

Mini-Rapport

Exploration de la notion de méta-apprentissage

*Dans quelle mesure un système apprenant peut
prendre conscience de ses performances et altérer son
comportement ?*

Yann Boniface, Alain Dutech, Nicolas Rougier
Matthieu Zimmer

20 avril 2012

1 Introduction

Notre intérêt s'est tourné vers l'article [Cleeremans Alex, 2007] et ses 2 types de réseaux proposés. Dans un premier temps, nous avons cherché à reproduire et expliquer les résultats donnés, et ensuite à des solutions pour tirer profit des paris réalisés.

Nous nous sommes également penchés sur [Pasquali Antoine, 2010] dont nous avons reproduit les expériences, mais leurs enjeux nous semblent encore vagues.

2 Dupliquer le premier réseau

2.1 Les bases

En premier lieu, rappelons la structure des réseaux et les résultats de l'article :

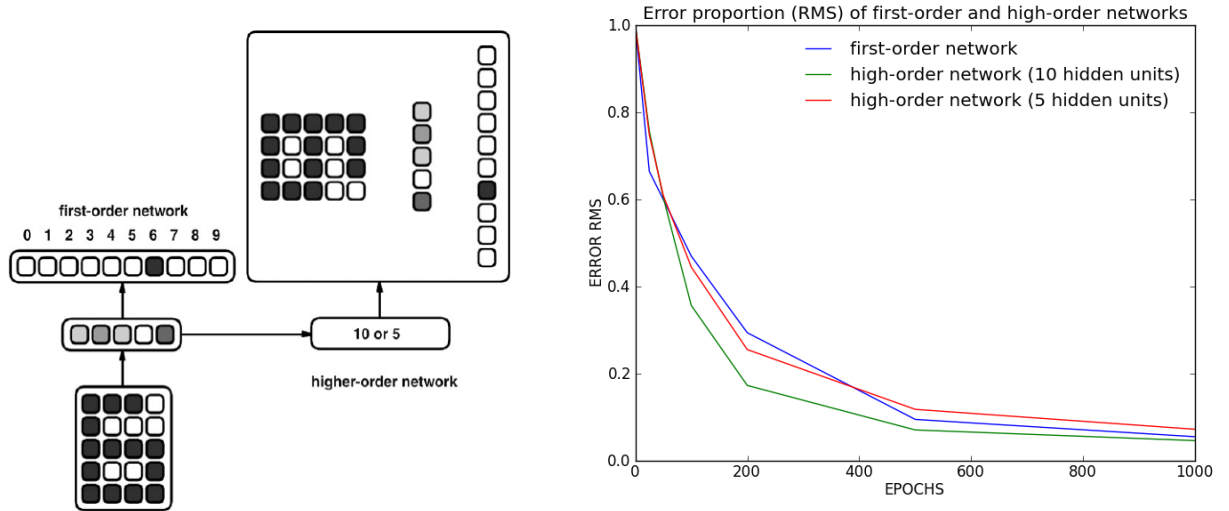


FIGURE 1 – [Cleeremans Alex, 2007] Architecture connexionniste avec méta-représentations

La formule RMS utilisée à une époque e est la suivante :

$$rms\ proportion_e = \frac{rms_e = \sqrt{\frac{1}{n} \sum_{i=1}^n (o_{i,e} - d_i)^2}}{\max(rms_{e'}, \forall e' \in epochs)}$$

with $\begin{cases} n : \text{number of neurons on the output layer} \\ o_{i,e} : \text{value obtained for the } i^{th} \text{ neuron at the } e^{th} \text{ epoch} \\ d_i : \text{value desired for the } i^{th} \text{ neuron} \end{cases}$

Nous avons décomposé les erreurs pour mieux comprendre le fonctionnement de l'architecture.

