

- apprentissage supervisé :

$$\mathcal{D} = \{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots\}$$

⇒ recherche d'un prédicteur minimisant la différence entre labels prédits et réels

⇒ nécessite l'annotation de longues bases de données

- apprentissage non supervisé

inference de l'info à partir d'observation uniquement
 a/ réduction de données b/ extraction de connaissances
 c/ clustering d/ estimation de densité

- Reinforcement learning

apprentissage par échantillon supervisé

- un petit nombre d'annotations (few shots)
- one shot learning

- les annotations sont trop bruitées ou imprecises

apprentissage semi-supervisé

petit nombre de données annotées

+ grand nombre de données non annotées

apprentissage actif

ensemble de données dont seulement certains sont annotés

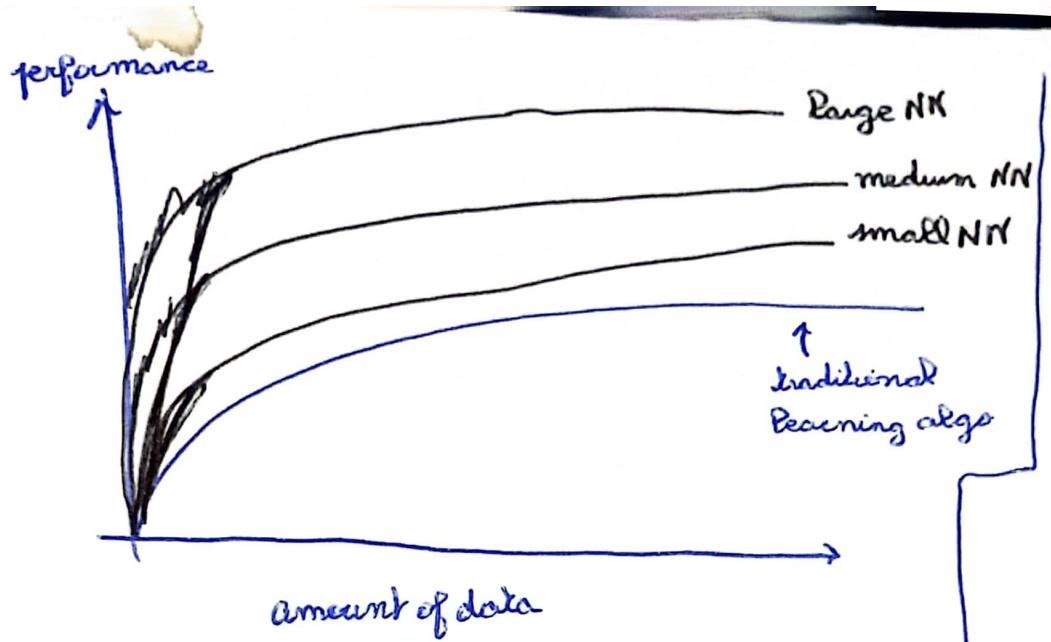
⇒ le modèle doit déterminer au travers d'un certain critère à définir quelles sont les données vont potentiellement lui fournir le plus d'infos pour évoluer vers de meilleurs perfos

⇒ un expert peut annoter ces données

Finance

* préduire le cours d'une action

x



⇒ scale drives deep learning progress:

$\left\{ \begin{array}{l} \text{data} \\ \text{computation} \\ \text{algorithms} \end{array} \right.$

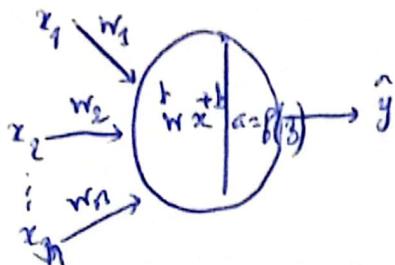
small learning data:

to get better performance:

skills at handling features

Perceptions

neuronne : biellet + poids synaptique/bias b
+ fonction d'activation



Linear - + MSE \Rightarrow regression
activation linéaire

sigmoid + cross entropy \Rightarrow regression
activation Pass logistique

algorithme de perception monocoche

1/ initialisation des poids $w_j^{(0)}$

2/ présentation du vecteur d'entrée $x^{(1)}, \dots, x^{(m)}$
et du vecteur de sortie $y^{(1)}, \dots, y^{(m)}$

3/ calcul de la sortie prédictive

$$\hat{y}^{(i)}\{h\} = f\left(\sum_{j=1}^n w_j^{(h)} x_j^{(i)}\right)$$

$$J(w) = \sum_{i=1}^m \text{perde}\left(\hat{y}^{(i)}\{h\}, y^{(i)}\right)$$

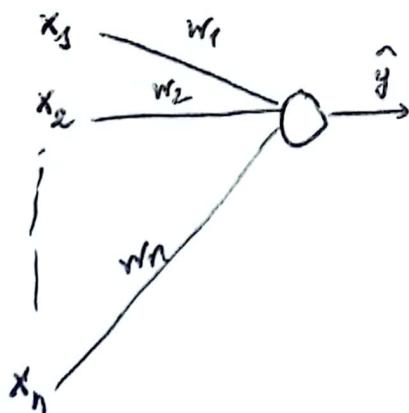
4/ mise à jour des poids : taux d'apprentissage

$$0 < \alpha < 1$$

$$w_j^{(h+1)} = w_j^{(h)} - \alpha \cdot \frac{\partial J}{\partial w_j}$$

5/ retourner en (2) jusqu'à la convergence

enche
d'entrée



Régression Linéaire

pour le i^e exemple :

$$\hat{y}^{(i)} = \theta_0 + \theta_1 x^{(i)}$$

~~equation~~

$$\mathcal{L}(y^{(i)}, \hat{y}^{(i)}) = (\hat{y}^{(i)} - y^{(i)})^2$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (\theta_0 x^{(i)} - y^{(i)})^2 : \boxed{\text{mean squared error}}$$

$$X \in \mathbb{R}^{m \times (d+1)}$$

$$y \in \mathbb{R}^m$$

fonction d'activation linéaire

+

MSE



régression linéaire

m : nb d'exemples

d : nb of features

$$\theta \in \mathbb{R}^{d+1}$$

On utilise un algo de descente de gradient

peut éviter l'inversion de matrice qui est

une opération très coûteuse

1. binary classification

examples:
cat vs no cat

$$y: (1) \quad (0)$$

$$\text{for } (X, y) \quad X \in \mathbb{R}^n \quad y \in \{0, 1\}$$

m learning examples:

$$((x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}))$$

$$m = m_{\text{train}}$$

m_{test} : number of learning examples

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix} \quad n$$

m columns

$$Y = [y^{(1)}, y^{(2)}, \dots, y^{(m)}] \in \mathbb{R}^{1 \times m}$$

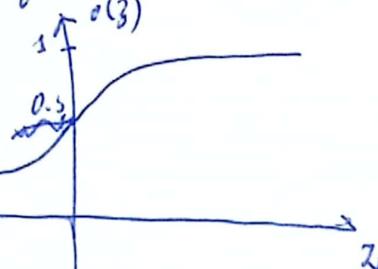
given x we want $\hat{y} = \text{PP}(y=1|x)$

$$x \in \mathbb{R}^n$$

parameters $w \in \mathbb{R}^n, b \in \mathbb{R}$

output: $c(w^T x + b) = \hat{y}$

sigmoid function



$$c(z) = \frac{1}{1 + e^{-z}}$$

$$\text{if } z \rightarrow +\infty \quad c(z) \rightarrow 1$$

$$\text{if } z \rightarrow -\infty \quad c(z) \rightarrow 0$$

One layer neuron networks (perceptrons)

dans des modèles linéaires:
on utilise la transposée de X
 $X \in \mathbb{R}^{n \times m}$

Cost function

$$\text{Given } \{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$$

$$\text{we want to get } \hat{y}^{(i)} \approx y^{(i)}$$

cross entropy loss

$$\mathcal{L}(\hat{y}, y) = - (y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$

if $y \geq 1 \Rightarrow \mathcal{L}(\hat{y}, y) = -\log(\hat{y}) \rightarrow$ we penalize the examples where $\hat{y} = 0$
 \Rightarrow we want \hat{y} to be large \Rightarrow we want \hat{y} large

if $y = 0 \Rightarrow \mathcal{L}(\hat{y}, y) = 0 - \log(1-\hat{y}) \rightarrow$ we penalize the examples where $\hat{y} = 1$
 \Rightarrow we want \hat{y} small

\Rightarrow Cost:

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})$$

1. Gradient descent

Want to find w and b to minimize $J(w, b)$

$J(w, b)$ is convex in this case.

initialize w and b as w_0 and b_0

Repeat {

$$w := w - \alpha \cdot \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \cdot \frac{\partial J(w, b)}{\partial b}$$

}

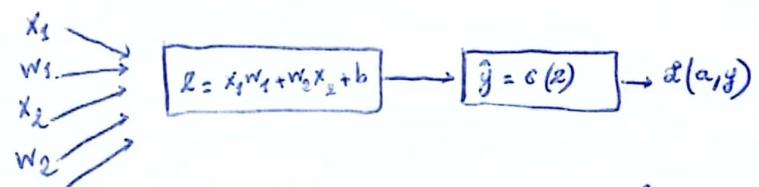
α : learning rate

for Regularization on a single learning example

$$z = w^T x + b$$

$$\hat{y} = c(z)$$

$$\mathcal{L}(\hat{y}, y) = - (y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$$



$$\frac{d\mathcal{L}(\hat{y}, y)}{d\hat{y}} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$\frac{d\mathcal{L}}{dz} = \frac{d\hat{y}}{dz} \cdot \frac{d\mathcal{L}}{d\hat{y}} = \hat{y} - y$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial z}, \frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z}, \kappa_1$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z}, \frac{\partial \mathcal{L}}{\partial w_2} = \frac{\partial \mathcal{L}}{\partial z}, \kappa_2$$

$$\frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial b} = \frac{\partial \mathcal{L}}{\partial z}$$

gradient descent for m training examples

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=1}^m \frac{\partial \mathcal{L}}{\partial w_1} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

algorithm

$$j=0 \quad dw_1 = 0 \quad dw_2 = 0 \quad db = 0$$

for i = 1 to m

$$z^{(i)} = w^T x^{(i)} + b$$

$$\alpha^{(i)} = \sigma(z^{(i)})$$

$$J_L = - \left\{ y^{(i)} \log(\alpha^{(i)}) + (1-y^{(i)}) \log(1-\alpha^{(i)}) \right\}$$

$$dz^{(i)} = \frac{dJ_L}{dz^{(i)}} = \alpha^{(i)} - y^{(i)}$$

$$\begin{cases} dw_1 += x_1^{(i)} \cdot dz^{(i)} \\ dw_2 += x_2^{(i)} \cdot dz^{(i)} \\ db += dz^{(i)} \end{cases}$$

$$J_L = m$$

$$dw_1 | = m$$

$$dw_2 | = m$$

$$db | = m$$

reducing this

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$dw = 0$$

for i = 1 to m

$$\begin{aligned} dw &+= x^{(i)} \cdot dz^{(i)} \\ db &+= x^{(i)} \cdot dz^{(i)} \end{aligned}$$

vectorizing Logistic Regression

propagation steps

$$\begin{aligned} z^{(i)} &= w^T x^{(i)} + b \\ \alpha^{(i)} &= \sigma(z^{(i)}) \end{aligned}$$

$$X = \begin{bmatrix} 1 & | & | & | \\ x^{(1)} & x^{(2)} & \cdots & x^{(m)} \end{bmatrix} \in \mathbb{R}^{n \times m}$$

$$\begin{aligned} Z &= [z^{(1)}, \dots, z^{(m)}] \\ &= w^T X + \underbrace{\begin{bmatrix} b, \dots, b \end{bmatrix}}_{(1, m)} \end{aligned}$$

$$Z = np.\text{dot}(w, X) + b$$

cell operation fail

$$Z = [w^T x^{(1)} + b, \dots, w^T x^{(m)} + b]$$

"broadcasting de b"

