scanu@insa-rouen.fr](mailto:scanu@insa-rouen.fr)



École d'été BasMatI 2018

Porquerolles, june 6, 2018

# Plan de l'exposé

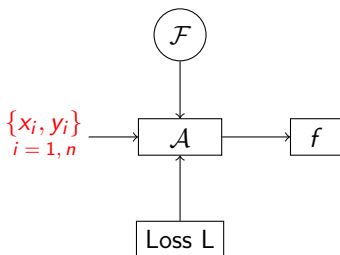
1 Apprendre a partir des données  $\{x_i, y_i\}$

2 le Cout  $L$  : premier algorithme MRE

3 Régularisation et Pénalités : MRS et  $\mathcal{F}$

4 Retour sur l'algorithme de minimisation  $\mathcal{A}$

5 Conclusion et perspectives



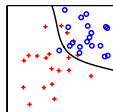
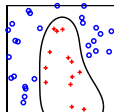
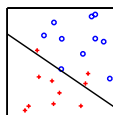
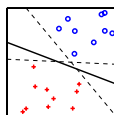
# Classification à 2 classes

données ▶ observation :  $x_i \in \mathbb{R}^p$   
▶ étiquettes  $y_i \in \{-1, 1\}$

individus  $n = 20$

variables  $p = 2$

classes  $k = 2$



Trouver un classifieur ( $f \in \mathcal{F}$ ) telle que

$f(x_i) > 0$  si  $x_i \in \text{Classe1}$  et en même temps  $f(x_i) < 0$  si  $x_i \in \text{Classe2}$

Loss  $L$

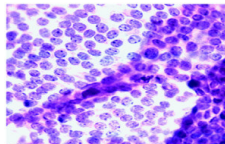
# Wisconsin Diagnostic Breast Cancer (WDBC)

- données
- ▶ observations : images
  - ▶ caractéristiques :  $x_i \in \mathbb{R}^p$
  - ▶ étiquettes  $y_i \in \{-1, 1\}$

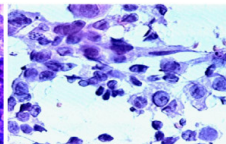
individus  $n = 569$

variables  $p = 30$

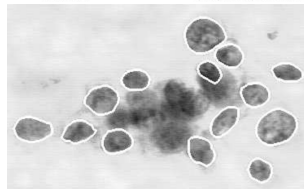
classes  $k = 2$



Smear with BENIGN diagnosis – uniform nucleus of cells, symmetrical, homogeneous, with areas within normal size



Smear with MALIGNANT diagnosis – nucleus of cells without uniformity, asymmetrical, not homogeneous (multiple sizes) and with areas above normal size



## Caractéristiques extraites à la main

Image  $\rightarrow$  caractéristiques  $x \rightarrow$  classe (bénin ou malin)  $y$

<http://archive.ics.uci.edu/ml/machine-learning-databases/breast-cancer-wisconsin/wdbc.names>

<ftp://ftp.cs.wisc.edu/math-prog/cpo-dataset/machine-learn/cancer/>

**Breast cancer diagnosis and prognosis via linear programming., Mangasarian, et al. Operations Research, 1995**

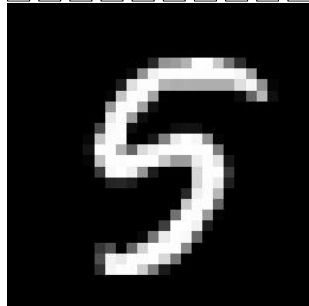
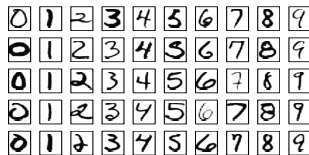
# MNIST

- données
- ▶ observations :  $x_i \in [0, 255]^p$   
256 niveaux de gris
  - ▶ étiquettes  $y_i \in \{0, 1, \dots, 9\}$

individus  $n = 60\,000$

variables  $p = 28 \times 28 = 784$

classes  $k = 10$



Images normalisée

Image (presque brute, juste normalisée)

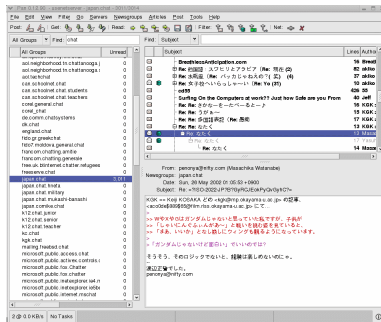
# Newsgroup

données ▶ observations : textes  
▶ étiquettes  $y_i \in \{1; \dots, 20\}$

individus  $n = 11\,314$

variables  $p = 4.8$  million

classes  $k = 20 \rightarrow 2$



Texte: variables = mots

Données creuses (*sparse*)