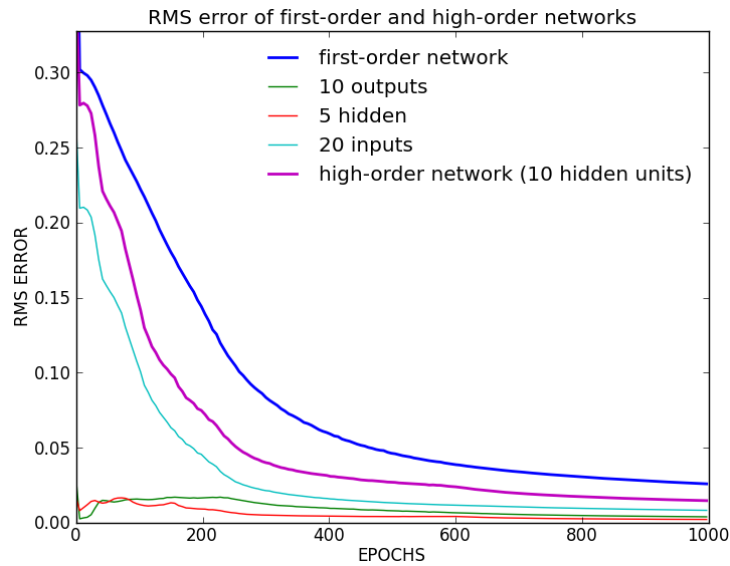


FIGURE 2 – **Erreur RMS sans proportion.** Le réseau supérieur à 5 unités cachées n'est plus considéré, et celui à 10 est découpé en 3 courbes pour représenter les 3 couches à reproduire. La courbe violette est donc la somme des 3 courbes des couches.

On peut en conclure 2 choses :

- la couche cachée et la couche de sortie ne pose aucun problème d'apprentissage
- les performances du second réseau dépendent principalement de sa capacité à reproduire les entrées

Il reste une question en suspens :

Pourquoi le second réseau apprendrait-il plus rapidement que le premier ?

Nous n'avons pas la réponse. Par contre, nous pouvons nous intéresser à l'erreur de classification.

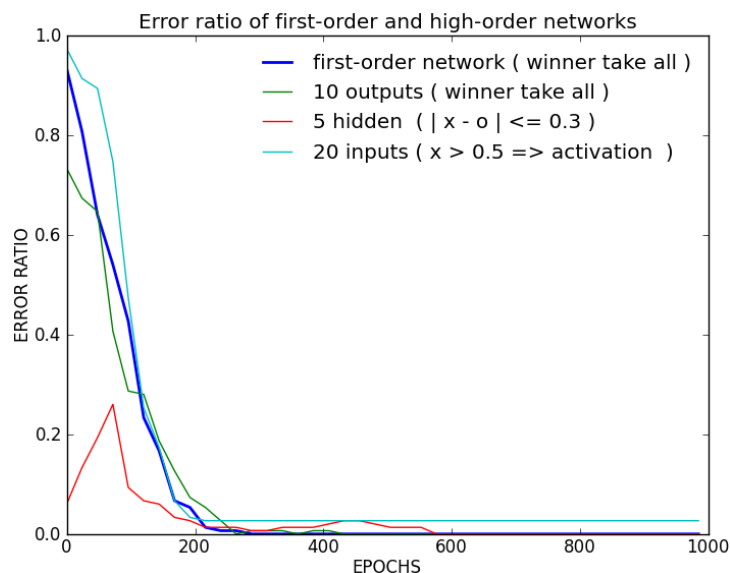


FIGURE 3 – **Erreur de classification.** Le réseau supérieur à 5 unités cachées n'est plus considéré, et celui à 10 est découpé en 3 courbes pour représenter les 3 couches à reproduire.

Il est alors beaucoup moins flagrant que le second réseau apprend mieux que le premier. Évidemment, parler de classification de la premier couche (courbe 20 inputs) et de la couche caché (courbe 5 hidden) est relativement dérisoire. La question de la rapidité d'apprentissage reste donc ouverte.

