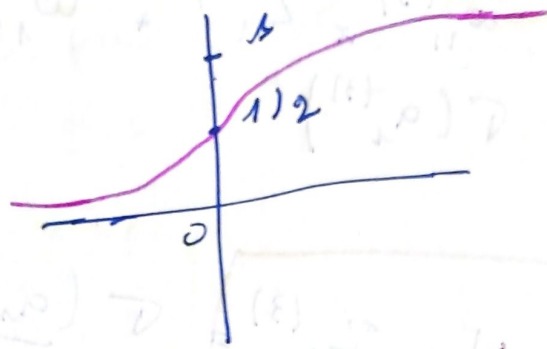


Exercices de Révision Sys

(5)

Q1:

$$1 - \sigma(x) = \frac{1}{1 + e^{-x}}$$



2 -

~~y(x)~~ = Feed forward:

- couche 1

$$\begin{cases} a_1^{(1)} = x_1 w_{11}^{(1)} + w_{10}^{(1)} \\ z_1^{(1)} = \sigma(a_1^{(1)}) \end{cases}$$

$$\begin{cases} a_2^{(1)} = x_2 w_{21}^{(1)} + w_{20}^{(1)} \\ z_2^{(1)} = \sigma(a_2^{(1)}) \end{cases}$$

$$\begin{cases} a_3^{(1)} = x_3 w_{31}^{(1)} + w_{30}^{(1)} \\ z_3^{(1)} = \sigma(a_3^{(1)}) \end{cases}$$

- couche 2:

$$\begin{cases} a_1^{(2)} = w_{11}^{(2)} z_1^{(1)} + w_{12}^{(2)} z_2^{(1)} + w_{10}^{(2)} \\ z_1^{(2)} = \sigma(a_1^{(2)}) \end{cases}$$

$$\begin{cases} a_2^{(2)} = w_{21}^{(2)} z_1^{(1)} + w_{22}^{(2)} z_2^{(1)} + w_{20}^{(2)} \\ z_2^{(2)} = \sigma(a_2^{(2)}) \end{cases}$$

couche 3:
→ Sortie

$$z_1^{(3)} = w_{11}^{(3)} \times z_1^{(2)} + w_{12}^{(3)} \times z_2^{(2)} + w_{10}^{(3)}$$

$$z_1^{(3)} = \sigma(a_1^{(3)})$$

$$y(w, x) = z_1^{(3)} = \sigma(a_1^{(3)})$$

3 - Back propagation

$$l(B) = - \sum_{i=1}^N \log t_i y(w, x) + (1 - t_i) \log (1 - y(w, x))$$

$$y(w, x) = \sigma(a_1^{(3)})$$

$$\frac{\partial l(B)}{\partial a_1^{(3)}} = \sigma(a_1^{(3)}) - t$$

- dérivation en chaîne

$$\begin{aligned} \frac{\partial l}{\partial w_{11}^{(3)}} &= \frac{\partial a_1^{(3)}}{\partial z_1^{(3)}} \times \frac{\partial z_1^{(3)}}{\partial a_1^{(2)}} \times \frac{\partial a_1^{(2)}}{\partial z_1^{(1)}} \times \frac{\partial z_1^{(1)}}{\partial a_1^{(1)}} \times \frac{\partial a_1^{(1)}}{\partial w_{11}^{(1)}} \\ &= w_{11}^{(3)} \times \sigma'(a_1^{(3)}) \times w_{11}^{(2)} \times \sigma'(a_1^{(2)}) \times x_1 \end{aligned}$$

$$\frac{\partial l}{\partial w_{11}^{(3)}} = \frac{\partial l}{\partial a_1^{(3)}} \times \frac{\partial a_1^{(3)}}{\partial w_{11}^{(3)}}$$

$$= (\sigma(a_1^{(3)}) - t) \times w_{11}^{(3)} \times \sigma'(a_1^{(2)}) \times w_{11}^{(2)} \times \sigma'(a_1^{(1)}) \times x_1$$