$$A = [\alpha^{(1)}, \dots, \alpha^{(m)}] = \sigma(Z)$$

$$A = \text{np.linalg.sigmoid}(Z)$$

$$dz^{(i)} = \alpha^{(i)} - y^{(i)}$$

$$dZ = [dz^{(1)}, \dots, dz^{(m)}]$$

$$dZ = A - y$$

$$db = \frac{1}{m} \cdot \sum_{i=1}^m dz^{(i)}$$

$$\Rightarrow db = \frac{1}{m} \cdot np.\text{sum}(dZ)$$

$$dw = \frac{1}{m} \cdot X \cdot dZ^T$$

$$dw = 1/m \cdot np.\text{dot}(X, dZ^T)$$

w : a n by 1 vector
 X : a n by m vector , y : a 1 by n vector

$$J = 0$$

$$dw = 0 \quad \# \text{ a } n \text{ by } 1 \text{ vector}$$

$$db = 0$$

$$Z = np \cdot \text{dot}(w^T, X) + b$$

$$A = \sigma(Z)$$

$$dZ = A - Y$$

$$dw = 1/m \cdot X \cdot dZ^T$$

$$db = 1/m \cdot np \cdot \text{sum}(dZ)$$

~~~~~

$$w = w - \alpha \cdot dw$$

$$b = b - \alpha \cdot db$$



This is a single iteration of  
a gradient descent

$$\frac{\partial J}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left[ \frac{1}{m} \sum_{i=1}^m -y^{(i)} \log(h_\theta(x^{(i)})) - (1-y^{(i)}) \log(1-h_\theta(x^{(i)})) \right]$$

$$= \frac{1}{m} \sum_{i=1}^m + y^{(i)} \frac{\partial}{\partial \theta_j} \log(h_\theta(x^{(i)})) - (1-y^{(i)}) \frac{\partial}{\partial \theta_j} \log(1-h_\theta(x^{(i)}))$$

$$z = \epsilon(x) \quad u = \epsilon(z) \text{ et } v = \log(u)$$

pour un exemple

$$\frac{\partial}{\partial \theta_j} (\log(h_\theta(x))) = \frac{\partial}{\partial \theta_j} (\log(\epsilon(\epsilon^\top x)))$$

$$= \frac{\partial v}{\partial u} \cdot \frac{\partial u}{\partial z} \cdot \frac{\partial z}{\partial \theta_j}$$

$$= \frac{1}{\epsilon(z)} \epsilon(z) \cdot (1-\epsilon(z)) x_j$$

$$= (1-h_\theta(x)) \cdot x_j$$

$$\text{et } \frac{\partial}{\partial \theta_j} (\log(1-h_\theta(x))) = -h_\theta(x) \cdot x_j$$

$$\text{ainsi } \frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

$$Z = w_1 x_1 + w_2 x_2 + b$$

$$\alpha = \frac{1}{1 + e^{-z}}$$

$$\mathcal{L} = -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(\alpha^{(i)}) + (1-y^{(i)}) \log(1-\alpha^{(i)})$$

$$w = w - \alpha \cdot \frac{\partial \mathcal{L}}{\partial w} \quad b = b - \alpha \cdot \frac{\partial \mathcal{L}}{\partial b}$$

$$\frac{\partial \mathcal{L}}{\partial w_1}$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \frac{\partial \mathcal{L}}{\partial a} \times \frac{\partial a}{\partial z} \times \frac{\partial z}{\partial w_1}$$

$$\frac{\partial a}{\partial z} = -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)}}{a^{(i)}} - \frac{(1-y^{(i)})}{1-a^{(i)}}$$

$$\frac{\partial a}{\partial z} = + \frac{e^{-z}}{(1+e^{-z})^2} = a(z) (1-a(z))$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = w_1 x_1$$

$$\frac{\partial \mathcal{L}}{\partial w_1} = \left( -\frac{1}{m} \sum_{i=1}^m \frac{y^{(i)}}{a^{(i)}} - \frac{(1-y^{(i)})}{1-a^{(i)}} \right) x_1 \\ (a'(z) = (1-a''(z))) \times x_1$$

$$= w_1 - \frac{x_1}{m} \sum_{i=1}^m \frac{y^{(i)}(1-a^{(i)})}{a^{(i)}} - (1-y^{(i)}) a^{(i)}$$

$$= -\frac{m}{m} \sum_{i=1}^m (y^{(i)} - a^{(i)}) x_1^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial w_2} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - a^{(i)}) x_2^{(i)}$$

$$\frac{\partial \mathcal{L}}{\partial b} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - a^{(i)})$$