2.2 Les rouages

Notre attention s'est également portée sur les mécanismes qui permettaient la réalisation de cette architecture. Nous avons alors pu remarquer que les neurones de la couche cachée du premier réseau se stabilisaient très rapidement (autour de la 50^{ième} époque en moyenne), le tout permettant au second réseau d'avoir des entrées très peu variables, favorisant et permettant donc son apprentissage.

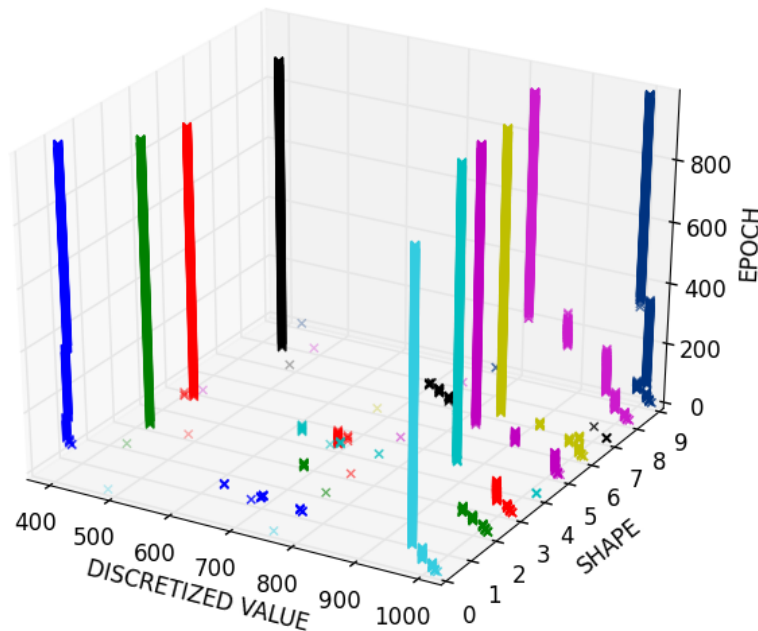


FIGURE 4 – Valeurs discrétisés de la couche caché du premier réseau. Chaque couleurs représentant un des 10 chiffres en entrée. Les courbes deviennent rapidement stables.

2.3 Ré-haussement

Pour revenir à l'importance des entrées, nous avons réalisé que le second réseau n'était capable de dupliquer le premier qu'uniquement parce que ces entrées sont triviales. Ainsi nous avons augmenté le nombre d'entrées en passant sur des chiffres manuscrits [Semeion Research Center, 1994], tout en augmentant proportionnellement le nombre de neurones des couches cachées :

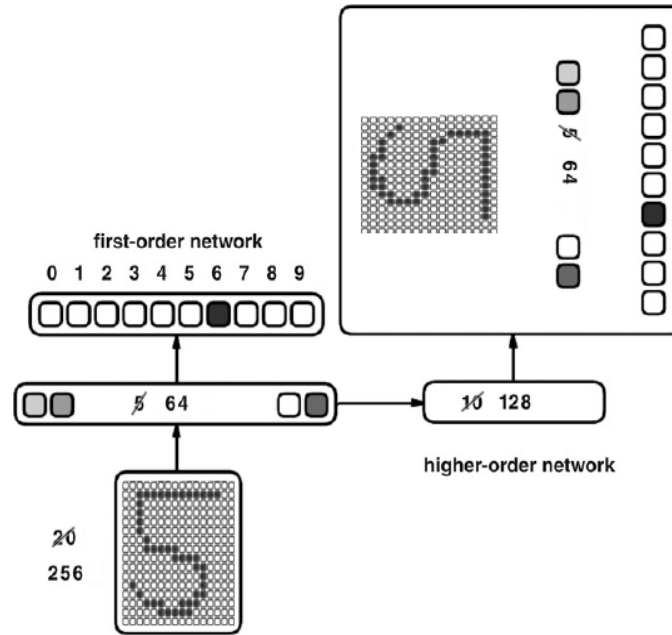
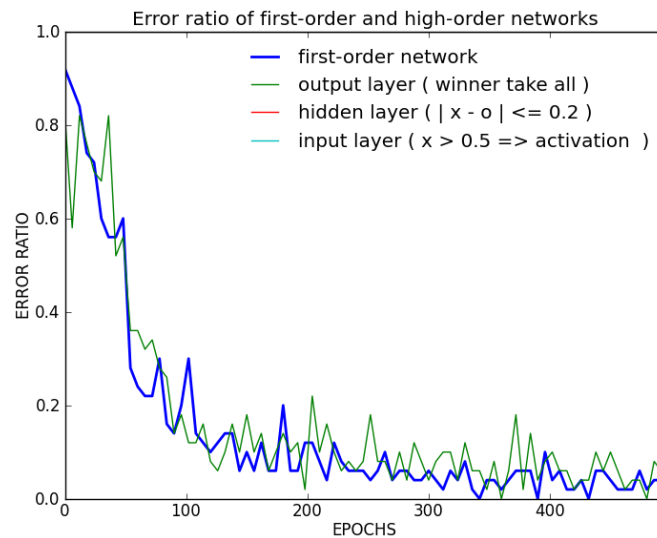


