# MULTI-LAYER NETWORKS

## (FOR VECTORS)

by

FRANCESCO A. N. PALMIERI
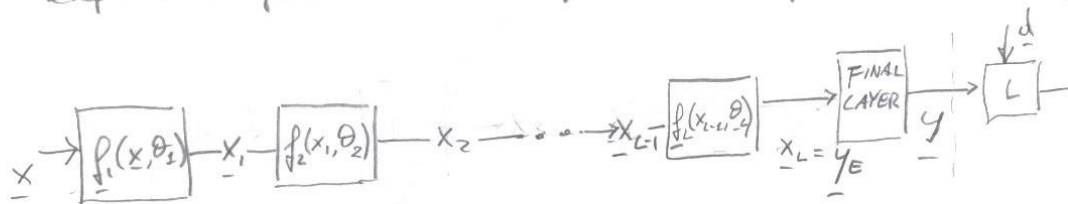
UNIVERSITA' DELLI STUDI
"LUIGI VANVITELLI"

Artificial

Neural networks have been originally inspired by biological networks where one of the main features is massive connectivity. Many paradigms that try to explain the signal processing capabilities of biological brains are based our progressive unfolding of the information coming from the various sensor modalities. Think of visual processing where, after a first transformation provided by the retina circuits, the visual cortex extracts relevant features such as lines bars and various patterns. This seems to be a strategy that we find both in the relatively simple brains of insects and in much more complex mammal neural architectures.

Therefore, even if we have shown in a previous chapter that two-layer network can already provide universal functional capabilities, there must be intrinsic advantages in using multiple layers. The reason for the means of multiple-layer architectures is not yet fully understood. Constrained connectivity may be one of reasons why we need multiple layers because of the need to provide sufficient mixing of the input components. Localized feature extraction may be another one.

Much experience has been gained in the last fourty years in researching the appropriate architecture for each problem, and almost invariably, multi-layer neural networks have

provided effective solutions. We should not forget that a working system has to provide sufficient generalisation and massive contained connection may provide after learning the necessary resolution.

In this chapter we consider general multi-layer architectures and provide explicit gradient computation for learning.



The figure above shows the cascade of L progressive layers that "prepare" the embedding $y_E$ for the final layer. Recall from the previous chapter, that the final layer can be

be of three different types : (1) linear for regression (2) linear + sigmoids for multiple binary classification; (3) linear + softmax for n-ary classification.
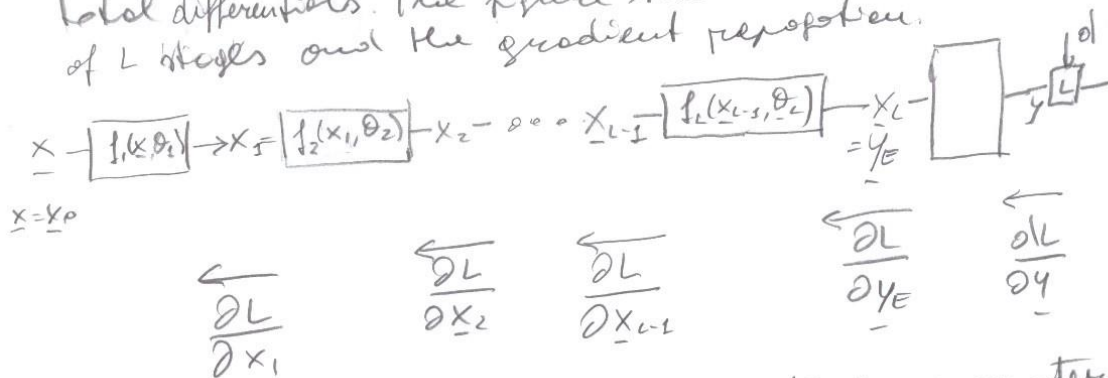
In all cases the multiple layers have to transform the input vector $\underline{x}$ into the embedding vector $\underline{y}_E$ in a way which is appropriate for the problem to be solved at the output $y$. Formally the transformation to the embedding space $\mathcal{Y}_\sigma$ is the nested function

$$y_E = f_L\left(f_{L-1}\left(f_{L-2}\left(\ldots f_2\left(f_1(\underline{x},\underline{\theta}_1),\underline{\theta}_2\right),\ldots \theta_{L-2}\right),\underline{\theta}_{L-1}\right),\underline{\theta}_L\right)$$

