

# MAP 569

## Machine Learning II

Christophe Giraud

PC2

Supervised Classification

# Supervised Classification

# Alternative modeling

## Parametric modeling

Modeling of the  $X_i$  by a mixture of Gaussians  $\implies$  LDA.

## Semi-parametric modeling

Modeling of the distribution on  $Y$  given  $X \implies$  logistic regression.

## Non-parametric modeling

For a given set  $\mathcal{H}$  of classifiers, take the empirical risk minimizer

$$\hat{h}_{\mathcal{H}} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h), \quad \text{where} \quad \hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{y_i \neq h(x_i)}$$

## Empirical risk minimizer

For some observations  $(x_i, y_i)_{i=1, \dots, n}$  and a set  $\mathcal{H}$  of classifiers

$$\hat{h}_{\mathcal{H}} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h), \quad \text{where} \quad \hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\mathbb{R}^+}(-y_i h(x_i))$$

## In practice

- 1  $\mathcal{H}$  non convex,
- 2  $\hat{R}_n(h)$  non convex.

Prohibitive computational complexity !

## Empirical risk minimizer

For some observations  $(x_i, y_i)_{i=1, \dots, n}$  and a set  $\mathcal{H}$  of classifiers

$$\hat{h}_{\mathcal{H}} = \operatorname{argmin}_{h \in \mathcal{H}} \hat{R}_n(h), \quad \text{where} \quad \hat{R}_n(h) = \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\mathbb{R}^+}(-y_i h(x_i))$$

## In practice

- ❶  $\mathcal{H}$  non convex,
- ❷  $\hat{R}_n(h)$  non convex.

**Prohibitive computational complexity !**

## Two issues

- 1  $\mathcal{H}$  non convex,
- 2  $\hat{R}_n(h)$  non convex.

## Convexification

For

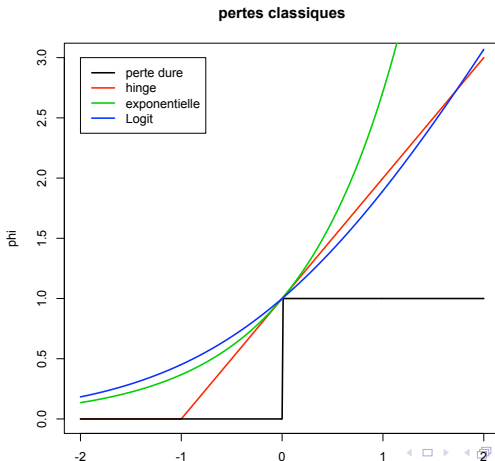
- $\mathcal{F}$  a convex set of functions  $f : \mathcal{X} \rightarrow \mathbb{R}$
- and  $\varphi : \mathbb{R} \rightarrow \mathbb{R}^+$  convex and non-decreasing

we define

$$\hat{h}_{\varphi, \mathcal{F}} = \text{sign}(\hat{f}_{\varphi, \mathcal{F}}) \quad \text{with} \quad \hat{f}_{\varphi, \mathcal{F}} = \underset{f \in \mathcal{F}}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n \varphi(-y_i f(x_i))$$

## Some popular $\varphi$

- **Hinge loss** :  $\varphi(x) = (1 + x)_+$
- **Exponential loss** :  $\varphi(x) = e^x$
- **Logit loss** :  $\varphi(x) = \log_2(1 + e^x)$



## Some popular $\mathcal{F}$

- **Linear classifier** :  $\mathcal{F} = \{\langle w, \cdot \rangle : \|w\| \leq R\}$  (exercise 1)
- **Convex hull of some basic classifiers**  $\{h_1, \dots, h_M\}$  :

$$\mathcal{F} = \left\{ f = \sum_{j=1}^M \theta_j h_j : \theta \in \Theta \right\}$$

with  $\Theta$  a convex subset of  $\mathbb{R}^M$ . (later)

- **Ball of a RKHS  $\mathcal{W}$**  : for  $R > 0$

$$\mathcal{F} = \{f \in \mathcal{W} : |f|_{\mathcal{W}} \leq R\}.$$

(next week)



# Support Vector Machine

## SVM

SVM corresponds to

- $\varphi(x) = (1 + x)_+$
- $\mathcal{F} = \{\langle w, \cdot \rangle : \|w\| \leq R\}, \text{ with } R > 0.$

## SVM : Lagrangian version

The classifier  $\hat{h}_{\varphi, \mathcal{F}}$  is defined by  $\hat{h}_{\varphi, \mathcal{F}}(x) = \text{sign}(\langle \hat{w}, x \rangle)$  with

$$\hat{w} = \underset{w \in \mathbb{R}^p}{\operatorname{argmin}} \left\{ \frac{1}{n} \sum_{i=1}^n (1 - y_i \langle w, x_i \rangle)_+ + \lambda \|w\|^2 \right\}$$

