

I. Quelques résultats

1. Origine de la sigmoïde :

$$\begin{cases} P(X=1) = p \\ P(X=0) = 1-p \end{cases}$$

$$\begin{cases} p(y|x=1) = \mathcal{N}(y; m_1; \sigma^2) \\ p(y|x=0) = \mathcal{N}(y; m_2; \sigma^2) \end{cases}$$

$$\Rightarrow P(X=1|y) = \frac{P(X=1) \cdot p(y|x=1)}{P(X=1) p(y|x=1) + P(X=0) p(y|x=0)}$$

$$\ln \frac{P(X=1|y)}{1 - P(X=1|y)} = ax + b$$

$$= \frac{p \cdot \mathcal{N}(y; m_1; \sigma^2)}{p \cdot \mathcal{N}(y; m_1; \sigma^2) + (1-p) \cdot \mathcal{N}(y; m_2; \sigma^2)}$$

$$= \frac{1}{1 + \frac{(1-p)}{p} \times e^{-\frac{1}{2\sigma^2} (y^2 - 2m_2 y + m_2^2 - y^2 + 2m_1 y - m_1^2)}}$$

$$= \frac{1}{1 + \frac{(1-p)}{p} \times e^{-\frac{1}{2\sigma^2} (y^2 - 2m_2 y + m_2^2 - y^2 + 2m_1 y - m_1^2)}}$$

$$= \frac{1}{1 + e^{\ln\left(\frac{1-p}{p}\right) - (a'y + b')}} = \frac{1}{1 + e^{-(a'y + b')}} = \text{Sig} (a'y + b)$$

$$= \text{Sig} (a'y + b)$$

2. Derivées

(2)

- $g(x) = \text{sign}(x)$

$$\frac{dg}{dx} = g(x) (1 - g(x))$$

- $f(z_1, \dots, z_K) ; \quad h(x) = f(g_1(x), \dots, g_K(x))$

$$\frac{dh}{dx} = \sum_{i=1}^K \frac{dg_i}{dx} \times \frac{df}{dg_i}$$

- $f_i(z_1, \dots, z_K) = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}} = \text{softmax}(z_1, \dots, z_K)$

$$* \frac{df_i}{dz_i} = \frac{e^{z_i} \times \sum_{k=1}^K e^{z_k} - e^{z_i} \times e^{z_i}}{\sum_k e^{z_k} \times \sum_k e^{z_k}}$$

