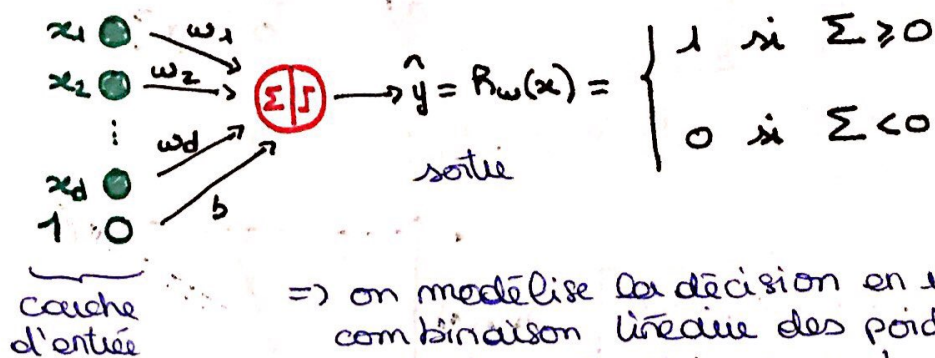


# Du perceptron à la régression logistique (cas de la classification binaire)

Architecture du perceptron (0 couches cachées)

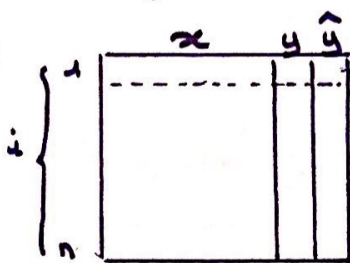


⇒ on modélise la décision en utilisant la combinaison linéaire des poids et des entrées suivie d'une fonction de seuil.

$$\hat{y} = R_w(x) = \text{Threshold}(x^T w)$$

donc on va chercher les meilleurs poids  $w$  et  $b$  pour obtenir la meilleure prédiction (proche de la vraie sur le jeu d'entraînement).

Algorithme

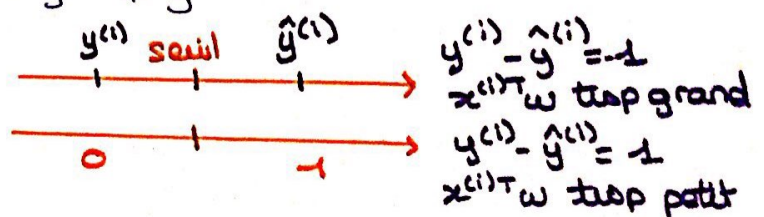


pour chaque paire  $x^{(i)}, y^{(i)}$

- calcul de la prédiction  $\hat{y}^{(i)}$
- comparer  $y^{(i)}$  à  $\hat{y}^{(i)}$

Idée =

- si  $y^{(i)} = \hat{y}^{(i)} \Rightarrow$  bonne prédiction pas de maj.
- si  $y^{(i)} \neq \hat{y}^{(i)}$



Formalisation en termes de coût.

$$L(y^{(i)}, \hat{y}^{(i)}) = -(y^{(i)} - \hat{y}^{(i)}) x^{(i)T} w$$

Règle de mise à jour des poids

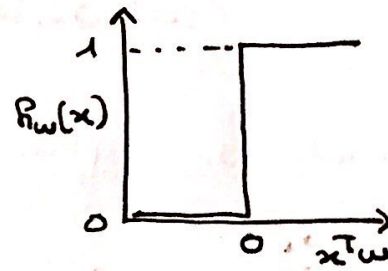
(descente de gradient)

$$w_d \leftarrow w_d - \alpha \frac{\partial L(y^{(i)}, \hat{y}^{(i)})}{\partial w_d}$$

on diminue les poids ou on ↑ les poids en fonction de "combien on s'est trompé"

Mais la dérivée de la fonction seuil en 0 n'existe pas :  
donc cela peut causer des pb pendant l'entraînement.

=> au lieu de prédire la classe, on va prédire une probabilité qui continue entre 0 et 1.

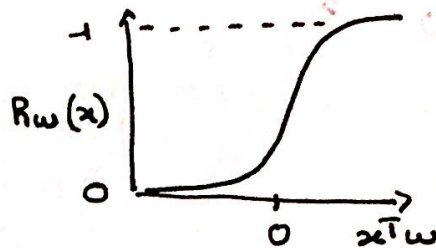


Pour cela, on n'utilise pas la fonction seuil mais la fonction sigmoïde.

$$\hat{y} = R_w(x) = \text{logistic}(x^T w)$$

$$= P(Y=1|x)$$

$$= \frac{1}{1 + e^{-x^T w}}$$



Fonction de perte associée

Binary cross entropy

$$L(y^{(i)}, \hat{y}^{(i)}) = -y^{(i)} \log(R_w(x^{(i)})) - (1-y^{(i)}) \log(1 - R_w(x^{(i)}))$$

Mise à jour des poids.

$$w_d \leftarrow w_d - \alpha \frac{\partial L(y^{(i)}, R_w(x^{(i)}))}{\partial w_d} = w_d + \alpha (y^{(i)} - R_w(x^{(i)})) x_d$$

identique à celle du perceptron.

on a juste changé la loss et la règle de décision.