

On recurrent neural networks

Giulio Galvan

12th August 2015

Contents

1	Artificial neural networks	5
1.1	A family of models	5
1.2	Feed forward neural networks	8
1.2.1	Learning with FFNNs	9
1.2.2	Gradient	10
1.3	Recurrent neural networks	12
1.3.1	Learning with RNNs	14
1.3.2	Gradient	15
1.3.3	The vanishing and exploding gradient problem	17
1.4	Activation functions and gradient	22
1.5	On expressiveness	25
2	An historical review of proposed solutions	27
2.1	Long short-term memory	27
2.2	Gated recurrent units	29
A	Notation	31

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, as can be seen in figure 1.1; the link which connects neurons i and j , is associated with a weight w_{ji} . Perceptrons share the same structure for all models, what really distinguish a particular model in the family of artificial neural networks is how the perceptron units are arranged and connected together, for example whether there are cycles or not, and how data inputs are *fed* to the network.

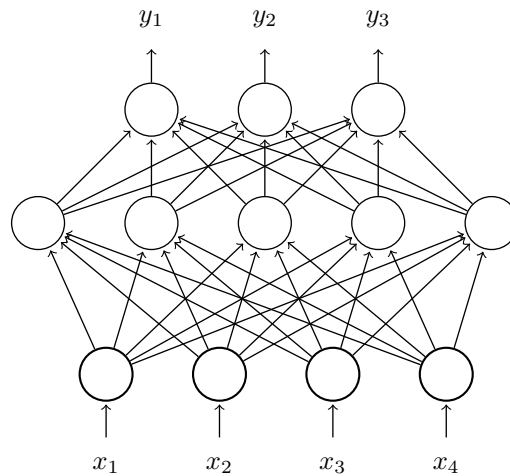


Figure 1.1: Artificial neural network example

As you can see in figure 1.2 each neuron is *fed* with a set of inputs which are the weighted outputs of other neurons and/or other external inputs. Formally

the output of a perceptron ϕ_j is defined as:

$$\phi_j \triangleq \sigma(a_j) \quad (1.1)$$

$$a_j \triangleq \sum_l w_{jl} \phi_l + b_j \quad (1.2)$$

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 a bias term. It's worth noticing that in this formulation the input ϕ_l can be the outputs of other neurons or provided external inputs.

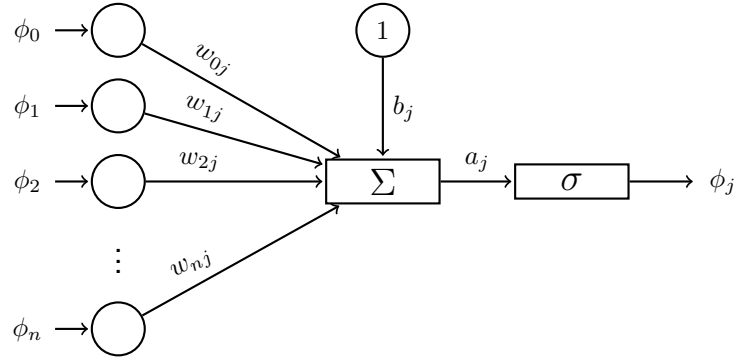


Figure 1.2: Neuron model

So, given a set of inputs $\{x\}_i$ which are *fed* to some of the units of the net which we call *input units* the output of the network $\{y\}_i$ is given by the some of units of the network, the 'upper' ones which we call *output units*. All remaining units, i.e the ones which are neither input nor output units are called *hidden units* because their value is not *observed* from the outside. The mapping from input to output is captured by equations 1.1 and 1.2.

The activation function The σ function is called *activation* function and should determine whether a perceptron unit is *active* or not. When artificial neural networks were first conceived, trying to mimic the brain structure, such function was a simple threshold function, trying to reproduce the behaviour of brain neurons: a neuron is *active*, i.e it's output ϕ_j is 1, if the sum of input stimuli $\sum_l w_{jl} \phi_l + b_j$ is greater than a given threshold τ .

$$\sigma_\tau(x) = \begin{cases} 1 & \text{if } x > \tau. \\ 0 & \text{otherwise.} \end{cases} \quad (1.3)$$

Such function, however, is problematic when we are to compute gradients because it's not continuous, so one of the following function is usually chosen:

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (1.4)$$

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (1.5)$$

These functions behave similarly to the threshold function, but, because of their *smoothness*, present no problems in computing gradients. Another function which is becoming a very popular choice is the *rectified linear unit*:

$$\text{ReLU}(x) = \begin{cases} x & \text{if } x > 0. \\ 0 & \text{otherwise.} \end{cases} \quad (1.6)$$

ReLU activation function is rather different from previous activation functions, some of these difference, in particular with respect to gradients will be analysed in later sections.

It's worth noticing that the activation function it's the only component which make artificial neuron networks a non linear model. Were we to choose a *linear* function as activation function we will end up with a simple linear model since the outputs of the network would be ultimately composed only of sums of products.

The bias term Let's consider the old threshold function σ_τ , and ask ourselves what the bias term is for, what does changing this term bring about. Suppose neuron j has no bias term, the neuron value would be $a_j = \sum_l w_{jl}\phi_l$; if $a_j > \tau$ that neuron is active otherwise it is not. Now, let's add the bias term to a_j ; we obtain that neuron j is active if $a_j > \tau - b_j$. So the effect of the bias term is to change the activation threshold of a given neuron. Using bias terms in a neural network architecture gives us the ability to change the activation threshold for each neuron; that's particularly important considering that we can learn such bias terms. We can do these same considerations in an analogous way for all the other activation functions.

Layered view of the net It's often useful to think of a neural network as series of layers, one on top of each other, as depicted in figure 1.3. 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 it's output is the output of the net.

When we describe a network in this way is also useful to adopt a different notation: we describe the weights of the net with a set of matrices W^k one for each layer, and neurons are no more globally indexed, instead with refer to a neuron with a relative index with respect to the layer; this allows to write easier equations in matrix notation ¹. In this notation W_{ij}^k is the weight of the link connecting neuron j of layer $k - 1$ to neuron i of level k

¹In the rest of the book we will refer to the latter notation as *layer notation* and to the previous one as *global notation*

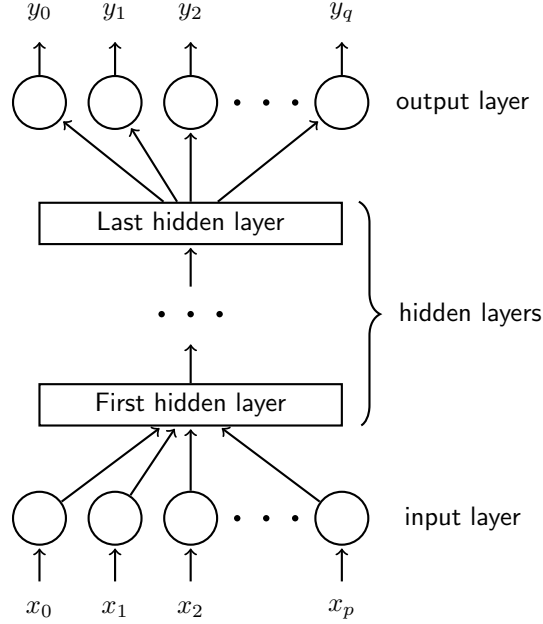


Figure 1.3: Layered structure of an artificial neural network

1.2 Feed forward neural networks

A feed forward 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 earlier layers are not possible.

