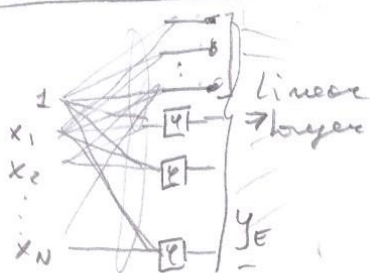


## OTHER TWO-LAYER ARCHITECTURES

The examples seen in the previous sections are only some of the possible architectures that can be built. The representation theorem suggests that there exist many other options to build the first layer. We will try to discuss a few of them in the following sections.

### SKIP-CONNECTION ARCHITECTURE



$$y_E = \begin{bmatrix} c_1^T x + b_1 \\ c_2^T x + b_2 \end{bmatrix} = \begin{bmatrix} y_{E1} \\ y_{E2} \end{bmatrix}$$

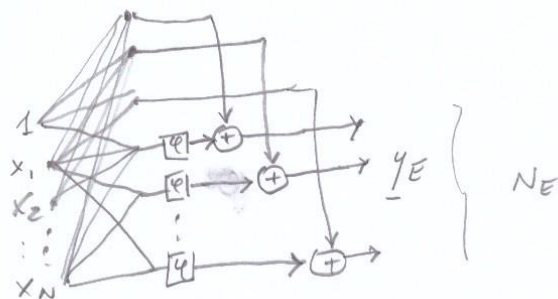
In this architecture the embedding space is augmented with skip-connections, i.e. with a strictly linear part. The activation functions can be sigmoids, RELUs, or any other function.

The idea here is to provide also a strictly linear part into the embedding space. Training of this architecture will ensure that classifiers and regressors are at least linear.

The nonlinear part would be used by the output layer if necessary. Note that the linear and the nonlinear parts can have different sizes.

## RESIDUAL ARCHITECTURE

The residual architecture also builds on the idea of generalizing the linear layer. In this case the output of the linear part is added to the output of the nonlinear part.



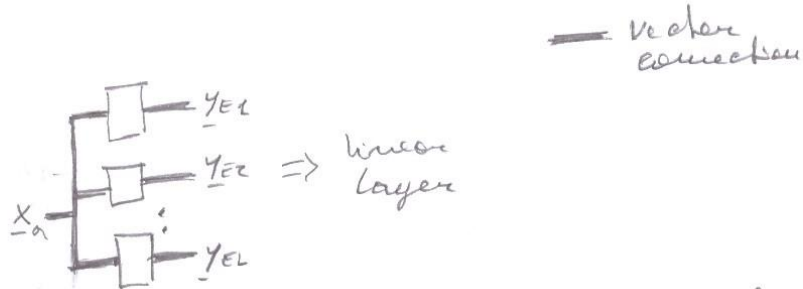
$$\underline{y}_E = (\underline{C}_1 \underline{x} + b_1) + \varphi(\underline{C}_2 \underline{x} + b_2)$$

Clearly here the sizes of the linear and nonlinear parts have to be the same ( $N_E$ ). The terminology "residual" comes from the fact that  $\underline{y}_E$  may be considered the difference between a linear (and/or nonlinear) hypothesis. (differences or miss can be corrected) Clearly this architecture is to be considered a bit less powerful than the skip connection parallel structure. Usually residual layers are used mostly on multi-layer architectures. A word of caution should be said about using structures with more less constrained parameters: knowing how to determine these

parameters and generalization should be checked. Therefore any choice of neural network structure should be done with careful attention to overfitting. The real challenge is therefore in finding the proper parametrization for the problem to be solved.

## PARALLEL ARCHITECTURES

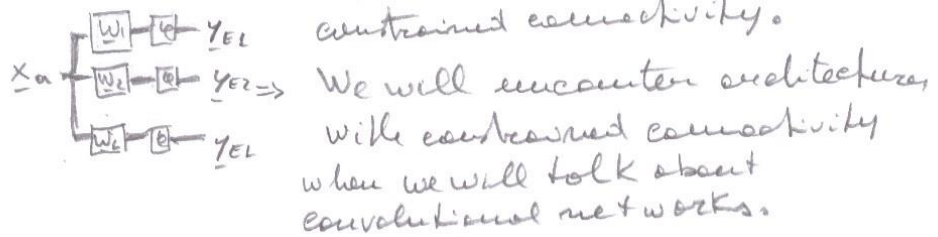
The SKIP-CONNECTION architecture is a special case of parallel structures



Each element of the parallel must be linear or nonlinear.  
The embedding space  $\underline{y}_E = \begin{bmatrix} y_{E1} \\ y_{E2} \\ \vdots \\ y_{EL} \end{bmatrix}$  is composed by the

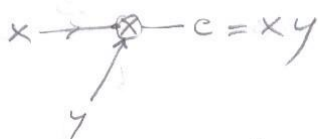
concatenation of the various vector-outputs that may specialize on different features.

When all the parallel structures are of the same type, for example linear + activation functions, the network is equivalent to a single network with constrained connectivity.



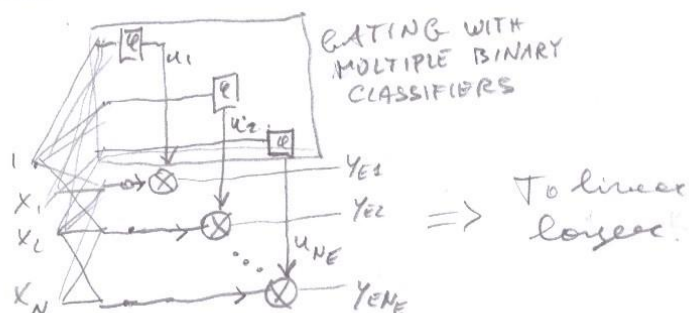
RDM-PRODUCT ARCHITECTURES

Until now we have considered only neural network architectures that use linear combinations and real or nonlinearities. Larger functional capabilities can be obtained if we include also product units, i.e. nodes where the inputs are multiplied.