FIGURE 5 – Architecture avec méta-représentation pour chiffres manuscrits

Les performances du second réseau se sont alors écroulées :

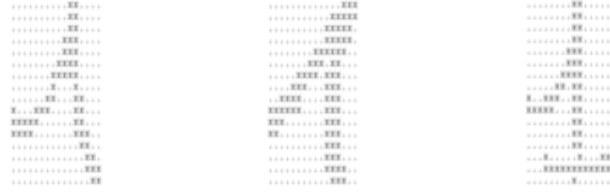
FIGURE 6 – **Erreur de classification** de l'architecture sur des chiffres manuscrits. L'erreur du second réseau est toujours divisé en 3.

FIXME Il n'est alors plus capable que de reproduire la couche de sortie. Petit bémol tout de même, le critère de classification pour la couche caché est peut-être trop rigide, il suffit qu'un neurone diffère de 0.2 (de la valeur voulue) pour considérer que c'est un échec.

Tentons d'expliquer cette chute :

La première chose à comprendre est que l'augmentation du nombre de neurone n'a rien à voir avec cette perte. En effet, nous avons fait les simulations avec les chiffres triviales

(seulement 10 formes différentes) agrandi en 16×16 et tout se passe très bien. Cette chute provient du fait que les entrées proposées sont très nombreuses. Il s'opère ainsi une première étape de classification dans la couche caché du premier réseau. Donc ces formes différentes peuvent être représenté par la même couche cachée.



Sauf que le second réseau n'ayant en entrée que la version factorisé (la couche caché), il lui est impossible de retrouver la forme initiale des 1.

2.4 Représentations

Enfin, nous avons remarqué qu'en bloquant l'apprentissage entre la couche cachée et les entrées du premier réseau, puis en changeant de tâche, le réseau était capable de réapprendre la nouvelle tâche, ce qui prouve bien la présence d'une représentation des entrées dans la couche cachée.