Definition 1 (Feed forward neural network). A feed forward neural network is tuple:

$$\text{FFNN} \triangleq \langle \mathbf{p}, \mathcal{W}, \mathcal{B}, \sigma(\cdot), F(\cdot) \rangle$$

- $\mathbf{p} \in \mathbb{N}^U$ is the vector whose elements $p(k)$ are the number of neurons of layer k ; U is the number of layers
- $\mathcal{W} \triangleq \{W_{p(k+1) \times p(k)}^k, k = 1, \dots, U-1\}$ is the set of weight matrices of each layer
- $\mathcal{B} \triangleq \{\mathbf{b}^k \in \mathbb{R}^{p(k)}, k = 1, \dots, U\}$ 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

Remark 1. Given a FFNN:

- The number of output units is $p(U)$

- The number of input units is $p(1)$
- The total number of weights is $\mathcal{N}(\mathcal{W}) \triangleq \sum_{k=1}^U p(k+1)p(k)$
- The total number of biases is $\mathcal{N}(\mathcal{B}) \triangleq \sum_{k=1}^U p(k)$

Definition 2 (Output of a FFNN). Given a FFNN and an input vector $\mathbf{x} \in \mathbb{R}^{p(1)}$ the output of the net $\mathbf{y} \in \mathbb{R}^{p(U)}$ is defined by the following:

$$\mathbf{y} = F(\mathbf{a}^U) \quad (1.7)$$

$$\phi^i \triangleq \sigma(\mathbf{a}_i), \quad i = 2, \dots, U \quad (1.8)$$

$$\mathbf{a}^i \triangleq W^{i-1} \cdot \phi^{i-1} + \mathbf{b}^i \quad i = 2, \dots, U \quad (1.9)$$

$$\phi^1 \triangleq \mathbf{x} \quad (1.10)$$

1.2.1 Learning with FFNNs

A widespread application of neural networks is that of machine learning. In the following we will model an optimisation problem which rely on FFNNs. To model an optimisation problem we first need to define a data-set D as

$$D \triangleq \{\bar{\mathbf{x}}^{(i)} \in \mathbb{R}^p, \bar{\mathbf{y}}^{(i)} \in \mathbb{R}^q, i = 1, \dots, N\} \quad (1.11)$$

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

$$L_D(\mathcal{W}, \mathcal{B}) \triangleq \frac{1}{N} \sum_{i=1}^N L(\bar{\mathbf{x}}^{(i)}, \bar{\mathbf{y}}^{(i)}) \quad (1.12)$$

$L(\mathbf{x}, \mathbf{y}) : \mathbb{R}^{p(U)} \rightarrow \mathbb{R}$ is an arbitrary loss function computed on the i^{th} example. Note that \mathbf{y} is the output of the networks, so it depends on $(\mathcal{W}, \mathcal{B})$

The problem is then to find a FFNN which minimise L_D . As we have seen feed forward neural networks allow for large customisation: the only variables in the optimisation problem are the weights and the biases, 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. 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 optimisation problem becomes:

$$\min_{\mathcal{W}, \mathcal{B}} L_D(\mathcal{W}, \mathcal{B}) \quad (1.13)$$

1.2.2 Gradient

Consider a FFNN $= \langle \mathbf{p}, \mathcal{W}, \mathcal{B}, \sigma(\cdot), F(\cdot) \rangle$, let $L : \mathbb{R}^{p(U)} \rightarrow \mathbb{R}$ a loss function and $g(\cdot) : \mathbb{R}^{\mathcal{N}(\mathcal{W}) + \mathcal{N}(\mathcal{B})} \rightarrow \mathbb{R}$ be the function defined by

$$g(\mathcal{W}, \mathcal{B}) \triangleq L(F(a^U(\bar{\mathbf{x}}^{(i)}(\mathcal{W}, \mathcal{B}))), \bar{\mathbf{y}}^{(i)}) \quad (1.14)$$

Equation 1.14, though it seems rather scary, express a very simple thing: we consider a single input example $\bar{\mathbf{x}}^{(i)}$, we run it through the network and we confront it's output $F(a^U(\bar{\mathbf{x}}^{(i)}))$ with it's label $\bar{\mathbf{y}}^{(i)}$ using the loss function L ; the function $g(\mathcal{W}, \mathcal{B})$ it's simply the loss function computed on the i^{th} example which of course depends only on the weights and biases of network since the examples are constant.

$$\frac{\partial g}{\partial W^i} = \nabla L \cdot J(F) \cdot \frac{\partial \mathbf{a}^U}{\partial W^i} \quad (1.15)$$

$$= \frac{\partial g}{\partial \mathbf{a}^U} \cdot \frac{\partial \mathbf{a}^U}{\partial W^i} \quad (1.16)$$

We can easily compute $\frac{\partial g}{\partial \mathbf{a}^U}$ once we define $F(\cdot)$ and $L(\cdot)$, note that the weights are not involved in such computation. Let's derive an expression for $\frac{\partial \mathbf{a}^U}{\partial W^i}$. We will start deriving such derivative using global notation. Let's consider a single output unit u and a weight w_{lj} linking neuron j to neuron l .

$$\frac{\partial a_u}{\partial w_{lj}} = \frac{\partial a_u}{\partial a_l} \cdot \frac{\partial a_l}{\partial w_{lj}} \quad (1.17)$$

$$= \delta_{ul} \cdot \phi_j \quad (1.18)$$

where we put

$$\delta_{ul} \triangleq \frac{\partial a_u}{\partial a_l}$$

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

$$P(l) = \{k : \exists \text{ a link between } l \text{ and } k \text{ with weight } w_{lk}\} \quad (1.19)$$

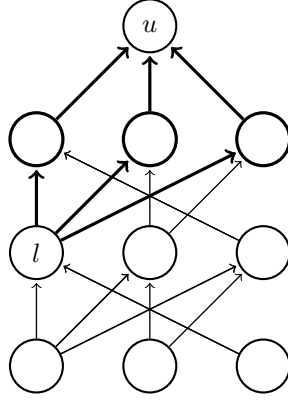
We can compute δ_{ul} simply applying the chain rule, if we write it down in bottom-up style, as can be seen in figure 1.4, we obtain:

$$\delta_{ul} = \sum_{k \in P(l)} \delta_{uk} \cdot \sigma'(a_k) \cdot w_{kl} \quad (1.20)$$