The input $\underline{x} = \underline{x}_0$ has size $N$, $\underline{x}_1 \to N_1$, $\underline{x}_2 \to N_2 \ldots X_L \to N_L = N_E$

The layer functions can be of different types and each $f_i$ may depend on a set of parameters $\theta_i$ to be learned. More generally some of the layer functions may be fixed with no free parameters as we will see in the following examples.

In all cases the parameters, to be learned, follow a gradient descent algorithm. The gradient is computed on the output loss function $L(y, d)$ must be "back propagated" following the rule of total differentials. The figure shows the cascade of $L$ stages and the gradient propagation.

$$x - \boxed{f_1(x, \theta_1)} \rightarrow x_1 - \boxed{f_2(x_1, \theta_2)} - x_2 - \cdots - x_{L-1} - \boxed{f_L(x_{L-1}, \theta_L)} \rightarrow x_L = y_E \quad \boxed{} \quad y \boxed{L} \overset{d}{}$$

$$x = x_0$$

$$\overleftarrow{\frac{\partial L}{\partial x_1}} \qquad \overleftarrow{\frac{\partial L}{\partial x_2}} \qquad \overleftarrow{\frac{\partial L}{\partial x_{L-1}}} \qquad \overleftarrow{\frac{\partial L}{\partial y_E}} \qquad \overleftarrow{\frac{\partial L}{\partial y}}$$

A parametric block $f_i(x_{i-1}, \theta_i)$ to update its parameters $\theta_i$ needs the gradient $\frac{\partial L}{\partial \theta_i}$, that can be computed from the back propagated gradient $\frac{\partial L}{\partial x_i}$ at its output $x_i$ as

$$\frac{\partial L}{\partial \theta_i} = \left(\frac{\partial x_i}{\partial \theta_i}\right)^T \frac{\partial L}{\partial x_i} \qquad i = 1, \ldots, L$$

Back propagation through layers consists in computing the Jacobian $\left(\frac{\partial x_{i+1}}{\partial x_i}\right)$ at each stage and propagate backward as follows:

$$\frac{\partial L}{\partial x_i} = \left(\frac{\partial x_{i+1}}{\partial x_i}\right)^T \frac{\partial L}{\partial x_{i+1}} \qquad i = L-1, L-2, \ldots, 1$$

Each layer may be of different type and also
be nonparametric, i.e. with no parameters to be
learned. In such a case we just backpropagate
the gradient through it.

In the following we examine different kinds
of blocks specifying the Jacobian for back propagation
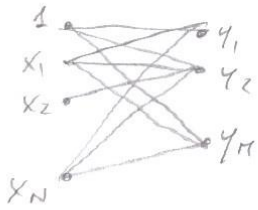and also the gradient for parametric learning

## LINEAR LAYER

We have already discussed linear transformations
both in the chapter on linear regression and
classifiers and in the two-layer architectures.
Linear layers (fully connected) can be inserted
at any stage in the multi-layer architecture.
They are usually followed by non linearities,
but we discuss them here separately as
a potential building block of our cascade.
Recall that by a linear transformation we mean
always an affine transformation, because we
include biases to avoid to have to worry about
means and baselines.