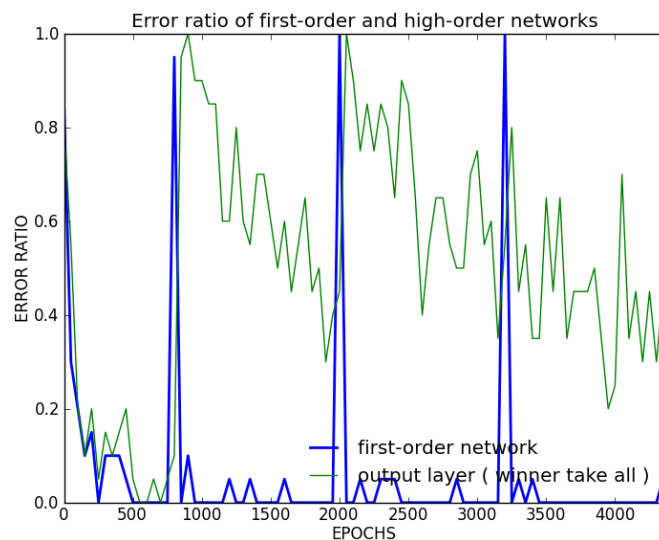


FIGURE 7 – **Erreur de classification** de l'architecture sur des chiffres manuscrits. Apprentissage bloqué à l'époque 800. Changement de tâche aux époques : 800, 2000, 3200

2.5 Remarque

Il faut cependant remarquer qu'un simple perceptron est suffisant pour réaliser la tâche du premier réseau (même sur les chiffres manuscrits), et donc, qu'il est possible que dans le cas d'un problème non linéairement séparable cette architecture soit invalidée.

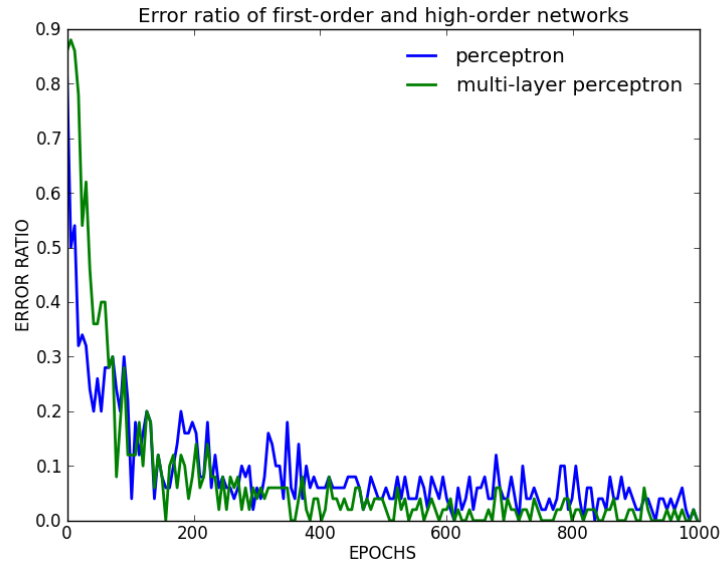


FIGURE 8 – **Erreur de classification** d'un perceptron et d'un perceptron multi-couches sur la base de chiffres manuscrits.

3 Parier sur le premier réseau

3.1 Les bases

Rappelons également la structure des réseaux et les résultats :