The derivatives with respect to biases are compute in the same way:

$$\frac{\partial a_u}{\partial b_l} = \frac{\partial a_u}{\partial a_l} \cdot \frac{\partial a_l}{\partial b_l} \quad (1.21)$$

$$= \delta_{ul} \cdot 1 \quad (1.22)$$

Figure 1.4: Nodes and edges involved in $\frac{\partial a_u}{\partial a_l}$

In layered notation we can rewrite the previous equations as:

$$\frac{\partial a^U}{\partial W^i} = \frac{\partial a^U}{\partial a^{i+1}} \cdot \frac{\partial^+ a^{i+1}}{\partial W^i} \quad (1.23)$$

$$\frac{\partial^+ a^{i+1}}{\partial W_j^i} = \begin{bmatrix} \phi_j^i & 0 & \cdots & \cdots & 0 \\ 0 & \phi_j^i & \cdots & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & \cdots & \cdots & \phi_j^i \end{bmatrix} \quad (1.24)$$

$$\frac{\partial a^U}{\partial a^i} \triangleq \Delta^i = \begin{cases} \Delta^{i+1} \cdot \text{diag}(\sigma'(\mathbf{a}^{i+1})) \cdot W^i & \text{if } i < U \\ Id & \text{if } i == U \\ 0 & \text{otherwise.} \end{cases} \quad (1.25)$$

where

$$\text{diag}(\sigma'(\mathbf{a}^i)) = \begin{bmatrix} \sigma'(a_1^i) & 0 & \cdots & \cdots & 0 \\ 0 & \sigma'(a_2^i) & \cdots & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \cdots & \cdots & \cdots & \sigma'(a_{p(i)}^i) \end{bmatrix} \quad (1.26)$$

$$\frac{\partial a^U}{\partial b^i} = \frac{\partial a^U}{\partial a^i} \cdot \frac{\partial a^i}{\partial b^i} \quad (1.27)$$

$$= \Delta^i \cdot Id \quad (1.28)$$

Backpropagation Previous equations are the core of the famous *backpropagation* algorithm which was first introduced by Rumelhart et al. [5]. The algorithm consists in two *passes*, a *forward pass* and a *backward pass* which give the name

to the algorithm. The *forward pass* start from the first layer, compute the hidden units values and the proceed to upper layers using the value of the hidden units \mathbf{a}^i of previous layers which have already been computed. The *backward pass* instead, start from the upmost layer and computes Δ^i which can be computed, as we can see from equation 1.25 , once it's known Δ^{i+1} , which has been computed in the previous step, and \mathbf{a}^i which has been computed in the *forward pass*.

Backpropagation algorithm has time complexity $\mathcal{O}(\mathcal{N}(\mathcal{W}))$.

1.3 Recurrent 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 $j < i$, as can be seen in figure 1.5. This is a key difference, as we will see in later sections, because it introduces *memory* in the network changing, somehow, the expressiveness of the model.

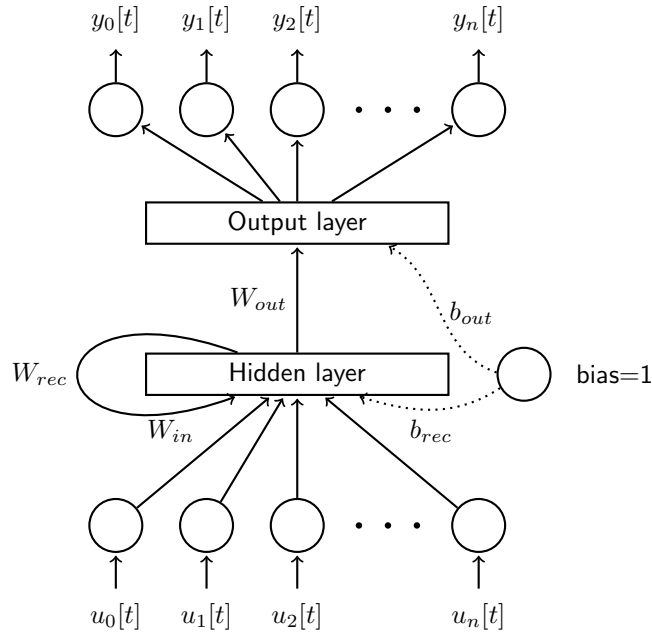


Figure 1.5: RNN model

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; in other words, 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 temporal inputs; that's not true, in fact we can apply RNNs 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 at all after the first time step.

Definition 3 (Recurrent neural network). A recurrent neural network is tuple

$$\text{RNN} \triangleq \langle \mathcal{W}, \mathcal{B}, \sigma(\cdot), F(\cdot) \rangle$$

- $\mathcal{W} \triangleq \{W^{in}, W^{rec}, W^{out}\}$ where
 - 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
- $\mathcal{B} \triangleq \{\mathbf{b}^{rec}, \mathbf{b}^{out}\}$ where
 - $\mathbf{b}^{rec} \in \mathbb{R}^r$ 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

Remark 2. Given a RNN

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

Definition 4 (Output of a RNN). Given a RNN and an input sequences $\{\mathbf{x}\}_{t=1, \dots, T}$, $\mathbf{x}_t \in \mathbb{R}^p$ the output sequence $\{\mathbf{y}\}_{t=1, \dots, T}$, $\mathbf{y}_t \in \mathbb{R}^o$ of the net is defined by the following:

$$\mathbf{y}_t \triangleq F(W^{out} \cdot \mathbf{a}_t + \mathbf{b}^{out}) \quad (1.29)$$

$$\mathbf{a}_t \triangleq W^{rec} \cdot \phi_{t-1} + W^{in} \cdot \mathbf{x}_t + \mathbf{b}^{rec} \quad (1.30)$$

$$\phi_t \triangleq \sigma(\mathbf{a}_t) \quad (1.31)$$

$$\phi_0 \triangleq \vec{0} \quad (1.32)$$

As we can understand from definition 4, there is only one recurrent layer, whose weights are the same for each time step, so one could 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 that of a feed forward neural network. As we can observe in figure 1.6, the unfolding of the network through time consist of putting identical version of the same recurrent layer one 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; another important unlikeness is of course the presence of additional timed inputs for each unfolded layer.

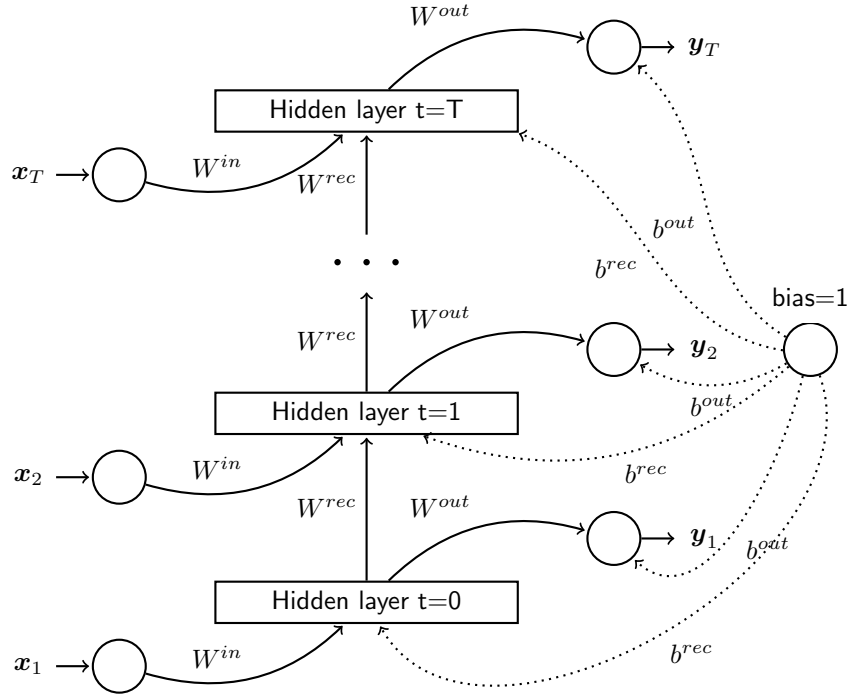


Figure 1.6: Unfolding of a RNN

1.3.1 Learning with RNNs

We can model an optimisation 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 data-set D :

$$D \triangleq \{ \{ \bar{\mathbf{x}} \}_{t=1, \dots, T}, \bar{\mathbf{x}}_t \in \mathbb{R}^p, \{ \bar{\mathbf{y}} \}_{t=1, \dots, T}, \bar{\mathbf{y}}_t \in \mathbb{R}^o; i = 1, \dots, N \} \quad (1.33)$$

