

Réseaux de neurones

Alexandre Allauzen

Université Paris-Sud / LIMSI-CNRS

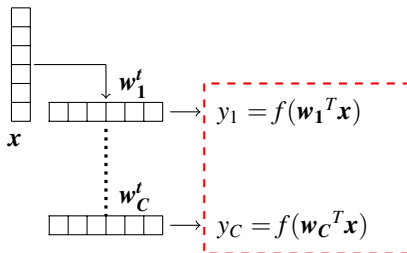
Plan

1 Classification multi-classe

2 Réseaux récurrents

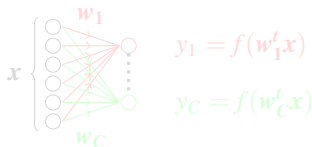
Classification multi-classe

Classification binaire vs C classes (Maxent)



$$f(a_j = \mathbf{w}_j^T \mathbf{x}) = \frac{e^{a_j}}{\sum_{j'=1}^K e^{a_{j'}}} = \frac{e^{a_j}}{Z(\mathbf{x})}$$

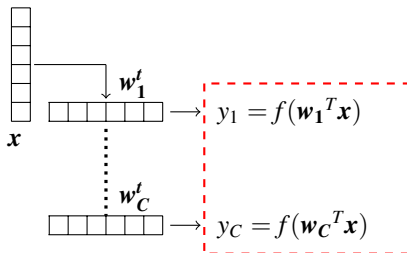
Réseau de neurones associé



- \mathbf{x} : input layer
- \mathbf{y} : output layer
- chaque y_j a ses paramètres $(\mathbf{w}_j) \rightarrow W$
- f : softmax

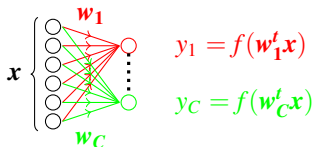
Classification multi-classe

Classification binaire vs C classes (Maxent)



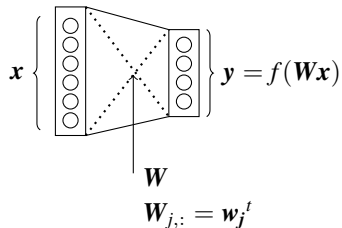
$$f(a_j = \mathbf{w}_j^T \mathbf{x}) = \frac{e^{a_j}}{\sum_{j'=1}^K e^{a_{j'}}} = \frac{e^{a_j}}{Z(\mathbf{x})}$$

Réseau de neurones associé



- \mathbf{x} : input layer
- \mathbf{y} : output layer
- chaque y_j a ses paramètres $(\mathbf{w}_j) \rightarrow \mathbf{W}$
- f : softmax

Deux couches neuronales

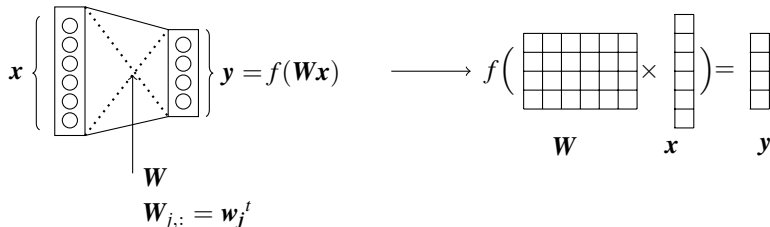


- f est souvent non-linéaire
- f s'applique composante par composante
- e.g le softmax:

$$y_j = P(c = j | \mathbf{x}) = \frac{e^{\mathbf{w}_j^T \mathbf{x}}}{\sum_{j'} e^{\mathbf{w}_{j'}^T \mathbf{x}}} = \frac{e^{\mathbf{W}_{j,:} \mathbf{x}}}{\sum_{j'} e^{\mathbf{W}_{j',:} \mathbf{x}}}$$

- tanh, sigmoid, relu, ...