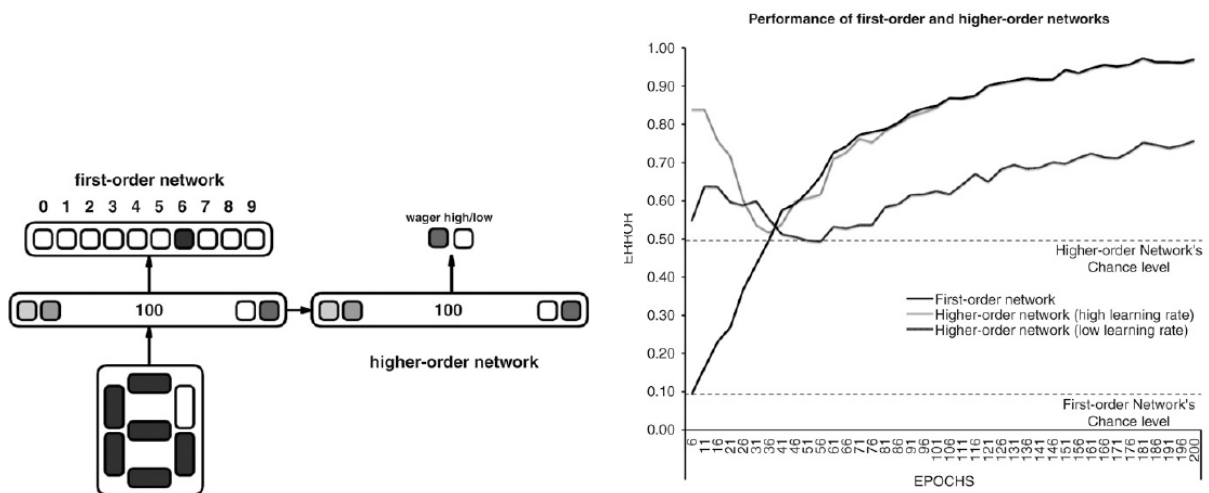


FIGURE 9 – [Cleeremans Alex, 2007] Architecture connexionniste avec paris

La première chose que nous avons faites à été d'améliorer les performances du second réseau en modifiant quelques paramètres (initialisation des poids sur $[-1 ; -1]$, momentum à 0.5).

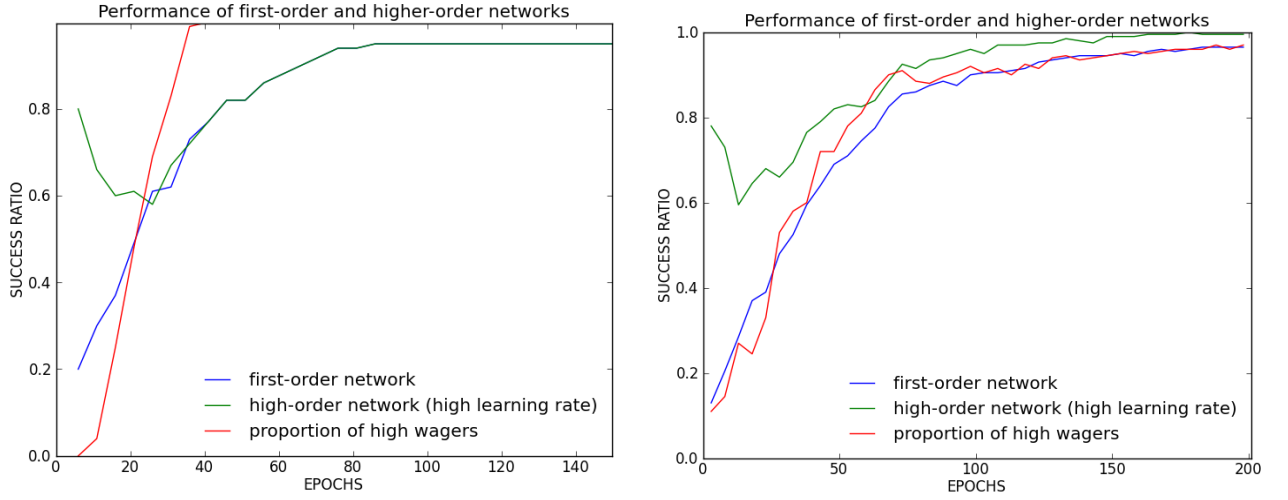


FIGURE 10 – Performance de l'architecture. À gauche la base, à droite la version améliorée. On ne considère plus le réseau avec apprentissage faible, mais on a ajouté le taux de paris hauts.

Contrairement à celui de l'article, il ne se contentera plus simplement de parier haut à chaque coups (après 40 époques). Il aura une longueur d'avance sur le premier réseau sur toute la durée de l'apprentissage. On pourra donc tirer profit de cette avance.

