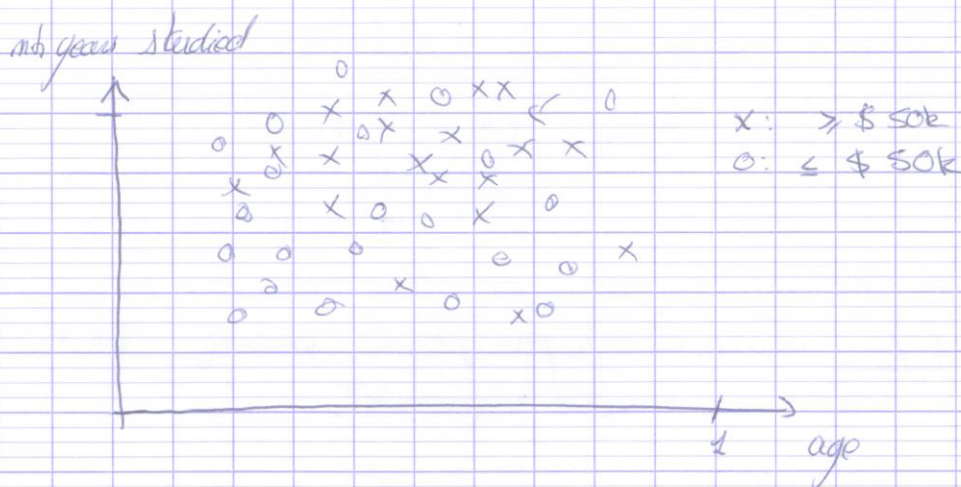


# Lecture 5: Learning and optimization



$$X = [0, 15]^2 ; Y = \{-1, +1\}$$

Aim: find a good classificat<sup>o</sup> rule  $h: X \rightarrow Y$

Hypothesis: data comes from a probability law  $P_{(X,Y)}$

Risk:  $R(h) = P_{(X,Y)}(h(X) \neq Y)$

$\rightarrow h^* = \underset{h}{\operatorname{argmin}} R(h)$  "Bayes classifier"

$$h^*(x) = \begin{cases} +1 & \text{if } P(Y=+1 | X=x) \geq \frac{1}{2} \\ -1 & \text{otherwise} \end{cases}$$

you don't know  $P_{(X,Y)}$  but you have  $(X_1, Y_1), \dots, (X_m, Y_m)$

$$\hat{h}_m = \underset{h}{\operatorname{argmin}} \hat{R}_m(h) \quad \text{where } \hat{R}_m(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{h(X_i) \neq Y_i\}}$$

(•  $h$ -nearest-neighbors)  $\mathbb{E}[\hat{R}_m(h)] = R(h)$

① Assume that the classes are linearly separable  
 $X = \mathbb{R}^p$ ,  $H = \{ \operatorname{sgn}(w^T x) ; w \in \mathbb{R}^p \}$

$$\ln(1 + e^{-ywx}) = \frac{-ywx e^{-ywx}}{1 + e^{-ywx}} = \frac{-ywx}{e^{ywx} + 1}$$

$$\omega \in \mathbb{R}^p, \quad R(\omega) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y_i \omega^T x_i \leq 0\}$$

Goal:  $\hat{\omega}_m \in \text{argmin}_{\omega} R(\omega)$

pb:  $\min_{\omega \in \mathbb{R}^p} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{y_i \omega^T x_i \leq 0\}$   $R_R$

$g_R(\omega)$  not diff. (unsmooth), not convex.

a) Algo: Rosenblatt's perceptron (Book of Mark, Friedman, Tibshirani)  
 → good for linear separability.

b) Logistic regression (convexificat° of the Risk!)