$$X = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \in \mathbb{R}^{m \times 1}$$

m : nombre de données

n : nombre de variables

$$W = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix} \in \mathbb{R}^{n \times 1}$$

$$Z = X \cdot W + b$$

$$Z \in \mathbb{R}^{m \times 1}$$

$$A = \begin{bmatrix} c(z_1) \\ \vdots \\ c(z_n) \end{bmatrix} \in \mathbb{R}^{n \times 1} = c(Z)$$

$$\mathcal{L} = -\frac{1}{m} \sum_{i=1}^m y \times \log(A) + (1-y) \times \log(1-A)$$

product terme à terme

## vectorisation de la descente du gradient

$$w = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

$$w = w - \alpha \cdot \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \vdots \\ \frac{\partial L}{\partial w_n} \end{bmatrix}$$

$$b = b - \alpha \cdot \frac{\partial L}{\partial b}$$

## vectorisation des gradients

$$\frac{\partial L}{\partial w} = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{bmatrix} = \frac{1}{m} \cdot \left[ \begin{array}{l} (a^{(1)} - y^{(1)}) x_1^{(1)} + (a^{(2)} - y^{(2)}) x_1^{(2)} + \dots + (a^{(m)} - y^{(m)}) x_1^{(m)} \\ (a^{(1)} - y^{(1)}) x_2^{(1)} + \dots + (a^{(m)} - y^{(m)}) x_2^{(m)} \end{array} \right]$$

$$= \frac{1}{m} \cdot X^T \cdot (A - Y)$$

$$\frac{\partial L}{\partial b} = \frac{1}{m} \sum_{i=1}^m (a^{(i)} - y^{(i)})$$

$$= \frac{1}{m} \cdot \sum A - Y$$

## Un neurone formel

### Définition 1

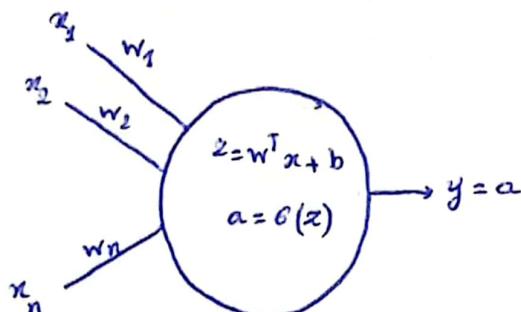
un neurone formel est une fonction paramétrée par  $(n+s)$  paramètres  $w_1, \dots, w_n, o$

$$y: \mathbb{R}^n \times \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}$$

$$(x, w, b) \rightarrow g(x, w, b) = \sigma \left( \sum_{i=1}^n w_i x_i + b \right)$$

$\sigma$ : appelée fonction d'activation

chaque  $w_i$ : poids synaptique associé au signal d'entrée  $x_i$



on a à notre disposition K points  $x^{(k)} \in \mathbb{R}^n$  et

$$y^{(k)} \in \mathbb{R}$$

### Apprentissage du neurone

estimation par les moindres carrés des paramètres du neurone

### Exemple 1

$$\theta: \mathbb{R} \rightarrow \mathbb{R}$$

$$z \rightarrow \frac{s - e^z}{1 + e^z}$$

$$r(\beta): \mathbb{R}^{n+s} \rightarrow \mathbb{R}^K$$

$$\beta \rightarrow \begin{pmatrix} r_1(\beta) \\ \vdots \\ r_K(\beta) \end{pmatrix}$$

$$\begin{pmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \\ b \end{pmatrix} \quad r_i(\beta) = y^{(i)} - g(x^{(i)}, w, o)$$

$$(\beta) \quad \min_{\beta} \{ \beta | r(\beta) = \frac{1}{2} \|r(\beta)\|^2 \} \quad \beta \in \mathbb{R}^{n+1}$$

2.2) ce n'est pas un problème linéaire

2.3)  $\sigma \theta: \mathbb{R} \rightarrow \mathbb{R}$

$$z \rightarrow \frac{s - e^z}{1 + e^z}$$

$$r_i(\beta) = y^{(i)} - g(x^{(i)}, w, o)$$

$$= y^{(i)} - \sigma \left( \sum_{k=1}^n w_k x_k^{(i)} + b \right)$$

$$= y^{(i)} - b - \sum_{k=1}^n w_k x_k^{(i)}$$

$$r(\beta) = \begin{pmatrix} y^{(1)} \\ \vdots \\ y^{(K)} \end{pmatrix} - \begin{pmatrix} 1 & x_1^{(1)} & \cdots & x_n^{(1)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(K)} & \cdots & x_n^{(K)} \end{pmatrix} \beta$$

Le problème devient linéaire

## Régression Logistique cas multivariée

Reclames avec  $P > \epsilon$