3.2 Feedbacks

À partir de cette différence de performances, nous avons imaginé plusieurs architectures, qui améliorent plus ou moins les performances de reconnaissance du réseau sur des chiffres manuscrits :

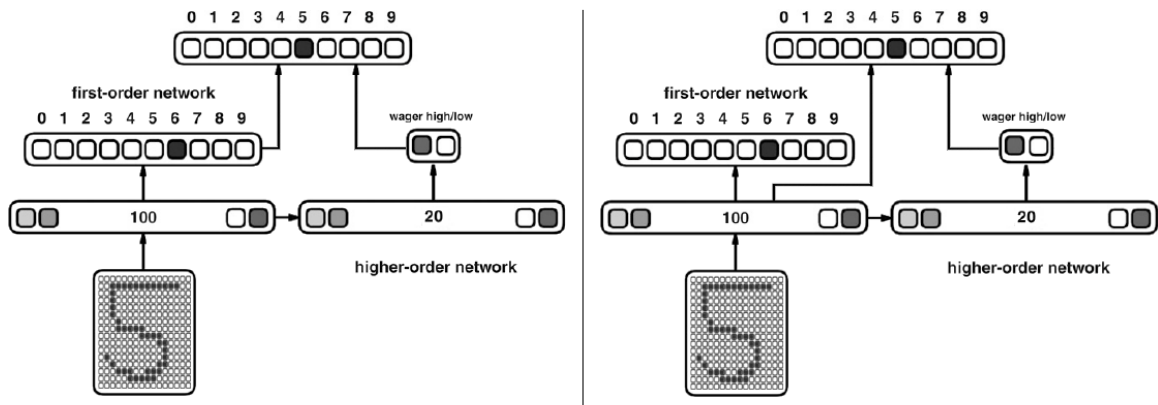


FIGURE 11 – Architecture avec 3^{ème} réseau

Dans ces 2 architectures, nous nous contentons de connecter un 3^{ème} réseau qui doit tirer des conclusions à partir d'informations sur les 2 premiers.

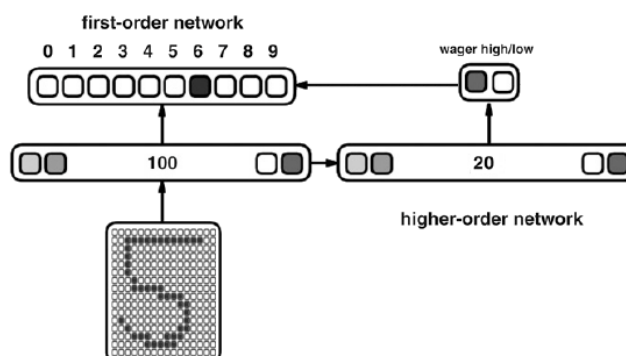


FIGURE 12 – Architecture par fusion

Ici, nous mélangeons un apprentissage par descente de gradient (sur le premier et second réseau) et un apprentissage perceptron (entre les 2 couches de sorties).

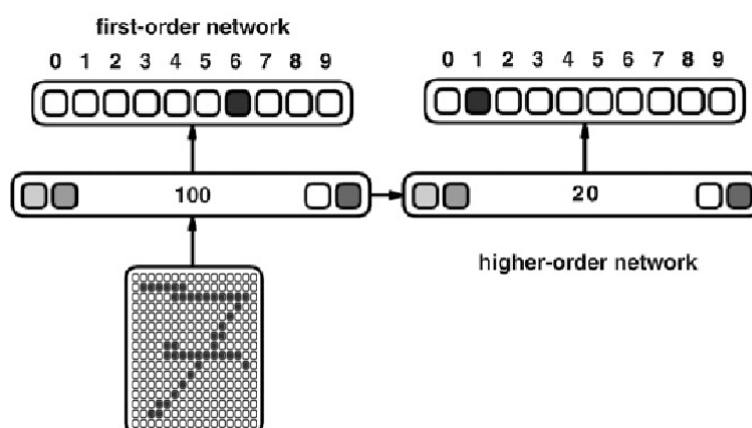


FIGURE 13 – Architecture par intuitions

Cette architecture est légèrement différente dans le sens où elle n'enregistre plus de pari mais l'indice du $n^{\text{ième}}$ neurone le plus actif contenant la bonne réponse. Exemple : le réseau supérieur sort 2 -> la réponse est le 3^{ième} neurone le plus actif de la couche de sortie du premier réseau.