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

$$L_D(\mathcal{W}, \mathcal{B}) \triangleq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T L_t(\bar{\mathbf{x}}_t^{(i)}, \mathbf{y}_t^{(i)}) \quad (1.34)$$

L_t is an arbitrary loss function at time step t .

The definition takes into account the output for each temporal step, depending on the task at hand, it could be relevant or not to consider intermediate outputs; that's not a limitation, in fact we could define a loss which is computed only on the last output vector, at time T , and adds 0 for each other time step.

1.3.2 Gradient

Consider a RNN $= \langle \mathcal{W}, \mathcal{B}, \sigma(\cdot), F(\cdot) \rangle$. Let $L : \mathbb{R}^o \rightarrow \mathbb{R}$ a loss function:

$$L(\mathcal{W}, \mathcal{B}) \triangleq \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T L_t(\bar{\mathbf{x}}_t^{(i)}, \mathbf{y}_t^{(i)})$$

Let $g_t(\cdot) : \mathbb{R}^{\mathcal{N}(W)+\mathcal{N}(B)} \rightarrow \mathbb{R}$ be the function defined by

$$g_t(\mathcal{W}, \mathcal{B}) \triangleq L(F(\mathbf{y}^t(\mathcal{W}, \mathcal{B})))$$

and

$$g(\mathcal{W}, \mathcal{B}) \triangleq \sum_{t=1}^T g_t(\mathcal{W}, \mathcal{B})$$

$$\frac{\partial g}{\partial W^{rec}} = \sum_{t=1}^T \nabla L_t \cdot J(F) \cdot \frac{\partial \mathbf{y}^t}{\partial \mathbf{a}^t} \cdot \frac{\partial \mathbf{a}^t}{\partial W^{rec}} \quad (1.35)$$

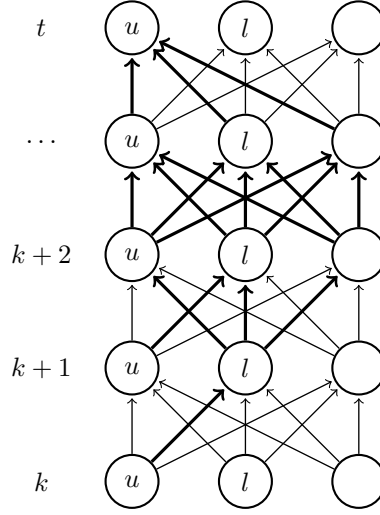
$$= \sum_{t=1}^T \frac{\partial g_t}{\partial \mathbf{a}^t} \cdot \frac{\partial \mathbf{a}^t}{\partial W^{rec}} \quad (1.36)$$

As we noticed for FNNs it's easy to compute $\frac{\partial g_t}{\partial \mathbf{a}^t}$ once we define $F(\cdot)$ and $L(\cdot)$, note that the weights are not involved in such computation. Let's see how to compute $\frac{\partial \mathbf{a}^t}{\partial W^{rec}}$.

Let's consider a single output unit u , and a weight w_{lj} , we have

$$\frac{\partial a_u^t}{\partial w_{lj}} = \sum_{k=1}^t \frac{\partial a_u^t}{\partial a_l^k} \cdot \frac{\partial a_l^k}{\partial w_{lj}} \quad (1.37)$$

$$= \sum_{k=1}^t \delta_{lu}^{tk} \cdot \phi_j^{t-1} \quad (1.38)$$

Figure 1.7: Nodes involved in $\frac{\partial a_u^t}{\partial a_l^k}$

where

$$\delta_{lu}^{tk} \triangleq \frac{\partial a_u^t}{\partial a_l^k} \quad (1.39)$$

Let's observe a first difference from ffn case: since the weights are shared in each unfolded layer, in equation 1.37 we have to sum over time.

Let $P(l)$ be the set of parents of neuron l , defined as the set of parents in the unfolded network.

$$\delta_{lu}^{tk} = \sum_{h \in P(l)} \delta_{hu}^{tk} \cdot \sigma'(a_h^{t-1}) \cdot w_{hl} \quad (1.40)$$

In figure 1.7 we can see the arcs which are involved in the derivatives in the unfolded network.

In matrix notation we have:

$$\frac{\partial \mathbf{a}^t}{\partial \mathbf{W}^{rec}} = \sum_{k=1}^t \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \cdot \frac{\partial^+ \mathbf{a}^k}{\partial \mathbf{W}^{rec}} \quad (1.41)$$

$$\frac{\partial^+ \mathbf{a}^k}{\partial \mathbf{W}_j^{rec}} = \begin{bmatrix} \phi_j^k & 0 & \dots & \dots & 0 \\ 0 & \phi_j^k & \dots & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & \dots & \dots & \dots & \phi_j^k \end{bmatrix} \quad (1.42)$$

$$\triangleq \Delta^{tk} \quad (1.43)$$

$$\Delta^{tk} = \Delta^{t(k+1)} \cdot \text{diag}(\sigma'(\mathbf{a}^k)) \cdot W^{rec} \quad (1.44)$$

$$= \prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot W^{rec} \quad (1.45)$$

The derivatives with respect to W^{in} and \mathbf{b}^{rec} have the same structure. The derivatives with respect to $W^{out}, \mathbf{b}^{out}$ are straightforward:

$$\frac{\partial \mathbf{g}}{\partial W^{out}} = \sum_{t=1}^T \frac{\partial g_t}{\partial \mathbf{y}^t} \cdot J(F) \cdot \frac{\partial \mathbf{y}^t}{\partial W^{out}} \quad (1.46)$$

$$\frac{\partial \mathbf{g}}{\partial \mathbf{b}^{out}} = \sum_{t=1}^T \frac{\partial g_t}{\partial \mathbf{y}^t} \cdot J(F) \cdot \frac{\partial \mathbf{y}^t}{\partial \mathbf{b}^{out}} \quad (1.47)$$

Backpropagation through time (BPTT) *Backpropagation through time* is an extension of the *backpropagation* algorithm we described for FNNs, we can think of BPTT simply as a standard BP in the unfolded network. The same considerations done for BP also apply for BPTT, the difference is of course in how derivatives are computed, equation 1.45. Time complexity is easily derived noticing that in the unfolded network there are $n \cdot T$ units, where n is the number of units of the RNN. This yields time complexity $\mathcal{O}(\mathcal{N}(\mathcal{W}) \cdot T)$. Please see [8] for more details.