Question 2 :

1. dans le graphe, on peut diviser les données en 2 parties sachant $x_2 = -2$

on sait que $y(x) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

et on a $x_2 + 2 = 0$

donc

$$2 + 0 \times x_1 + x_2 = 0$$

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$$

par identification on trouve

$$\begin{cases} \beta_0 = 2 \\ \beta_1 = 0 \\ \beta_2 = -1 \end{cases} \Rightarrow \beta(2, 0, 1)$$

2. Minimiser la fct coût : Régularisation Ridge

on a : MSE

$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^i - h(\beta))^2 + \lambda \sum_{j=1}^2 \beta_j^2$$

si nous utilisons une fonction de coût type Ridge les valeurs β_0 et β_1 ne changeront pas car

la pénalité Ridge n'affecte pas β_0 et aussi β_1 est le plus petit possible. Pour β_2 ~~il va~~

en ce qui concerne β_2 il va décroître en fonction du coefficient de pénalité λ : $\beta_2 = 1 - \delta$

$$\delta \geq 0$$

Question 4:

8

$$1. \quad \varepsilon \sim N(0, \sigma^2)$$

$$\varepsilon = t - \hat{t}$$

$$\hat{t} = \beta_0 + \beta_1 x_1 + \dots + \beta_d x_d$$

$$t = \hat{t} + \varepsilon$$

$$= \beta_0 + \beta^T x + \varepsilon$$

and $\varepsilon \sim (0, \sigma^2)$

$$p(t | x, \beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\varepsilon^2}{2\sigma^2}\right)$$

$$\varepsilon = t - \hat{t}$$

$$= t - \beta_0 - \beta^T x$$

$$p(t | x, \beta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t - \beta_0 - \beta^T x)^2}{2\sigma^2}\right)$$

$$\rightarrow \ell(\beta) = \prod_{i=1}^N p(t = t^i | x^i, \beta)$$

$$= \prod_{i=1}^N \left(\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2} (t - \beta_0 - \beta^T x^i)^2\right] \right)$$

$$\log \ell(\beta) = \sum_{i=1}^N \left(\log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) + \underbrace{\left(-\frac{1}{2\sigma^2}\right) (t - \beta_0 - \beta^T x^i)^2}_{A} \right)$$

$$\downarrow$$

$$\sqrt{2\pi\sigma^2}$$

A

$$2. \frac{\partial \log L(\beta)}{\partial \sigma} = 0$$

$$= \sum \left(\frac{\left(\frac{1}{\sqrt{2\pi}\sigma^2} \right)'}{\left(\frac{1}{\sqrt{2\pi}\sigma^2} \right)^2} + \frac{-1 \times (-2\sigma^2)'}{(2\sigma^2)^2} \right) \times A$$

$$(\sqrt{f(n)})' = \frac{f'(n)}{2\sqrt{f(n)}}$$

$$\log(f(n))' = \frac{f'(n)}{f(n)}$$

$$= \sum \left(\frac{-\left(\frac{1}{\sqrt{2\pi}\sigma^2} \right)'}{\left(\frac{1}{\sqrt{2\pi}\sigma^2} \right)^2} + \frac{4\sigma}{(2\sigma^2)^2} \times A \right)$$

$$= \sum \left(-\frac{(2\pi\sigma^2)'}{2\sqrt{2\pi}\sigma^2} + \frac{4}{\sigma^3} \times A \right)$$

$$\frac{\partial \log L(\beta)}{\partial \sigma} = \sum \left(-\frac{4\pi\sigma}{2\sqrt{2\pi}\sigma^2} + \frac{1}{\sigma^3} \times (t - \beta_0 - \beta^T x^i)^2 \right)$$

$$= \sum \left(-\frac{2\pi}{\sqrt{2\pi}} + \frac{1}{\sigma^3} \times (t - \beta_0 - \beta^T x^i)^2 \right)$$

