

Naturellement → apprentissage supervisé

$$\underbrace{\{x^{(i)}, t^{(i)}\}_{i=1}^N}$$

$$x^{(i)} \in \mathbb{R}^D$$

→ régression linéaire

→ descente de gradient

→ réduction d'équations normales

→ overfitting

$$h_{\beta}(x) = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}$$

↳ décomposition biais variance

→ Régularisation

→ Best subset selection

→ Ridge

→ LASSO

Classification $t^{(i)} \in \{1, \dots, K\}$

↳ classification for neural network

→ 1 centre 1

→ 1 centre tous

→ discriminer à K classes

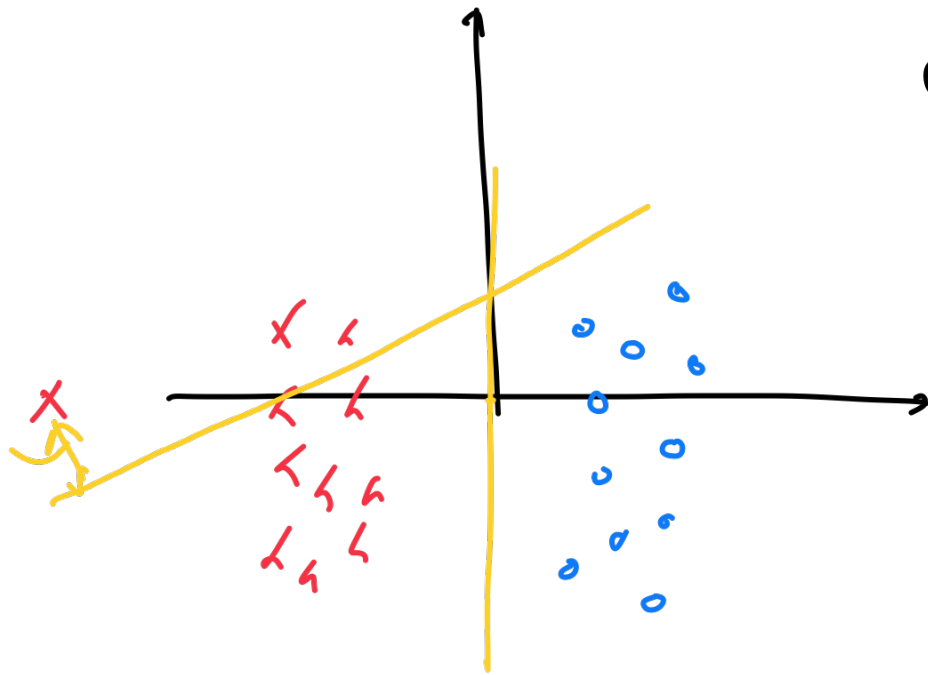
$$X \in \mathbb{R}^{N \times D+2}$$

$$B = \begin{bmatrix} - & \beta_1^T & - \\ & \vdots & \\ - & \beta_K^T & - \end{bmatrix}$$

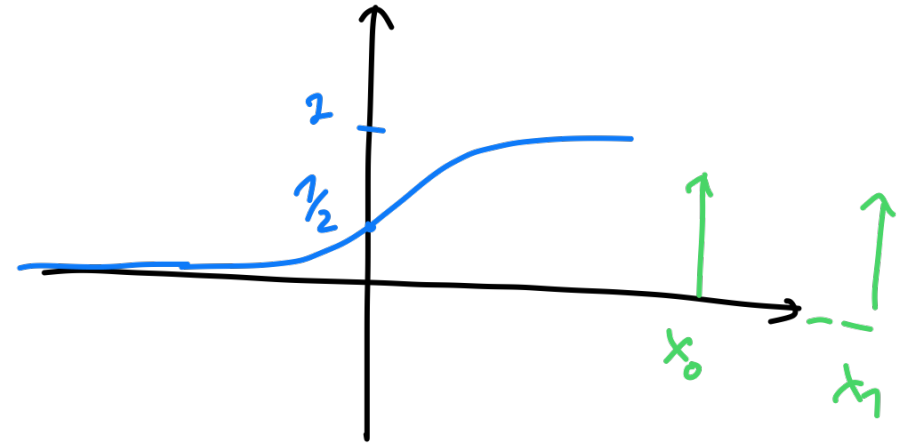
$$T \approx X B^T$$

$$T = \underbrace{\begin{bmatrix} 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \\ \vdots & & & & \end{bmatrix}}_K$$