<https://archive.ics.uci.edu/ml/datasets/Twenty+Newsgroups>

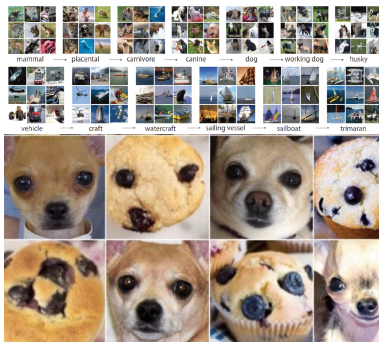
# ImageNet Large Scale Visual Recognition Challenge 2012

données ▶ observations : images  
▶ étiquettes  $y_i \in \{1; \dots, 1000\}$

individus  $n = 1.2$  million

variables  $p = 3 \times 224 \times 224 = 150\,000$

classes  $k = 1000$



Images butes d'un seul objet

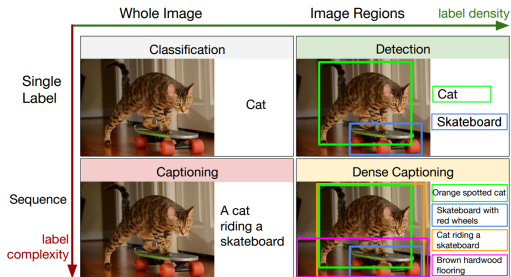
Étiquetées par des humains

# Les données

	jouet	Cancer	MNIST	NewsGroup	ImageNet
$\approx$ chronologie		1970	1990	2000	2010
individus $n =$	20	569	60 000	11 314	1.2 M $\rightarrow$ 14.2 M
variables $p =$	2	30	784	4.8 M	150 000
classes $k =$	2	2	10	20 $\rightarrow$ 2	1000

l'équation des données d'aujourd'hui (2020) :  $n = p = k = \text{millions}$ .

- Visual Genome <http://visualgenome.org/>
- génomique :  $n = 100\,000$  individus  $p = 30$  millions





# Plan de l'exposé

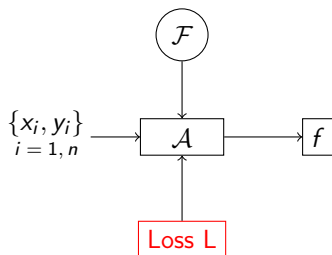
1 Apprendre a partir des données  $\{x_i, y_i\}$

2 le Cout  $L$  : premier algorithme MRE

3 Régularisation et Pénalités : MRS et  $\mathcal{F}$

4 Retour sur l'algorithme de minimisation  $\mathcal{A}$

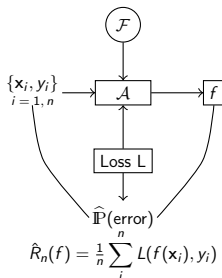
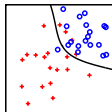
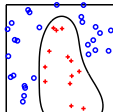
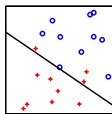
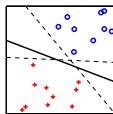
5 Conclusion et perspectives



## Cout et risque

Trouver un algorithme qui minimise le nombre d'erreur (le cout 0/1)

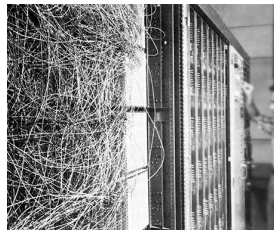
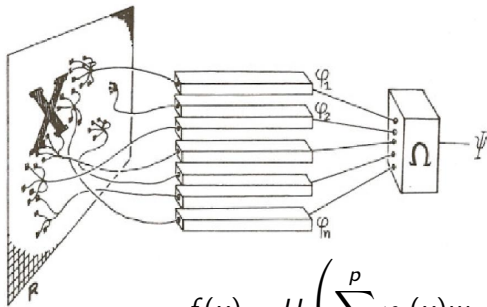
$$L(f(\mathbf{x}_i), y_i) = \begin{cases} 0 & \text{si } f(\mathbf{x}_i) = y_i \\ 1 & \text{sinon} \end{cases}$$