$$= f_i - f_i^2 = f_i (1 - f_i)$$

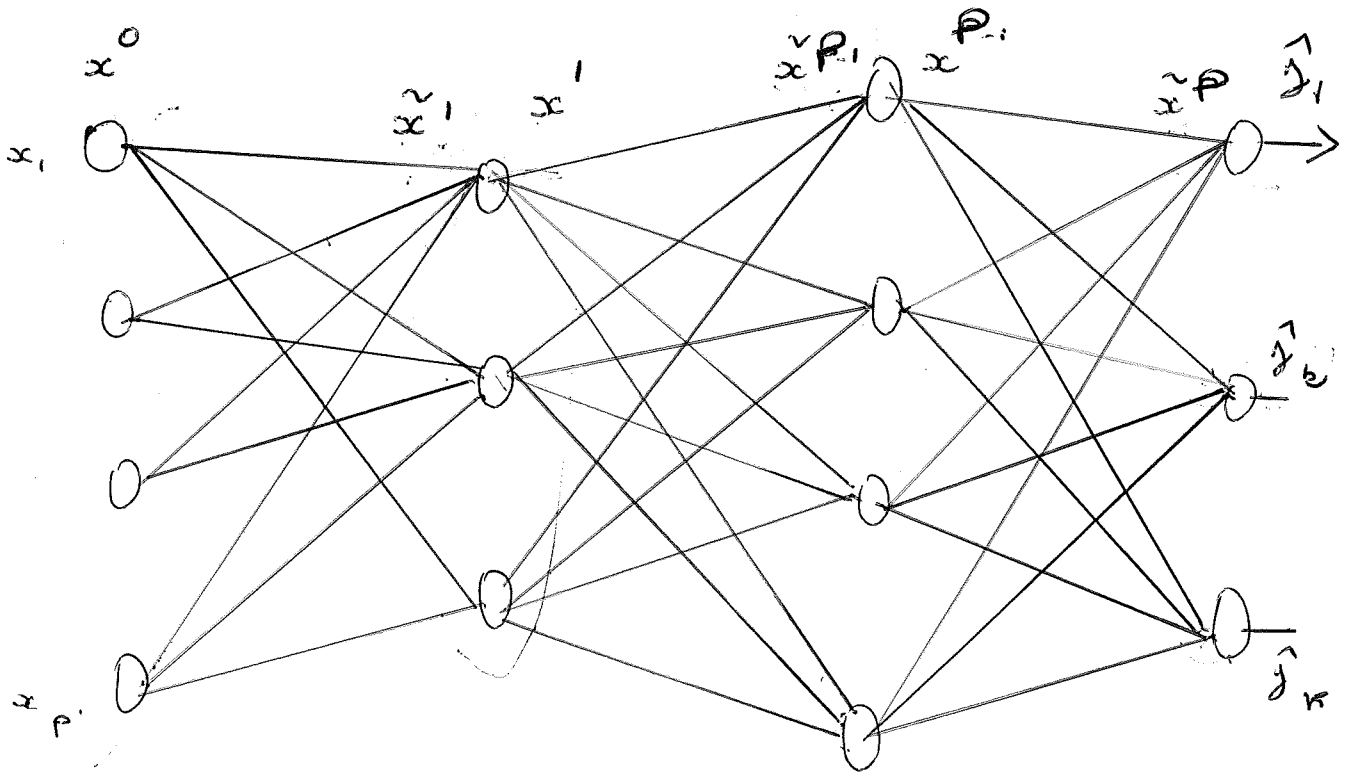
$$\frac{df_i}{dz_j} = \frac{-e^{z_i} e^{z_j}}{\sum_{k=1}^K e^{z_k} \times \sum_{k=1}^K e^{z_k}} = -f_i \times f_j$$

B. Réseau de neurones pour classif multiclasse (3)

objectif: modéliser $P(Y=k | X=x) = f_k(x)$

$$\Rightarrow f_k(x) \rightarrow f_{k,0}(x) \text{ tq } \sum_{k=1}^K f_{k,0}(x) = 1$$

Réseau de neurones



$$\tilde{x}_j^P = w_j^P \cdot x^{P-1} + b_j^P \quad w_j^P = \text{vecteur de poids}$$

$$x_j^P = g(\tilde{x}_j^P) \quad \text{avec } g(x) = \text{sign}(x)$$

$$\hat{j}_k = \text{softmax}_k(\tilde{x}_1^P, \dots, \tilde{x}_K^P)$$

• Objectif: estimer $\theta = (w^P, b^P)$ pour $p = 1: P$. (4)

à partir de $(x^{(1)}, \dots, x^{(n)})$
 $y^{(1)}, \dots, y^{(n)}$

• Démarche: Maximisons la vraisemblance.

Pour une donnée (x, y) : $p(x, y) = p(x) \cdot p(y|x)$

$$\log p(x, y) = \log p(x) + \log p_{\theta}(y|x)$$

\Rightarrow maximisons $\sum_{i=1}^n \log p_{\theta}(y^{(i)} | x^{(i)})$

calculons $\nabla_{\theta} \log p(y|x)$ pour (x, y)

$$y = i \Leftrightarrow y = [0, 0, \underset{\substack{\uparrow \\ \text{ième composant}}}{1}, \dots, 0]$$