1.3.3 The vanishing and exploding gradient problem

The training of recurrent neural networks is afflicted by the so called *exploding* and *vanishing* gradient problem. Such problem is directly linked to the notion of memory. When we talk about memory what are we really talking about is the dependency of neurons output at a given time t from previous time steps, that is how ϕ^t depends on ϕ^k with $t > k$. This dependency is captured by the expression $\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k}$. Obviously when such expression equals zero it means neuron output at time t is not affected by output at time k . The terms $\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k}$ are usually referred as *long term* contribution when $k \ll t$ or *short term* contributions otherwise. We talk of *exploding* gradient problem when *long term* components grow exponentially, on the contrary of *short term* ones, causing the overall norm of the gradient to explode. We refer to *vanishing* gradient problem when, vice versa, *long term* components diminish exponentially.

Vanishing gradient means long term components approach zero value, this in turn leads to the fact that the output of the net won't depend on inputs of distant temporal steps, i.e the output sequence is determined only by recent temporal input: we say that the network doesn't have memory. Evidently this can have catastrophic effects on the classification error. Imagine we would like

to classify an input sequence as positive whether or not it contains a given character. It would seem a rather easy task, however, if the neural network we are training suffers from the *vanishing* gradient issue, it could perform the classification using only the most recent temporal inputs. What if the character was at the beginning of the sequence? Of course the prediction would be wrong.

Exploding gradient seems to be a different kind of a problem, it does not affect the ability of the network to use information from distant temporal step, on the contrary we have very strong information about where to go using the gradient direction. One could argue that some components, namely, the long term ones, have gradient norm exponentially bigger than short term ones. I fail to see why this could be a problem: the output of the first time steps influences all successive steps, so changing the output of a long term neuron do imply big changes in the output of more distant in time neurons. The only situation in which *exploding* gradient can be a problem is when using learning algorithms like gradient descent with constant step: if we are to compute a step in the gradient direction with a fixed step and gradient has too big norm we may make a too big step.

Let's now return to the nature of the problem and try to explaining the mechanics of it. We have seen in the previous section that

$$\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} = \prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot W^{rec} \quad (1.48)$$

Intuitively we can understand why such problems arises, more evidently in *long term* components, just by looking at equation 1.55. We can notice each temporal contribution is the product of $l = t - k - 1$ jacobian matrix, so in *long term* components l is large and depending on the values of the matrices in the product we can go exponentially fast towards 0 or infinity.

Hochreiter Analysis: Weak upper bound In this paragraph we report some useful consideration made by Hochreiter, please see [?] for more details.

Let's put:

$$\|A\|_{max} \triangleq \max_{i,j} |a_{ij}|$$

$$\sigma'_{max} \triangleq \max_{i=k,\dots,t-1} \{\|\text{diag}(\sigma'(\mathbf{a}^i))\|_{max}\}$$

Since

$$\|A \cdot B\|_{max} \leq r \cdot \|A\| \cdot \|B\|_{max} \quad \forall A, B \in \mathbb{R}_{r \times r} \quad (1.49)$$

it holds:

$$\left\| \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \right\|_{max} = \left\| \prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot W^{rec} \right\|_{max} \quad (1.50)$$

$$\leq \prod_{i=t-1}^k r \cdot \|\text{diag}(\sigma'(\mathbf{a}^i))\|_{max} \cdot \|W^{rec}\|_{max} \quad (1.51)$$

$$\leq (r \cdot \sigma'_{max} \cdot \|W^{rec}\|_{max})^{t-k-1} \quad (1.52)$$

$$= \tau^{t-k-1} \quad (1.53)$$

where

$$\tau \triangleq r \cdot \sigma'_{max} \cdot \|W^{rec}\|_{max}$$

So we have exponential decay if $\tau < 1$. We can match this condition if $\|W^{rec}\|_{max} \leq \frac{1}{r \cdot \sigma'_{max}}$. As pointed out by Hochreiter in his work we can match this condition, in the case of sigmoid activation function by choosing $\|W^{rec}\|_{max} < \frac{1}{0.25 \cdot r}$.

Let's note that we would actually reach this upper bound for some i, j only if all the path cost have the same sign and the activation function takes always maximal value.

An upper bound with singular values Lets decompose W^{rec} using the singular value decomposition. We can write

$$W^{rec} = S \cdot D \cdot V^T \quad (1.54)$$

where S, V^T are squared orthogonal matrices and $D \triangleq \text{diag}(\mu_1, \mu_2, \dots, \mu_r)$ is the diagonal matrix containing the singular values of W^{rec} . Rewriting equation 1.55 using this decomposition leads to

$$\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} = \prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot S \cdot D \cdot V^T \quad (1.55)$$

Recalling that for any orthogonal matrix U it holds

$$\|U \mathbf{x}\|_2 = \|\mathbf{x}\|_2$$

and using

$$\|\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_r) \cdot \mathbf{x}\|_2 \leq \lambda_{max} \|\mathbf{x}\|_2$$

we get

$$\left\| \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \cdot \mathbf{e}_i \right\|_2 = \left\| \left(\prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot S \cdot D \cdot V^T \right) \cdot \mathbf{e}_i \right\|_2 \quad (1.56)$$

$$\leq (\sigma'_{max} \cdot \mu_{max})^{t-k-1} \quad (1.57)$$

The previous equation provide us a sufficient condition, $\sigma'_{max} \cdot \mu_{max} < 1$, as in Hochreiter's analysis, for exponential decay of long term components. In this case however the bound depends on the singular value of the recurrent weights rather than on the maximal weight of the matrix itself.

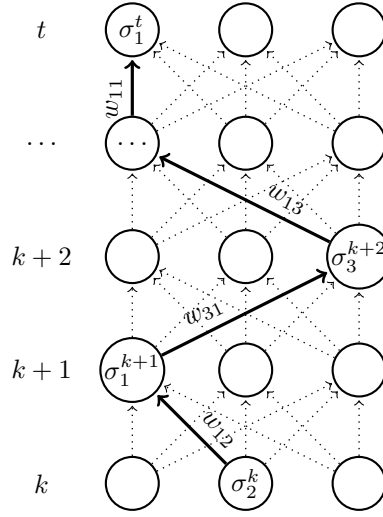


Figure 1.8: The cost for a path from neuron 2 at time k to neuron 1 at time t is $w_{12}w_{31}w_{13} \dots w_{11} \cdot \sigma_2^k \sigma_1^{k+1} \sigma_3^{k+2} \dots \sigma_1^{t-1}$

Explaining the problem using the network's graph Let's now dig a bit deeper and rewrite equation 1.55 with respect to a couple of neurons i and j .

$$\frac{\partial a_i^t}{\partial a_j^k} = \sum_{q \in P(j)} \sum_{l \in P(q)} \dots \sum_{h: i \in P(h)} w_{qj} \dots w_{jh} \cdot \sigma'(a_j^k) \sigma'(a_q^{k+1}) \dots \sigma'(a_i^{t-1}) \quad (1.58)$$

Observing the previous equation we can argue that each derivatives it's the sum of p^{t-k-1} terms; each term represents the path cost from neuron i to neuron j in the unfolded network, obviously there are p^{t-k-1} such paths. If we bind the cost $\sigma'(a_i^t)$ to neuron l in the t^{th} layer in the unfolded network we can read the path cost simply surfing the unfolded network multiply the weight of each arc we walk through and the cost of each neuron we cross, as we can see from figure 1.8.