⇒ On utilise la fonction softmax

$$P[y = i | x; \theta] = \frac{e^{\theta_i^T x}}{\sum_{j=1}^k e^{\theta_j^T x}}$$

avec  $y^{(i)} = [0 \dots 0 \underset{i}{\text{1}} \dots 0]$ :  $i$  est la  $i$ -ème coordonnée

⇒ Pour chaque vecteur de données  $x$  est  $x^{(i)}$

On calcule un vecteur de probas

L'sigmoid de ces probabilités est la proba

$$\begin{aligned} P_\theta(x^{(i)}) &= \begin{bmatrix} P(y^{(i)}=1 | x^{(i)}; \theta) \\ P(y^{(i)}=0 | x^{(i)}; \theta) \end{bmatrix} \\ &= \begin{bmatrix} \frac{1}{1 + \sum_{j=1}^k e^{\theta_j^T x^{(i)}}} \\ \begin{matrix} e^{\theta_1^T x^{(i)}} \\ \vdots \\ e^{\theta_k^T x^{(i)}} \end{matrix} \end{bmatrix} \end{aligned}$$

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^k \mathbb{I}\{y^{(i)}=j\} \log \left( \frac{e^{\theta_j^T x^{(i)}}}{\sum_{p=1}^k e^{\theta_p^T x^{(i)}}} \right)$$

$$\nabla_{\theta_j} J(\theta) = -\frac{1}{m} \sum_{i=1}^m x^{(i)} (\mathbb{I}\{y^{(i)}=j\} - P(y^{(i)}=j | x^{(i)}; \theta))$$

### exemple

Classification :  $y \in \{-1, 1\}$

$$\hat{y}(x) = f(\sum w_i x_i + b)$$

$$f(z) = \tanh z = \text{sigmoid}(z) - \text{sigmoid}(-z)$$

$$\mathcal{L}(y, \hat{y}) = (y - \hat{y})^2$$

### exemple 2

$$\hat{y}(w) = f\left(\sum_{i=1}^n w_i x_i\right) \quad \boxed{\text{f sigmoid}}$$

$$\mathcal{J}(w) = P(\hat{y}(w) \neq y) \quad f(z) = \frac{1}{1+e^{-z}}$$

$$\text{1er cas: } \mathcal{L}(\hat{y}, y) = -y \log(\hat{y}) - (1-y) \log(1-\hat{y})$$

$$\mathcal{L}(w) = \prod_{i=1}^m \mathcal{L}(w, x^{(i)}, y^{(i)}) \quad \text{constante}$$

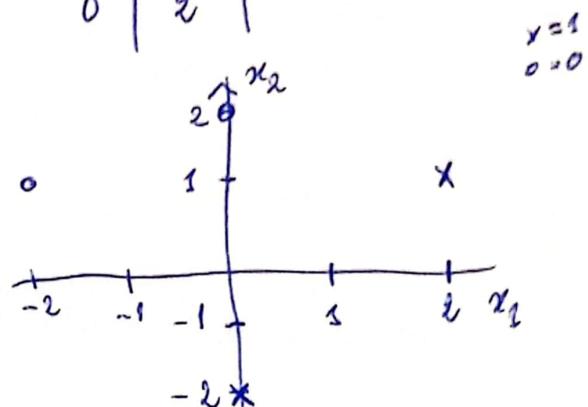
$$\nabla_w \mathcal{J}(w) = \mathcal{L}'(\hat{y}(w), y) \cdot \nabla_{w_i} (\hat{y}) \\ = \left( -\frac{y}{\hat{y}(w)} - \frac{(1-y)}{1-\hat{y}(w)} \right) \cdot \mathcal{L}'\left(\sum_{i=1}^n w_i x_i\right) \otimes x_i$$

$$\text{avec } \mathcal{L}'(z) = \frac{e^{-z}}{(1-e^{-z})^2} = f'(z)(1-f(z))$$

$$\text{2ème cas: } \mathcal{L}(\hat{y}, y) = (y - \hat{y})^2: \text{MSE}$$

$$\nabla_{w_i} \mathcal{J}(w) = \mathcal{L}'(\hat{y}(w), y) \cdot \nabla_{w_i} (\hat{y}) \\ = (\hat{y} - y) \cdot \mathcal{L}'\left(\sum_{i=1}^n w_i x_i\right) \cdot x_i$$

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 2     | 1     | 1   |
| 0     | -2    | 1   |
| -2    | 1     | 0   |
| 0     | 2     | 0   |



2) pour que la perception sans bias peut représenter les données avec la sigmoid activation

$$w = \begin{pmatrix} 1 \\ -1 \end{pmatrix}$$

$$\hat{y}^{(i)} = f(w^\top x^{(i)}) \xrightarrow{z \rightarrow +\infty} y^{(i)}$$

$$J = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \xrightarrow{c \rightarrow \infty} 0$$

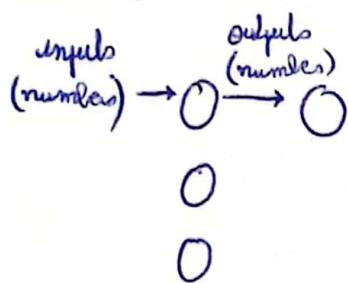
avec  $c = 1$

$$w^\top x^{(i)} \quad \begin{cases} +1 & \text{si } i = 1 \\ +2 & \text{si } i = 2 \\ -1 & \text{si } i = 3 \\ -2 & \text{si } i = 4 \end{cases}$$

# Deep Learning

## Neural Networks

- \* speech recognition
- \* images
- \* text (NLP)

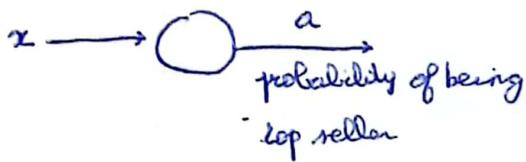


example:

$x = \text{price} : \text{input}$

$$a = f(x) = \frac{1}{1 + e^{-(wx+b)}} : \text{output}$$

activation



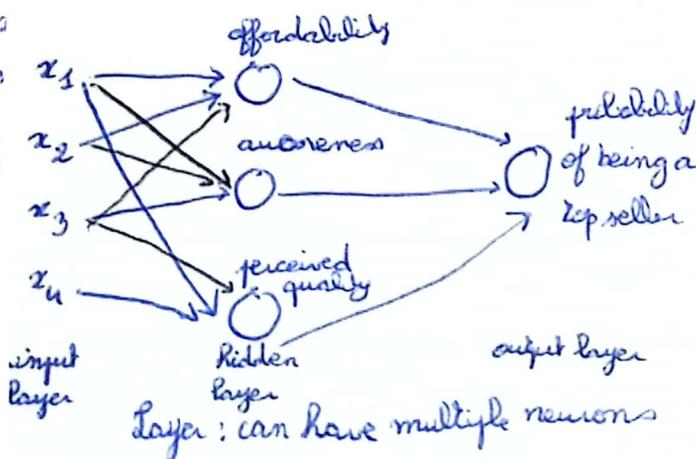
more complex example

$x_1$ : price

$x_2$ : shipping fast

$x_3$ : marketing