La minimisation du risque empirique  $\hat{R}_n(f)$

$$\min_{f \in \mathcal{F}} \hat{R}_n(f) = \frac{1}{n} \sum_i^n L(f(\mathbf{x}_i), y_i)$$

# Perceptron (Rosenblatt, 1957)



$$f(x) = H \left( \sum_{j=1}^p \varphi_j(x) w_j + w_0 \right)$$

## Règle du Perceptron

Tant que on n'a pas convergé :

- 1 tirer un exemple  $(x_i, y_i)$
- 2 calculer la prédiction  $f(x_i)$
- 3 adapter les poids  
 $w \leftarrow w + \frac{\rho}{2} (y_i - f(x_i)) \varphi(x_i)$

## Du point de vue optimisation

$$\hat{R}_n(w) = \sum_{i=1}^n \max(0, -y_i w^T \varphi(x_i))$$

Méthode de gradient stochastique

## La régression logistique (Cox, 1958)

Le modèle logistique :  $y_i$  réalisation d'une VA  $Y \sim \text{Bernoulli}$  de paramètre

$$p(x_i) = \frac{\exp^{\varphi(x_i)^T w}}{1 + \exp^{\varphi(x_i)^T w}}$$

Le cout logistique : le max de vraisemblance

$$\hat{R}_n(w) = -\log \mathcal{L}(w) = \sum_{i=1}^n -y_i \left( \varphi(x_i)^T w \right) + \log \left( 1 + \exp^{\varphi(x_i)^T w} \right)$$

Cette fonction cout est non linéaire mais

- différentiable
- convexe
- consistante

Du point de vue de l'optimisation

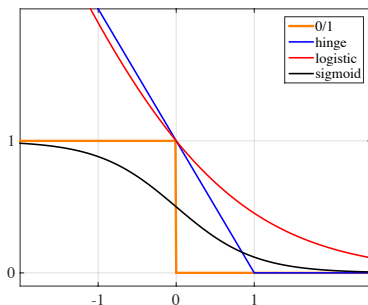
Méthode itérative du second ordre (Newton)

# Exemples de différents coûts utilisés en apprentissage

Principe de minimisation du risque empirique

$$\hat{R}_n(f) = \frac{1}{n} \sum_i^n L(f(x_i), y_i)$$

Loss $L$	Convex	Diff.
logistic $\log(1 + \exp^{f(x_i)})$	✓	✓
sigmoïd $\frac{1}{1 + \exp(y_i f(x_i))}$	✗	✓
hinge $\max(0, 1 - y_i f(x_i))$	✓	✗
cout 0/1 $\begin{cases} 0 & \text{si } f(x_i) = y_i \\ 1 & \text{sinon} \end{cases}$	✗	✗



## Comparaison des méthodes (cas favorable)

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) \quad \hat{f} = \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_n(f)$$

méthodes	Gradient stochastique	1er ordre Gradient	2nd ordre Q Newton
pour une itération	$p$	$np$	$p(n+p)$
nombre d'itération	$\frac{\kappa\gamma}{\varepsilon}$	$\kappa \log \frac{1}{\varepsilon}$	$\log \log \frac{1}{\varepsilon}$
temps pour $\varepsilon$	$\frac{\kappa\gamma p}{\varepsilon}$	$\kappa np \log \frac{1}{\varepsilon}$	$p(n+p) \log \log \frac{1}{\varepsilon}$

$\varepsilon$  : la précision de l'algorithme  $\|f_{\mathcal{A}} - \hat{f}\|^2 \leq \varepsilon$

$\kappa$  : le conditionnement de  $H$  la hessienne

$\gamma = \operatorname{trace}(H^{-1} \nabla \hat{R}_n)$

Les méthodes du second ordre sont les plus efficaces

... si  $L$  est différentiable et si le critère est la précision  $\varepsilon$

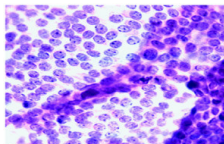
# Wisconsin Diagnostic Breast Cancer (WDBC)

- données
- ▶ observations : images
  - ▶ caractéristiques :  $x_i \in \mathbb{R}^p$
  - ▶ étiquettes  $y_i \in \{-1, 1\}$

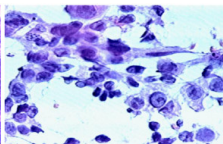
individus  $n = 569$

variables  $p = 30$

classes  $k = 2$



Smear with BENIGN diagnosis – uniform nucleus of cells, symmetrical, homogeneous, with areas within normal size



Smear with MALIGNANT diagnosis – nucleus of cells without uniformity, asymmetrical, not homogeneous (multiple sizes) and with areas above normal size

Toutes les méthodes minimisent l'erreur en moins d'une seconde  
(meilleurs résultats avec des méthodes non linéaires)

Propriété d'universalité de l'apprentissage

Que ce passe t'il lorsque  $p$  est grand ?