We can further characterize each path cost noticing that we can separate two components, one that depends only on the weights $w_{qj} \dots w_{jh}$ and the other that depends both on the weights and the inputs $\sigma'(a_j^k) \sigma'(a_q^{k+1}) \dots \sigma'(a_i^{t-1})$.

The ReLU case ReLU case is a bit special, because of it's derivative. ReLU's derivative is a step function, it can assume only two values: 1 when the neuron is active, 0 otherwise. Returning to the path graph we introduced earlier we can say that a path is *enabled* if each neuron in that path is active. In fact if we encounter a path which cross a non active neuron it's path cost will be 0; on the contrary for an *enabled* path the cost will be simply the product of weight of the arcs we went through, as we can see in figure 1.9

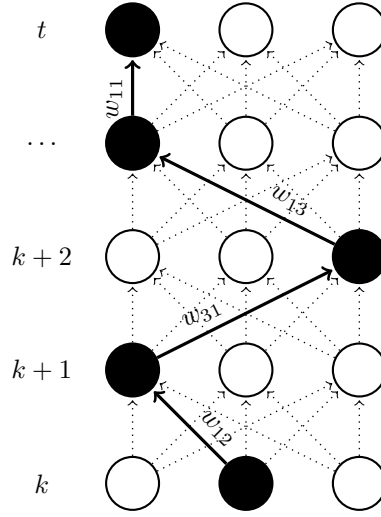


Figure 1.9: The cost for an enabled path from neuron 2 at time k to neuron 1 at time t is $w_{12}w_{31}w_{13}\dots w_{11}$

So $|(\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k})_{ij}|$ ranges from 0, when no path is enabled to, $|((W^{rec})^{t-k-1})_{ij}|$ when all paths are enabled and all path cost have the same sign, which is consistent with what we found in Hochreiter analysis. We can argue then that ReLU has an advantage over sigmoidal activation functions, for instance, because gradient depends only on the W^{rec} matrix: ReLU function *only* enables or disables the paths but doesn't change their costs as sigmoids do

Poor solutions Let's pretend we have found, with some learning technique, an assignment for all the weights which causes the gradient to have zero norm. We could be happy with it and claim to have 'solved' the problem. However, by chance, we discover that $\frac{\partial \mathbf{a}^T}{\partial \mathbf{a}^k}$ has zero norm for all time steps $k < \tau$. So, the output of the network doesn't depend on the inputs of the sequence for those time steps. In other words we have found a possibly optimal solution for the truncated sequence $x_{[\tau:T]}$. The solution we have found is an optimal candidate to be a bad local minimum.

As a final observation on this matter it's worth noticing how a bad initialization of W^{rec} can lead to poor solutions or extremely large convergence time just because such initialization imply $\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k}$ approaching zero norm for $t \gg k$. Moreover, even if we somehow provide an initialization matrix which is unaffected by this curse, it's certainly possible that we reach such a bad matrix during learning phase. Several techniques have been proposed to overcome this problem, they will be the topic of later chapters.

1.4 Activation functions and gradient

Activation functions play a central role in the artificial neural networks model, they are responsible for the non linearity of the model. In the history of neural networks several activation functions have been proposed and used, in the following some of them are taken into consideration underling the difference between them, with a special focus on their derivative expression. A special class of activation function, is that of *squashing* functions.

Definition 5. A function $f(\cdot) : \mathbb{R} \rightarrow [a, b]$ with $a, b \in \mathbb{R}$ is said to be a *squashing* function if it is not decreasing and

$$\lim_{x \rightarrow +\infty} f(x) = b \quad (1.59)$$

$$\lim_{x \rightarrow -\infty} f(x) = a \quad (1.60)$$

Step function, ramp function and all sigmoidal functions are all examples of squashing functions.

Remark 3. An important property of a *squashing* function $\sigma(\cdot)$ is that

$$\lim_{\alpha \rightarrow +\infty} \sigma(\alpha \cdot (x - \tau)) = \begin{cases} b \cdot \sigma_{\tau}(x) & \text{if } x > \tau \\ a + \sigma_{\tau}(x) & \text{otherwise} \end{cases} \quad (1.61)$$

being σ_{τ} the usual step function. This property is extensively used in several proofs of the universal approximator property of neural networks. Roughly speaking we can say that *squashing* functions act as step functions at the limit; please note that this property has a practical use since inputs of activation functions are the weighted sum of some neurons output, so activation function inputs can be arbitrarily big or small.

Sigmoid

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (1.62)$$

$$\text{sigmoid}'(x) = \text{sigmoid}(x) \cdot (1 - \text{sigmoid}(x)) \quad (1.63)$$

As we can see from figure 1.10, the sigmoid derivative has only one maximum in 0 where it assume value 0.25. Receding from 0, in both direction leads to regions where the the derivative take zero value, such regions are called *saturation* regions. If we happen to be in such regions, for a given neuron, we cannot learn anything since that neuron doesn't contribute to the gradient.

Tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (1.64)$$

$$\tanh'(x) = 1 - \tanh^2(x) \quad (1.65)$$

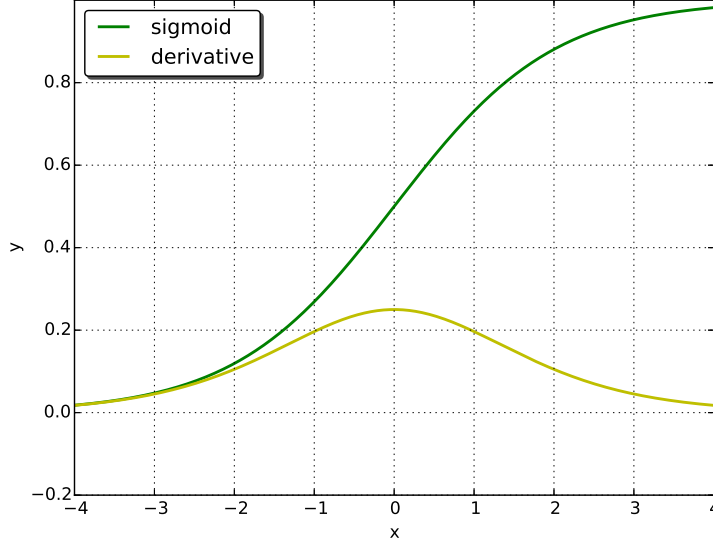


Figure 1.10: sigmoid and its derivative

As we can see from figure 1.11 \tanh (and its derivative) have a behavior similar to the sigmoid one; Again we have two saturation region towards infinity: that's typical of all squashing functions.