## Reminder on convex optimization

Let  $f, -g_1, \dots, -g_n$  be  $\mathcal{C}^1$  convex functions and

$$\hat{x} = \underset{g_i(x) \geq 0}{\operatorname{argmin}} f(x).$$

### Karush-Kuhn-Tucker necessary conditions

Set

$$L(x, \lambda) = f(x) - \sum_{i=1}^n \lambda_i g_i(x).$$

There exists  $\hat{\lambda}$  such that

- ❶  $\nabla_x L(\hat{x}, \hat{\lambda}) = 0$
- ❷  $\min(\hat{\lambda}_i, g_i(\hat{x})) = 0$  for  $i = 1, \dots, n$

### Strong duality

$$\hat{\lambda} = \underset{\lambda \geq 0}{\operatorname{argsup}} \inf_x L(x, \lambda)$$

# Geometric interpretation

We have shown that

$$\hat{f}_{\varphi, \mathcal{F}}(x) = \langle \hat{w}, x \rangle, \text{ with } \hat{w} = \sum_{i=1}^n \hat{\beta}_i x_i$$

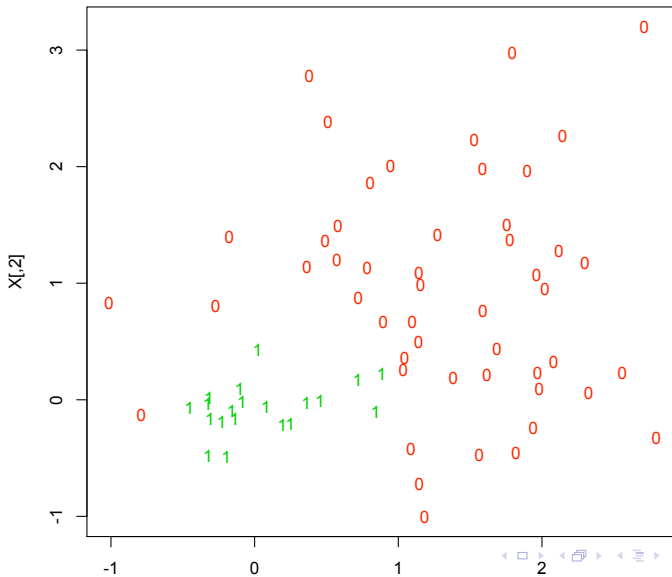
where

## KKT

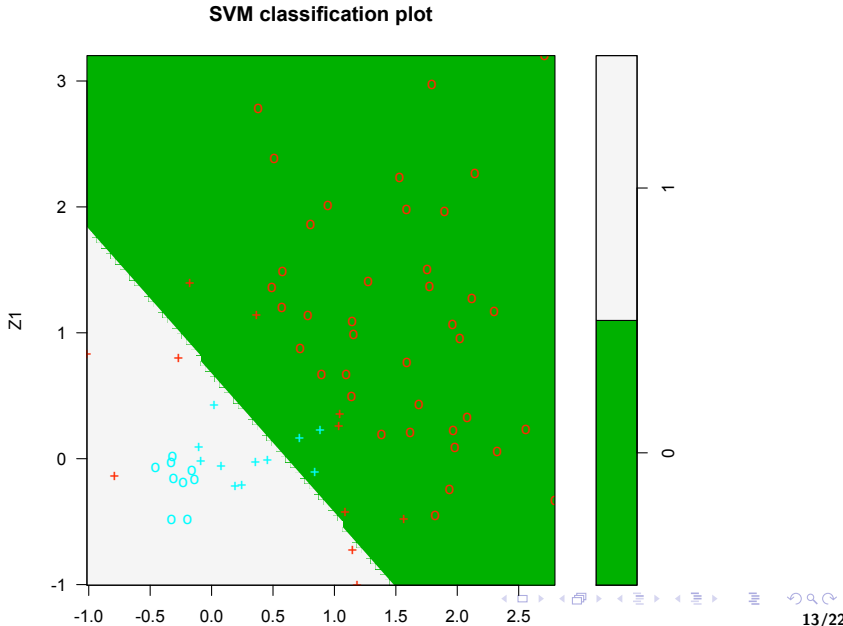
- if  $y_i \hat{f}(x_i) > 1$  then  $\hat{\beta}_i = 0$ ,
- if  $y_i \hat{f}(x_i) < 1$  then  $\hat{\beta}_i = y_i / (2\lambda n)$ ,
- if  $y_i \hat{f}(x_i) = 1$  then  $0 \leq \hat{\beta}_i y_i \leq 1 / (2\lambda n)$ ,

Geometric interpretation ?

