

I. Rappels: estimation de paramètres dans des modèles à variables latentes

Soit x v.a observée $\in \mathbb{R}$
 h v.a cachée $\in \mathbb{R}$

• On se donne $p_\theta(x, h)$ dépend d'un paramètre θ qu'on aimerait apprendre à partir de $(x^{(1)}, \dots, x^{(n)})$ iid

$$\begin{aligned} \hat{\theta}_{MV} &= \arg \max_{\theta} \prod_{i=1}^n p_\theta(x^{(i)}) = \arg \max_{\theta} p_\theta(x^{(1)}, \dots, x^{(n)}) \\ &= \arg \max_{\theta} \sum_{i=1}^n \log p_\theta(x^{(i)}) \end{aligned}$$

avec $p_\theta(x) = \int p_\theta(x, h) dh$

est l'estimateur du MV : bonnes propriétés asymptotiques

• minimise la DKL $\int p(x) \log \frac{p(x)}{p_\theta(x)} dx$

où $p(x)$ est la vraie loi des observations

ie $\hat{\theta}_{MV} \xrightarrow[n \rightarrow \infty]{\text{probab}} \arg \min_{\theta} DKL(p, p_\theta)$

\Rightarrow difficile à appliquer si $\int p_\theta(x, h) dh$ inconnue

• Algorithme EM : itératif (pour ! ∞) (2)

partant de $\theta^{(i)}$:

$$E: E_{\theta^{(i)}} [\log p_{\theta}(x, h) | x] = Q(\theta^{(i)}, \theta)$$

$$M: \theta^{(i+1)} = \underset{\theta}{\text{Arg max}} Q(\theta^{(i)}, \theta)$$

$$Rq: E: = \int \log p_{\theta}(x, h) \cdot p_{\theta^{(i)}}(h|x) dh$$

$$\left[\text{Th de Transfert} \quad E(f(x, y)) = \iint f(x, y) p(x, y) dx dy \right]$$

On a tout de même besoin de la loi

$$\text{à posteriori: } p_{\theta^{(i)}}(h|x) = \frac{p_{\theta^{(i)}}(x, h)}{p_{\theta^{(i)}}(x)}$$

Vocab: $p(h) \rightarrow$ prior

$p(x|h) \rightarrow$ likelihood conditionnelle

$p(x) \rightarrow$ likelihood

$p(h|x) \rightarrow$ posterior

Problème: on a besoin de $p(h|x)$ pas évident à calculer.

Soit $q_S(h|x)$ une loi quelconque donnée de paramètre S .

$$\begin{aligned} \text{DKL}(q_S(h|x), p_\theta(h|x)) &= \int \log\left(\frac{q_S(h|x)}{p_\theta(h|x)}\right) q_S(h|x) dh \\ &= \int \log\left(\frac{q_S(h|x)}{p_\theta(h|x)}\right) q_S(h|x) + \int \log(p_\theta(x)) q_S(h|x) dh \\ &= \int \log\left(\frac{q_S(h|x)}{p_\theta(h|x)}\right) q_S(h|x) + \log p_\theta(x) \end{aligned}$$

$$\Rightarrow \log p_\theta(x) = - \int \log\left(\frac{q_S(h|x)}{p_\theta(h|x)}\right) q_S(h|x) + \underbrace{\text{DKL}(q_S(h|x), p_\theta)}_{\geq 0}$$

$$\Rightarrow \log p_\theta(x) \geq - \int \log\left(\frac{q_S(h|x)}{p_\theta(h|x)}\right) q_S(h|x) = \text{ELBO}(\theta, S)$$

Rq: à à $\theta = \theta^{(i)}$ fixé

Alors $q_S^{(i)}(h|x) = p_{\theta^{(i)}}(h|x)$ maximise $\text{ELBO}(\theta^{(i)}, S)$ à $\theta \in \Theta$

$$\hat{a} \quad q_S(h|x) = p_{\theta^{(i)}}(h|x)$$

$$\text{ELBO}(\theta) = \int \log p_\theta(h|x) p_{\theta^{(i)}}(h|x)$$

À maximiser / à $\hat{a} \in \Theta$

On retrouve EM.

- on se donne une classe de q_s paramétrée par s en faisant en sorte $q_s(h|x)$ "ne soit pas trop éloignée de" $p_\theta(h|x)$ et on optimise $ELBO(\theta, s)$ / r à (θ, s) (gradient) de manière à pousser $p_{q_s}(x)$ vers la gauche.