It may be argued that multiplication could be implemented miming the logs and then exponentiating. However there are some architectures where the product nodes may be very useful and act as "gates".

It should be pointed out that these structures are not the non-product architectures used to model random variables distributions. They simply indicate the inclusion of product nodes.

GATED ARCHITECTURES

The embedding vector is obtained as

$$\underline{y}_E = f(\underline{C}^T \underline{x} + \underline{d}) \odot (\underline{C}^T \underline{x} + \underline{b})$$

$$\underline{C}_i = [\underline{g}_1, \underline{g}_2, \dots, \underline{g}_{N_i}]$$

$$\underline{C} = [\underline{C}_1, \underline{C}_2, \dots, \underline{C}_{N_E}]$$

Here each binary classifier "activates" a linear output. If one classifier works in situation  $x \geq 1$  or  $x \leq 0$ , it turns on and off the connected linear filter at the bottom. If instead the specific classifier works in its non-linear range, it implements a quadratic transformation

More specifically if the  $j$ th classifier gives (at a specific  $x$ )  $u_j \approx 1$ ,  $y_{E_j} \approx \underline{e}_j^T x + b_j$ , otherwise if  $u_j \approx 0$   $y_{E_j} \approx 0$ .

Vice versa if the  $j$ th classifier works in its linear range, we have  $u_j \approx \underline{g}_j^T x + d_j$

and

$$y_{E_j} \approx (\underline{g}_j^T x + d_j)(\underline{e}_j^T x + b_j)$$

$$= \underline{g}_j^T x \underline{e}_j^T x + d_j \underline{e}_j^T x + b_j \underline{g}_j^T x + d_j b_j$$

which is a quadratic function of  $x$ .

Once we have set  $N_E$ , the learning algorithm will provide the parameters  $\underline{g}_j, \underline{d}_j, \underline{e}_j, b_j$  that will characterize the embedding space needed for the final recognizer or classifier. (\*)

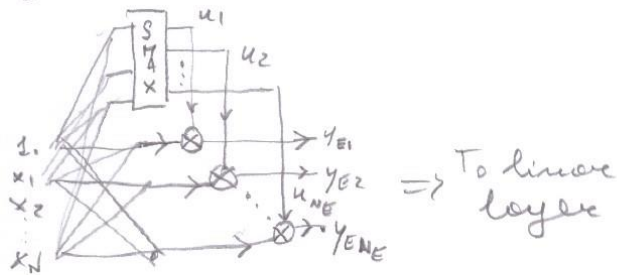
This is a very powerful architecture for its capabilities of handling linear and non-linear embeddings.

This gating unit is utilized in the LSTM (LONG-SHORT TIME MEMORY) architecture, which is a very powerful recurrent network. The LSTM will be discussed in one of the following chapters.

(\*) Recall the gaussian classifier in the model-based chapter of this book that requires quadratic log-likelihoods. Therefore this architecture may be suitable to build a classifier with clusters that have different covariance matrices. ... (?)

Wu

Another gated architecture is the following one TWO.30



$$y_E = \sum_{\max} (\underline{c}^T \underline{x} + d) \odot (\underline{c}^T \underline{x} + b)$$

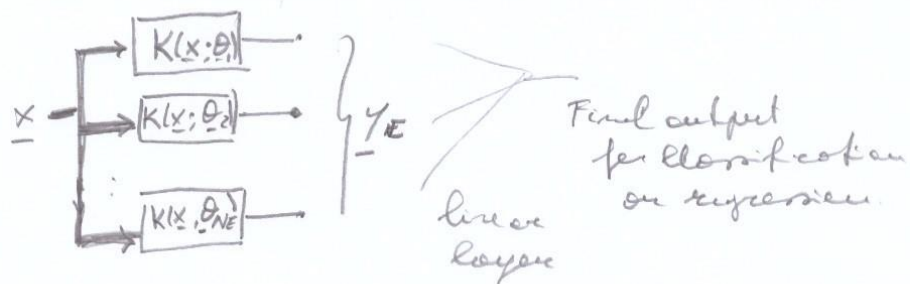
Here the gating is done by an  $N_E$ -class linear classifier. The difference with the previous architecture is that here the target is one class out of  $N_E$  (Recall that  $\sum_{i=1}^{N_E} u_i = 1$ ). Therefore if  $\underline{x}$  is such that  $u_j \approx 1$  and the others are  $\approx 0$ ,  $y_{Ej} \approx \underline{c}_j^T \underline{x} + b_j$ . Essentially in every decision region of the linear classifier, one output of  $y_E$  will be active and is a linear function of  $\underline{x}$ .

Note that both in this and the previous gated architectures,  $N_E$  can be made very large giving to the system the opportunity of building many useful basis functions  $\{y_{Ei}(\underline{x}), i=1, \dots, N_E\}$ .

## KERNEL-BASED ARCHITECTURES

The representation theorem focuses on finding the proper embedding space with basis functions that can be linearly combined to produce the system output.