GENERIC LINEAR BLOCK
To facilitate notations consider first a generic
linear block $\underline{x} \rightarrow \underline{y}$

Forward propagation can be written as

$$y = \underline{\underline{W}}^T \underline{x} + \underline{b} = \begin{bmatrix} \underline{w}_1^T \\ \underline{w}_2^T \\ \vdots \\ \underline{w}_M^T \end{bmatrix} \underline{x} + \underline{b}$$

or in the augmented form

$$\underline{y} = \underbrace{\begin{bmatrix} \underline{b} & \begin{matrix} \underline{w}_1^T \\ \underline{w}_2^T \\ \vdots \\ \underline{w}_M^T \end{matrix} \end{bmatrix}}_{\underline{\underline{W}}_a^T} \underbrace{\begin{bmatrix} 1 \\ \underline{x} \end{bmatrix}}_{\underline{x}_a}$$

$$\underline{\underline{W}}_a = \begin{bmatrix} \underline{b}^T \\ \underline{w}_1 \ \underline{w}_2 \dots \underline{w}_M \end{bmatrix}$$

The Jacobian is defined and computed as

$$\frac{\partial \underline{y}}{\partial \underline{x}} \triangleq \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \dots & \frac{\partial y_1}{\partial x_N} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \dots & \frac{\partial y_2}{\partial x_N} \\ \frac{\partial y_M}{\partial x_1} & \frac{\partial y_M}{\partial x_2} & & \frac{\partial y_M}{\partial x_N} \end{pmatrix} = \begin{bmatrix} \underline{w}_1^T \\ \underline{w}_2^T \\ \vdots \\ \underline{w}_M^T \end{bmatrix} = \underline{\underline{W}}^T$$

Note that it does not depend on the biases.

Gradient propagation through the block follows the rule

$$\frac{\partial L}{\partial \underline{x}} = \left( \frac{\partial \underline{y}}{\partial \underline{x}} \right)^T \frac{\partial L}{\partial \underline{y}} = \underline{\underline{W}} \frac{\partial L}{\partial \underline{y}}$$

To train $\underline{\underline{W}}$ and $\underline{b}$ we need the gradients

$$\frac{\partial L}{\partial \underline{\underline{W}}} \quad \text{and} \quad \frac{\partial L}{\partial \underline{b}}$$

that can be computed from $\frac{\partial L}{\partial \underline{y}}$ as

$$\frac{\partial L}{\partial \underline{W}} = \left[\frac{\partial L}{\partial \underline{w_1}} \frac{\partial L}{\partial \underline{w_2}} \cdots \frac{\partial L}{\partial \underline{w_H}}\right] = \left[\frac{\partial L}{\partial y_1}\frac{\partial y_1}{\partial \underline{w_1}}, \frac{\partial L}{\partial y_2}\frac{\partial y_2}{\partial \underline{w_2}}, \cdots, \frac{\partial L}{\partial y_H}\frac{\partial y_H}{\partial \underline{w_H}}\right]$$

$$= \left[\frac{\partial L}{\partial y_1}\underline{x}, \frac{\partial L}{\partial y_2}\underline{x}, \cdots, \frac{\partial L}{\partial y_H}\underline{x}\right] = \left(\frac{\partial L}{\partial \underline{y}}\right)^T \otimes \underline{x}$$

$$\frac{\partial L}{\partial \underline{b}} = \begin{bmatrix} \frac{\partial L}{\partial b_1} \\ \frac{\partial L}{\partial b_2} \\ \frac{\partial L}{\partial b_H} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_1}\frac{\partial y_1}{\partial b_1} \\ \frac{\partial L}{\partial y_2}\frac{\partial y_2}{\partial b_2} \\ \vdots \\ \frac{\partial L}{\partial y_H}\frac{\partial y_H}{\partial b_H} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_1}\cdot 1 \\ \frac{\partial L}{\partial y_2}\cdot 1 \\ \vdots \\ \frac{\partial L}{\partial y_H}\cdot 1 \end{bmatrix} = \frac{\partial L}{\partial \underline{y}}$$

In the augmented representation

$$\frac{\partial L}{\partial \underline{\underline{W}}_a} = \left(\frac{\partial L}{\partial \underline{y}}\right)^T \otimes \underline{x}_a$$

The above results are immediately applied to the ith linear stage

The forward propagation is

$$\underline{x}_i = f_i(\underline{x}_{i-1}, \underline{\theta}_i) = \underline{\underline{W}}_i^T \underline{x}_{i-1} + \underline{b}_i \quad \text{or in the augmented representation}$$

$$\underline{x}_i = \underline{\underline{W}}_{ia}^T \underline{x}_{i-1a}$$

the Jacobian is $\frac{\partial \underline{x}_i}{\partial \underline{x}_{i-1}} = \underline{\underline{W}}_i^T$ and the gradient backpropagation

$$\boxed{\frac{\partial L}{\partial \underline{x}_{i-1}} = \underline{\underline{W}}_i \frac{\partial L}{\partial \underline{x}_i}}$$