$$= \sum \left(-\sqrt{2\pi} + \frac{1}{\sigma^3} \times (t - \beta_0 - \beta^T x^i)^2 \right) \quad (1)$$

$$\frac{\partial^2 \log L(\beta)}{\partial \sigma^2} = \sum \left(0 + \frac{-(\sigma^3)'}{(\sigma^3)^2} \times (t - \beta_0 - \beta^T x^i)^2 \right)$$

$$= \sum \left(\frac{-3\sigma^2}{\sigma^6} \times (t - \beta_0 - \beta^T x^i)^2 \right)$$

$$= \sum \left(-\frac{3}{\sigma^4} \times (t - \beta_0 - \beta^T x^i)^2 \right) = 0 \quad (2)$$

resoudre $\frac{\partial \log(\mathcal{L}(\beta))}{\partial \sigma} = 0$

$$\Rightarrow -\sqrt{2\pi} + \frac{1}{\sigma^3} (t - \beta_0 - \beta^T x)^2 = 0$$

$$\Rightarrow \frac{1}{\sigma^3} = \frac{\sqrt{2\pi}}{(t - \beta_0 - \beta^T x)^2}$$

$$\Rightarrow \sigma^3 = -\frac{(t - \beta_0 - \beta^T x)^2}{\sqrt{2\pi}}$$

$$\Rightarrow \sigma = \left(-\frac{1}{\sqrt{2\pi}} (t - \beta_0 - \beta^T x)^2 \right)^{1/3}$$

Question 5.2)

argmin $\frac{1}{2} \sum (t^i - \beta x^{(i)})^2$ avec

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}(\beta)}{\partial \beta_0} = \sum (t^i - \beta x^{(i)}) = 0 \\ \frac{\partial \mathcal{L}(\beta)}{\partial \beta_j} = \sum (t^i - \beta x^{(i)}) x^{(i)} = 0 \end{array} \right.$$

5.3

Question 13:

11

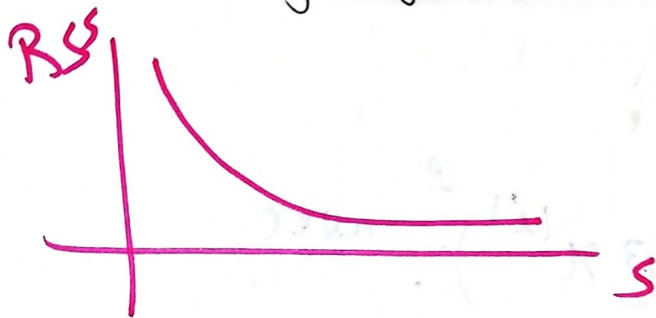
$$1 - \sum_{i=1}^n (t^{(i)} - \beta_0 - \sum \beta_j x_{ij})^2 \quad \sum |\beta_j| \leq s$$

$$s_i \leq \uparrow \leadsto \text{RSS}$$

$$\lambda |\beta| \leq s$$

$\{ s_i \lambda \uparrow \sim \beta \downarrow \Rightarrow \text{forte pénalisation} \leadsto \text{RSS} \uparrow$
 $\{ s_i \downarrow \sim \beta \uparrow \Rightarrow \text{faible pénalisation} \leadsto \text{RSS} \downarrow$

- Faux, RSS d'entraînement diminue lorsque $s \uparrow$
 car moins de pénalisation permet un meilleur ajustement aux données