One of the early ideas about building such space, is to postulate a number of parametric "kernels". The following figure shows such an architecture



Each embedding function is a kernel function

$$y_{Ei} = K(\underline{x}; \underline{\theta}_i)$$

that depends on a number of parameters  $\underline{\theta}_i$ .  
A scalar output <sup>(single regression)</sup> would be the superposition

$$y(\underline{x}) = \sum_{i=1}^{N_E} c_i K(\underline{x}, \underline{\theta}_i) + b$$

or <sup>multiple</sup> multiple outputs for classification

### RADIAL-BASIS FUNCTIONS

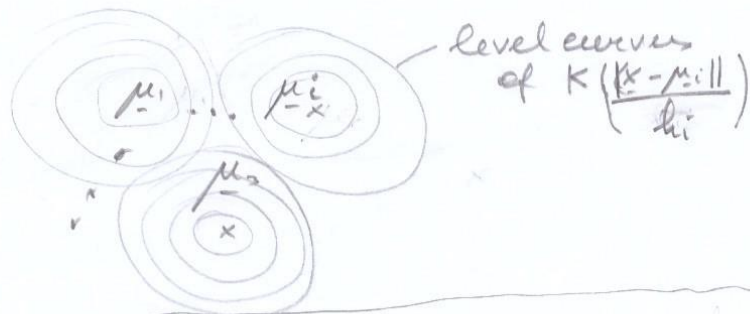
The choice of the kernel functions  $K(\underline{x}, \underline{\theta}_i)$  is done essentially to "cover" the input space, with functions that are usually chosen to have



a radial symmetry

$$K(\underline{x}; \underline{\theta}_i) = K\left(\frac{\|\underline{x} - \underline{\mu}_i\|}{h_i}\right)$$

where  $\|\underline{x} - \underline{\mu}_i\|$  is the square of the euclidean distance between  $\underline{x}$  and a "cluster point"  $\underline{\mu}_i$ , and  $h_i$  controls the spread.



NOTE:

The idea of using radial functions has come from the traditional problem of density estimation in statistics. <sup>(Parzen density estimator)</sup> More specifically if we are given a set of samples  $\underline{x}_1, \underline{x}_2, \dots, \underline{x}_L$  we could build density function of  $\underline{x}$  by placing a function around each point  $\underline{x}_i$  and summing up.

$$\hat{p}(\underline{x}) = \sum_{i=1}^L \alpha_i K\left(\frac{\|\underline{x} - \underline{x}_i\|}{h_i}\right)$$

To be a valid pdf we have to impose

$$\int_{\mathcal{X}} \hat{p}(\underline{x}) d\underline{x} = 1 \Rightarrow \sum_{i=1}^L \alpha_i \int_{\mathcal{X}} K\left(\frac{\|\underline{x} - \underline{x}_i\|}{h_i}\right) d\underline{x} = 1$$

Kernel functions can be normalized and the coefficients,  $\alpha_i$  regulate the height of each function.

The most common choice for  $K(\cdot)$  is the Gaussian function

$$K\left(\frac{\|x-\mu\|^2}{h^2}\right) = e^{-\frac{\|x-\mu\|^2}{2h^2}}$$

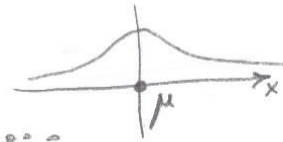
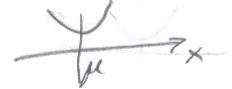
even if other choices are possible, such as the multiquadrate

$$K(\|x-\mu\|^2, \beta) = \sqrt{\|x-\mu\|^2 + \beta^2}$$

or the inverse multiquadrate

$$K(\|x-\mu\|^2, \beta) = \frac{1}{\sqrt{\|x-\mu\|^2 + \beta^2}}$$

$\mu$  is the center and  $\beta$  controls the spread



There are other radial functions...

In any case more generally the  $\mu_i$  are our cluster points and  $h_i$  controls the size of each the spherical region around  $\mu_i$ .

The superposition of these spherical functions can build the necessary likelihoods or log-likelihoods for regression or classification.

Usually the determination of the cluster points  $\mu_i$  and the size  $h_i$  is the result of our "unsupervised" phase. We have not talked about unsupervised learning yet as the topic will be more specifically addressed in one of the following chapters.

Essentially, we can assume that the clusters will be determined only from  $\{x^{(u)}, u=1, \dots, u_2\}$  in the training set, with an algorithm such as K-means or similar. (we will address this issue in the following)  
 The supervised learning phase determines then the best linear combination for classification or regression.

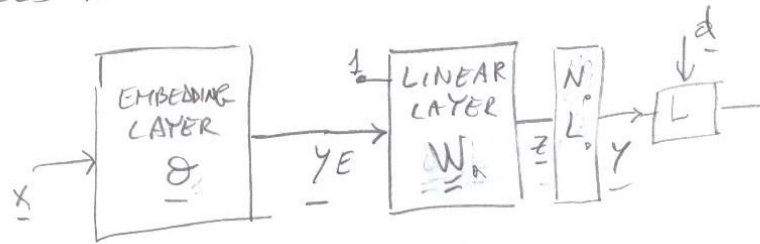
Usually, the size of the embedding space  $N_E$  is kept much smaller than the number of examples and the size of the clusters  $h_i$  is chosen to make the overall function sufficiently smooth. Sometimes  $h_i = \frac{d_{max}}{\sqrt{N_E}}$  ← max distance between clusters  
 ↑ # of clusters.  
 --- more ---  
 --- example ---

Note that the clusters can be built to have different elliptical shapes. Obviously there would be many more parameters to be learned.

~~(scribble)~~

## LEARNING IN TWO-LAYER ARCHITECTURES

The general two-layer architecture of Figure (x), repeated here for convenience, contains free parameters to be learned in both layers



The N.L. block is not present in regression and  
 It is made up of  
 It is made up of  $\sigma$  neurons in the multiple binary classifier and a softmax in the multi-class classifier.