# Plan de l'exposé

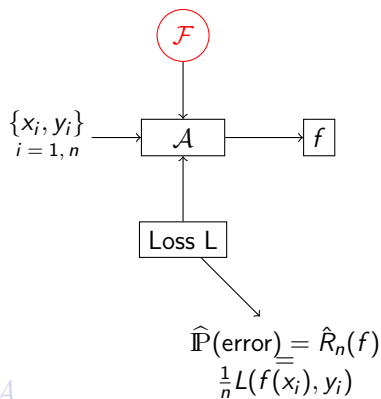
1 Apprendre à partir des données  $\{x_i, y_i\}$

2 le Coût  $L$  : premier algorithme MRE

3 Régularisation et Pénalités : MRS et  $\mathcal{F}$

4 Retour sur l'algorithme de minimisation  $\mathcal{A}$

5 Conclusion et perspectives





## Que se passe t'il lorsque $p$ augmente

Augmenter le nombre de caractéristiques pour obtenir l'**universalité** :

### Les caractéristiques comme une fonction

$p$	$\longrightarrow$	$\infty$	
$\mathbf{x}_i$	$\longrightarrow$	$k(\mathbf{x}_i, \bullet)$	une fonction noyau
$\mathbb{R}^p$	$\longrightarrow$	$\mathcal{H}$	un RKHS
$\mathcal{F}$	$\longrightarrow$	$\mathcal{H}$	
$w^\top \mathbf{x}_i$	$\longrightarrow$	$f(\mathbf{x}_i) = \langle f(\bullet), k(\mathbf{x}_i, \bullet) \rangle_{\mathcal{H}}$	

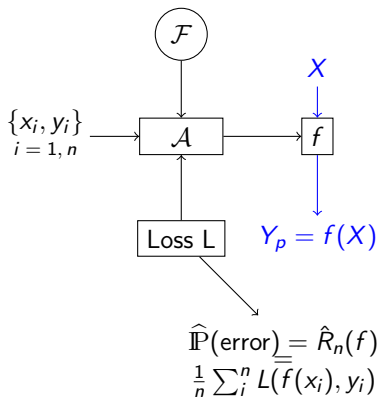
Mécaniquement :

plus  $p$  augmente, plus  $\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i), y_i)$  diminue

Jusqu'à apprendre « par cœur », et mal généraliser :  $\hat{R}_n(f) = 0$

## Apprentissage et généralisation

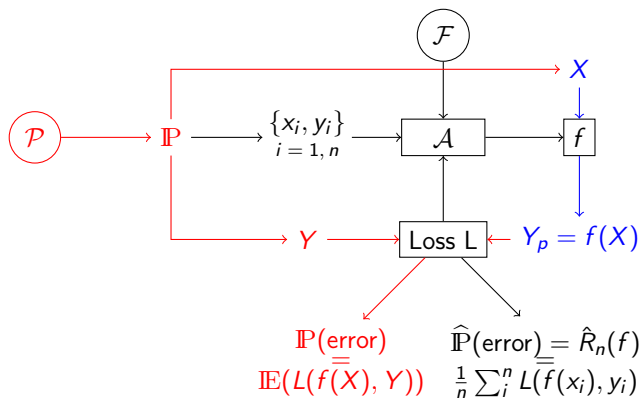
$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$



# Apprentissage et généralisation

$$R(f) = \mathbb{E}(L(f(X), Y))$$

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$



# Théorie de l'apprentissage statistique (COLT)

$$R(f) = \mathbb{E}(L(f(X), Y)) \qquad \hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i)$$

L'idéal :  $f^* = \operatorname{argmin}_f R(f) \Rightarrow R(f^*) \leq R(\hat{f})$   
Ce qu'on a :  $\hat{f} = \operatorname{argmin}_f \hat{R}_n(f) \Rightarrow \hat{R}_n(\hat{f}) \leq \hat{R}_n(f^*)$

$$R(\hat{f}) - R(f^*) < \underbrace{|R(\hat{f}) - \hat{R}_n(\hat{f})|}_{\leq \varepsilon} + \underbrace{|\hat{R}_n(f^*) - R(f^*)|}_{\leq \varepsilon}$$

$$\mathbb{P}\left(\sup_{f \in \mathcal{V}(\mathcal{F})} |R(\hat{f}) - \hat{R}_n(\hat{f})| \leq \varepsilon\right) \geq \delta$$