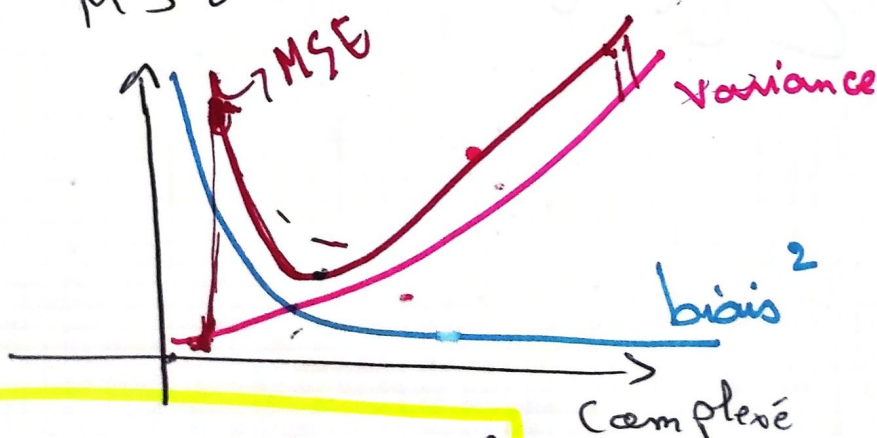


2 - Faux

3 - F

4 - V

5 : $MSE = \text{biais}^2 + \text{variance}$



variance
 bjs : croissante

$$\lambda \uparrow \leadsto \text{MSE} \uparrow$$

5 - Faux

6 - Faux

7 - ~~Faux~~ Faux

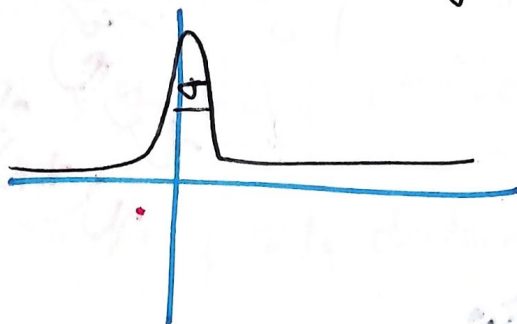
8 - ~~Faux~~ Vrai

variance entre β va
diminuer si on a une forte
pénalité

dans le cas
global

variance
entre β

$\lambda \uparrow \leadsto$ Variance \downarrow ; $\beta \downarrow$



Question 15 :

1 - partie différentiable et non :

* différentiable

$$\frac{1}{N} \sum (t^i - \beta_0 + \sum_{j=1}^D \beta_j x_j^i)^2 + \lambda_2 \sum \beta_j^2$$

* non-différentiable :

$$\lambda_1 \sum_{j=1}^D |\beta_j|$$

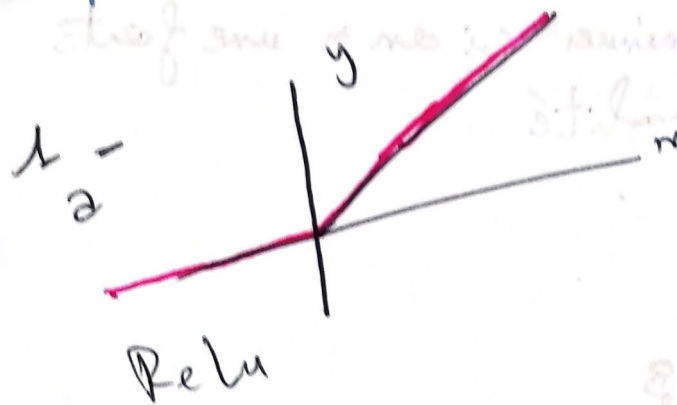
2) Voir figure 5

Lasso

Ridge

Question 7

13



$$y = \begin{cases} 0 & \text{si } x \leq 0 \\ x & \text{si } x > 0 \end{cases}$$

$$Relu = \begin{cases} x & \text{si } x \geq 0 \\ 0 & \text{si } x < 0 \end{cases}$$