Logistic →  $P(Y=1 | X=x) = \text{logit}^{-1}(\omega^T x)$

$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right), \quad \text{logit}^{-1}(u) = \frac{e^u}{1+e^u}$

$$P(Y=-1 | X=x) = 1 - \frac{e^{\omega^T x}}{1+e^{\omega^T x}} = \frac{1}{1+e^{\omega^T x}}$$

$$P(Y=1 | X=x) = \frac{1}{1+e^{-\omega^T x}}$$

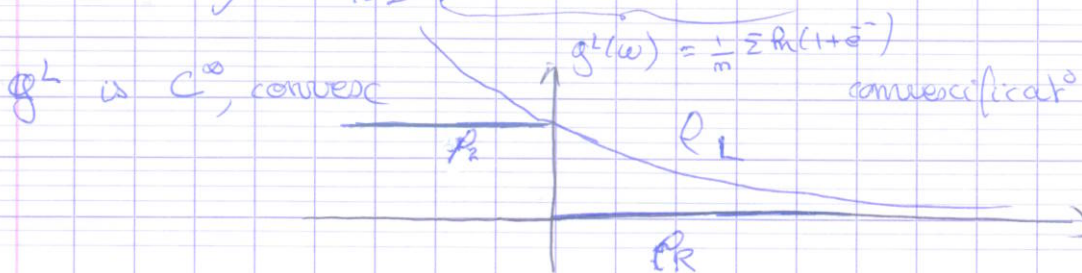
$$g \in \{-1, +1\}, \quad P(Y=g | X=x) = \frac{1}{1+e^{-g\omega^T x}}$$

ML estimator of  $\omega$ :

$$\hat{\omega}_m^L = \underset{\omega}{\text{argmax}} \prod_{i=1}^m P(Y=y_i | X=x_i)$$

$$= \underset{\omega}{\text{argmax}} \sum_{i=1}^m \log\left(\frac{1}{1+e^{-g_i \omega^T x_i}}\right)$$

$$= \underset{\omega}{\text{argmin}} \sum_{i=1}^m \underbrace{\ln(1+e^{-g_i \omega^T x_i})}_{\ell^L}$$



$$\nabla g^2(w) = \sum_{i=1}^n \frac{-x_i y_i}{1 + e^{g(w)x_i}}$$

## ② Convex optimization

Reference: Lectures notes by Sebastian Bubeck (Princeton)

• Tutorial Part 3 CML from Boston, Noredal.  
Curtis

GD:  $w_0 \in \mathbb{R}^p$ ,  $w_{k+1} = w_k - \gamma_k \nabla \ell(w_k)$   
where  $\ell: \mathbb{R}^p \rightarrow \mathbb{R}$  we for to be minimized

Example:  $\ell(w) = \frac{1}{2} \|w\|_2^2$

Q: for what values of  $\gamma_k = \gamma$  does the GD  $w_k$  to 0? What is the optimal value for  $\gamma$ ?

$$\nabla \ell(w) = w$$

$$w_0 \in \mathbb{R}^p \quad w_{k+1} = w_k - \gamma w_k = (1 - \gamma) w_k$$

$$\text{i.e. } w_k = (1 - \gamma)^k w_0$$

$$\text{if } \gamma > \frac{2}{L}, \quad \|w_k\|_2 \rightarrow \infty$$

$$\text{if } \gamma = \frac{2}{L}, \quad \|w_k\|_2 \rightarrow 0$$

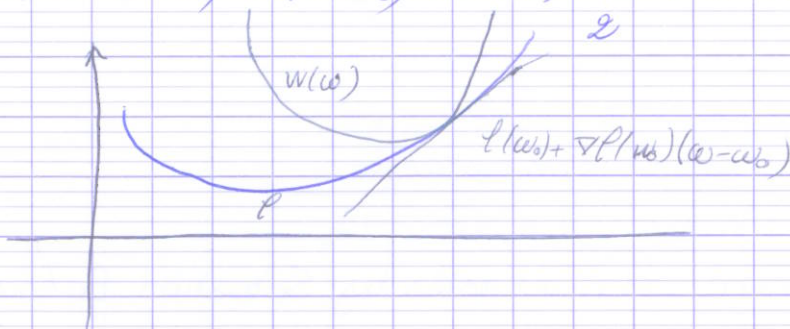
$$\text{if } 0 < \gamma < \frac{2}{L}, \quad w_k \rightarrow 0 \quad \text{as } k \rightarrow \infty, \quad \text{the best } \gamma \text{ is } \gamma = \frac{1}{L}$$

Here  $w \mapsto \frac{L}{2} \|w\|_2^2$  has an  $L$ -Lipschitz gradient

In the sequel, we assume that function  $\ell$  has a  $L$ -Lipschitz gradient  $\|\nabla \ell(x) - \nabla \ell(y)\| \leq L \|x - y\|$  (for example, true with  $g^2$  with  $L = \frac{1}{4}$ )