Nous avons aussi essayé quelques modèles où le réseau supérieur servait de superviseur à l'apprentissage du premier réseau. Par exemple, s'il parie haut, l'apprentissage du premier réseau sera faible, sinon il sera accentué.

4 La suite

Ce que nous continuons d'étudier :

- validation sur des expériences plus complexes (qui ne peuvent être résolue directement par un perceptron)

- relation entre la taille de la couche cachée du premier réseau et le taux de paris avantageux
- nouveau critère de classification pour la couche caché dans la première architecture sur les chiffres manuscrits
- approfondir les intérêts du second article [Pasquali Antoine, 2010]
- de nouvelles architectures axées méta-apprentissage où le réseau d'ordre supérieur contrôle le premier (taux d'apprentissage, momentum, entrées à approfondir, ...) telle une conscience

Références

- [Cleeremans Alex, 2007] Cleeremans Alex, Timmermans Bert, P. A. (2007). Consciousness and metarepresentation : A computational sketch. *doi :10.1016/j.neunet.2007.09.011*.
- [Pasquali Antoine, 2010] Pasquali Antoine, Timmermans Bert, C. A. (2010). Know thyself : Metacognitive networks and measures of consciousness.
- [Semeion Research Center, 1994] Semeion Research Center, o. S. o. C. (1994). Semeion handwritten digit data set. 1593 handwritten digits from around 80 persons.

A Algorithme

Ensemble d'inscrutions executé pour une époque et un tirage d'entrées/sorties.
Cet algorithme se déroule en plusieurs temps :

- a On calcule la sortie du premier réseau en ignorant le second (valeur à 0)
- b On calcule la sortie du second réseau
- On calcule la sortie du premier réseau avec les sorties de b)

```

first_order.calc_hidden_layer(samples.inputs)
high_order.calc_output_layer(first_order.hidden_layer)
first_order.calc_output_layer(first_order.hidden_layer, [0, ..., 0])

```

```

h_output ← ampli(high_order.output_layer)
right_houtput ← [0, 0]
if good_answer(first_order) then
    right_houtput[1] ← 1
else
    right_houtput[0] ← 1
end if
first_order.calc_output_layer(first_order.hidden_layer, h_output)

```

```
calc_stats()
```

```

high_order.train(first_order.hidden_layer, right_houtput)
first_order.train(samples.inputs, samples.outputs, h_output)

```

L'interêt se porte sur

```

function train(inputs, outputs, add)
    for i = 0 → output_neurons.length do
        y_output[i] ← g'(output_neurons[i].a) × (outputs[i] - output_neurons.state)
    end for

    for i = 0 → hidden_neurons.length do
        w_sum ←  $\sum_{j=0}^{output\_neurons.length} output\_neurons[j].weights[i] \times y_{output[j]}$ 
        y_hidden[i] ← g'(hidden_neurons[i].a) × w_sum
    end for
    update_weights_hidden_layer(y_hidden)

    for i = 0 → output_neurons.length do
        output_neurons[i].update_weights_gradient(y_output[i], hidden_neurons, add)
        output_neurons[i].update_weights_perceptron(outputs[i], hidden_neurons, add)
    end for
end function

function update_weights_gradient(error, intputs, add)
    calc_output(inputs + add)

    for j = 0 → inputs.length do
        dw ← weights[j] - last_weights[j]
        p ← error × inputs[j]
    end for
end function

```

```

         $weights[j] \leftarrow weights[j] + learning\_rate \times p + momentum \times dw$ 
    end for
end function

function update_weights_perceptron(goal, inputs, add)
    calc_output(inputs + add)

    for  $j = inputs.length \rightarrow inputs.length + add.length$  do
         $dw \leftarrow weights[j] - last\_weights[j]$ 
         $p \leftarrow (goal - state) \times add[inputs.length - j]$ 
         $weights[j] \leftarrow weights[j] + \frac{learning\_rate \times p + momentum \times dw}{add.length}$ 
    end for
end function

```