The gradients for $\underline{\underline{W}}_i$ and $\underline{b}_i$ for learning are

$$\boxed{\frac{\partial L}{\partial \underline{\underline{W}}_i} = \left(\frac{\partial L}{\partial \underline{x}_i}\right)^T \otimes \underline{x}_{i-1}}, \quad \boxed{\frac{\partial L}{\partial \underline{b}} = \frac{\partial L}{\partial \underline{x}_i}}$$

In the augmented format

$$\boxed{\frac{\partial L}{\partial \underline{\underline{W}}_{ia}} = \left(\frac{\partial L}{\partial \underline{x}_i}\right)^T \otimes \underline{x}_{i-1a}}$$

$$x_i = f_i(x_{i-1}) = \varphi(x_{i-1})$$

$$N_i = N_{i-1}$$

There are no parameters to be learned, therefore we have to provide only the propagation rule for gradients.

The Jacobian is diagonal

$$\frac{\partial x_i}{\partial x_{i-1}} = \text{diag}\left(\varphi'(x_{i-1})\right)$$

and the gradient $\frac{\partial L}{\partial x_i}$ is backpropagated as

$$\frac{\partial L}{\partial x_{i-1}} = \left(\frac{\partial x_i}{\partial x_{i-1}}\right)^T \frac{\partial L}{\partial x_i} = \text{diag}\left(\varphi'(x_{i-1})\right)\frac{\partial L}{\partial x_i}$$

The non linearities could be sigmoids, RELUs or many others listed in Appendix X with their derivatives. To be specific consider as an example, logistics and RELUS.

For Logistics $\varphi(\xi) = \frac{1}{1+e^{-\xi}}$ , $\varphi'(\xi) = \varphi(\xi)(1-\varphi(\xi))$

the vector of derivatives can be written as

$$\varphi'(\xi) = \varphi(\xi) \odot (1 - \varphi(\xi)) \qquad 1 = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

Therefore

$$\frac{\partial L}{\partial x_{i-1}} = \text{diag}\left(x_i \odot (1 - x_i)\right)\frac{\partial L}{\partial x_i}$$

For RELUS $\varphi(\xi) = r(\xi)$ , $\varphi'(\xi) = u(\xi)$
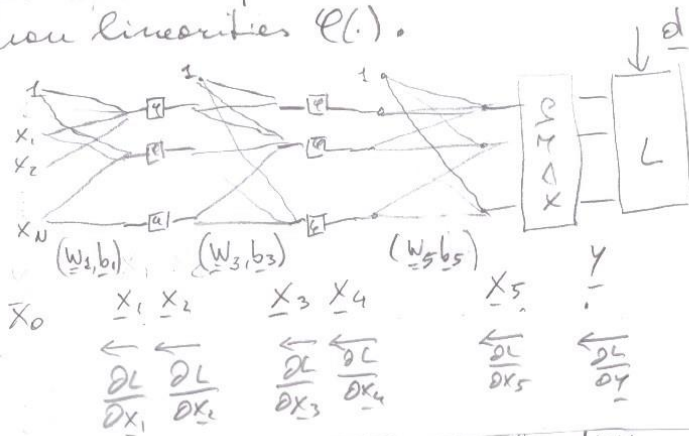
the vector of derivatives can be written as

$$\varphi'(\xi) = u(\xi)$$

therefore

$$\frac{\partial L}{\partial x_{i-1}} = \text{diag}\left(u(x_{i-1})\right)\frac{\partial L}{\partial x_i}$$

the RELUs in the gradient backward flow act as gates because they do not propagate the output gradient if the input is negative.

_EXAMPLE_  Consider the following multi-layer $ML.9$ architecture for classification. Assume logistic non-linearities $\varphi(\cdot)$.