Deux couches neuronales

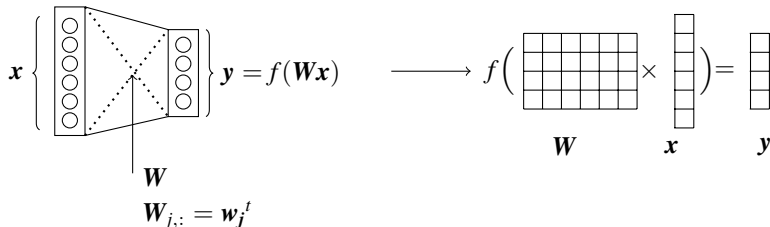


- f est souvent non-linéaire
- f s'applique composante par composante
- e.g le softmax:

$$y_j = P(c = j|x) = \frac{e^{w_j^T x}}{\sum_{j'} e^{w_{j'}^T x}} = \frac{e^{W_{j,:} x}}{\sum_{j'} e^{W_{j',:} x}}$$

- tanh, sigmoid, relu, ...

Deux couches neuronales

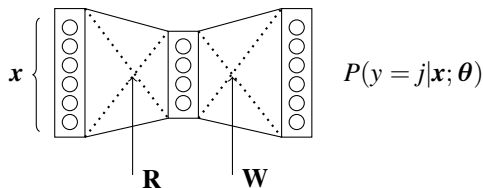


- f est souvent non-linéaire
- f s'applique composante par composante
- e.g le softmax:

$$y_j = P(c = j | x) = \frac{e^{w_j^T x}}{\sum_{j'} e^{w_{j'}^T x}} = \frac{e^{W_{j,:} x}}{\sum_{j'} e^{W_{j',:} x}}$$

- tanh, sigmoid, relu, ...

Classification de document



- $x : (|\mathcal{V}|, 1)$
- $\mathbf{R} : (K, |\mathcal{V}|)$
- $\mathbf{d} : (K, 1)$
- $\mathbf{W} : (C, K)$
- $\mathbf{y} : (C, 1)$

$$\mathbf{d} = \mathbf{R} \times x$$

$$\mathbf{y} = \text{softmax}(\mathbf{W} \times \mathbf{d})$$

Fonction objectif et Complexité

La fonction softmax est “coûteuse” à calculer:

- La multiplication matricielle de sortie est en $(C \times K)$,
- puis la normalisation (somme et division sur C éléments)

Log-loss (conditional log-likelihood, cross-entropy)

Les données : $\mathcal{D} = (\mathbf{x}_{(i)}, c_{(i)})_{i=1}^N, c_{(i)} \in \{1, 2, \dots, C\}$

$$\mathcal{L}(\boldsymbol{\theta}) = \sum_{i=1}^N l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)}) = \sum_{i=1}^N \left(- \sum_{c=1}^C \mathbb{I}\{c = c_{(i)}\} \log(P(c|\mathbf{x}_{(i)})) \right) \quad (1)$$

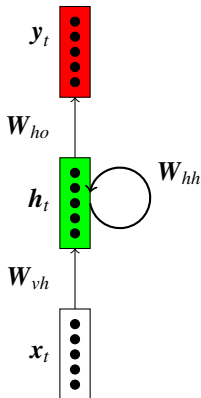
$$l(\boldsymbol{\theta}, \mathbf{x}_{(i)}, c_{(i)}) = - \sum_{k=1}^C \mathbb{I}\{k = c_{(i)}\} \log(y_k) \quad (2)$$

Plan

1 Classification multi-classe

2 Réseaux récurrents

Réseau récurrent



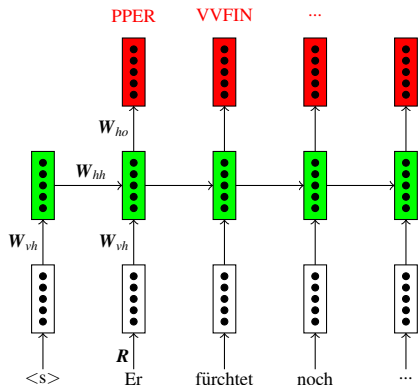
Un modèle de séquence.

Pour chaque instant t :

- maintient une représentation interne (état caché) de l'historique: h_t
- Mise à jour avec une observation x_t et l'état précédent h_{t-1}
- La prédiction y_t dépend de l'état caché (h_t)
- x_t vient des embeddings.