The cost function to be minimized is

$$\mathcal{E}(\underline{\theta}, \underline{W}_n) = \frac{1}{n} \sum_{u=1}^n L(\underline{y}^{[u]}, \underline{d}^{[u]})$$

As explained in one of the previous chapters we need to compute the gradients

$$\nabla_{\underline{W}_n} \mathcal{E}(\underline{\theta}, \underline{W}_n) \text{ and } \nabla_{\underline{\theta}} \mathcal{E}(\underline{\theta}, \underline{W}_n)$$

for making the updates

$$\underline{W}_n^{[k]} = \underline{W}_n^{[k-1]} - \mu \nabla_{\underline{W}_n} \mathcal{E}(\underline{\theta}^{[k-1]}, \underline{W}_n^{[k-1]})$$

$$\underline{\theta}^{[k]} = \underline{\theta}^{[k-1]} - \mu \nabla_{\underline{\theta}} \mathcal{E}(\underline{\theta}^{[k-1]}, \underline{W}_n^{[k-1]})$$

To simplify notation let us drop the example index  $n$  and concentrate on forward propagation and gradients on a single example

## THE OUTPUT LAYER

TWO.36

Recall the structure of the linear layer

$$\underline{z} = \underline{W}_a^T \underline{y}_E = \begin{bmatrix} b_1 & \underline{w}_1^T \\ b_2 & \underline{w}_2^T \\ \vdots & \vdots \\ b_M & \underline{w}_M^T \end{bmatrix} \begin{bmatrix} 1 \\ \underline{y}_E \end{bmatrix} = \begin{bmatrix} b \\ \underline{w}^T \end{bmatrix} \begin{bmatrix} 1 \\ \underline{y}_E \end{bmatrix}$$

with the augmented notation

$$\underline{W}_a = (\underline{w}_1^a \ \underline{w}_2^a \ \dots \ \underline{w}_M^a) = \begin{bmatrix} b_1 & b_2 & \dots & b_M \\ \underline{w}_1 & \underline{w}_2 & \dots & \underline{w}_M \end{bmatrix} = \begin{bmatrix} \underline{b}^T \\ \underline{W} \end{bmatrix}$$

that contains in the first row the biases.

Alternatively to the augmented notation we could explicitly sum the biases to the output

$$\underline{z} = \underline{W}^T \underline{y}_E + \underline{b}$$

This is totally equivalent to the augmented notation but it will allow easier analysis of the gradient propagation.

Now the loss  $L(\underline{y}, \underline{d})$  has to be differentiated with respect to  $\underline{W}$ ,  $\underline{b}$ , and  $\underline{\theta}$ .

For the second layer we need the gradients

$$\frac{\partial L}{\partial \underline{w}_j} = \frac{\partial L}{\partial z_j} \frac{\partial z_j}{\partial \underline{w}_j} \quad ; \quad \frac{\partial L}{\partial b_j} = \frac{\partial L}{\partial z_j} \frac{\partial z_j}{\partial b_j} \quad j=1, \dots, M$$

since  $\frac{\partial z_j}{\partial b_j} = 1$  for the bias  $\frac{\partial L}{\partial \underline{b}} = \frac{\partial L}{\partial \underline{z}}$ ,

Similarly, since  $\frac{\partial z_j}{\partial \underline{w}_j} = \underline{y}_E$ ,  $\frac{\partial L}{\partial \underline{W}_j} = \frac{\partial L}{\partial z_j} \underline{y}_E$ .

Now we need the gradient of  $L$  w.r.t.  $\underline{z}$ . For a generic  $z_j$ , using total gradient, we have

$$\frac{\partial L}{\partial z_j} = \sum_{i=1}^M \frac{\partial L}{\partial y_i} \frac{\partial y_i}{\partial z_j} = \left( \frac{\partial L}{\partial y_1} \ \frac{\partial L}{\partial y_2} \ \dots \ \frac{\partial L}{\partial y_M} \right) \begin{pmatrix} \frac{\partial y_1}{\partial z_j} \\ \frac{\partial y_2}{\partial z_j} \\ \vdots \\ \frac{\partial y_M}{\partial z_j} \end{pmatrix} = \left( \frac{\partial L}{\partial \underline{y}} \right)^T \frac{\partial \underline{y}}{\partial z_j}$$

In a comprehensive matrix form (gradient matrix flow) for  $\underline{w}$  we have

$$\frac{\partial L}{\partial \underline{w}} = \begin{pmatrix} \frac{\partial L}{\partial w_1} & \frac{\partial L}{\partial w_2} & \dots & \frac{\partial L}{\partial w_H} \end{pmatrix} = \left[ \left( \frac{\partial L}{\partial y} \right)^T \frac{\partial y}{\partial z_1} \quad y_E \left( \frac{\partial L}{\partial y} \right)^T \frac{\partial y}{\partial z_2} \quad y_E \right] \cdot \left( \frac{\partial L}{\partial y} \right)^T \frac{\partial y}{\partial z_H} y_E$$

$$= \left( \frac{\partial L}{\partial y} \right)^T \left( \frac{\partial y}{\partial z} \right) \otimes y_E$$

where  $\frac{\partial y}{\partial z}$  is the Jacobian matrix

$$\left( \frac{\partial y}{\partial z} \right) = \begin{bmatrix} \left( \frac{\partial y_1}{\partial z} \right)^T \\ \left( \frac{\partial y_2}{\partial z} \right)^T \\ \left( \frac{\partial y_H}{\partial z} \right)^T \end{bmatrix} = \begin{pmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_1}{\partial z_2} & \dots & \frac{\partial y_1}{\partial z_H} \\ \frac{\partial y_2}{\partial z_1} & \frac{\partial y_2}{\partial z_2} & & \frac{\partial y_2}{\partial z_H} \\ \frac{\partial y_H}{\partial z_1} & \frac{\partial y_H}{\partial z_2} & & \frac{\partial y_H}{\partial z_H} \end{pmatrix}$$