$x_4$ : material



affordability      awareness      perceived quality      } activations

$x_1$ ,  $x_2$ ,  $x_3 = x_4$       } input layer

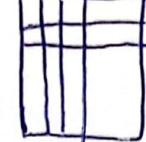
## perceptrons multicouche

### in practice

every neuron will have access to every feature from the previous layer

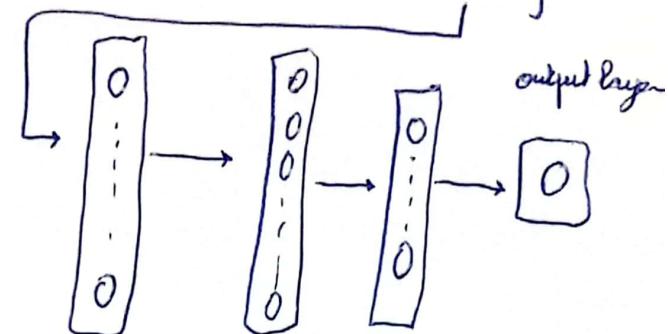
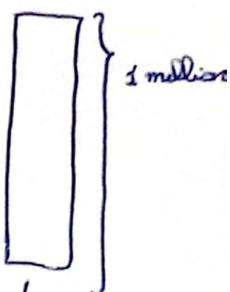
example: image recognition

3000 pixels



1000 pixels

unroll the vector

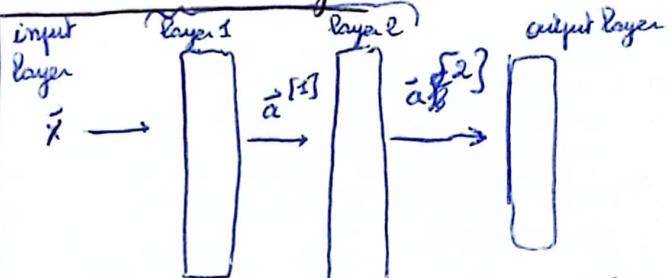


### neurons

first layer: neurons are looking for short lines

second layer: group together short lines in order to look for parts of faces

neural network layer      } hidden layers



the dimension of  $a^{[i]}$  = number

of neurons in layer i

paramètre du réseau = 3 \* 2 Hidden Layers + 1 output layer

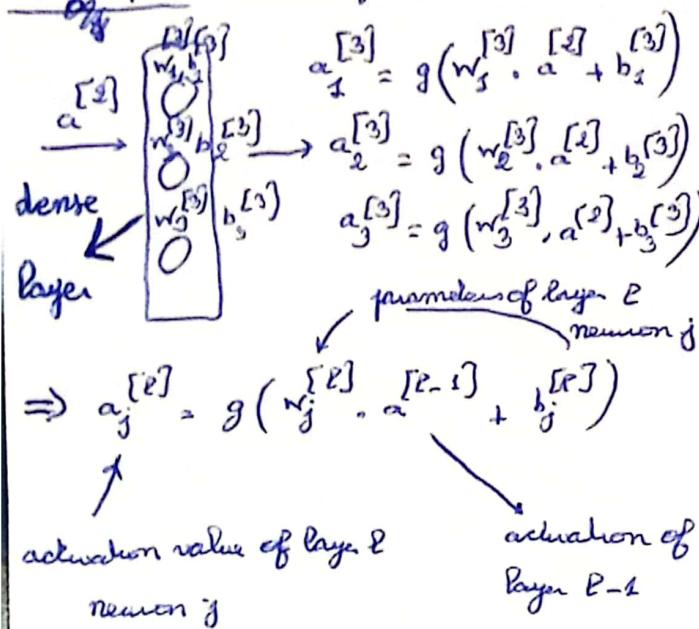
binary choice in output layer:

we need one neuron

## more complex neural networks

4 layers : 3 hidden layers + output layer  
⇒ we don't count the input layer

example : activation function = sigmoid



NB :

$$\vec{x} \rightarrow \vec{a}_1 \rightarrow \vec{a}_2 \rightarrow \vec{a}_3 :$$

forward propagation

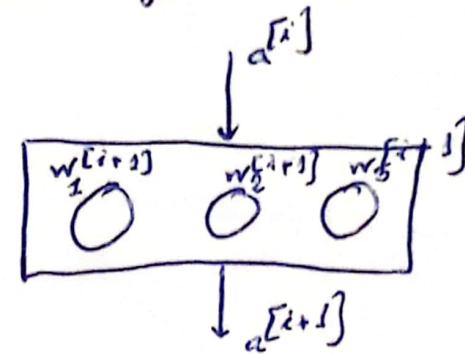
Exemple basique :

regression linéaire simple = un seul neuronne + sans activation

étapes

- 1) forward propagation : de la couche d'entrée à la couche de sortie
- 2) calcul de l'erreur de sortie après la propagation des données
- 3) calcul des gradients d'erreurs pour la couche de sortie ⇒ corriger les poids synaptiques de la couche de sortie
- 4) calcul des gradients d'erreurs pour corriger les poids synaptiques des neurones des couches cachées
- 5) mise à jour des poids synaptiques

1/ dense layer

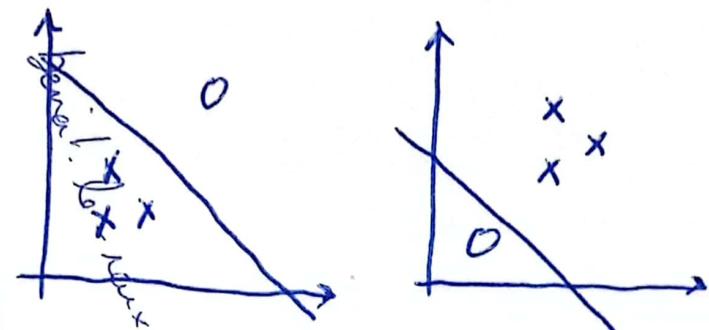


1/ multiplier par le poids

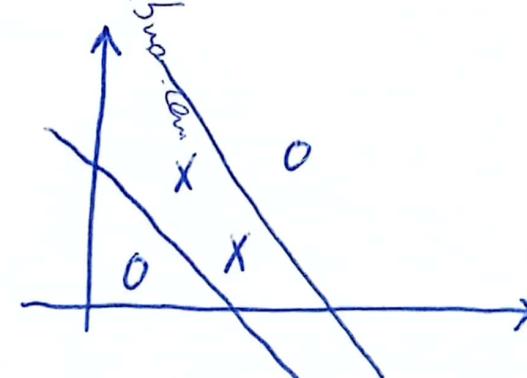
2/ ajouter le bias

3/ fonction d'activation

Remarque un noeud crée un seul banc à de décision



cette classification peut se faire avec un seul noeud



cette classification ne peut pas se faire avec un seul neurone en dense layer

Mais doit se faire en deux

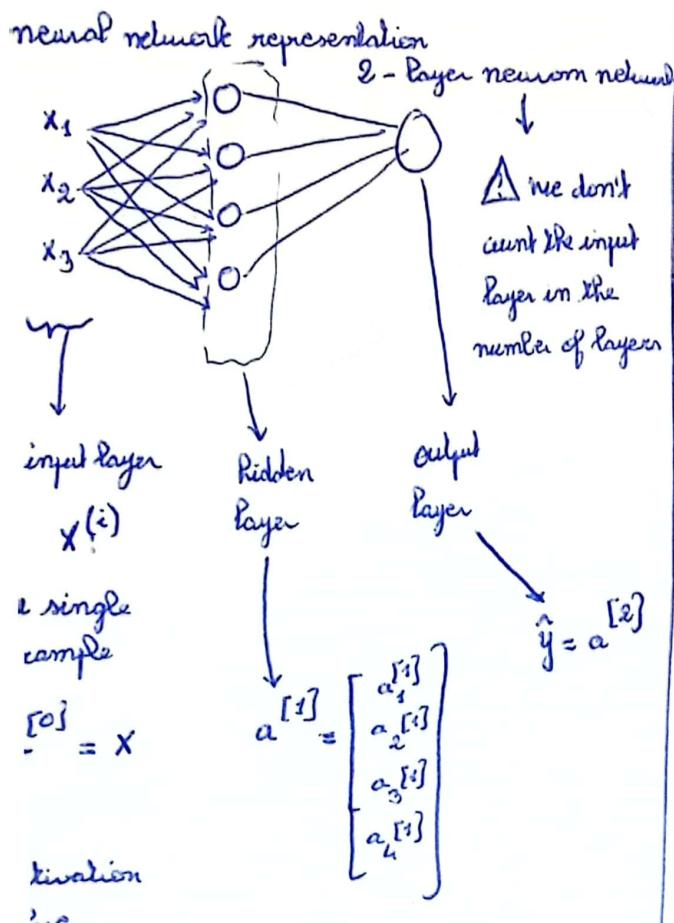
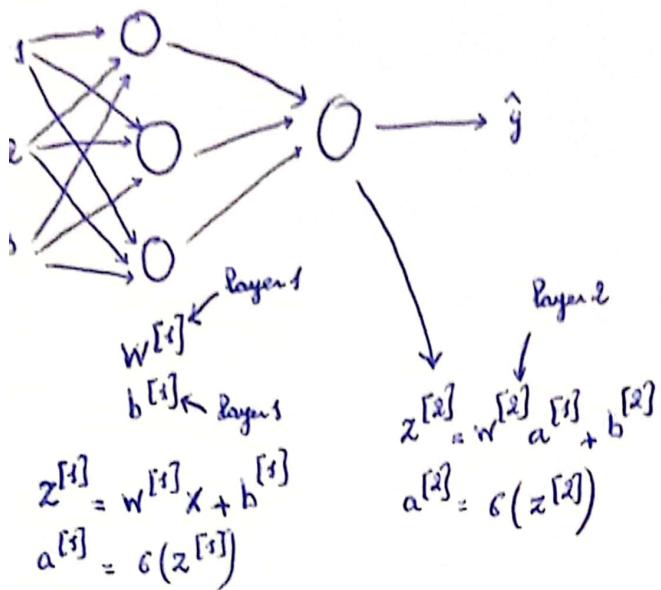
dans ce cas on a besoin :

- 1 dense layer with 2 nodes
- 1 dense layer with 1 node

$$z^i = Wx + b^{(i)}$$

## Neural network

I | only one example



the dimension of  $a^{[i]}$  = number of neurons in layer  $i$

the dimension of  $a^{[0]}$  = number of features

parameters of Layer 1:  $w^{[1]}, b^{[1]}$   
parameters of Layer 2:  $w^{[2]}, b^{[2]}$

in our case:

in Layer 1: we have 4 neurons  
+ 3 inputs

$$\Rightarrow w^{[1]}. \text{shape} = (4, 3)$$

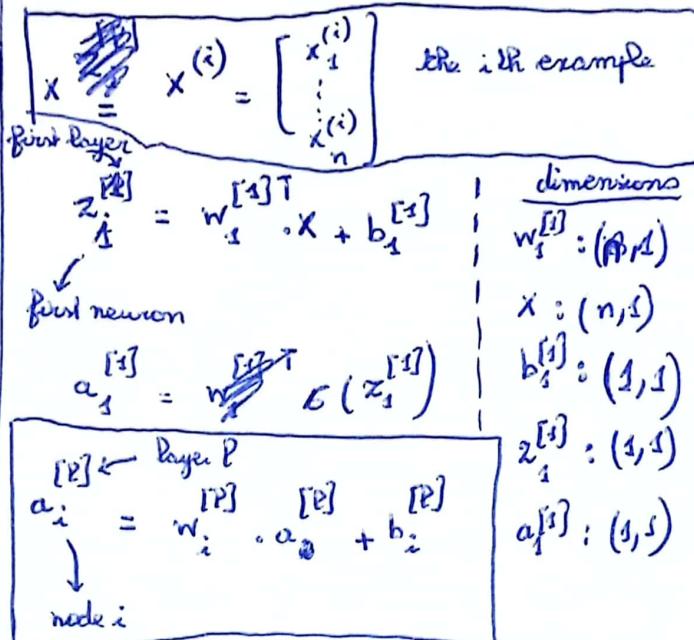
$$b^{[1]}. \text{shape} = (4, 1) \quad ;$$

in Layer 2: we have 1 neuron  
4 inputs

$$w^{[2]}. \text{shape} = (1, 4)$$

$$b^{[2]}. \text{shape} = (1, 1)$$

what does this NN compute



$$z^{[1]}_1 = w^{[1]T} x + b^{[1]}$$

$$z^{[2]}_2 = w^{[2]T} x + b^{[2]}$$

$$z^{[3]}_3 = w^{[3]T} x + b^{[3]}$$

$$z^{[4]}_4 = w^{[4]T} x + b^{[4]}$$

$$\Leftrightarrow w^{[i]} = \begin{bmatrix} -w_1^{[i]T} \\ -w_2^{[i]T} \\ -w_3^{[i]T} \\ -w_4^{[i]T} \end{bmatrix}$$