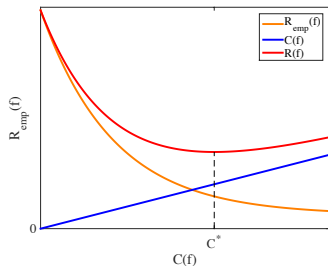
Convergence uniforme localisée (Mendelson, Bousquet, Tsybakov, Massart)

$$\varepsilon(n, C(\mathcal{F})) \propto \mathcal{O}\left(\sqrt{\frac{C(\mathcal{F})}{n} \log\left(\frac{n}{C(\mathcal{F})}\right)}\right)$$

$C(\mathcal{F})$  = mesure de la complexité de  $\mathcal{F}$

## L'apprentissage : un problème d'optimisation bi objectif

$$\left\{ \begin{array}{ll} \min_{f \in \mathcal{F}} \hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) & \text{attache aux données} \\ \min_{f \in \mathcal{F}} C(\mathcal{F}) & \text{mesure de la complexité de } \mathcal{F} \end{array} \right.$$



Le principe de minimisation du risque structurel

## Exemple de minimisation du risque structurel : les SVM

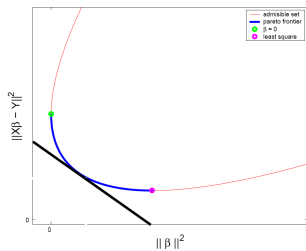
$$C(\mathcal{F}) = \|\mathbf{f}\|^2 \text{ et } \mathcal{F} = \mathcal{H} \text{ un RKHS}$$

L'apprentissage comme un problème d'estimation fonctionnel

$$\left\{ \begin{array}{ll} \min_{f \in \mathcal{H}} \hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i(f(x_i) + b)) & \text{attache aux données} \\ \min_{f \in \mathcal{H}} C(\mathcal{F}) = \|\mathbf{f}\|^2 & \text{mesure de la complexité} \end{array} \right.$$

### 3 formulations equivalentes

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) + \lambda \|f\|^2$$

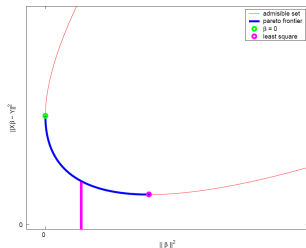


grâce à la convexité

### 3 formulations equivalentes

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) + \lambda \|f\|^2$$

$$\left\{ \begin{array}{l} \min_f \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) \\ \text{avec } \|f\|^2 \leq k \end{array} \right.$$



grâce à la convexité

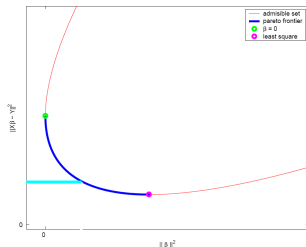


### 3 formulations equivalentes

$$\min_{f \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) + \lambda \|f\|^2$$

$$\left\{ \begin{array}{l} \min_f \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) \\ \text{avec } \|f\|^2 \leq k \end{array} \right.$$

$$\left\{ \begin{array}{l} \min_f \|f\|^2 \\ \text{avec } \frac{1}{n} \sum_{i=1}^n L(f(x_i), y_i) \leq k' \end{array} \right.$$



grâce à la convexité

C'est un programme quadratique (QP)

# Primal Dual

Primal

$$\left\{ \begin{array}{l} \min_{f, b, \xi \in \mathbb{R}^n} \quad \frac{1}{2} \|f\|^2 + C \sum_{i=1}^n \xi_i \\ \text{avec} \quad y_i (f(\mathbf{x}_i) + b) \geq 1 - \xi_i \\ \quad \quad \xi_i \geq 0 \quad i = 1, n \end{array} \right.$$

- QP à  $\#f + n + 1$  inconnues
- et  $2n$  contraintes

Dual

$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{avec} \quad \mathbf{y}^\top \alpha = 0 \\ \text{ert} \quad 0 \leq \alpha_i \leq C \quad i = 1, n \end{array} \right.$$

- QP à  $n$  inconnues
- $G$  matrice de Gram  $\mathcal{O}(pn^2)$
- $2n$  contraintes de boîte
- complexité empirique  $\sim \mathcal{O}(n^{1.5})$

$$\sum_{j=1}^d f(\mathbf{x}_j) + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