The network implements the function
$$y = f(\underline{x}, \underline{\theta})$$ where $0 \leq y_i \leq 1$ $i=1,..,M$ represent the estimated posterior associated to the sample $\underline{x}$.
The parameters to be learned are $\{(\underline{W_1}, \underline{b_1}), (\underline{W_3}, \underline{b_3}), (\underline{W_5}, \underline{b_5})\}$
The forward flow is

$$\underline{x_0} = \underline{x} \rightarrow \underline{x_1} = \underline{W_1^T} \underline{x_0} + \underline{b_1} \rightarrow \underline{x_2} = \varphi(\underline{x_1}) \rightarrow \underline{x_3} = \underline{W_3^T}\underline{x_2} + \underline{b_3} \rightarrow \underline{x_4} = \varphi(\underline{x_3}) \rightarrow \underline{x_5} = \underline{W_5^T}\underline{x_4} + \underline{b_5}$$

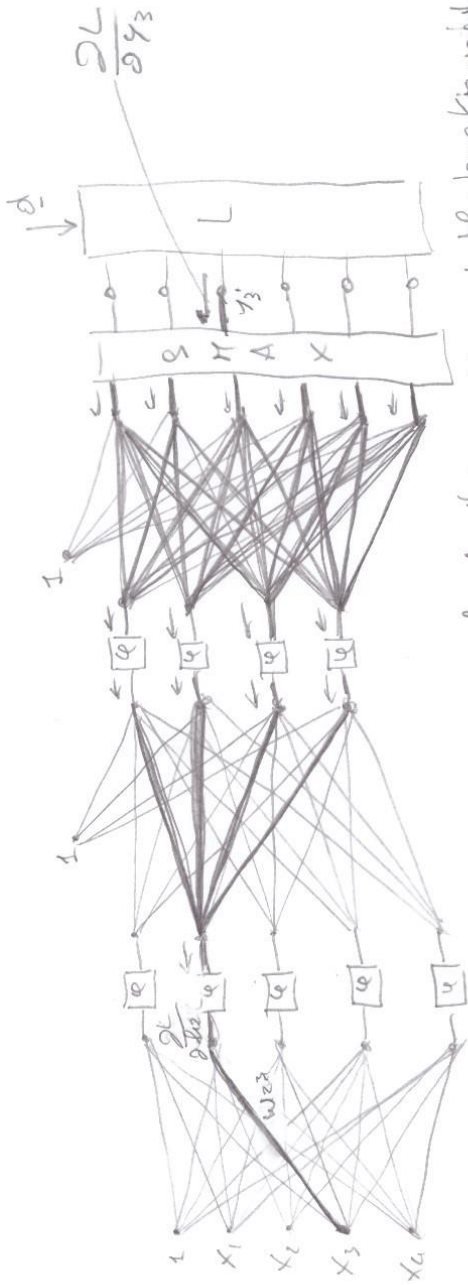$$\rightarrow \underline{y} = S_{max}(\underline{x_5})$$

Using the cross-entropy as the loss function, the gradients vectors backpropagated are:

$$\frac{\partial L}{\partial \underline{y}} = -\frac{\underline{d}}{\underline{y}} \rightarrow \frac{\partial L}{\partial \underline{x_5}} = \left(\text{diag}(\underline{y}) - \underline{y}\,\underline{y}^T\right)\frac{\partial L}{\partial \underline{y}} \rightarrow \frac{\partial L}{\partial \underline{x_4}} = \underline{W_5}\frac{\partial L}{\partial \underline{x_5}} \rightarrow$$

$$\frac{\partial L}{\partial \underline{x_3}} = \text{diag}\left(\underline{x_4} \odot (1-\underline{x_4})\right)\frac{\partial L}{\partial \underline{x_4}} \rightarrow \frac{\partial L}{\partial \underline{x_2}} = \underline{W_3}\frac{\partial L}{\partial \underline{x_3}} \rightarrow \frac{\partial L}{\partial \underline{x_1}} = \text{diag}\left(\underline{x_3} \odot (1-\underline{x_3})\right)\frac{\partial L}{\partial \underline{x_2}}$$