$$B = (X^T X)^{-1} X^T T$$



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



→ amélioration : on va considérer l'ajout d'une fonction d'activation

$$\sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D)}}$$

Regression
logistique

$$\rightarrow \begin{cases} p(t=0|x) = \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D) \\ p(t=1|x) = 1 - \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D) \end{cases}$$

$$\{x^{(i)}, t^{(i)}\}$$

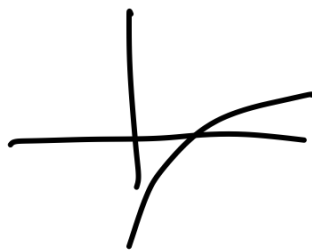
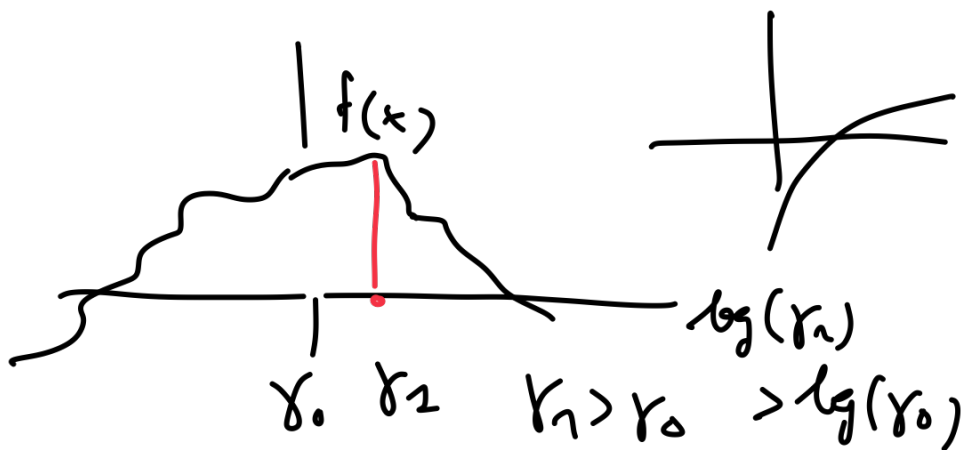
pour une paire $\{x^{(i)}, t^{(i)}\}$

$$p(t=t^{(i)} | x=x^{(i)}) = \underbrace{p(t=0|x^{(i)})}^{1-t^{(i)}} \underbrace{(p(t=1|x^{(i)}))}^{t^{(i)}}$$

$$p(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\}) = \prod_{i=1}^N p(t=t^{(i)} | x=x^{(i)}) \quad \begin{array}{l} \text{(si on suppose} \\ \{x^{(i)}, t^{(i)}\} \\ \text{indépendance} \end{array}$$

$$p(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\}) = \prod_{i=1}^N \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})^{1-t^{(i)}} \times (1 - \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^{t^{(i)}}$$

$$\beta = \underset{\beta}{\operatorname{argmax}} p(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\})$$



$$\beta^* = \underset{\beta}{\operatorname{argmax}} \log \left(\prod_{i=1}^N \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D)^{1-t^{(i)}} (1 - \sigma(\beta_0 + \beta_1 x_1 + \dots + \beta_D x_D))^{t^{(i)}} \right) \quad \text{LLE}$$

$$\underset{\beta}{\operatorname{argmax}} \sum_{i=1}^N (1-t^{(i)}) \log(\sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})) + t^{(i)} \log(1 - \sigma(\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))$$

$$\underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (1-t^{(i)}) \log(\sigma) + t^{(i)} \log(1 - \sigma)$$

