

# On recurrent neural networks

Giulio Galvan

23rd March 2015



# Contents

<b>1</b>	<b>Artificial neural networks</b>	<b>5</b>
1.1	A family of models . . . . .	5
1.2	Feed foward neural networks . . . . .	7
1.2.1	On expressivness of ffnn . . . . .	7
1.2.2	Learning with ffnn . . . . .	7
1.2.3	Gradient . . . . .	8
1.3	Recursive neural networks . . . . .	9
1.3.1	On expressivness of rnn . . . . .	10
1.3.2	Learning with ffnn . . . . .	11



# Chapter 1

## Artificial neural networks

### 1.1 A family of models

An artificial neural network is a network of connected units called neurons or perceptrons, each arc which connects two neurons  $i$  and  $j$  is associated with a weight  $w_{ji}$ . Perceptrons share the same structure for all models, what really define the model in the family is how the perceptrons are arranged and how they are connected, for example if there are cycles or not. As you can see in figure 1.1 each neuron is *fed* with a set of inputs which are the weighted outputs of other neurons. Formally the output of a perceptron  $\phi_j$  is defined as:

$$\phi_j \triangleq \sigma(a_j) \tag{1.1}$$

$$a_j \triangleq \sum_l w_{jl}\phi_l + b_j \tag{1.2}$$

where  $w_{jl}$  is the weight of the connection between neuron  $l$  and neuron  $j$ ,  $\sigma(\cdot)$  is a non linear function and  $b_j \in \mathbb{R}$  is called bias.

A neural network looks like: INSERT FIGURE FULLY CONNECTED

It is sometimes useful to think of a neural network as a series of layers, one on top of each other, as depicted in figure ???. The first layer is called the input layer and its units are *fed* with external inputs, the upper layers are called *hidden layers* because their outputs are not observed from outside except the last one which is called *output layer* because its output is the output of the net.

When we describe a network in this way it is also useful to adopt a different notation: we describe the weights of the net with a set of matrixes  $W_i$  one for each layer, and neurons are no more linearly indexed, instead we refer to a neuron with a relative index with respect to the layer; this allows to write easier equations in matrix notation.