# Primal Dual

Primal

$$\left\{ \begin{array}{l} \min_{f, b, \xi \in \mathbb{R}^n} \quad \frac{1}{2} \|f\|^2 + C \sum_{i=1}^n \xi_i \\ \text{avec} \quad y_i (f(\mathbf{x}_i) + b) \geq 1 - \xi_i \\ \quad \quad \xi_i \geq 0 \quad i = 1, n \end{array} \right.$$

- QP à  $\#f + n + 1$  inconnues
- et  $2n$  contraintes

Dual

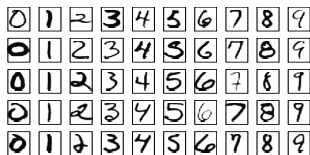
$$\left\{ \begin{array}{l} \min_{\alpha \in \mathbb{R}^n} \quad \frac{1}{2} \alpha^\top G \alpha - \mathbf{e}^\top \alpha \\ \text{avec} \quad \mathbf{y}^\top \alpha = 0 \\ \text{ert} \quad 0 \leq \alpha_i \leq C \quad i = 1, n \end{array} \right.$$

- QP à  $n$  inconnues
- $G$  matrice de Gram  $\mathcal{O}(pn^2)$
- $2n$  contraintes de boîte
- complexité empirique  $\sim \mathcal{O}(n^{1.5})$

$$\sum_{j=1}^d f(\mathbf{x}_j) + b = \sum_{i=1}^n \alpha_i y_i (\mathbf{x}^\top \mathbf{x}_i) + b$$

## Les résultats sur MNIST

- MNIST<sup>1</sup>, data = « image-label »
- $n = 60\,000$ ;  $p = 700$ ;  $k$  classes = 10
- Kernel error rate (2002) = 0.56 %,
- Best error rate (2012) = 0.23 % .



Les SVM sélectionnent les individus mais :

$n$  ne doit pas être trop grand

$p$  grand et  $n$  grand, donc on repasse dans le primal

<sup>1</sup><http://yann.lecun.com/exdb/mnist/index.html>

## Une autre approche quand $p$ est grand

Retour au modèle linéaire :

$$f(\mathbf{x}_i) = \mathbf{w}^\top \mathbf{x}_i$$

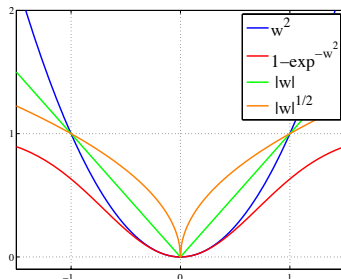
Dans le primal : sélection de variable

$$\left\{ \begin{array}{ll} \min_{f \in \mathbb{R}^p} & \hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{w}^\top \mathbf{x}_i, y_i) \quad \text{attache aux données} \\ & C(\mathbb{R}^p) = \|\mathbf{w}\|_0 \leq k \quad \text{nombre de variables } k < p \end{array} \right.$$

Comme pour le cout 0/1, on utilise une autre pénalité (relaxation)

# Exemples de différentes pénalités utilisées en apprentissage

Pénalité	Convexe	Diff.
$\omega(w) = w^2$ $\Omega(\mathbf{x}) = \ \mathbf{x}\ ^2$	✓	✓
$\omega(w) = 1 - \exp(-w^2)$	✗	✓
$\omega(w) =  w $ $\Omega(\mathbf{x}) = \ \mathbf{x}\ _1$	✓	✗
$\omega(w) = \sqrt{ w }$ $\Omega(\mathbf{x}) = \ \mathbf{x}\ _{1/2}$	✗	✗



# Gradient proximal pour gérer la non différentiabilité

selon la différentiabilité des critères

$$\min_f L(f(x_i), y_y) + \lambda\Omega(f)$$

$L$  et  $\Omega$  différentiables  $\rightarrow$  gradient

$$f^{(t+1)} = f^{(t)} - \rho \nabla (L(f(x_i), y_y) + \lambda\Omega(f))$$

$\Omega$  non différentiable  $\rightarrow$  gradient proximal une méthode du premier ordre pour résoudre certains programme quadratiques

$$\begin{aligned} g^{(t+1)} &= f^{(t+1)} - \rho \nabla (L(f(x_i), y_y)) \\ f^{(t+1)} &= \text{prox}_{\rho\lambda\Omega}(g^{(t+1)}) \end{aligned}$$