BINARY CROSS ENTROPY

Modèles de classification : Modèles discriminants
probabilistes

$$p(t|x)$$

ex: regression

Modèles génératifs

$$p(x|t)$$

$$p(t|x) = \frac{p(x|t)p(t)}{p(x)} \leftarrow \sum_c p(x|t_c)p(t_c)$$

$$\underline{t(x)} = \underset{k}{\operatorname{argmax}} p(t=k|x^{(i)}) = \underset{k}{\operatorname{argmax}} p(x|t=k)p(t=k)$$

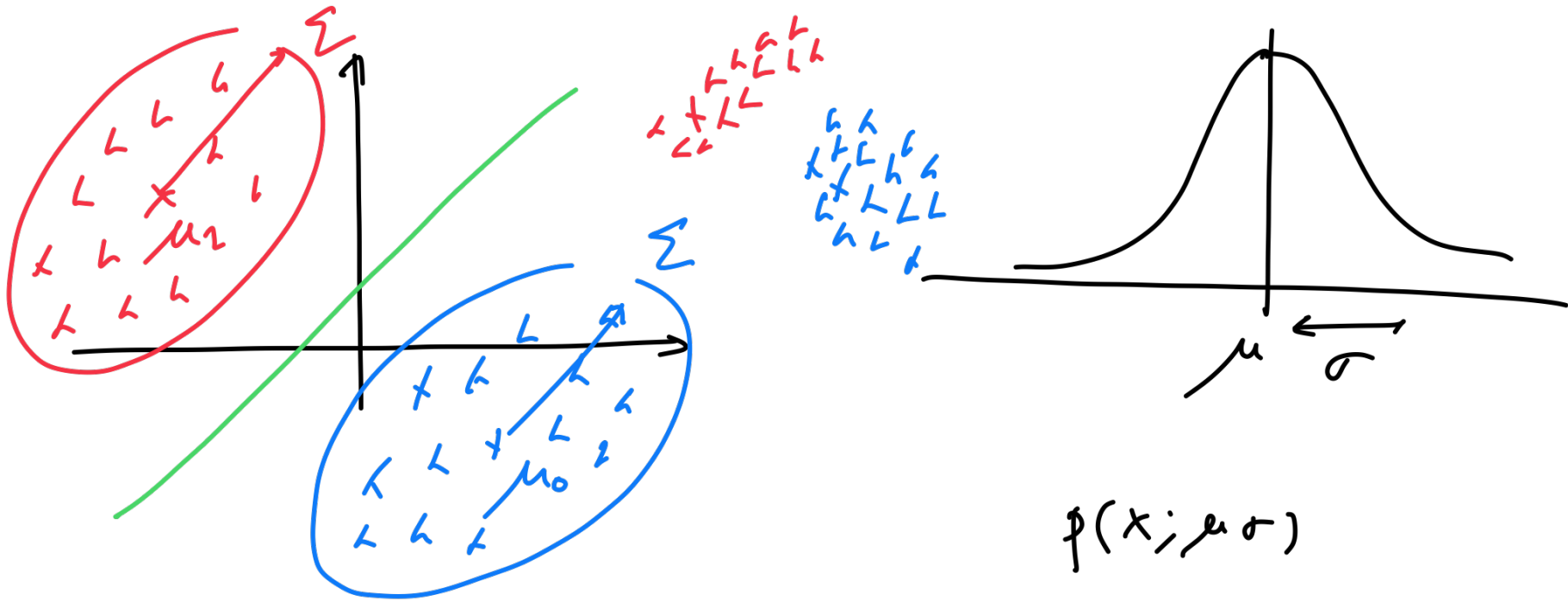
Un des modèles génératifs les plus populaires consiste à faire $p(x|t)$ donnée par une distribution normale multivariée.

$$p(\vec{x} | t=0) = \frac{1}{(2\pi)^{D/2} |\Sigma_0|^{1/2}} \exp\left(-\frac{1}{2}(\vec{x} - \mu_0)^T \Sigma_0^{-1} (\vec{x} - \mu_0)\right)$$

$$p(\vec{x} | t=1) = \frac{1}{(2\pi)^{D/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(\vec{x} - \mu_1)^T \Sigma_1^{-1} (\vec{x} - \mu_1)\right)$$