Consequence:  $\forall \omega_0 \in \mathbb{R}^p, \forall \omega \in \mathbb{R}^p$

$$l(\omega) \leq l(\omega_0) + \nabla l(\omega_0)^T (\omega - \omega_0) + \frac{L}{2} \|\omega - \omega_0\|^2$$



### Maximizat° - Minimizat°

- $\omega_0 \in \mathbb{R}^p$

- $g_k: \mathbb{R}^p \rightarrow \mathbb{R}$

$$g_k(\omega_k) = l(\omega_k)$$

$$\nabla g_k(\omega_k) = \nabla l(\omega_k)$$

$$\forall \omega, g_k(\omega) \geq l(\omega)$$

- $\omega_{k+1} = \underset{\omega \in \mathbb{R}^p}{\operatorname{argmin}} g_k$

If you choose  $g_k(\omega) = l(\omega_k) + \nabla l(\omega_k)^T (\omega - \omega_k) + \frac{L}{2} \|\omega - \omega_k\|^2$

$$\nabla g_k(\omega) = \nabla l(\omega_k) + L(\omega - \omega_k) = 0$$

$$\Leftrightarrow \omega_{k+1} = \omega_k - \frac{1}{L} \nabla l(\omega_k)$$

[Proposit° 2.2 in Optimizat° with 1<sup>st</sup> order Surrogate Fets by Daniel]

[Nesterov - complex optimizat°]

Assume that  $l$  is convex, and that for some  $R > 0$ ,

$$\|\omega - \omega^*\|_2 \leq R$$

for all  $\omega \in \mathbb{R}^p$  st  $l(\omega) \leq l(\omega_0)$ ,  $\omega_k = \underset{\omega \in \mathbb{R}^p}{\operatorname{argmin}} l$

$$\text{Then } l(\omega_k) - l(\omega^*) \leq \frac{2LR^2}{k+2} \quad \forall k \geq 1$$

Proof: [Nesterov] Define  $h_k = g_k - \ell$

Lemma: Then  $\ell(w_k) \leq \min_{w \in \mathbb{R}^D} \left\{ \ell(w) + \frac{L}{2} \|w - w_{k-1}\|^2 \right\}$

$$\ell(w_k) \leq \min_{\alpha \in [0,1]} \ell(\alpha w^* + (1-\alpha)w_{k-1}) + \frac{L}{2} \alpha^2 \|w^* - w_{k-1}\|^2$$

convexity of  $\ell$   $\leq \min_{\alpha \in [0,1]} \alpha \ell(w^*) + (1-\alpha) \ell(w_{k-1}) + \frac{L\alpha^2}{2} \|w^* - w_{k-1}\|^2$

But (lemma) the sequence  $\ell(w_k)$  is decreasing in  $k$ , and thus:

$$\ell(w_k) - \ell(w^*) \leq \min_{0 \leq \alpha \leq 1} (1-\alpha)(\ell(w_{k-1}) - \ell(w^*)) + \frac{L\alpha^2}{2} R^2$$

2 cases:

• if  $\ell(w_{k-1}) - \ell(w^*) \geq LR^2$

then choosing  $\alpha = 1$  yields  $\ell(w_k) - \ell(w^*) \leq \frac{LR^2}{2}$

• otherwise take  $\alpha = \frac{\ell(w_{k-1}) - \ell(w^*)}{LR^2}$

and denote  $r_k = \ell(w_k) - \ell(w^*)$

Then:

$$r_k \leq r_{k-1} \left( 1 - \frac{r_{k-1}}{2LR^2} \right)$$

Back to upper bound

$$\frac{1}{r_k} \geq \frac{1}{r_{k-1}} \cdot \frac{1}{1 - \frac{r_{k-1}}{2LR^2}}$$

$$\geq \frac{1}{r_{k-1}} \left( 1 + \frac{r_{k-1}}{2LR^2} \right)$$

$$= \frac{1}{r_{k-1}} + \frac{1}{2LR^2}$$

$$\geq \frac{1}{r_1} + \frac{k-1}{2LR^2}$$

$$r_k - r_{k-1} \leq -\frac{r_{k-1}^2}{2LR^2}$$

$$y' = -\frac{y^2}{2LR^2}$$

$$\rightarrow y = \frac{2LR^2}{2k-1}$$

$\rightarrow y = \frac{2LR^2}{2k-1}$  linearizes at  $r_k := \ell(w_k) - \ell(w^*)$  has limit (Cauchy).