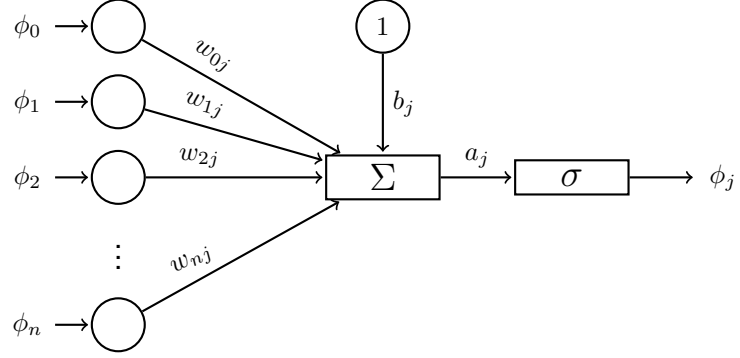


Figure 1.1: Neuron model

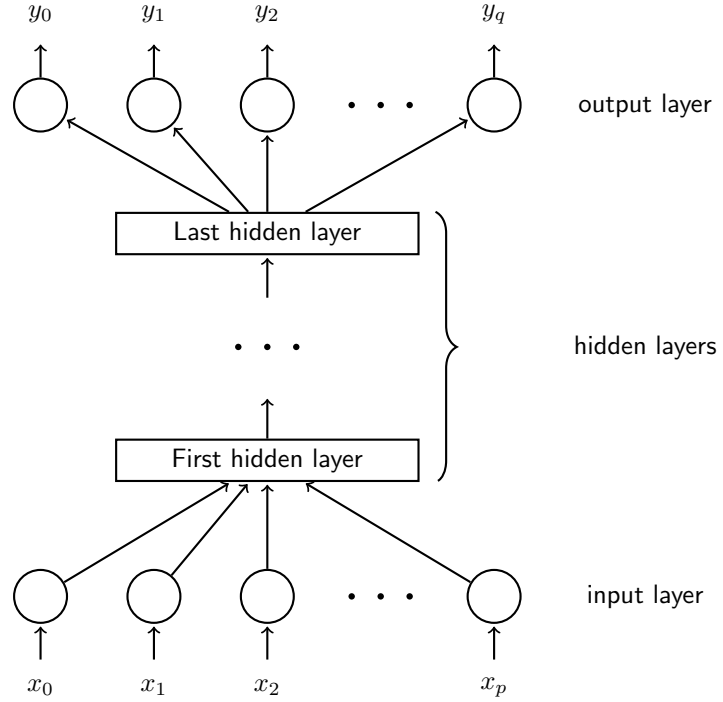


Figure 1.2: Feed forward neural network model

## 1.2 Feed foward neural networks

A feed foward neural network is an artificial neural network in which there are no cycles, that is to say each layer output is *fed* to the next one and connections to any earlier layer are not possible.

**Definition 1** (Feed foward neural network). A feed foward neural network is tuple

$$FFNN \triangleq \langle \mathbf{p}, W, b, \sigma(\cdot), f(\cdot) \rangle$$

- $\mathbf{p} \in \mathbb{N}^U$  is the vector whose elements  $p(k)$  are the number of neurons of layer  $k$ ;  $U$  is the number of layers
- $W \triangleq \{W_{p(k+1) \times p(k)}^k, k = 1, \dots, U\}$  is the set of weight matrixes of each layer
- $b \triangleq \{\mathbf{b}^k \in \mathbb{R}^{p(k)}, k = 1, \dots, U\}$  is the set of bias vectors
- $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is the activation function
- $f(\cdot) : \mathbb{R}^{p(U)} \rightarrow \mathbb{R}^{p(U)}$  is the output function

The total number of weights si given by  $\mathcal{N}(W) \triangleq \sum_{k=1}^U p(k+1)p(k)$ , the number of biases by  $\mathcal{N}(b) \triangleq \sum_{k=1}^U p(k)$

Given a  $FFNN$  and an input vector  $\mathbf{x}$  the output  $\mathbf{y}$  of the net is defined by the following:

$$\mathbf{y} = f(\phi^U) \tag{1.3}$$

$$\phi^i \triangleq \sigma(\mathbf{a}_i), \quad i = 1, \dots, U \tag{1.4}$$

$$\mathbf{a}^i \triangleq W^{i+1} \cdot \phi^{i-1} + \mathbf{b}^i \quad i = 1, \dots, U \tag{1.5}$$

$$\phi^0 \triangleq \mathbf{x} \tag{1.6}$$

### 1.2.1 On expressivness of ffnn

### 1.2.2 Learning with ffnn

As we have seen in the previous section a model sush as  $FFNN$  can approximate arbitrary well any smooth function, so a natural application of feed foward neural networks is machine learning. To model an optimization problem we first need to define a dataset  $D$  as

$$D \triangleq \{x^{(i)} \in \mathbb{R}^p, y^{(i)} \in \mathbb{R}^q, i \in [1, N]\} \tag{1.7}$$

Then we need a loss function  $L_D : \mathbb{R}^{\mathcal{N}(W)+\mathcal{N}(b)} \rightarrow \mathbb{R}_{\geq 0}$  over  $D$  defined as

$$L_D(W, b) \triangleq \frac{1}{N} \sum_{i=1}^N L^i(W, b) \tag{1.8}$$

$L^i$  is an arbitrary loss function on the  $i^{th}$  example.

The problem is then to find a *FFNN* which minimize  $L$ . As we have seen feed forward neural network allow for large customization: the only variables in the optimization problem are the weights, the other parameters are said *hyper-parameters* and are determined *a priori*. Usually the output function is chosen depending on the output, for instance for multi-way classification is generally used the softmax function, for regression a simple identity function.