Therefore we see

$$\frac{\partial L}{\partial \underline{b}} = \frac{\partial L}{\partial z} = \begin{pmatrix} \frac{\partial L}{\partial z_1} \\ \frac{\partial L}{\partial z_2} \\ \frac{\partial L}{\partial z_H} \end{pmatrix} = \begin{pmatrix} \frac{\partial y^T}{\partial z_1} \frac{\partial L}{\partial y} \\ \frac{\partial y^T}{\partial z_2} \frac{\partial L}{\partial y} \\ \frac{\partial y^T}{\partial z_H} \frac{\partial L}{\partial y} \end{pmatrix} = \begin{pmatrix} \frac{\partial y_1}{\partial z_1} & \frac{\partial y_2}{\partial z_1} & \dots & \frac{\partial y_H}{\partial z_1} \\ \frac{\partial y_1}{\partial z_2} & \frac{\partial y_2}{\partial z_2} & & \frac{\partial y_H}{\partial z_2} \\ \frac{\partial y_1}{\partial z_H} & \frac{\partial y_2}{\partial z_H} & & \frac{\partial y_H}{\partial z_H} \end{pmatrix} \begin{pmatrix} \frac{\partial L}{\partial y_1} \\ \frac{\partial L}{\partial y_2} \\ \dots \\ \frac{\partial L}{\partial y_H} \end{pmatrix}$$

$$= \left( \frac{\partial y}{\partial z} \right)^T \left( \frac{\partial L}{\partial y} \right)$$

For the first layer, we have computed all the  $N_b$  free parameters in a vector  $\underline{\theta}$  and we need the <sup>two</sup> <sub>38</sub> gradient vectors

$$\frac{\partial L}{\partial \underline{\theta}} = \begin{bmatrix} \frac{\partial L}{\partial \theta_1} \\ \frac{\partial L}{\partial \theta_2} \\ \frac{\partial L}{\partial \theta_{N_b}} \end{bmatrix} = \begin{bmatrix} \left(\frac{\partial L}{\partial y_E}\right)^T \frac{\partial y_E}{\partial \theta_1} \\ \left(\frac{\partial L}{\partial y_E}\right)^T \frac{\partial y_E}{\partial \theta_2} \\ \vdots \\ \left(\frac{\partial L}{\partial y_E}\right)^T \frac{\partial y_E}{\partial \theta_{N_b}} \end{bmatrix}$$

$$\left(\frac{\partial L}{\partial y_E}\right)^T = \left(\frac{\partial L}{\partial y_{E1}} \quad \frac{\partial L}{\partial y_{E2}} \quad \dots \quad \frac{\partial L}{\partial y_{EN_E}}\right)$$

$$= \left[ \sum_{i=1}^M \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial y_{E1}}, \sum_{i=1}^M \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial y_{E2}}, \dots, \sum_{i=1}^M \frac{\partial L}{\partial z_i} \frac{\partial z_i}{\partial y_{EN_E}} \right]$$

$$= \begin{pmatrix} \frac{\partial L}{\partial z_1} & \frac{\partial L}{\partial z_2} & \dots & \frac{\partial L}{\partial z_M} \end{pmatrix} \begin{pmatrix} \frac{\partial z_1}{\partial y_{E1}} & \frac{\partial z_1}{\partial y_{E2}} & \dots & \frac{\partial z_1}{\partial y_{EN_E}} \\ \frac{\partial z_2}{\partial y_{E1}} & \frac{\partial z_2}{\partial y_{E2}} & \dots & \frac{\partial z_2}{\partial y_{EN_E}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_M}{\partial y_{E1}} & \frac{\partial z_M}{\partial y_{E2}} & \dots & \frac{\partial z_M}{\partial y_{EN_E}} \end{pmatrix}$$

$$= \left[ \left(\frac{\partial L}{\partial y}\right)^T \frac{\partial y}{\partial z_1}, \left(\frac{\partial L}{\partial y}\right)^T \frac{\partial y}{\partial z_2}, \dots, \left(\frac{\partial L}{\partial y}\right)^T \frac{\partial y}{\partial z_M} \right] \begin{pmatrix} \frac{\partial z_1}{\partial y_{E1}} & \frac{\partial z_1}{\partial y_{E2}} & \dots & \frac{\partial z_1}{\partial y_{EN_E}} \\ \frac{\partial z_2}{\partial y_{E1}} & \frac{\partial z_2}{\partial y_{E2}} & \dots & \frac{\partial z_2}{\partial y_{EN_E}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial z_M}{\partial y_{E1}} & \frac{\partial z_M}{\partial y_{E2}} & \dots & \frac{\partial z_M}{\partial y_{EN_E}} \end{pmatrix} \quad ||$$

$$= \begin{bmatrix} \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_1}{\partial z_1} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_2}{\partial z_1} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_4}{\partial z_1} \end{bmatrix} \begin{bmatrix} \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_1}{\partial z_2} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_2}{\partial z_2} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_4}{\partial z_2} \end{bmatrix} \begin{bmatrix} \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_1}{\partial z_H} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_2}{\partial z_H} \\ \left(\frac{\partial L}{\partial Y}\right)^T \frac{\partial y_4}{\partial z_H} \end{bmatrix} \begin{bmatrix} \frac{\partial z_1}{\partial y_{E1}} & \frac{\partial z_1}{\partial y_{E2}} & \frac{\partial z_1}{\partial y_{ENE}} \\ \frac{\partial z_2}{\partial y_{E1}} & \frac{\partial z_2}{\partial y_{E2}} & \frac{\partial z_2}{\partial y_{ENE}} \\ \frac{\partial z_H}{\partial y_{E1}} & \frac{\partial z_H}{\partial y_{E2}} & \frac{\partial z_H}{\partial y_{ENE}} \end{bmatrix}$$

$$= \left(\frac{\partial L}{\partial Y}\right)^T \left(\frac{\partial Y}{\partial Z}\right) \left(\frac{\partial Z}{\partial Y_E}\right) \quad (\text{row vector})$$