By case 1:  $R_1 \leq \frac{LR^2}{2}$

$$\frac{1}{R_k} \geq \frac{2}{LR^2} \cdot \frac{k-1}{2LR^2} = \frac{k+3}{2LR^2}$$

$$\Leftrightarrow R_k \leq \frac{2LR^2}{k+3}$$

□

Lemma 1:  $g_k \geq f$   $h_k = g_k - f$

$$\bullet |h_k(w)| \leq \frac{L}{2} \|w - w_k\|^2$$

$$\bullet f(w') \leq f(w) + \frac{L}{2} \|w - w_k\|^2$$

proof:

$$f(w_k) \leq g_{k-1}(w_k) \leq g_{k-1}(w) = f(w) + h_{k-1}(w)$$

$$g_{k-1}(w) = f(w_{k-1}) + \nabla f(w_{k-1})(w - w_{k-1}) + \frac{L}{2} \|w - w_{k-1}\|^2$$

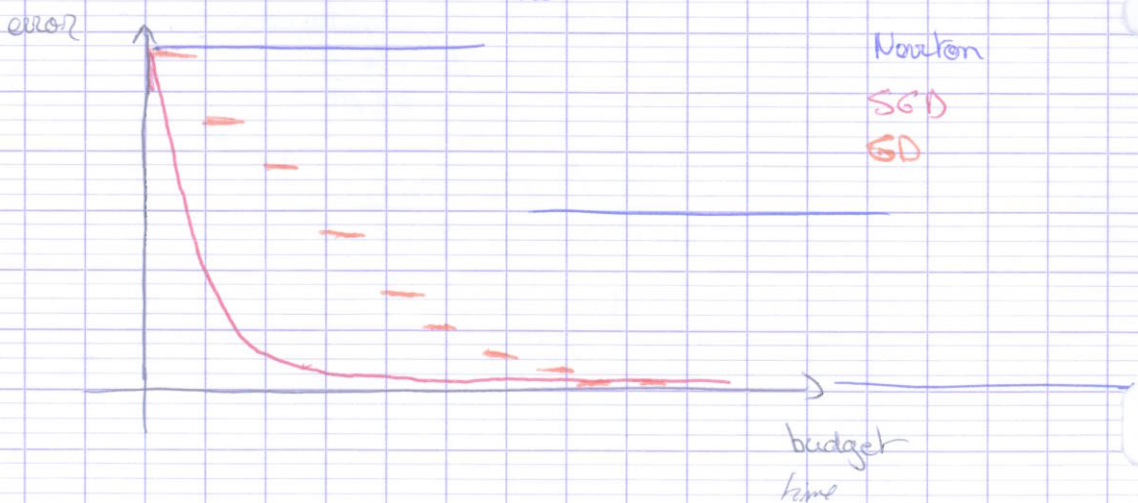
$$\leq f(w) + \frac{L}{2} \|w - w_{k-1}\|^2$$

Lemma 2.2: The sequence  $(f(w_k))_k$  is mon-increasing

proof  $f(w_k) \leq g_k(w_k)$

$$\leq g_k(w_{k-1}) = f(w_{k-1})$$

□



$$|h_k(w)| = |g_k(w) - f(w)|$$

$$= |f(w_k) + \nabla f(w_k)(w - w_k) + \frac{L}{2} \|w - w_k\|^2 - f(w)|$$

1)  $f(w) \leq g_k(w)$

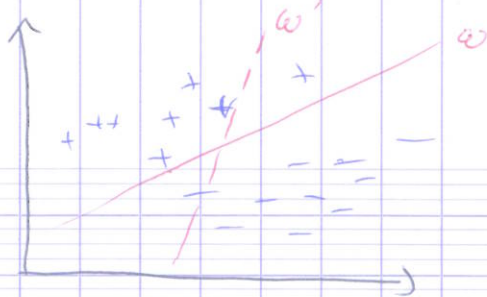
2)  $f(w_k) - f(w) + \nabla f(w_k)(w - w_k) \leq 0$