GDA \rightarrow Analyse Gaussienne discriminante

Si $\Sigma_0 = \Sigma_1 = \Sigma \rightarrow$ LDA Analyse linéaire discriminante



$$p(x; \mu, \sigma)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

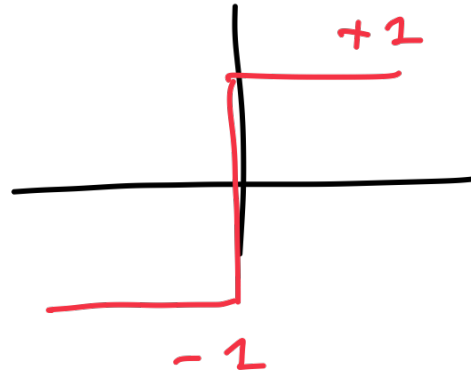
les paramètres du modèle peuvent être
calculés explicitement

$$p(t=0) = \frac{N_0}{N} \quad p(t=1) = \frac{N_1}{N}$$

$$\mu_0 = \frac{1}{N_0} \sum_{i \in C_0} x^{(i)} \quad \mu_1 = \frac{1}{N_1} \sum_{i \in C_1} x^{(i)}$$

$$\Sigma = \frac{1}{N} \left(\sum_{i \in C_0} (x^{(i)} - \mu_0)(x^{(i)} - \mu_0)^T + \sum_{i \in C_1} (x^{(i)} - \mu_1)(x^{(i)} - \mu_1)^T \right)$$

$$\sigma(x) = \begin{cases} +1 & \text{if } x \geq 0 \\ -1 & \text{if } x < 0 \end{cases}$$



→ fonction d'activ. de gradient nul presque partout
 on va considérer une nouvelle fonction de coût

$$\mathcal{L} = - \sum_{i \in \text{misclassified}} t^{(i)} (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})$$