$$y = \frac{x + |x|}{2}$$

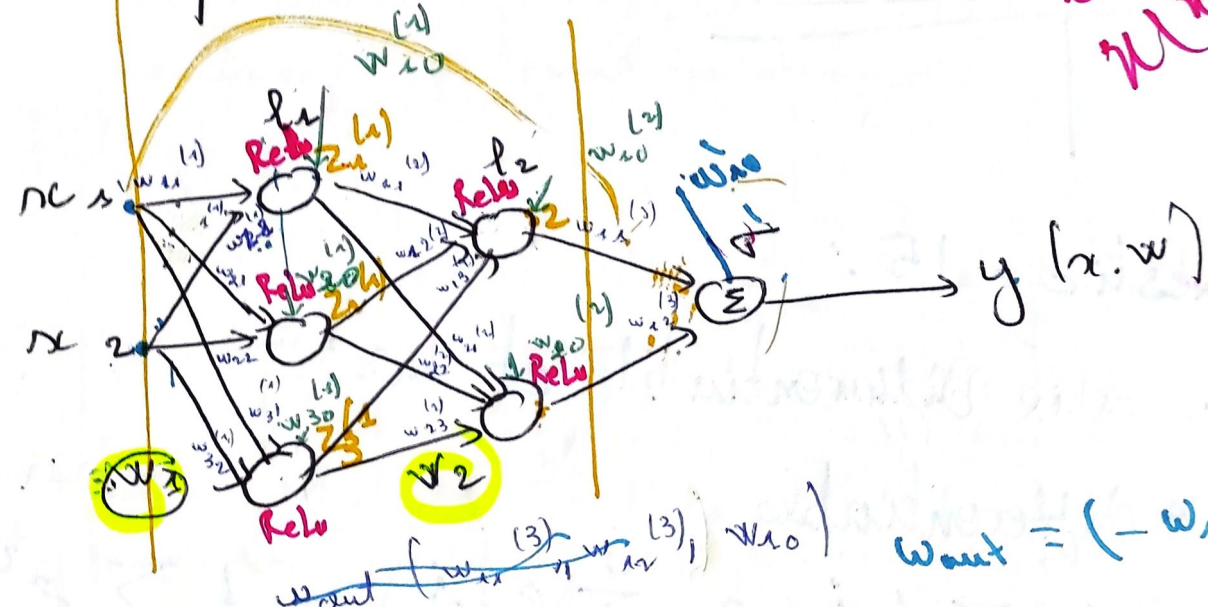


$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\sigma(x) = \begin{cases} 1 & x > 0 \\ 1/2 & x = 0 \\ 0 & x < 0 \end{cases}$$

Handwritten notes in pink:
 $\sigma(x) = \frac{1}{1 + e^{-x}}$
 $\sigma(x) = \begin{cases} 1 & x > 0 \\ 1/2 & x = 0 \\ 0 & x < 0 \end{cases}$

2 -



$$w_{out} = (-w_{10} -)$$

$$w_{out} = (w_{11}^{(3)}, w_{12}^{(3)}, w_{10}^{(3)})$$

$$y(x, w) = \sigma \left(\underbrace{Relu(x \cdot w_2^T + w_2)}_{\text{la sortie précédente}} \times \underbrace{w_{out}}_{\text{poids actuelle}} \right)$$

1. $Relu(w_0, x)$

$x = Relu(w_2 \cdot Relu(w_0, x)) \Rightarrow y = \sigma(x \cdot w_{out})$

Handwritten notes in pink:
 la sortie précédente
 poids actuelle
 c'est-à-dire w_0

- couche 1:

$$\text{Relu}(w_1 \cdot x)$$

- couche 2:

$$\text{Relu}(w_2 \cdot \text{Relu}(w_1 \cdot x))$$

→ **Sortie**

$$\sigma(\text{Relu}(w_2 \cdot \text{Relu}(w_1 \cdot x)))$$

explication:

- w_1, w_2 : contiennent intercepts
- w_1 poids de couche 1
- w_2 — — — — —
- w_2 poids de dernière sortie

$$\text{d'où: } y(x, w_1, w_2) = \sigma(\text{Relu}(w_2 \cdot \text{Relu}(w_1 \cdot [1, x]^T)))$$

Question 8:

least squares classifier : classificateurs des moindres carrés

(21)

$x = [1, x]$ car
les matrices w_1 et
 w_2 contiennent les
intercepts