True since  $f$  is convex (convex definit° dérivée direct. + def° convexe)

$$\Rightarrow -\frac{L}{2} \leq f(w_k) - f(w) + \nabla f(w_k)(w - w_k) \leq 0$$

$$\Rightarrow |h_k(w)| \leq \frac{L}{2} \|w - w_k\|^2$$

Other solut<sup>o</sup>:



$$d(\omega, x) = \left| \frac{\omega^T x}{\|\omega\|} \right|, \quad x = x_{\omega^T} + x_{\omega^\perp}$$

$$d(\omega, x) = \left| \left( \frac{\omega}{\|\omega\|} \right)^T x_{\omega} \right|$$

$$= \left| \frac{\omega^T x}{\|\omega\|} \right|$$

$\rightarrow$  Pb  $\max_{\omega} \min_{\substack{1 \leq i \leq m \\ \Delta \in \{y_i\}} \omega^T x_i \geq 1} d(\omega, x)$

$\Leftrightarrow \max_{\omega} \frac{1}{\|\omega\|}$

s.t.  $\forall i$   
 $y_i \omega^T x_i \geq 1$

$\Leftrightarrow \min \frac{1}{2} \|\omega\|^2$

s.t.  $\forall i$   
 $y_i \omega^T x_i \geq 1$

SVM (support vector machine)

simpler if  $\omega$  low dim.

Lagrangian:

$$L(\omega, \alpha_i) = \frac{1}{2} \|\omega\|^2 - \sum_{i=1}^m \alpha_i (y_i \omega^T x_i - 1)$$

$$\nabla L(\omega) = \omega - \sum_{i=1}^m \alpha_i y_i x_i \Rightarrow \hat{\omega} = \sum_{i=1}^m \alpha_i y_i x_i$$

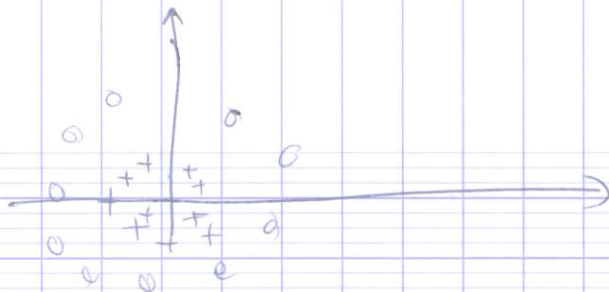
and  $\forall i, \alpha_i (y_i \omega^T x_i - 1) = 0$

Dual formulation of SVM:

$$\max_{\alpha_1, \dots, \alpha_m} \frac{1}{2} \left\| \sum_{i=1}^m \alpha_i y_i x_i \right\|^2 - \sum_{i=1}^m \alpha_i \left( y_i \left( \sum_{j=1}^m \alpha_j y_j x_j \right)^T x_i \right)$$

simpler  $\rightarrow$  if  $\omega$  high dim  $\approx -\alpha K\alpha + L\alpha$

$K, L$  depends on  $(x_i)$  only through the scalar product  $(x_i^T x_j)_{i,j}$



In the case  $\uparrow$ :

$$\mathbb{R}^p \quad x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}, \quad \tilde{x} = \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$$

$$\Phi: X \rightarrow X' \subset \mathbb{R}^{p'} \quad p' \gg p$$

$x \mapsto \Phi(x)$

$\hookrightarrow$  you only need to compute the  $(\langle \Phi(x_i); \Phi(x_j) \rangle)_{i,j}$  in order to find the SVM object

$\rightarrow$  you just need to give a fct

$$k: X \times X \rightarrow \mathbb{R}$$

$$(x, x') \mapsto k(x, x')$$

$$\langle \Phi(x), \Phi(x') \rangle$$

For example:  $k(x, x') = e^{-\frac{\|x-x'\|^2}{2\sigma^2}}$

Remark:

$$\omega_m^S = \underset{\omega}{\operatorname{argmin}} \quad \underbrace{\frac{\|\omega\|^2}{2}}_{\text{penalisation}} + \frac{c}{m} \sum_{i=1}^m (1 - y_i \omega^T x_i)_+$$

$$f^S(\omega) = (1 - \omega)_+$$

Not done: boosting.