↑
↑

$$\beta = \argmin_{\beta} \mathcal{L}(\beta)$$

PERCEPTRON

→ Tant qu'il existe des points incorrectement classés

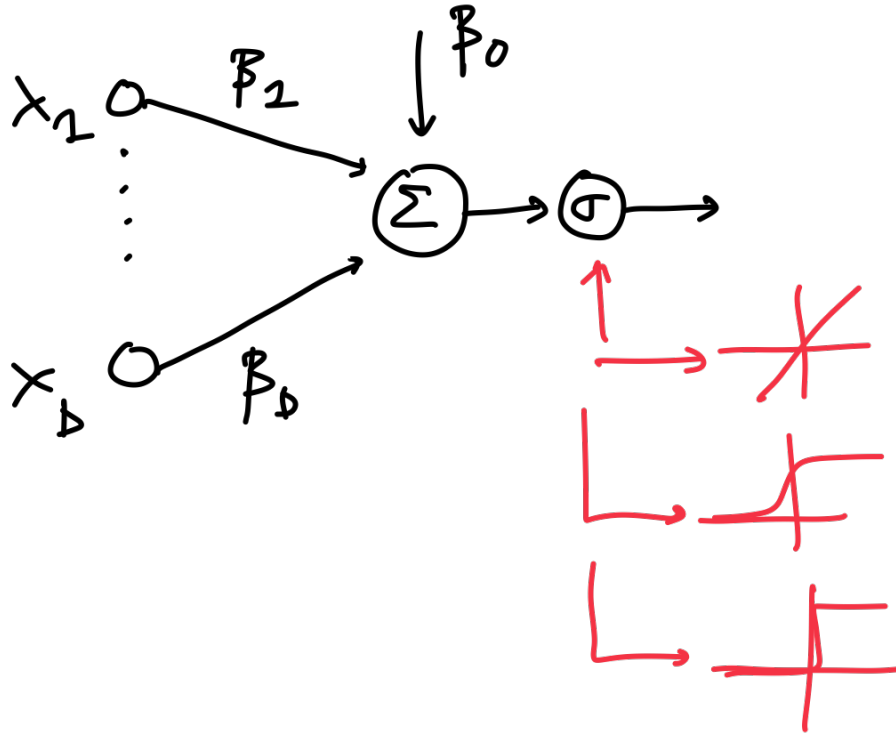
→ sélectionner un de ces points

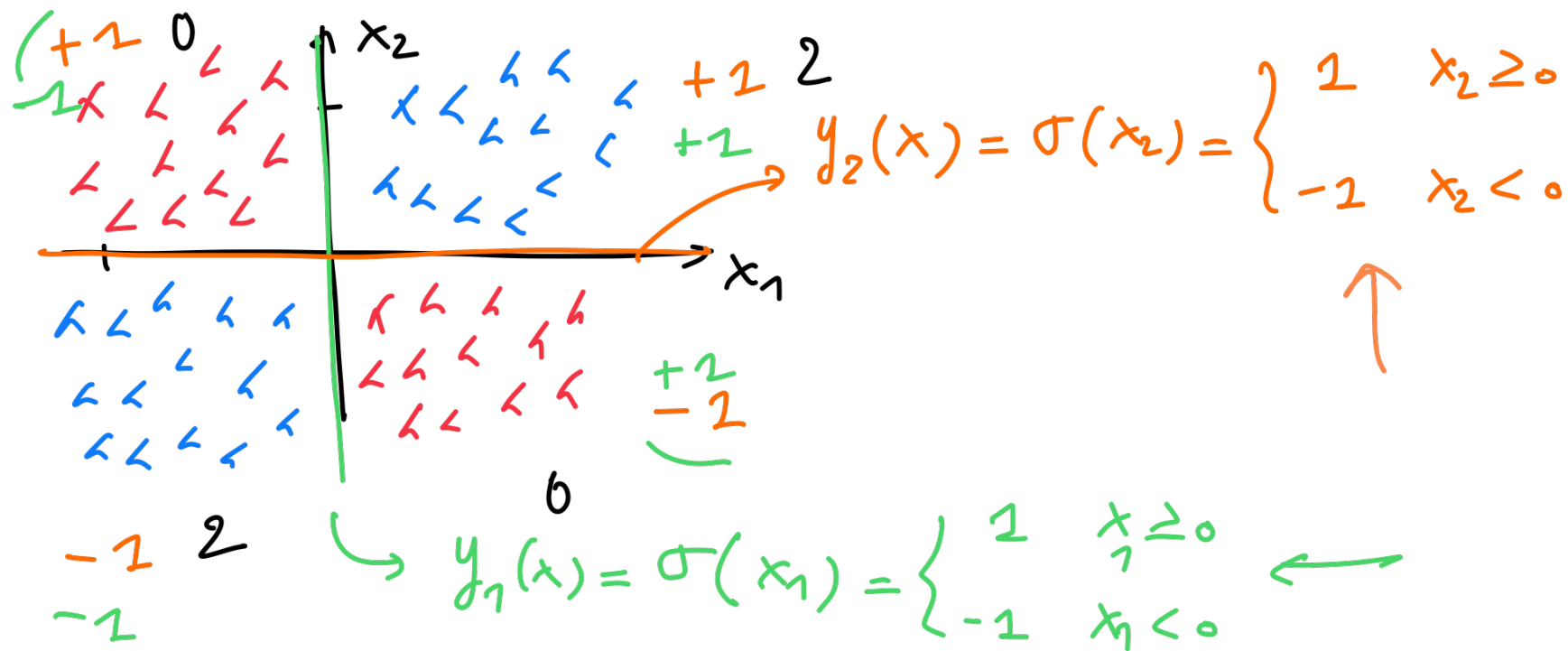
→ Mettre à jour β

$$\beta \leftarrow \beta + \eta \cdot t^{(i)} \tilde{x}^{(i)}$$

→ Si les données sont linéairement séparables,
l'algorithme converge en 1 nombre fini d'itérations.