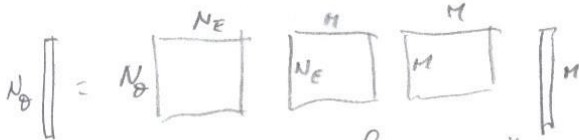
Defining the Jacobian

$$\frac{\partial Y_E}{\partial \theta} = \begin{bmatrix} \frac{\partial y_{E1}}{\partial \theta_1} & \frac{\partial y_{E1}}{\partial \theta_2} & \frac{\partial y_{E1}}{\partial \theta_{NB}} \\ \frac{\partial y_{E2}}{\partial \theta_1} & \frac{\partial y_{E2}}{\partial \theta_2} & \frac{\partial y_{E2}}{\partial \theta_{NB}} \\ \frac{\partial y_{ENE}}{\partial \theta_1} & \frac{\partial y_{ENE}}{\partial \theta_2} & \frac{\partial y_{ENE}}{\partial \theta_{NB}} \end{bmatrix}$$

We have

$$\frac{\partial L}{\partial \theta} = \left[ \left(\frac{\partial L}{\partial Y}\right)^T \left(\frac{\partial Y}{\partial Z}\right) \left(\frac{\partial Z}{\partial Y_E}\right) \left(\frac{\partial Y_E}{\partial \theta}\right) \right]^T$$

$$= \left(\frac{\partial Y_E}{\partial \theta}\right)^T \left(\frac{\partial Z}{\partial Y_E}\right)^T \left(\frac{\partial Y}{\partial Z}\right)^T \frac{\partial L}{\partial Y}$$



Since the linearity of the next layer

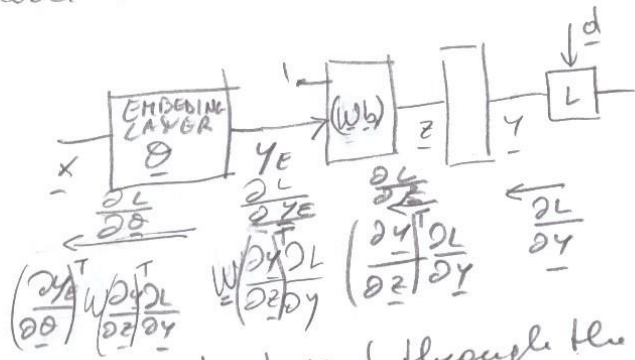
$$\left(\frac{\partial Z}{\partial Y_E}\right)^T = \begin{bmatrix} \frac{\partial z_1}{\partial y_{E1}} & \frac{\partial z_1}{\partial y_{E2}} & \frac{\partial z_1}{\partial y_{ENE}} \\ \frac{\partial z_2}{\partial y_{E1}} & \frac{\partial z_2}{\partial y_{E2}} & \frac{\partial z_2}{\partial y_{ENE}} \\ \frac{\partial z_H}{\partial y_{E1}} & \frac{\partial z_H}{\partial y_{E2}} & \frac{\partial z_H}{\partial y_{ENE}} \end{bmatrix} = \begin{bmatrix} w_1^+ \\ w_2^+ \\ w_H^+ \end{bmatrix} = W$$



Therefore

$$\frac{\partial L}{\partial \theta} = \left( \frac{\partial y^E}{\partial \theta} \right)^T W \left( \frac{\partial y}{\partial z} \right)^T \frac{\partial L}{\partial y}$$

Note how the backward flow evolves over the network architecture



Note that going backward through the linear block the bias part carries no gradient.

This is the essence of the backpropagation algorithm, that will be extended to multi-layer architecture, that essentially consists in propagating in the backward direction the gradient flow.

BACKWARD FLOW ON THE OUTPUT LAYER  
 Now we examine <sup>in more detail the 2-class</sup> the 3 types of problems: regression, multiple binary classification and multi-class classification.

Regression In regression the loss function is

$$L(y, d) = \sum_{i=1}^n \psi(d_i - y_i) = \sum_{i=1}^n \psi(\epsilon_i)$$

and  $y = z$ . Therefore

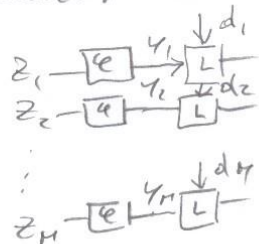
$$\frac{\partial L}{\partial y} = \begin{bmatrix} \frac{\partial \psi(\epsilon_1)}{\partial y_1} \\ \vdots \\ \frac{\partial \psi(\epsilon_n)}{\partial y_n} \end{bmatrix} = - \begin{bmatrix} \psi'(\epsilon_1) \\ \psi'(\epsilon_2) \\ \vdots \\ \psi'(\epsilon_n) \end{bmatrix} = -\psi'(\underline{\epsilon}) \quad \text{and} \quad \frac{\partial y}{\partial z} = \mathbf{I}_n$$