the updates on $(W_1, b_1), (\underline{W_3}, \underline{b_3}), (W_5, b_5)$ are based on the gradients

$$\frac{\partial L}{\partial \underline{W_1}} = \left(\frac{\partial L}{\partial \underline{x_1}}\right)^T \otimes \underline{x} \;,\; \frac{\partial L}{\partial \underline{b_1}} = \frac{\partial L}{\partial \underline{x_1}} \;,\; \frac{\partial L}{\partial \underline{W_3}} = \left(\frac{\partial L}{\partial \underline{x_3}}\right)^T \otimes \underline{x_2} \;;\; \frac{\partial L}{\partial \underline{b_3}} = \frac{\partial L}{\partial \underline{x_3}} \;;\; \frac{\partial L}{\partial \underline{W_5}} = \left(\frac{\partial L}{\partial \underline{x_5}}\right)^T \otimes \underline{x_4}$$

The figure shows an example in which it is evidenced the backpropagation flow for a specific weight $W_{23}$ in the first layer and gradient $\frac{\partial L}{\partial y_i}$ in the output layer. The gradient is sent to specific weight and nonlinearities to $W_{23}$. The network weight and nonlinearities to $W_{23}$ backward screen the effect of $\frac{\partial L}{\partial y_3}$ is translated in one update on $W_{23}$ proportional to

$$\frac{\partial L}{\partial \ell_2} \cdot X_3 \cdot$$

$$\frac{\partial L}{\partial y_3}$$

The following steps are the standard gradient updates for the backpropagation algorithm. The procedure is presented in the "batch" version. "Minibatch" and "stochastic" versions are easy modifications.

Our training set $\Sigma = \{(x[n], d[n]), n = 1, ..., n_\Sigma\}$ will be used many times for learning the system parameters. Each time the training set is used, we say that we perform an "EPOQUE". Typically the parameters obtained at the end of each epoque, are used as initial condition for the following one.

## INITIALIZATION

At the beginning of each training session, the network parameters are usually initialized to random values. Means and variances are usually determined heuristically and may depend on the problem at hand. Loss and gradients are set to zero

$$L_0 = 0 \qquad g_{\partial 0} = \underline{0}$$

## EPOQUE K

Keep from the previous epoque only the network parameters $\theta_{K-1}$; Set epoque loss and gradients to zero: $L[0] = 0$, $g_\partial[0] = \underline{0}$.

EXAMPLE $n = 1 : n_\Sigma$
- Get example $(x[n], d[n])$ from the training set;
- Propagate forward $x[n]$ computing all the activations in the network.
- Compute the instantaneous loss $L[n] = L(y[n], d[n])$ at the output
- Compute all instantaneous gradients $\frac{\partial L[n]}{\partial \theta[n]}$ for every parameter backpropagating $\frac{\partial L[n]}{\partial y[n]}$.
- Update the loss with a running mean

$$L[n] = \frac{n-1}{n} L[n-1] + \frac{1}{n} L[n]$$

- update the gradient for each parameter with a running mean

$$g_\theta[n] = \frac{n-1}{n} g_\theta[n-1] + \frac{1}{n} \frac{\partial L[n]}{\partial \theta[n]}$$
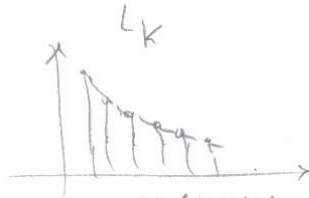
END OF EPOQUE K.

If the loss $L = L[nz]$ has reached a
satisfactory value $\rightarrow$ EXIT
Otherwise update the parameters

$$\theta_K = \theta_{K-1} - \mu \, g_\theta[nz]$$

and start a new epoque.

Note that the total loss after each epoque could
be computed separately using the parameters obtained
at the end of that epoque. We presented the
algorithm in a combined way to present the training
set to the network only once for each epoque. With the gradients
the only disadvantage is a useless computation of
the gradients in the last epoque.