Tous les modèles courants jusqu'ici peuvent être représentés par le diagramme suivant



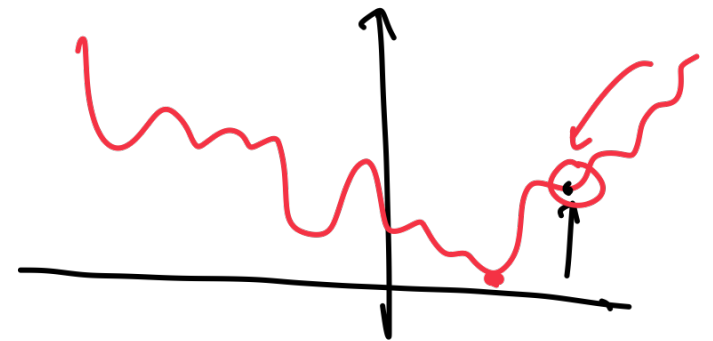
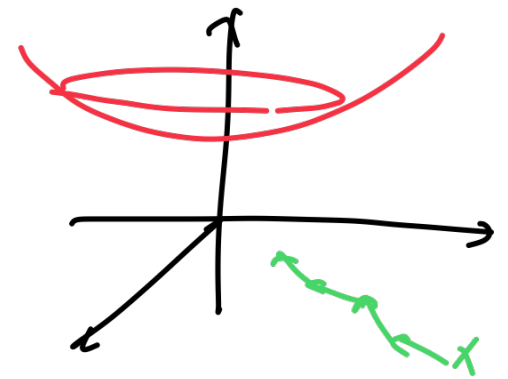
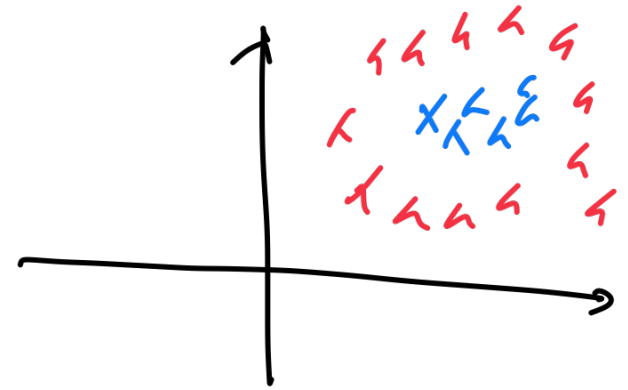


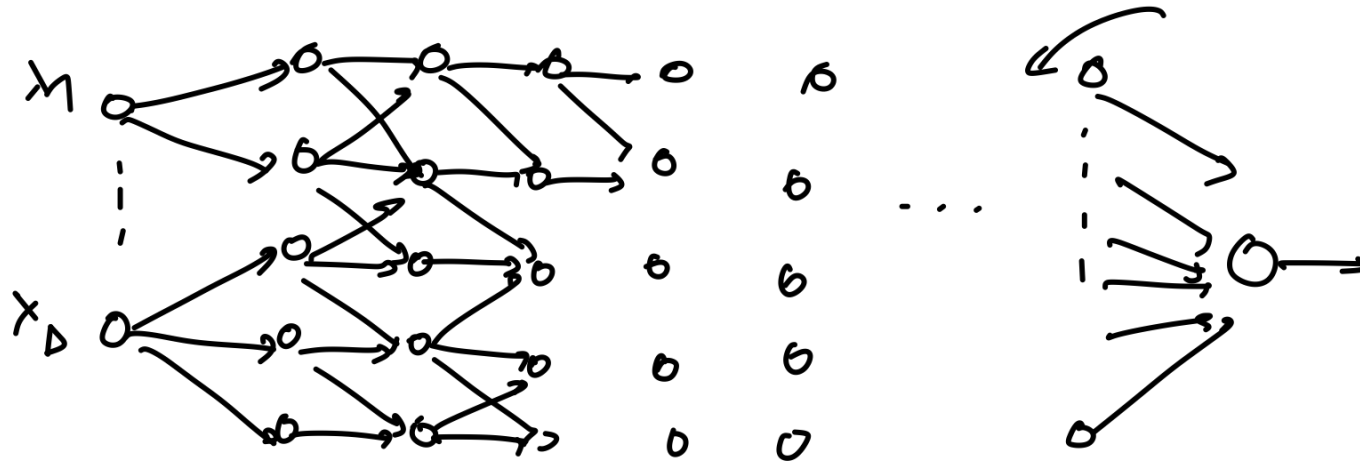
$$\begin{aligned}
 y(x) &= | \underbrace{y_1(x)}_{\text{green}} + \underbrace{y_2(x)}_{\text{orange}} | = \sigma^{(2)} \left(\underbrace{1}_{w_1=2} \cdot \underbrace{\sigma^{(1)}(x_1)}_{\uparrow} + \underbrace{\sigma^{(2)}(x_2)}_{w_2=1} \right) \\
 &= f(x)
 \end{aligned}$$