The derivatives  <sup>$\psi'$</sup>  of the loss functions have been discussed in previous chapters. For example for the quadratic loss  $\psi(\epsilon_i) = \epsilon_i^2$  we have simply

$$\frac{\partial L}{\partial \epsilon} = -2\epsilon$$

### MULTIPLE BINARY CLASSIFIER

In the multiple binary classifier we have



$$\left. \begin{array}{l} y_i = \psi(z_i), i=1, \dots, M \\ 0 \leq d_i \leq 1 \end{array} \right\}$$

The total loss is the sum of the losses on each output

$$L(y, d) = \sum_{i=1}^M L(y_i, d_i)$$

We could use standard loss functions  <sup>$\psi(\cdot)$</sup>  just like in regression

$$L(y, d) = \sum_{i=1}^M \psi(d_i - y_i) = \sum_{i=1}^M \psi(\epsilon_i)$$

so that

$$\frac{\partial L}{\partial y} = -\psi'(\epsilon)$$

that for the quadratic loss  $\Rightarrow \frac{\partial L}{\partial y} = -2\epsilon$ .  
More appropriately, using the binary crossentropy

$$L(y, d) = \sum_{i=1}^M H(d_i, y_i) = \sum_{i=1}^M (d_i \log \frac{1}{y_i} + (1-d_i) \log \frac{1}{1-y_i})$$

already discussed in the chapters on linear binary classifiers. The gradient is easily computed (see Appendix)

$$\frac{\partial L}{\partial \mathbf{y}} = \begin{bmatrix} \frac{y_1 - d_1}{y_1(1-y_1)} \\ \frac{y_2 - d_2}{y_2(1-y_2)} \\ \vdots \\ \frac{y_n - d_n}{y_n(1-y_n)} \end{bmatrix}$$

The gradient <sup>matrix</sup> across the nonlinearities is diagonal and is

$$\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \text{diag} \left( \frac{\partial y_1}{\partial z_1}, \frac{\partial y_2}{\partial z_2}, \dots, \frac{\partial y_n}{\partial z_n} \right) = \text{diag} \left( \varphi'(z_1), \varphi'(z_2), \dots, \varphi'(z_n) \right)$$

A table of sigmoidal functions is reported in Appendix X. For the popular logistic function

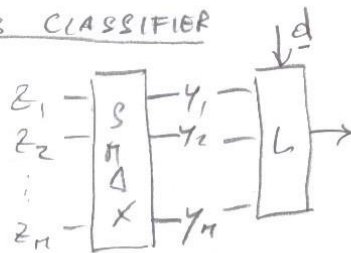
$$\varphi'(z_i) = \varphi(z_i)(1 - \varphi(z_i)), \text{ and}$$

$$\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = \text{diag} (y_1(1-y_1), y_2(1-y_2), \dots, y_n(1-y_n))$$

Note that using the cross entropy the denominators in  $\frac{\partial L}{\partial \mathbf{y}}$  cancel out and we have simply

$$\frac{\partial L}{\partial \mathbf{z}} = \left( \frac{\partial \mathbf{y}}{\partial \mathbf{z}} \right)^T \frac{\partial L}{\partial \mathbf{y}} = \mathbf{y} - \mathbf{d}$$

## MULTI-CLASS CLASSIFIER



Two.43

(Recall that  $\underline{d}$  may be the one-hot coded class or a smooth desired distribution.)

The total loss can be computed just as in regression, using the sum of individual losses, as

$$L(\underline{y}, \underline{d}) = \sum_{i=1}^M \psi(d_i - y_i) = \sum_{i=1}^M \psi(\epsilon_i)$$

that for the quadratic loss-function gives

$$L(\underline{y}, \underline{d}) = -2 \underline{\epsilon}$$

More appropriately using the cross entropy the loss is

$$L(\underline{y}, \underline{d}) = H(\underline{d}, \underline{y}) = \sum_{i=1}^M d_i \log \frac{d_i}{y_i}$$

and the gradient is (see Appendix)

$$\frac{\partial L}{\partial \underline{y}} = - \begin{bmatrix} \frac{d_1}{y_1} \\ \frac{d_2}{y_2} \\ \vdots \\ \frac{d_M}{y_M} \end{bmatrix} = - \frac{\underline{d}}{\underline{y}} \quad / \text{division element by element.}$$

The Jacobian across the Softmax function is (see Appendix)

$$\frac{\partial \underline{y}}{\partial \underline{z}} = \text{diag}(\underline{y}) - \underline{y} \underline{y}^T$$

Recall that both  $\underline{d}$  and  $\underline{y}$  are distributions, i.e.  $\sum_{i=1}^M d_i = 1$  and  $\sum_{i=1}^M y_i = 1$ . Therefore when we use the cross-entropy

Two.44

$$\frac{\partial L}{\partial \underline{z}} = \left( \frac{\partial \underline{y}}{\partial \underline{z}} \right)^T \frac{\partial L}{\partial \underline{y}} = -\text{diag}(\underline{y}) \frac{\underline{d}}{\underline{y}} + \underline{y} \underline{y}^T \frac{\underline{d}}{\underline{y}}$$

$$= -\underline{d} + \begin{pmatrix} y_1^2 & y_1 y_2 & \dots & y_1 y_M \\ y_2 y_1 & y_2^2 & \dots & y_2 y_M \\ \vdots & \vdots & \ddots & \vdots \\ y_M y_1 & y_M y_2 & \dots & y_M^2 \end{pmatrix} \begin{pmatrix} \frac{d_1}{y_1} \\ \frac{d_2}{y_2} \\ \vdots \\ \frac{d_M}{y_M} \end{pmatrix}$$

$$= -\underline{d} + \begin{pmatrix} y_1 d_1 + y_1 d_2 + \dots + y_1 d_M \\ y_2 d_1 + y_2 d_2 + \dots + y_2 d_M \\ \vdots \\ y_M d_1 + y_M d_2 + \dots + y_M d_M \end{pmatrix}$$

$$= -\underline{d} + \underline{y} = \underline{y} - \underline{d}$$

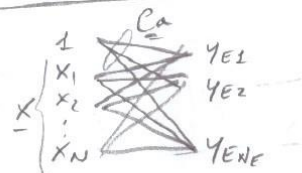
It simplifies  
just as in  
the multiple  
binary classifier.