$$p_{\theta}(y|x) = \prod_{k=1}^K p_{\theta}(y=k|x)^{y_k}$$

$$\Rightarrow \log p_{\theta}(y|x) = \sum_{k=1}^K y_k \log p_{\theta}(y=k|x)$$

$$= \sum_{k=1}^K y_k \log \hat{y}_k$$

$$= -L(\hat{y}_1, \dots, \hat{y}_K)$$

minimisons $\frac{d}{d\theta} = \sum_{k=1}^K y_k \log \hat{y}_k = y / \theta \hat{\theta}$

$$\frac{d\mathcal{L}}{dw_j^P} = \frac{d\mathcal{L}}{d\tilde{x}_j^P} \times \underbrace{\frac{d\tilde{x}_j^P}{dw_j^P}}_{\text{trivial}}$$

(5)

$$\frac{d\mathcal{L}}{d\tilde{x}_j^P} = \sum_{k=1}^K - \frac{y_k}{\hat{y}_k} \times \frac{d\hat{y}_k}{d\tilde{x}_j^P}$$

AVEC $\hat{y}_k = \text{softmax}_k(\tilde{x}_1^P, \dots, \tilde{x}_K^P)$

$$= \sum_{k \neq j} - \frac{y_k}{\hat{y}_k} \cdot \frac{d\hat{y}_k}{d\tilde{x}_j^P} - \frac{y_j}{\hat{y}_j} \cdot \frac{d\hat{y}_j}{d\tilde{x}_j^P}$$

$$= \sum_{k \neq j} + \frac{y_k}{\cancel{\hat{y}_k}} \times \cancel{\hat{y}_k} \times \hat{y}_j - \frac{y_j}{\cancel{\hat{y}_j}} \cancel{\hat{y}_j} (1 - \hat{y}_j)$$

$$= \sum_{k \neq j} \hat{y}_k \hat{y}_j - y_j + \hat{y}_j \hat{y}_j$$

$$= \sum_{k=1}^K \hat{y}_k \hat{y}_j - y_j = \hat{y}_j \sum_{k=1}^K \hat{y}_k - y_j$$

$$\left\| \frac{d\mathcal{L}}{d\tilde{x}_j^P} = \hat{y}_j - y_j = c_j^P \right.$$

$$\Rightarrow \left\| \begin{aligned} \frac{d\mathcal{L}}{dw_j^P} &= c_j^P \times x^{P-1} \\ \frac{d\mathcal{L}}{db_j^P} &= c_j^P \end{aligned} \right.$$

• $\frac{d\mathcal{L}}{dw_j^P}$ / P quelconque ?

(6)

$$\frac{d\mathcal{L}}{d\tilde{x}_j^P} = \sum_P \frac{d\mathcal{L}}{d\tilde{x}_P^{P+1}} \times \frac{d\tilde{x}_P^{P+1}}{d\tilde{x}_j^P} \times \frac{dx_j^P}{d\tilde{x}_j^P}$$

$$c_j^P = \sum_P c_P^{P+1} \times w_{jP}^{P+1} \times x_j^P (1-x_j)^P$$

$$\Rightarrow \frac{d\mathcal{L}}{dw_j^P} = \frac{d\mathcal{L}}{d\tilde{x}_j^P} \times \frac{d\tilde{x}_j^P}{dw_j^P}$$

$$\left\{ \begin{array}{l} \frac{d\mathcal{L}}{dw_j^P} = c_j^P \times x^{P-1} \\ \frac{d\mathcal{L}}{db_j^P} = c_j^P \end{array} \right.$$

Algorithme de rétropropagation du gradient

• ex: Écrire le pseudo code (matriciellement) avec les n contraintes, que pour train-RBM.

II Préapprentissage de DNN

(7)

considérer les couches observées + cachées
(couche de classif exclue) comme 1 DNN.

- Apprendre de manière non supervisée le DNN via greedy Layer wise procédure, pour ~~init~~ initialiser les paramètres du DNN
- Fine tuning : apprentissage supervisé du DNN par rétropropagation du gradient.