II Deep. Belief Network

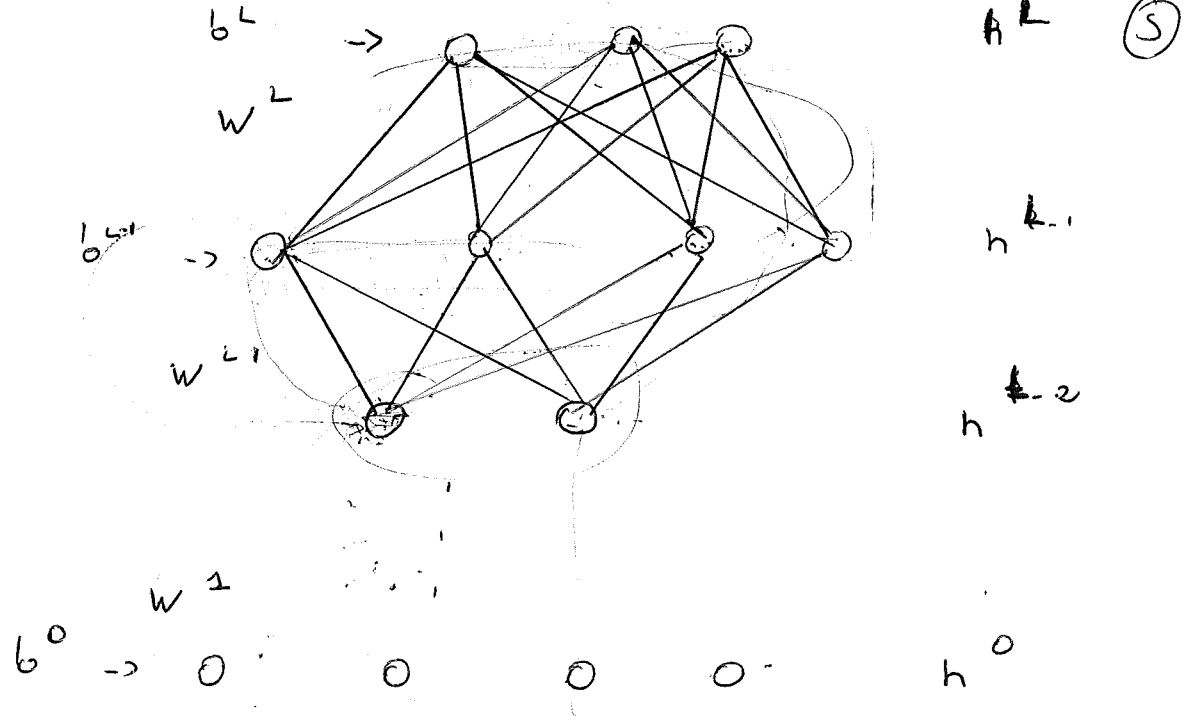
Soit ~~$v = \{v_i\}$~~ $V = H^0 = (V_1, \dots, V_P) \in \{0, 1\}^P$
 $H^P = (H_1^P, \dots, H_{q_P}^P) \in \{0, 1\}^{q_P}$
 pour $P = 1: L$

Déf: $p(h^0, h^1, \dots, h^L) = p(h^{L-1}, h^L) \prod_{i=1}^{L-1} p(h^{i-1} | h^{i+1})$

où i) $p(h^{L-1}, h^L)$ est un RBM
 $= \frac{1}{Z} e^{-E(h^{L-1}, h^L)}$

ii) $p(h^{i-1} | h^{i+1}) = \prod_{j=1}^{q_{i-1}} p(h_j^{i-1} | h^{i+1})$

iii) $p(h_j^{i-1} = 1 | h^i) = \text{sign}(b_j^{i-1} + \sum_k w_{jk}^i h_k^i)$



on a fabriqué un modèle génératif

$$p(v) = p(h^0) = \sum_{h^1, \dots, h^L} p(h^0, h^1, \dots, h^L)$$

et qui généralise le RBN (on prend $L=1 \Rightarrow \text{RBN}$)

Rq: $p(h^{p-1} | h^p) = p_{\text{RBN}}(h^{p-1} | h^p)$

mais (h^{p-1}, h^p) n'est pas un RBN
(sauf pour $p=L$)

car $p(h^p | h^{p-1}) \neq p_{\text{RBN}}(h^p | h^{p-1})$

(6)

• Estimer les paramètres d'un DBN à partir
 $(x^{(1)}, \dots, x^{(n)})$ iid

i) estimateur du MV: $p_{\theta}(x) = \sum_{\substack{h^1, \dots, h^L}} p(x, h^1, \dots, h^L)$

$$\theta = (b^i, w^i)_{i=1:L} + b^0$$

\Rightarrow somme sur h^1, \dots, h^L, \dots

ii) en $\Rightarrow p(h^1, \dots, h^L | x) = \frac{p(x, \dots, h^L)}{\sum_{h^1, \dots, h^L} p(x, \dots, h^L)}$

iii) Variational inference.

Idee: DBN comme un empilement de RBM.

$$q(h_j^1 | x) = \text{sign}(\tilde{b}_j^1 + \sum_i w_{ij}^1 x_i)$$

$$q(h_j^p | x) = q(h_j^p | h^{p-1}) = \text{sign}(\tilde{b}_j^p + \sum_i w_{ij}^p h_i^{p-1})$$

où h^{p-1} est $h^{p-1} = f(x)$ où on

a appliqué les transformations précédentes.

Une fois cette loi q variationnelle donnée, on aq (7)
l'optimisation de la ELBO se déduit :

Algorithme Greedy, Layer wise procédure

x : utiliser CD-1 pour apprendre

b^0, w^1, \tilde{b}^1

calculer la sortie \tilde{x} de x avec ce RBM

\tilde{x} utiliser CD-1 pour apprendre

b^1, w^2, \tilde{b}^2

\vdots

(Learning Deep Architectures for AI.)

Pseudo-code $x, \text{DBN} \rightarrow \text{liste de RBM}$

train-DBN()

for $i = 1 : \text{nb-couches}$

DBN[i] = train-RBM(x, \dots, \dots)

$x = \text{entree-sortie}(x, \text{DBN}[i])$

fin