ReLU

$$ReLU(x) = \begin{cases} x & \text{if } x > 0. \\ 0 & \text{otherwise.} \end{cases} \quad (1.66)$$

$$ReLU'(x) = \begin{cases} 1 & \text{if } x > 0. \\ 0 & \text{otherwise.} \end{cases} \quad (1.67)$$

ReLU is a bit different from other activation function seen so far: the main difference is that's it's not a squashing function. As we can see from figure 1.12, ReLU's derivative is the step function; it has only one *saturation* region $(-\infty, 0]$ and a region in which is always takes value 1, $(0, +\infty]$ This leads to the fact that we cannot learn to *turn on* a switched off neuron ($x < 0$), but we have no *saturation* region toward infinity. The constancy of the derivative in the *active* region is a distinctive property of ReLU which plays an important role in the *vanishing gradient* problem too, as it's explained in the relative section.

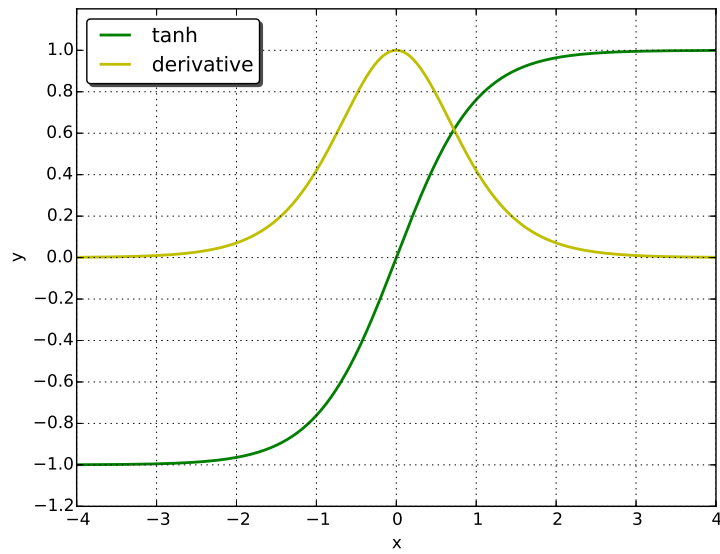


Figure 1.11: tanh and its derivative

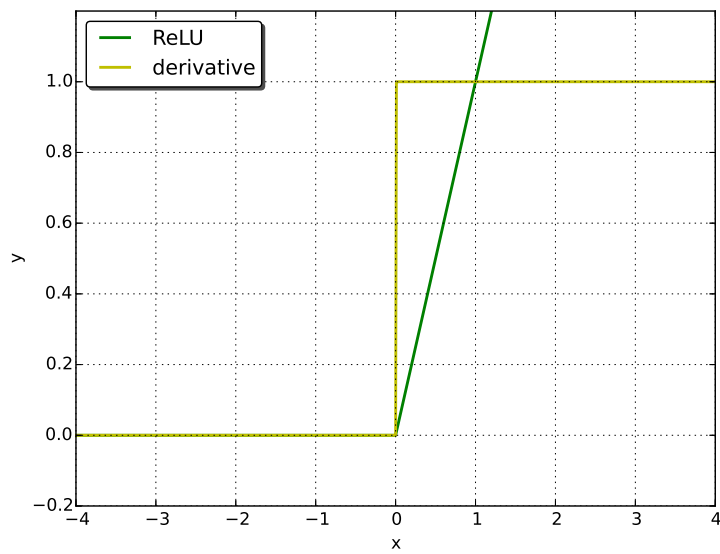


Figure 1.12: ReLU and its derivative

1.5 On expressiveness

In this section we will investigate the expressive power of neural networks, presenting some results that motivate the use of neural networks as learning models and underline the differences between the two paradigm of FNNs and RNNs.

One of the first import results regarding the expressive power of neural networks it's due to Hornik et al. [3] which basically states “*Multilayered feed forward networks with at least one hidden layer, using an arbitrary squashing function, can approximate virtually any function of interest to any desired degree of accuracy provided sufficiently many hidden units are available*”.

To give a more formal result we need first to define what *approximate to any degree of accuracy means*, this concept is captured in definition 6

Definition 6. A subset S of \mathbb{C}^n (continuous functions in \mathbb{R}^n) is said to be *uniformly dense on compacta in \mathbb{C}^n* if \forall compact set $K \subset \mathbb{R}^n$ holds: $\forall \epsilon > 0$, $\forall g(\cdot) \in \mathbb{C}^n \exists f(\cdot) \in S$ such that $\sup_{x \in K} \|f(x) - g(x)\| < \epsilon$

Hornik result is contained in theorem 1.

Theorem 1. For every squashing function σ , $\forall n \in \mathbb{N}$, feed forward neural networks with one hidden layer are a class of functions which is *uniformly dense on compacta in \mathbb{C}^n* .

Theorem 1 extends also to Borel measurable functions, please see [3] for more details.

A survey of other approaches, some of which constructive, which achieve similar results can be found in [6] At the moment I don't know of any results concerning ReLU activation function.

This results implies that FNN are *universal approximators*, this is a strong argument for using such models in machine learning. It's important to notice, however, that the theorem holds if we have *sufficiently many* units. In practice the number of units will be bounded by the machine capabilities and by computational time, of course greater the number of units greater will be the learning time. This will limit the expressiveness of the network to a subset of all measurable functions.

Let's now turn our attention to RNNs and see how the architectural changes, namely the addition of backward links, affect the expressive power of the model. It suffice to say that RNNs are as powerful as Turing machine. Siegelman and Sontag [7] proved the existence of a finite neural network, with sigmoid activation function, which simulates a universal Turing machine. Hyötyniemi [4] proved, equivalently, that Turing machine are recurrent neural network showing how to build a network, using instead ReLU activation function, that performs step by step all the instruction of a computer program. Hyötyniemi work is particularly interesting because it shows how to construct a network that simulate an algorithm written in a simple language equivalent to a Turing machine. For

each instruction type (increment,decrement,conditional branch,...) a particular setting of weights and neuron is devised allowing the net to simulate step by step the behaviour of the program. In the program equivalent network there are a unit for each program variable and one or two, depending on the instruction type, units for each program instruction. This is very interesting from an expressiveness point of view since it bounds the number of units we ought to use with the length of the algorithm we are trying to reproduce.

For better understanding the implications of this fact, imagine how many complex function you can express with short algorithms, for example fractals (approximations). It's worth underling the difference with feed forward neural networks where a large number of units seems to be required. This seems to suggest that FFNNs and RNNs differ mainly in a manner of representation, where FNNs use space to define a somehow explicit mapping from input to output, RNNs use time to implicitly define an algorithm responsible for such mapping.

This seems extremely good news, since we could simulate turing machines, hence all algorithms we can think of, using a recurrent neural network with a relatively small number of units; recall that for FFNN we had to suppose infinitely many units to obtain the universal approximator property. Of course there is a pitfall: we can simulate any turing machine but we have to allow sufficiently many time steps and choose a termination criterion. This is of course impractical, and we don't use RNNs in this way. Usually the number of time steps is choose to be equal to the input sequence length. This of course restrict the class of algorithm we can learn with RNNs. A particular class of algorithms suited to be learned by such models is that of algorithms consisting in one loop, inside which some invariants are enforced. The output, step by step, depends on the loop invariants.