Data :



les points "+" sont les vecteurs supports



# Strong duality

$$\begin{aligned}(\hat{\alpha}, \hat{\gamma}) &\in \operatorname{argmax}_{(\alpha, \gamma) \geq 0} \min_{\beta, \xi} \left\{ \lambda \langle K\beta, \beta \rangle - \langle K\beta, y.\alpha \rangle + \langle \alpha, 1 \rangle + \langle \xi, \frac{1}{n} - \alpha - \gamma \rangle \right\} \\&\in \operatorname{argmax}_{(\alpha, \gamma) \geq 0} \min_{\xi} \left\{ -\frac{1}{4\lambda} \langle K(y.\alpha), y.\alpha \rangle + \langle \alpha, 1 \rangle + \langle \xi, \frac{1}{n} - \alpha - \gamma \rangle \right\} \\&\in \operatorname{argmax}_{(\alpha, \gamma) \geq 0 \ \& \ \alpha + \gamma = \frac{1}{n}} \left\{ -\frac{1}{4\lambda} \langle K(y.\alpha), y.\alpha \rangle + \langle \alpha, 1 \rangle \right\} \\&\in \operatorname{argmax}_{0 \leq \alpha \leq \frac{1}{n} \ \& \ \gamma = \frac{1}{n} - \alpha} \left\{ -\frac{1}{4\lambda} \langle K(y.\alpha), y.\alpha \rangle + \langle \alpha, 1 \rangle \right\}\end{aligned}$$

# Selection of $\lambda$

## V-fold Cross-Validation

**Recipe :** split the data into  $V$  groups

- 1 learn  $\hat{h}_\lambda$  on  $V - 1$  "training" groups
- 2 test  $\hat{h}_\lambda$  on the remaining "test" group
- 3 iterate by permuting the "train" and "test" groups
- 4 keep  $\hat{h}_\lambda$  with the smallest average misclassification error on the  $V$  tests.

## Example : 5-fold CV

train	train	train	train	test
train	train	train	test	train
train	train	test	train	train
train	test	train	train	train
test	train	train	train	train



# Elastic Net

# Regression setting

## Linear model

$$y_i = \langle \beta^*, x_i \rangle + \epsilon_i \quad i=1, \dots, n$$

## Vectorial writing

$$Y = \mathbf{X}\beta^* + \epsilon$$

# Least squares

## Least-squares

$$\hat{\beta}^{LS} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \|Y - \mathbf{X}\beta\|^2$$

## Issues

- ❶ no unique solution if  $p > n$
- ❷ if  $\operatorname{cov}(\epsilon) = \sigma^2 I_n$  then the average error

$$\mathbb{E}[\|\hat{\beta}^{LS} - \beta^*\|^2] = \operatorname{Tr}((\mathbf{X}^T \mathbf{X})^{-1}) \sigma^2$$

can be huge. 😞

# Sparse regression

## Sparse regression paradigm

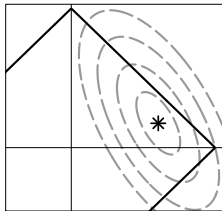
Only a few features matters :  $\beta^*$  is sparse or is close to a sparse vector.

## Lasso

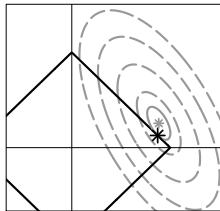
$$\hat{\beta}_\mu \in \underset{\beta \in \mathbb{R}^p}{\operatorname{argmin}} \mathcal{L}(\beta) \quad \text{with} \quad \mathcal{L}(\beta) = \|Y - \mathbf{X}\beta\|^2 + \mu|\beta|_{\ell^1}.$$

# Singularities of the $\ell^1$ -ball induce feature selection

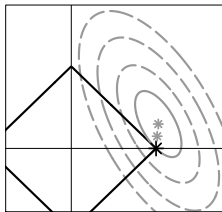
R=2



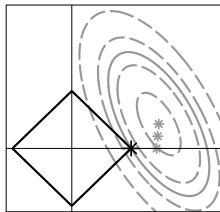
R= 1.4



R= 1.2



R= 0.82



# Elastic net

## Issue

If two (or more) important features are strongly correlated, then only one of the two will be selected. 😞

## Elastic Net

Recipe : add a  $\ell^2$  penalty

$$\hat{\beta}_{\lambda,\mu} \in \operatorname{argmin}_{\beta \in \mathbb{R}^p} \mathcal{L}(\beta) \quad \text{with} \quad \mathcal{L}(\beta) = \|Y - \mathbf{X}\beta\|^2 + \lambda\|\beta\|^2 + \mu|\beta|_{\ell^1}.$$

# Illustration

Intermediate between Ridge and Lasso