$$z^{[i]} = \begin{bmatrix} z^{[1]}_1 \\ z^{[2]}_2 \\ z^{[3]}_3 \\ z^{[4]}_4 \end{bmatrix}$$

$$\begin{aligned}
 z^{[3]} &= w^{[3]} x + b^{[3]} \\
 a^{[3]} &= \sigma(z^{[3]}) \\
 z^{[2]} &= w^{[2]} a^{[3]} + b^{[2]} \\
 a^{[2]} &= \sigma(z^{[2]}) \\
 &\vdots
 \end{aligned}$$

dimensions

$$x \in (n, 1)$$

$$z^{[i]} : (n^{[i]}, n^{[i]}, 1)$$

$$w^{[i]} : (n^{[i]}, n)$$

$$b^{[i]} : (n^{[i]}, 1)$$

$$a^{[i]} : (n^{[i]}, 1)$$

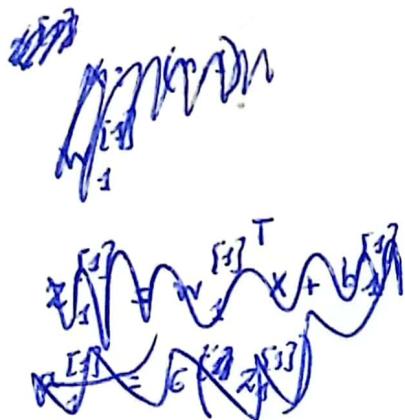
$$w^{[i]} : (n^{[i]}, n^{[i]})$$

$$b^{[i]} : (n^{[i]}, 1)$$

$$z^{[i]} : (n^{[i]}, 1)$$

$$a^{[i]} : (n^{[i]}, 1)$$

⚠ The dimension of  $a^{[i]} = \text{number of neurons in layer } i$



II - across multiple training examples

We have  $m$  training examples

$$x^{(1)}, \dots, x^{(m)}$$

$$x^{(i)} \rightarrow a^{[2]}(i) = \hat{y}^{(i)}$$

$$\vdots$$

$$x^{(m)} \rightarrow a^{[2]}(m) = \hat{y}^{(m)}$$

for  $i = 1 \dots m$  # loop over  $m$  training examples

$$z^{[1]}(i) = w^{[1]} x^{(i)} + b^{[1]}$$

$$a^{[1]}(i) = \sigma(z^{[1]}(i))$$

$$z^{[2]}(i) = w^{[2]} a^{[1]}(i) + b^{[2]}$$

$$a^{[2]}(i) = \sigma(z^{[2]}(i))$$

$w^{[1]}, w^{[2]}, b^{[1]}, b^{[2]}$  re-used per iteration & individual

~~w<sup>[1]</sup>(i)~~

dimensions

$$X: (n, m)$$

$$w^{[1]}: (n^{[1]}, n)$$

~~w<sup>[1]</sup>~~

$$A^{[1]} \text{ et } Z^{[1]}: (n^{[1]}, m)$$

$$w^{[2]}: (n^{[2]}, n^{[1]})$$

$$A^{[2]} \text{ et } Z^{[2]}: (n^{[2]}, m)$$

Vectorization

$$X = \begin{bmatrix} x^{(1)} \\ \vdots \\ x^{(m)} \end{bmatrix} \in \mathbb{R}_{n, m}$$

from  $x^{[i]} \rightarrow Z^{[i]} = \begin{bmatrix} z^{[1]}(1) \\ \vdots \\ z^{[1]}(m) \end{bmatrix}$

$$A^{[1]} = W^{[1]} X + b^{[1]}$$

$$A^{[1]} = \sigma(Z^{[1]})$$

$$Z^{[2]} = W^{[2]} X + b^{[2]}$$

$$A^{[2]} = \sigma(Z^{[2]})$$

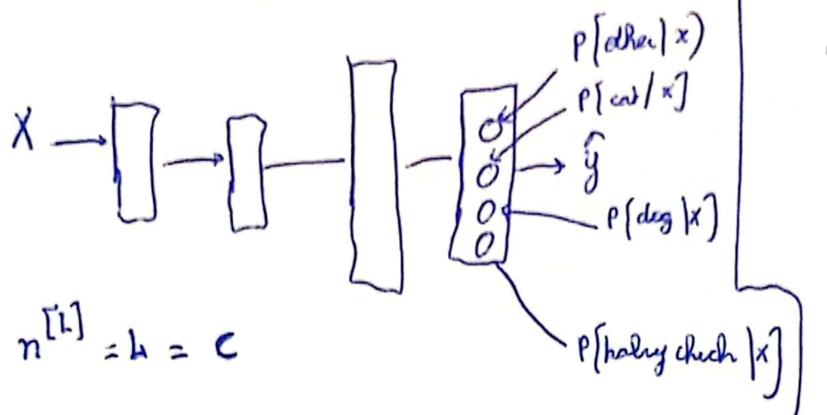
four lines per example

## softmax Regression for multiclassification

example: cat vs dogs vs baby chick vs other

1      2      3      4

$$c = \# \text{ classes} = 4$$



$$z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$$

we have



$$a^{[l]} = \frac{e^{(z^{[l]})}}{\sum_{i=1}^c z_i}$$

⚠  $y$  must be also  $(c, s)$

if we have a cat  $y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$   $\hat{y} = \begin{bmatrix} 0.3 \\ 0.6 \\ 0.1 \\ 0.4 \end{bmatrix}$

$$\mathcal{L}(\hat{y}, y) = - \sum_{j=1}^h y_j \log(\hat{y}_j)$$

$$= -y_2 \log(\hat{y}_2)$$

$$\mathcal{J}(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$



8/11

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \dots$$

⚠ if we have

$$y = [1, 2, 0, \dots]$$

on view

$$y = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \dots$$

we use one-hot encodings

in TensorFlow

$$y\_one\_hot = \text{tf\_one\_hot}(\text{label}, c, \text{axis}=0, \text{depth}, 1))$$

app didn't apply to our example

if we have our labels in this form

$$y = \begin{bmatrix} y_1^{(1)}, y_1^{(2)}, \dots, y_1^{(m)} \end{bmatrix} (c, m)$$

$$\hat{y} = \begin{bmatrix} \hat{y}_1^{(1)}, \hat{y}_1^{(2)}, \dots, \hat{y}_1^{(m)} \end{bmatrix}$$

on tensor flow:

to compute categorical cross entropy:

- about the shape:

$$y = \begin{bmatrix} \overbrace{\quad \quad y^{(1)} \quad \quad} \\ \overbrace{\quad \quad y^{(2)} \quad \quad} \\ \vdots \\ \overbrace{\quad \quad y^{(m)} \quad \quad} \end{bmatrix} (m, c)$$

- about logits

if in the forward propagation

1/ we finished with  $Z^{[L]}$  without calculating

$$A^{[L]}$$

$\Rightarrow$  we need to set `from_logits = True` in  
the parameters of the cross entropy

2/ we finished by  $A^{[L]}$

$\Rightarrow$  set `from_logits = False`

## Multiclass classification with neural networks

### ① Logistic Regression

(2 possible output values)

$$z = \vec{w} \cdot \vec{x} + b$$

$$\alpha_1 = g(\beta) = \frac{e^\beta}{1+e^{-\beta}} = P(y=1 | \vec{x})$$

$$\alpha_2 = 1 - \alpha_1 = P(y=0 | \vec{x})$$

#### cost function

$$\text{Loss} = -y \log(\alpha_1) - (1-y) \log(\alpha_2)$$

$$J(\vec{w}, b) \approx \text{average loss}$$

### ② softmax regression

$$y = 1, 2, 3, 4$$

$$z_1 = \vec{w}_1 \cdot \vec{x} + b_1$$

$$\alpha_1 = \frac{e^{z_1}}{\sum_{j=1}^N e^{z_j}} = P(y=1 | \vec{x})$$

$$z_2 = \vec{w}_2 \cdot \vec{x} + b_2$$

$$\alpha_2 = \frac{e^{z_2}}{\sum_{j=1}^N e^{z_j}} = P(y=2 | \vec{x})$$

$$z_3 = \vec{w}_3 \cdot \vec{x} + b_3$$

$$\alpha_3 = \frac{e^{z_3}}{\sum_{j=1}^N e^{z_j}} = P(y=3 | \vec{x})$$

$$z_4 = \vec{w}_4 \cdot \vec{x} + b_4$$

$$\alpha_4 = \frac{e^{z_4}}{\sum_{j=1}^N e^{z_j}} = P(y=4 | \vec{x})$$

### ③ softmax Regression generalization

$$z_j = \vec{w}_j \cdot \vec{x} + b_j \quad j=1, \dots, N$$

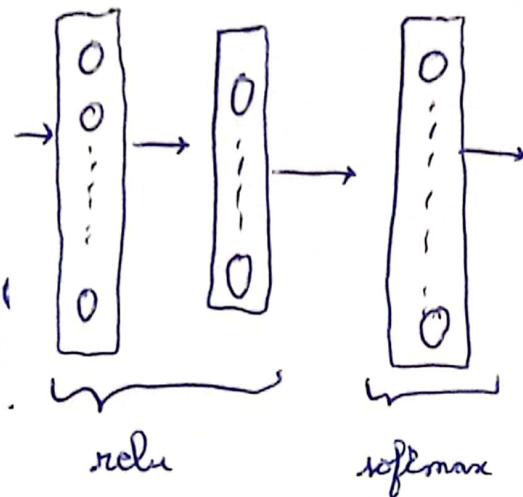
$$\alpha_j = \frac{e^{z_j}}{\sum_{k=1}^N e^{z_k}} = P(y=j | \vec{x})$$

$$\sum_{i=1}^N \alpha_i = 1$$

#### cost

$$\text{Loss}(\alpha_1, \dots, \alpha_N, y) = \begin{cases} -\log \alpha_1 & \text{if } y=1 \\ -\log \alpha_2 & \text{if } y=2 \\ \vdots \\ -\log \alpha_N & \text{if } y=N \end{cases}$$

## softmax with NN



## Logistic Regression

$$a_1^{[3]} = g(\beta_0^{[3]}) \quad \text{et} \quad a_2^{[3]} = g(\beta_0^{[3]})$$

## Softmax Regression

$$a^{[3]} = (a_1^{[3]}, \dots, a_{30}^{[3]}) = g(\beta_0^{[3]}, \dots, \beta_{30}^{[3]})$$

## numerical Roundoff errors

$$x = \frac{2.0}{10000} \rightarrow 0,00020000$$

$$x_2 = 1 + (1/10000) - (1 - 1/10000)$$

$$\hookrightarrow 0,000199999\dots$$

## improved implementation of Logistic Loss

① on calcule  $a$  et puis on calcule le loss

$$a = \frac{1}{1+e^{-y}} \quad \text{enfonction de } a$$

$$\text{Loss} = -y \log(a) - (1-y) \log(1-a)$$

② improved version : ne pas calculer  $a$  en intermédiaire

$$\text{Loss} = -y \log\left(\frac{1}{1+e^{-y}}\right) - (1-y) \log\left(1 - \frac{1}{1+e^{-y}}\right)$$

import tensorflow as tf

from tensorflow.keras import Sequential

from tensorflow.keras.layers import Dense

model = Sequential ([

Dense (units = 25, activation = 'relu'),

Dense (units = 10, activation = 'relu'),

Dense (units = 10, activation = 'softmax'),

])

from tensorflow.keras.losses import

SparseCategoricalCrossentropy

model.compile (loss = SparseCategoricalCrossentropy ())

model.fit (X, y, epochs = 100)

⇒ the model output is  $a_1, \dots, a_{30}$

improved version SparseCategorical

① model.compile (loss = ~~SparseCategoricalCrossEntropy~~ (from\_logits = True))

② set the output layer activation function to 'linear'

Logits =  $y$