The activation function  $\sigma(\cdot)$  is often chosen from:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (1.9)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (1.10)$$

$$\text{ReLU}(x) = \begin{cases} x, & \text{if } x > 0. \\ 0, & \text{otherwise.} \end{cases} \quad (1.11)$$

For what concerns the number of layers and the number of units per layers they are chosen relying on experience or performing some kind of hyper-parameter tuning, which usually consists on training nets with some different configurations of such parameters and choosing the best one.

Once we have selected the values for all hyper-parameters the optimization problem becomes:

$$\min_{W,b} L_D(W,b) \quad (1.12)$$

### 1.2.3 Gradient

We can compute partial derivatives with respect to a single weight  $w_{lj}$ , using simply the chain rule, as

$$\frac{\partial L}{\partial w_{lj}} = \frac{\partial L}{\partial a_l} \cdot \frac{\partial a_l}{\partial w_{lj}} = \delta_l \cdot \phi_j$$

where we put

$$\delta_l \triangleq \frac{\partial L}{\partial a_l} \quad (1.13)$$

So we can easily compute  $\delta_u = \frac{\partial L^{(i)}}{\partial a_u}$  for each output unit  $u$  once we choose a differentiable loss function; note that we don't need the weights for such a computation.

Let  $P(l)$  be the set of parents of neuron  $l$ , formally:

$$P(l) = \{k : \exists \text{ a link between } l \text{ and } k \text{ with weight } w_{lk}\} \quad (1.14)$$



Again, simply using the chain rule, we can write, for each non output unit  $l$ :

$$\delta_l = \sum_{k \in P(l)} \frac{\partial L^{(i)}}{\partial a_k} \cdot \frac{\partial a_k}{\partial a_l} = \sum_{k \in P(l)} \delta_k \cdot \frac{\partial a_k}{\partial \phi_l} \cdot \frac{\partial \phi_l}{\partial a_l} = \sum_{k \in P(l)} \delta_k \cdot w_{kl} \cdot \sigma'(a_l) \quad (1.15)$$

For output units instead we can compute  $\delta_u = \frac{\partial L^{(i)}}{\partial a_u}$  directly once we define the loss function.

For biases variables partial derivatives are simply given by:

$$\frac{\partial L}{\partial b_l} = \frac{\partial L}{\partial a_l} \cdot \frac{\partial a_l}{\partial w_l} = \delta_l \cdot 1$$

In the following we rewrite the previously derived equations in matrix notation. Let us recall that the weight matrix for the  $i^{th}$  layer is the  $p(i) \times p(i-1)$  matrix whose elements  $w_{l,k}$  are the weights of the arcs which link neuron  $k$  from level  $i-1$  to neuron  $l$  from level  $i$ , where  $p(i)$  is the number of neuron layer  $i$  is composed of.

We can rewrite equation 1.15 in matrix notation as:

$$\frac{\partial L}{\partial W^i} = \frac{\partial L}{\partial \mathbf{a}^i} \cdot \left( \frac{\partial \mathbf{a}^i}{\partial W^i} \right)^T = \Delta^i \cdot (\phi^{i-1})^T \quad (1.16)$$

where

$$\Delta^i \triangleq \frac{\partial L}{\partial \mathbf{a}^i} \quad (1.17)$$

$$\Delta^i = (W^{i+1})^T \cdot \Delta^{i+1} \circ \sigma(\Delta^i) \quad (1.18)$$

$$\frac{\partial L}{\partial b^i} = \frac{\partial L}{\partial \mathbf{a}^i} \cdot \left( \frac{\partial \mathbf{a}^i}{\partial b^i} \right)^T = \Delta^i \cdot Id \quad (1.19)$$

### 1.3 Recursive neural networks

Recurrent neural networks differ from feed forward neural networks because of the presence of recurrent connections: at least one perceptron output at a given layer  $i$  is *fed* to another perceptron at a level  $j < i$ . This is a key difference: as we will see in the next section, rnn are not only more powerfull than ffnn but as powerfull as turing machines.

