

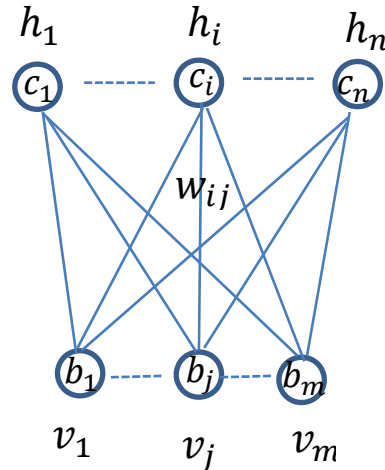
Examen Apprentissage Statistique

Master DAC (Informatique) et Spécialité Big Data (master de Math.)

10/02/2015

Durée 2 h, notes de cours autorisées

Les notations en **gras** désignent des vecteurs et les notations normales des scalaires.



Une machine de Boltzmann restreinte (RBM) est un réseau de neurones avec une couche d'unités visibles \mathbf{V} et une couche d'unités cachées \mathbf{H} (voir Figure). On va associer à ce réseau de neurones un modèle probabiliste. Les unités visibles et les unités cachées sont des variables aléatoires, notées respectivement $\mathbf{V} = (V_1, \dots, V_m)$ et $\mathbf{H} = (H_1, \dots, H_n)$. On considère par la suite des RBM binaires où V_i et H_i prennent des valeurs $v_j \in \{0,1\}$ et $h_i \in \{0,1\}$.

La distribution jointe de V et H est décrite par une distribution de Gibbs $p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})}$ où $E(\mathbf{v}, \mathbf{h})$ est une fonction énergie qui a la forme suivante :

$$E(\mathbf{v}, \mathbf{h}) = -\mathbf{h}^T \mathbf{W} \mathbf{v} - \mathbf{c}^T \mathbf{v} - \mathbf{b}^T \mathbf{h} = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{j=1}^m b_j v_j - \sum_{i=1}^n c_i h_i$$

$w_{ij} \in \mathbb{R}$ est le poids de la connexion entre la cellule visible V_j et la cellule cachée H_i , $b_j \in \mathbb{R}$ et $v_i \in \mathbb{R}$ sont des termes de biais associés à V_j et H_i , et Z est la fonction de partition¹ $Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$ qui sert à normaliser les probabilités de façon à ce qu'elles somment à 1.

Dans ce modèle, les variables cachées sont conditionnellement indépendantes connaissant les variables visibles et réciproquement :

$$p(\mathbf{h}|\mathbf{v}) = \prod_{i=1}^n p(h_i|\mathbf{v}) \text{ et } p(\mathbf{v}|\mathbf{h}) = \prod_{j=1}^m p(v_j|\mathbf{h})$$

¹ Ici la somme est calculée sur toutes les valeurs possibles des vecteurs \mathbf{v}' et \mathbf{h}' . Si ils sont binaires, de taille respective n et m , il y aura 2^m valeurs possibles pour \mathbf{v}' et 2^n pour \mathbf{h}' .

Une RBM va être entraînée à apprendre la distribution jointe $p(V, H)$. Pour cela on va maximiser la vraisemblance du modèle sur un ensemble d'exemples d'apprentissage, avec comme paramètres du modèle les w, b et c . Une fois entraînée, elle pourra servir à effectuer des inférences

1. Inférence

On suppose que la RBM a été entraînée sur un corpus d'apprentissage. Elle modélise alors une distribution $p(V, H)$ apprise sur ces données. Pour une donnée v , on peut calculer la probabilité $p(v)$ que cette donnée ait été générée par le modèle appris.

Pour cela on a besoin de calculer la probabilité marginale $p(v) = \sum_{\mathbf{h}} p(v, \mathbf{h})$ où la somme est prise sur toutes les valeurs possible de $\mathbf{h} \in \{0,1\}^n$.

1. On peut réécrire $p(v)$ sous la forme $p(v) = \sum_{\mathbf{h}} p(v|\mathbf{h})p(\mathbf{h})$. A quel type de modèle est ce que cela vous fait penser ? Combien y-a-t-il de composantes et est-ce que ce calcul direct vous paraît raisonnable en terme de complexité ?
2. On va regarder plus précisément comment se décompose $p(v)$. Montrer

$$p(v) = \frac{1}{Z} \sum_{h_1=0}^1 \dots \sum_{h_n=0}^1 e^{\sum_{j=1}^m b_j v_j} \prod_{i=1}^n e^{h_i(c_i + \sum_{j=1}^m w_{ij} v_j)}$$

puis

$$p(v) = \frac{1}{Z} \prod_{j=1}^m e^{b_j v_j} \prod_{i=1}^n (1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j})$$

3. Quelle est la complexité du calcul du numérateur $\prod_{j=1}^m e^{b_j v_j} \prod_{i=1}^n (1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j})$ et celle du dénominateur $Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$? Commenter.
4. Classification

Plusieurs méthodes sont possibles si on veut utiliser ce modèle pour faire de la classification.