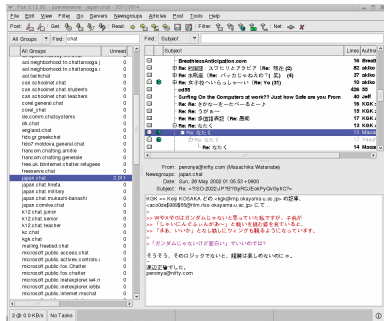
# Newsgroup

données ▶ observations : textes  
▶ étiquettes  $y_i \in \{1; \dots, 20\}$

individus  $n = 11\,314$

variables  $p = 4.8$  million

classes  $k = 20 \rightarrow 2$



modèle linéaire, cout logistique, pénalité L1

- computing time (2010) = 2 minutes
- Best error rate = 2.73 % .

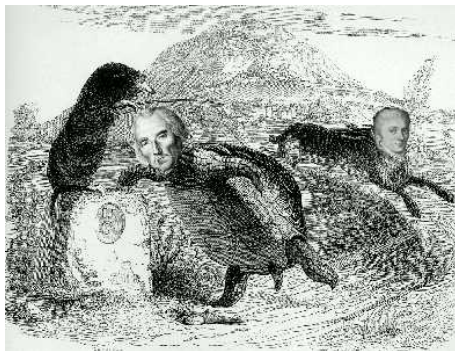
An interior-point method for large-scale l1-regularized logistic regression K Koh, SJ Kim, S Boyd - JMLR, 2007



# Le lièvre gaussien et la tortue laplacienne

Loi de Gauss  $\|w\|_2$

Loi de Laplace  $\|w\|_1$



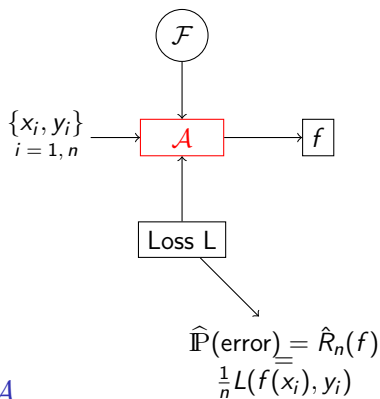
Gauss  $p^3$  vs Laplace (LP, QP)  $n^2 \times k^2$ ,  $k$  nombre de variables utiles

**Problème : comment gérer**

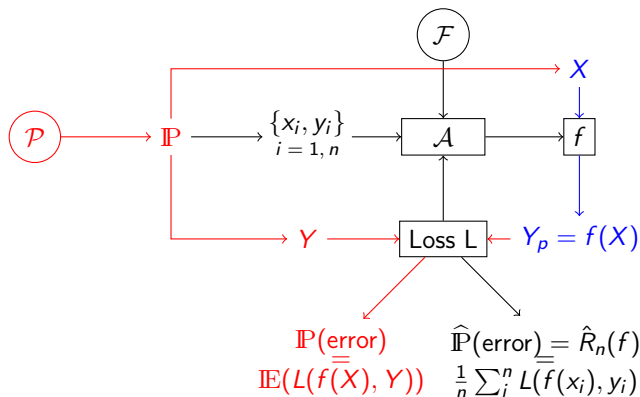
pour  $n$  et  $p$  grand.

# Plan de l'exposé

- 1 Apprendre a partir des données  $\{x_i, y_i\}_{i=1, n}$
- 2 le Cout  $L$  : premier algorithme MRE
- 3 Régularisation et Pénalités : MRS et  $\mathcal{F}$
- 4 Retour sur l'algorithme de minimisation  $\mathcal{A}$
- 5 Conclusion et perspectives



Quel algorithme choisir lorsque  $n$  et  $p$  sont grands



$$\mathbb{P}(|R(\hat{f}) - R(f^*)| \leq \varepsilon) \geq \delta$$

L'erreur statistique est de l'ordre de  $\varepsilon(n) \propto \mathcal{O}\sqrt{\frac{1}{n}} \Rightarrow n \propto \frac{1}{\varepsilon^2}$