Chapter 2

An historical review of proposed solutions

distinction between architectural and non architectural solutions, historical review

2.1 Long short-term memory

Long short-term memory (LSTM) were proposed (1997) by Hochreiter and Schmidhuber [2] as a novel network structure to address the vanishing gradient problem, which was first studied by Hochreiter (1991) in his diploma thesis, a milestone of deep learning.

The idea behind this structure is to enforce a constant error flow, that is to say, to have constant gradient norm, thus preventing the gradient to vanish. This is done by introducing special types of neurons called *memory cells* and *gate units*. As we can see by looking at figure 2.1, a memory cell is essentially a neuron with a self connection with unitary weight, whose input and output are managed by two multiplicative neurons: the gate units.

The memory cell and the gate units behave accordingly to the following:

$$u_j^t \triangleq \sigma(\sum_i w_{ui} \cdot u_j^{t-1}) \quad (2.1)$$

$$o_j^t \triangleq \sigma(\sum_i w_{oi} \cdot o_j^{t-1}) \quad (2.2)$$

$$a_j^t \triangleq g(\sum_i w_{ij} \cdot \phi_i^t) \quad (2.3)$$

$$m^t \triangleq a_j \cdot u_j^t + (1 \cdot m^{t-1}) \quad (2.4)$$

$$\phi_j \triangleq h(m^t) \cdot o_j^t \quad (2.5)$$

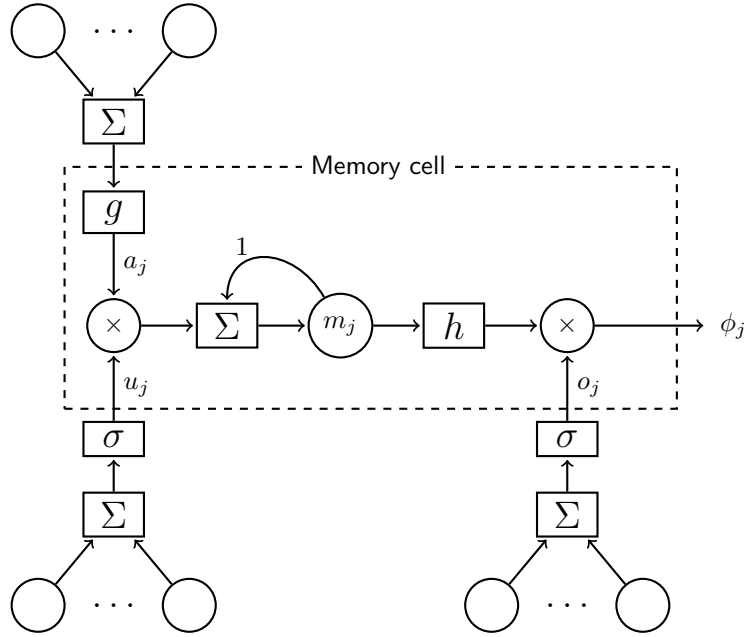


Figure 2.1: Memory cell and gate units of LSTM network

As we can see from equation 2.4, the value of the memory cell $m(t)$ remains constant as long as the input gate u doesn't "open" causing a "write" operation. Similarly the output o of the memory cell, which is connected with the other neurons of the network, is controlled by an output gate: the memory will have a non zero output only if the output gate opens, which we could call a "read" operation. As for constant error flow it is ensured because the memory cell has only a self-loop with unitary weight.

Memory cells, guarded by gate units can be employed in networks with different topology alongside traditionally input, output and hidden units. Another way to look at this kind of architecture is to think of memory cells as units able to store one bit of information, even for long periods of time, hence able to learn distant time correlations between inputs.

As we have seen these network units are specifically designed to store information, through the use of gates; these gates however are no different from other units, apart from the fact they are multiplicative units, hence without further precautions, the networks would incur in the same vanishing problem it aimed to resolve. In fact LSTM comes with a proper learning algorithm: essentially errors arriving at memory cells inputs are not propagated back in time, only the error within the memory cell gets propagated; in other words gradients are truncated taking into account only the self-connection of the memory cells and not its other input connections, hence providing constant error flow.

LSTM units have proven to be very successful in various tasks and even

at present times (2015), they continue to be employed. In recent implementations however, alongside small modifications, as the introduction of other gates, LSTM architecture is often used without the original learning algorithm which is often replaced by a standard stochastic gradient descend as done in [1].

2.2 Gated recurrent units

Appendix A

Notation

Let $F : \mathbb{R}^N \rightarrow \mathbb{R}^M$ be defined by

$$F(\mathbf{x}) = (f_1(\mathbf{x}), f_2(\mathbf{x}), \dots, f_M(\mathbf{x})) \text{ for some } f_i : \mathbb{R}^N \rightarrow \mathbb{R} \quad (\text{A.1})$$

Definition 7 (Derivative with respect to a vector). We define the derivative of $F(x(\mathbf{w}))$ with respect to a vector \mathbf{w} of p elements as the $M \times p$ matrix

$$\frac{\partial F}{\partial \mathbf{w}} \triangleq \begin{bmatrix} \frac{\partial f_1}{\partial \mathbf{w}_1} & \frac{\partial f_1}{\partial \mathbf{w}_2} & \dots & \dots & \frac{\partial f_1}{\partial \mathbf{w}_p} \\ \frac{\partial f_2}{\partial \mathbf{w}_1} & \frac{\partial f_2}{\partial \mathbf{w}_2} & \dots & \dots & \frac{\partial f_2}{\partial \mathbf{w}_p} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \frac{\partial f_M}{\partial \mathbf{w}_1} & \frac{\partial f_M}{\partial \mathbf{w}_2} & \dots & \dots & \frac{\partial f_M}{\partial \mathbf{w}_p} \end{bmatrix} \quad (\text{A.2})$$

Definition 8 (Derivative with respect to a matrix). We define the derivative of $F(x(W))$ with respect to a vector W , being W_j the j^{th} column of a $p \times m$ matrix W as the $M \times (p \cdot m)$ matrix:

$$\frac{\partial F}{\partial W} \triangleq \left[\begin{array}{c|c|c|c} \frac{\partial F}{W_1} & \frac{\partial F}{W_2} & \dots & \frac{\partial F}{W_m} \end{array} \right] \quad (\text{A.3})$$

Bibliography

- [1] Alex Graves. Generating sequences with recurrent neural networks. *CoRR*, abs/1308.0850, 2013.
- [2] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997.
- [3] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Netw.*, 2(5):359–366, July 1989.
- [4] Heikki Hyotyniemi. Turing machines are recurrent neural networks, 1996.
- [5] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Parallel distributed processing: Explorations in the microstructure of cognition, vol. 1. chapter Learning Internal Representations by Error Propagation, pages 318–362. MIT Press, Cambridge, MA, USA, 1986.
- [6] Franco Scarselli and Ah Chung Tsoi. Universal approximation using feed-forward neural networks: A survey of some existing methods, and some new results. *Neural Netw.*, 11(1):15–37, January 1998.
- [7] Hava T. Siegelmann and Eduardo D. Sontag. Turing computability with neural nets. *Applied Mathematics Letters*, 4:77–80, 1991.
- [8] Ronald J. Williams and Jing Peng. An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural Computation*, 2:490–501, 1990.