- 4.1. on peut entraîner une RBM par classe et ensuite présenter une donnée test à chacune des RBM. Si on suppose que l'on a 2 classes, on notera $p_1(V)$ et $p_2(V)$ les 2 distributions apprises.
 - Comment peut-on effectuer la classification d'une donnée de test x ?
 - Quelle est la complexité de ce classifieur et est-ce que cela vous paraît réaliste ?
- 4.2. Une autre façon est de considérer que la couche visible contient des cellules d'entrée x qui coderont les caractéristiques de la donnée et des cellules de classe y qui coderont les classes suivant un codage 1 parmi k si on a k classes.
 - Illustrer cette RBM par un dessin
 - Donner informellement une solution permettant après apprentissage de déterminer la classe la plus probable quand on présente un exemple de test².
 - En utilisant la relation $p(y|x) = \frac{\sum_{\mathbf{h}} p(x, y, \mathbf{h})}{\sum_{y', \mathbf{h}} p(x, y', \mathbf{h})}$ expliquer comment on peut faire la classification.
 - Quelle est la complexité de ce classifieur et est-ce que cela vous paraît réaliste ?

² Pour classer, il suffit de savoir quelle est la classe dont la probabilité est la plus élevée sans avoir besoin de connaître cette probabilité.

2. Apprentissage

Neurones stochastiques

La RBM peut être interprétée comme un réseau de neurones stochastiques avec (résultat admis):

$$p(H_i = 1|\mathbf{v}) = \sigma(\sum_{j=1}^m w_{ij}v_j + c_i) \quad (i)$$

et

$$p(V_j = 1|\mathbf{h}) = \sigma(\sum_{i=1}^n w_{ij}h_i + b_j) \quad (ii)$$

Où σ est la fonction logistique $\sigma(x) = \frac{1}{1+e^{-x}}$

Gradient de la log vraisemblance

Pour entraîner une RBM, on va maximiser la log vraisemblance du modèle sur les exemples d'apprentissage. On peut faire cela avec une montée de gradient. On note θ les paramètres du modèle (ici les paramètres w, b, c).

Si on considère un unique exemple \mathbf{v}^0 , la log vraisemblance s'écrit :

$$\log p(\mathbf{v}^0|\theta) = \log \frac{1}{Z} \sum_{\mathbf{h}} e^{-E(\mathbf{v}^0, \mathbf{h})} = \log \sum_{\mathbf{h}} e^{-E(\mathbf{v}^0, \mathbf{h})} - \log \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})}$$

1. Montrer

$$p(\mathbf{h}|\mathbf{v}) = \frac{e^{-E(\mathbf{v}, \mathbf{h})}}{\sum_{\mathbf{h}'} e^{-E(\mathbf{v}, \mathbf{h}')}}}$$

2. En utilisant ce résultat, montrer

$$\partial \log \frac{p(\mathbf{v}^0|\theta)}{\partial \theta} = - \sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^0) \frac{\partial E(\mathbf{v}^0, \mathbf{h})}{\partial \theta} + \sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial \theta}$$

La première expression est l'espérance par rapport à la distribution conditionnelle des cellules cachées étant donné l'exemple \mathbf{v}^0 , la seconde expression est l'espérance calculée par rapport à la distribution jointe $p(V, H)$ apprise par le modèle.

3. On va calculer ces dérivées pour le paramètre w_{ij} , montrer

$$\sum_{\mathbf{h}} p(\mathbf{h}|\mathbf{v}^0) \frac{\partial E(\mathbf{v}^0, \mathbf{h})}{\partial w_{ij}} = - p(H_i = 1|\mathbf{v}^0)v_j^0 \quad (iii)$$

$$\sum_{\mathbf{v}, \mathbf{h}} p(\mathbf{v}, \mathbf{h}) \frac{\partial E(\mathbf{v}, \mathbf{h})}{\partial w_{ij}} = - \sum_{\mathbf{v}} p(\mathbf{v})p(H_i = 1|\mathbf{v})v_j \quad (iv)$$

Le gradient de la vraisemblance par rapport à w_{ij} en un point s'écrit donc

$$\partial \log \frac{p(\mathbf{v}^0|w_{ij})}{\partial w_{ij}} = p(H_i = 1|\mathbf{v}^0)v_j^0 - \sum_{\mathbf{v}} p(\mathbf{v})p(H_i = 1|\mathbf{v})v_j$$

Calculer (iii) facile grâce à (i). Par contre calculer (iv) nécessite une somme sur toutes les valeurs possibles de \mathbf{v} , ce qui est prohibitif. On va utiliser une approximation qui repose sur 2 idées. La première consiste à remplacer l'espérance (iv) par la valeur du gradient prise en un point \mathbf{v} . La seconde consiste à échantillonner ce point selon la distribution apprise par la RBM. On va utiliser pour cet échantillonnage la procédure suivante.

4. Algorithme : Apprentissage RBM

Input : $\text{RBM}(V_1, \dots, V_m, H_1, \dots, H_n)$,

Output : approximation du gradient Δw_{ij} pour $i = 1 \dots n, j = 1, \dots, m$

Init $\Delta w_{ij} = 0$ pour $i = 1 \dots n, j = 1, \dots, m$

Répéter

Choisir \mathbf{x} dans l'ensemble d'apprentissage

$\mathbf{v}^0 = \mathbf{x}$

pour $i = 1 \dots n$ faire échantillonner $h_i^0 \sim p(h_i | \mathbf{v}^0)$

pour $j = 1 \dots m$ faire échantillonner $v_j^1 \sim p(v_j | \mathbf{h}^0)$

pour $i = 1 \dots n, j = 1 \dots m$ faire

$\Delta w_{ij} = p(H_i = 1 | \mathbf{v}^0) v_j^0 - p(H_i = 1 | \mathbf{v}^1) v_j^1$

- Expliquer comment on peut réaliser les deux étapes d'échantillonnage.
- A quelle famille d'algorithme de gradient est ce que cet algorithme s'apparente ?
- Donner la formule de gradient utilisée pour l'apprentissage.