$$p(t=0|x) = y_w(x)$$

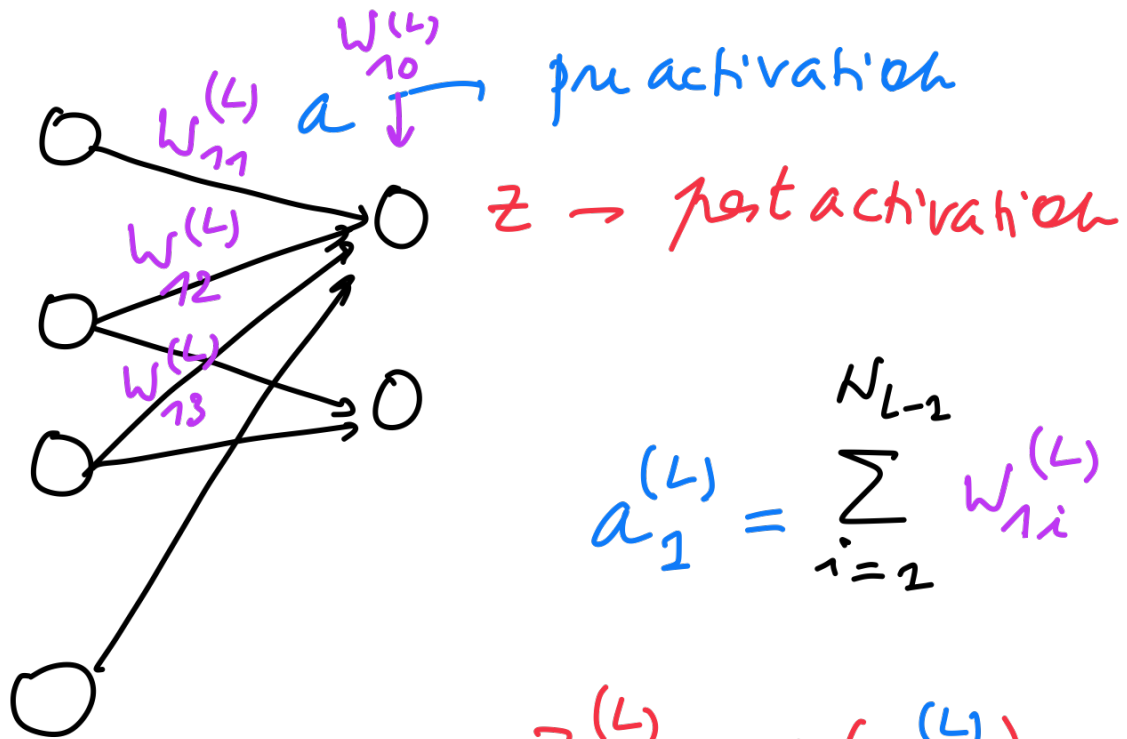
$$p(t=1|x) = 1 - y_w(x)$$

$$l(w) = - \sum_{i=1}^N (1 - t^{(i)}) \log(y_w) + t^{(i)} \log(1 - y_w)$$





Dans le cadre des réseaux de neurones le calcul du gradient se base sur l'idée de back propagation



$$a_1^{(L)} = \sum_{i=1}^{N_{L-1}} W_{1i}^{(L)} z_i^{(L-1)} + W_{10}^{(L)}$$

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

$$= \sigma \left(\sum_{i=1}^{N_{L-1}} W_{1i}^{(L)} z_i^{(L-1)} + W_{10}^{(L)} \right)$$