Les mêmes paramètres sont partagés au cours du temps

Un modèle déplié dans le temps



à chaque instant t

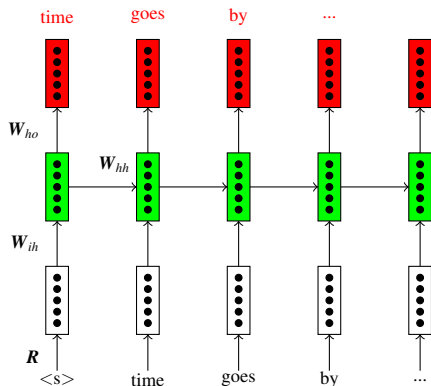
- Lire le mot $w_t \rightarrow \mathbf{x}_t$ de R
- mettre à jour l'état interne

$$\mathbf{h}_t = f(\mathbf{W}_{vh}\mathbf{x}_t + \mathbf{W}_{hh}\mathbf{h}_{t-1})$$
- (prédire à t une sortie à partir de \mathbf{h}_t :

$$\mathbf{y}_t = g(\mathbf{W}_{ho}\mathbf{h}_t)$$

- g est la fonction softmax

Variante : Prédiction de séquence



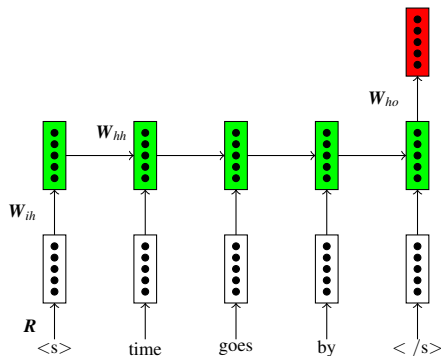
à chaque instant t

- lire $w_t \rightarrow x_t$ de R
- $h_t = f(W_{ih}x_t + W_{hh}h_{t-1})$
- prédire le mot/symbôle à $t + 1$ à partir h_t :

$$y_t = g(W_{ho}h_t)$$

- g : softmax

Variante : Classification de séquence



à chaque instant t

- lire $w_t \rightarrow x_t$ de R
- $h_t = f(W_{ih}x_t + W_{hh}h_{t-1})$

prédire la classe à partir de h_L

Apprentissage

Algorithme

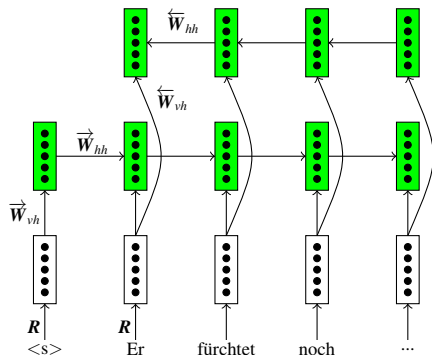
Back-Propagation through time

- à chaque instant t
 - calcul de la fonction objectif et de son gradient
 - Back-Propagation au travers du réseau “déplié”

Problèmes

- Vanishing/exploding gradient [?] → LSTM, Gradient clipping
- Mémoire à long-terme → Bi-recurrent network

Bi-recurrent network



à chaque instant t , de gauche à droite

- $w_t \rightarrow \mathbf{x}_t$
- $\vec{\mathbf{h}}_t = f(\vec{\mathbf{W}}_{vh}\mathbf{x}_t + \vec{\mathbf{W}}_{hh}\vec{\mathbf{h}}_{t-1})$

à chaque instant t , de droite à gauche

- $w_t \rightarrow \mathbf{x}_t$
- $\overleftarrow{\mathbf{h}}_t = f(\overleftarrow{\mathbf{W}}_{vh}\mathbf{x}_t + \overleftarrow{\mathbf{W}}_{hh}\overleftarrow{\mathbf{h}}_{t-1})$

prédiction:

$$\mathbf{y}_t = g(\mathbf{W}_{ho}[\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t])$$

$[\vec{\mathbf{h}}_t; \overleftarrow{\mathbf{h}}_t]$: représentation contextuelle de w_t