To keep track of convergence we look at the sequence
$\{L_K\}$



Also global statistics about coefficients and performance
metrics may be tracked during learning.

We have discussed the method with reference to
a **batch** approach. The strategy is easily modified
for **minibatches** simply considering as epoques
(sub epoques) subset of the training set. Recall at the
end of each subset, we keep only the network
parameters for the new batch. The **stochastic** version
simply updates the parameters after each example.

The backpropagation algorithm being a gradient search
on a nonconvex cost function, will converge
inevitably to a local minimum. To reach
a global minimum is generally a rather hopeless
task because in theory global convergence can
be obtained starting from sufficiently many initial
conditions. Therefore we have to rely on solutions that
represent satisfactory local minimum.
Obviously the solutions obtained must be verified
on the validation set for generalization.
In practice a neural network, to be useful
in an application, has to go through various cycles
of learning perhaps adjusting parameters to obtain
a good compromise between performances on the training set and
generalization.
Generalization is a crucial feature, as pointed
out many times before, and can be thought as a
sort of smoothness: if the training set samples are
exactly matched the obtained function may be
overfitting the data because it interpolates the training
points too finely. A smoother function may do a better
job in containing the variations outside the training set.
Therefore the problem of obtaining a smooth solution
may be posed as a regularization problem

$$\theta^o = \operatorname*{argmin}_{\theta}\left(\sum_{u=1}^{u_2} L\left(y[u], d[u]\right) + \lambda R(\theta)\right)$$

where $R(\theta)$ is a function of the parameters and $\lambda$ is

the lagrange multiplier, that regulates the effect of the additional term. For example $R(\underline{\theta})$ may be $R(\underline{\theta}) = \sum_{i=1}^{M_\theta} \theta_i^2$ (Tikhonov regularization).

The constraint usually lead to a decaying term that encourages smaller weights. Smaller weights that contribute the input of a sigmoidal function should give a smoother function. Regulating $\lambda$ to larger values favors smoother at the expenses of accuracy in the cost function. Vice-versa small $\lambda$ affect less the cost function encouraging smoother.

There are other regularization strategies, also on the gradient, to encourage smoother.

The Most popular strategies to encourage generalization are based on simple heuristics:

Heuristic 1 : EARLY STOPPING

This strategy suggests to stop training before reaching convergence. This provides solution that are less accurate, but that perhaps are smoother

Heuristic 2 : ENSEMBLING

We train different networks, using different initial conditions, and combine the prediction/results. In regression we could use the mean or the median. Max rule could be used for classification (the classification that gets the largest score). This is sometimes called the committee machine.

Another strategy could be to train the system on different subsets of the training set obtaining different networks. The predictions from the networks are then combined (averaged or maxed out)

This strategy in statistics is known as
BOOTSTRAP AGGREGATION or BAGGING.
The solutions are expected to be smoother.

### Heuristic 3 : DROPOUT

In this strategy during training some
hidden units are randomly set to zero at
each gradient backpropagation. The number
of censored units is usually 50% or less.

In this way the gradient estimates become
noisy and may make the system more
robust moving the evolution during training
in a more scattered pattern that may lead to
better minima and better generalization.
Dropout has allowed a major leap in
performance for convolutional networks in
solving image recognition tasks.

### Heuristic 4 : ADD NOISE

Dropout discussed above can be seen as applying
bernoulli noise to the hidden units during
training. The idea can be extended by applying
noise during gradient calculations. In regression
it can be shown to be equivalent to regularization.
Noise can be also added to the weights

### CONVERGENCE SPEED

The gradient updates are controlled by the stepsize
parameter $\mu[k]$ that can be made adaptive.

Clearly larger $\mu$ correspond to faster updates that may overshoot the valleys of the cost function. Vice versa small $\mu$ may be slow and more precise in catching the minimum, but they may be bad (high cost) local minimum.

Generally at the beginning of training $\mu[k]$ is kept large with progressive decay around convergence. A good compromise that tracks the evolution of the gradient is ADAM widely used and already reminded in one of our previous chapter. Also ADAM has some free parameters that are determined by trial-and-error.