## Quel algorithme choisir lorsque $n$ et $p$ sont grands

L'erreur statistique est de l'ordre de  $n \propto \frac{1}{\varepsilon^2}$

Principe : l'erreur statistique  $\approx$  l'erreur algorithmique

méthodes	Gradient stochastique	1er ordre Gradient	2nd ordre Q Newton	QP dual
pour une itération	$p$	$np$	$p(n + p)$	
nombre d'itération	$\frac{\kappa\gamma}{\varepsilon}$	$\kappa \log \frac{1}{\varepsilon}$	$\log \log \frac{1}{\varepsilon}$	
temps pour $\varepsilon$	$\frac{\kappa\gamma p}{\varepsilon}$	$\kappa np \log \frac{1}{\varepsilon}$	$pn \log \log \frac{1}{\varepsilon}$	$n^2$
temps pour $\varepsilon$				

## Quel algorithme choisir lorsque $n$ et $p$ sont grands

L'erreur statistique est de l'ordre de  $n \propto \frac{1}{\varepsilon^2}$

Principe : l'erreur statistique  $\approx$  l'erreur algorithmique

méthodes	Gradient stochastique	1er ordre Gradient	2nd ordre Q Newton	QP dual
pour une itération	$p$	$np$	$p(n + p)$	
nombre d'itération	$\frac{\kappa\gamma}{\varepsilon}$	$\kappa \log \frac{1}{\varepsilon}$	$\log \log \frac{1}{\varepsilon}$	
temps pour $\varepsilon$	$\frac{\kappa\gamma p}{\varepsilon}$	$\kappa np \log \frac{1}{\varepsilon}$	$pn \log \log \frac{1}{\varepsilon}$	$n^2$
temps pour $\varepsilon$	$\frac{\kappa\gamma p}{\varepsilon}$	$\kappa \frac{p}{\varepsilon^2} \log \frac{1}{\varepsilon}$	$\frac{p}{\varepsilon^2} \log \log \frac{1}{\varepsilon}$	$\frac{1}{\varepsilon^4}$

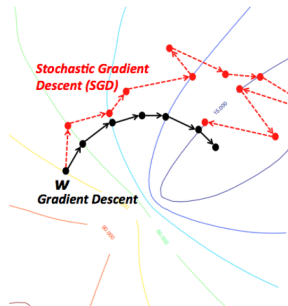
L'algorithme le plus rapide est

le gradient stochastique

Pour une comparaison détaillée voir par exemple <http://leon.bottou.org/projects/sgd>

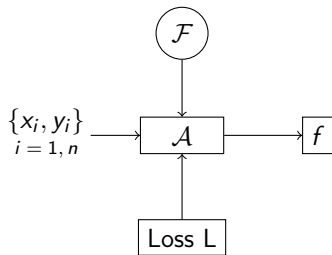
# Le beurre et l'argent du beurre

- Accélération de la convergence
  - ▶ méthode de Nesterov (1983)
  - ▶ momentum (heuristique)
- Moyenner
  - ▶ les exemples (mini batch)
  - ▶ les paramètres (Polyak and Juditsky, 1992)
  - ▶ les gradients (stochastic average gradient (SAG-A), Le Roux et al 2012)
  - ▶ réduction de la variance (Johnson, Zhang, 13)
- Moyenner et accélérer
  - ▶ (Dieuleveut, Flammarion & Bach, 2016)
- Adapter les pas
  - ▶ Adaptive Moment Estimation – ADAM (Kingma & Ba, 2015)



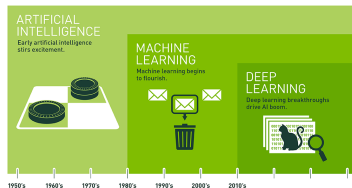
## Revenir au non linéaire avec $n$ et $p$ grands

- cout  $L$
- regulariser  $C(\mathcal{F})$
- gradient stochastique (++)  $\mathcal{A}$
- Modèle non linéaire à taille fixe : apprendre les caractéristiques  $f$



# Plan de l'exposé

- 1 Apprendre a partir des données  $\{x_i, y_i\}$
- 2 le Cout  $L$  : premier algorithme MRE
- 3 Régularisation et Pénalités : MRS et  $\mathcal{F}$
- 4 Retour sur l'algorithme de minimisation  $\mathcal{A}$
- 5 Conclusion et perspectives



Since an early flash of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.



## Conclusion

- la taille des données explose ( $n$  exemples et  $p$  variables)

$n$  et  $p$  petits régression logistique et Newton

$p$  grand non linéaire – régularisation SVM (QP dual)

$p$  grand linéaire – régression logistique sparse (L1 proximal)

$n$  grand (et  $p$  grand) approximation stochastique SGD

→ deep learning (cours de M. Cord)

→ sélectionner les variables et les individus

- données hétérogènes : forêts aléatoires