This difference in topology reflects also on the network's input and output domain, where in feed forward neural networks inputs and outputs were real valued vectors, recursive neural networks deal with sequences of vectors, that is to say that now time is also considered. One may argue that taking time (and sequences) into consideration is some sort of limitation because it restricts our model to deal only with a temporal inputs; that's not true, in fact we can apply rnn to non temporal data by considering space as the temporal dimension or we can feed the network with the same input for all time steps, or just providing no input after the first time step.

FIGURE RNN

**Definition 2** (Recurrent neural network). A recurrent neural network is tuple

$$RNN \triangleq \langle W_{in}, W_{out}, W_{rec}, \mathbf{b}_{out}, \mathbf{b}_{rec}, \sigma(\cdot), f(\cdot) \rangle$$

- $W^{in}$  is the  $r \times p$  input weight matrix
- $W^{rec}$  is the  $r \times r$  recurrent weight matrix
- $W^{out}$  is the  $o \times r$  output weight matrix
- $\mathbf{b}^{rec} \in \mathbb{R}^r$  is the bias vector for the recurrent layer
- $\mathbf{b}^{out} \in \mathbb{R}^o$  is the bias vector for the output layer
- $\sigma(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$  is the activation function
- $f(\cdot) : \mathbb{R}^o \rightarrow \mathbb{R}^o$  is the output function

$p$  is the size of input vectors,  $r$  is the number of hidden units,  $o$  is the size of output vectors. The total number of weights is given by  $\mathcal{N}(W) \triangleq rp + r^2 + ro$ , the number of biases by  $\mathcal{N}(b) \triangleq r + o$

Given a  $RNN$  and an input sequences  $\{\mathbf{x}\}_t$  the output sequence  $\{\mathbf{y}\}_t$  of the net is defined by the following:

$$\mathbf{y}_t \triangleq f(W^{out} \cdot \phi_t + \mathbf{b}^{out}) \quad (1.20)$$

$$\mathbf{a}_t \triangleq W^{rec} \cdot \phi_{t-1} + W^{in} \cdot \mathbf{x}_t + \mathbf{b}^{rec} \quad (1.21)$$

$$\phi_t \triangleq \sigma(\mathbf{a}_t) \quad (1.22)$$

As we can understand from definition 2, there is only one recurrent layer, whose weights are the same for each time step, so one can ask where does the deepness of the network come from. The answer lies in the temporal unfolding of the network, in fact if we unfold the network step by step we obtain a structure similar to the structure of a feed forward neural network. As we can observe in figure ??, the unfolding of the network through time consist of putting identical version of the same recurrent layer on top of each other and linking the inputs of one layer to the next one. The key difference from feed forward neural networks is, as we have already observed, that the weights in each layer are identical, and of course the additional timed inputs which are different for each unfolded layer.

UNFOLDING FIGURE

### 1.3.1 On expressiveness of rnn

Rnns are as powerful as turing machines

### 1.3.2 Learning with rnn

We can model an optimization problem in the same way we did for feed forward neural networks, the main difference is, again, that we now deal with temporal sequences so we need a slightly different loss function. Given a dataset  $D$ :

$$D \triangleq \{\mathbf{x}_t^{(i)} \in \mathbb{R}^p, \mathbf{y}_t^{(i)} \in \mathbb{R}^q; t = 1, \dots, T; i = 1, \dots, N\} \quad (1.23)$$

we define a loss function  $L_D : \mathbb{R}^{\mathcal{N}(W)+\mathcal{N}(b)} \rightarrow \mathbb{R}_{\geq 0}$  over  $D$  as

$$L_D(W, b) \triangleq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T L_t^i(W, b) \quad (1.24)$$

$L_t^i$  is an arbitrary loss function for the  $i^{th}$  example at time step  $t$ .



# Bibliography

- [1] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory, 1995.