# GRADIENTS FOR THE FIRST LAYER

Two.45

clearly the gradient  $\frac{\partial y_E}{\partial \theta}$  depends on the type of embedding layer. let us look at few options.

## LINEAR LAYER



$$y_E = C_a^T x_a = \begin{bmatrix} c_{10} & c_{11} & \dots & c_{1N} \\ c_{20} & c_{21} & \dots & c_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{N0} & c_{N1} & \dots & c_{NN} \end{bmatrix} \begin{bmatrix} 1 \\ x \\ \vdots \\ x \end{bmatrix}$$

$$C_a = \begin{bmatrix} c_{10} & c_{20} & \dots & c_{N0} \\ c_{11} & c_{21} & \dots & c_{N1} \\ \vdots & \vdots & \ddots & \vdots \\ c_{1N} & c_{2N} & \dots & c_{NN} \end{bmatrix} = \begin{bmatrix} c_{11}^a & c_{12}^a & \dots & c_{1N}^a \\ c_{21}^a & c_{22}^a & \dots & c_{2N}^a \\ \vdots & \vdots & \ddots & \vdots \\ c_{N1}^a & c_{N2}^a & \dots & c_{NN}^a \end{bmatrix} = \begin{bmatrix} c_{11}^{aT} \\ c_{21}^{aT} \\ \vdots \\ c_{N1}^{aT} \end{bmatrix} \begin{bmatrix} 1 \\ x \\ \vdots \\ x \end{bmatrix}$$

Since the following layer is linear

$$z = W_a^T y_E = W_a^T \begin{bmatrix} 1 \\ C_a^T x_a \end{bmatrix}$$

The entire mapping from  $x$  to  $z$  will be linear, and it may be useful to consider this option because  $C_a$  could be absorbed into the next layer matrix  $W_a$ .

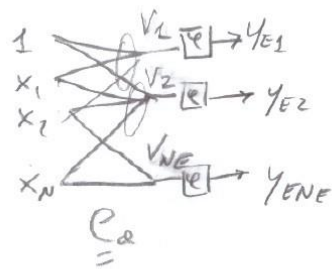
A possibility that could be interesting is to use the first layer as a compression layer, i.e. with  $N_E \ll N$  as in principal component analysis. In any case looking at the connectivity

$$\frac{\partial L}{\partial c_{ij}^a} = \frac{\partial L}{\partial y_{Ej}} \frac{\partial y_{Ej}}{\partial c_{ij}^a} = \frac{\partial L}{\partial y_{Ej}} x_{aj} \quad j=1, \dots, N_E$$

$$\frac{\partial L}{\partial C_a^T} = \begin{bmatrix} \frac{\partial L}{\partial c_{11}^a} & \frac{\partial L}{\partial c_{12}^a} & \dots & \frac{\partial L}{\partial c_{1N}^a} \\ \frac{\partial L}{\partial c_{21}^a} & \frac{\partial L}{\partial c_{22}^a} & \dots & \frac{\partial L}{\partial c_{2N}^a} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial L}{\partial c_{N1}^a} & \frac{\partial L}{\partial c_{N2}^a} & \dots & \frac{\partial L}{\partial c_{NN}^a} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_{E1}} x_{a1} & \frac{\partial L}{\partial y_{E2}} x_{a2} & \dots & \frac{\partial L}{\partial y_{EN}} x_{aN} \end{bmatrix} = \left( \frac{\partial L}{\partial y_E} \right)^T \otimes x_a$$

# LINEAR LAYER WITH NONLINEARITIES

TWO.46



$\phi(\cdot)$  could be ReLU  
Sigmoid

$$\underline{c}_a = \begin{bmatrix} c_{10} & c_{20} & \dots & c_{NE0} \\ c_{11} & c_{12} & \dots & c_{1NE} \\ c_{21} & c_{22} & \dots & c_{2NE} \\ \dots & \dots & \dots & \dots \\ c_{NE1} & c_{NE2} & \dots & c_{NENE} \end{bmatrix} = \begin{bmatrix} \underline{c}_1^a & \underline{c}_2^a & \dots & \underline{c}_{NE}^a \end{bmatrix}$$

$$\underline{y}_E = \begin{bmatrix} \phi(v_1) \\ \phi(v_2) \\ \dots \\ \phi(v_{NE}) \end{bmatrix} = \phi(\underline{v}) = \phi(\underline{c}_a^+ \underline{x}_a)$$

$$\frac{\partial L}{\partial \underline{c}_j^a} = \frac{\partial L}{\partial y_{Ej}} \frac{\partial y_{Ej}}{\partial c_j^a} = \frac{\partial L}{\partial y_{Ej}} \frac{\partial y_{Ej}}{\partial v_j} \frac{\partial v_j}{\partial c_j^a}$$

$$= \frac{\partial L}{\partial y_{Ej}} \phi'(v_j) x_a$$

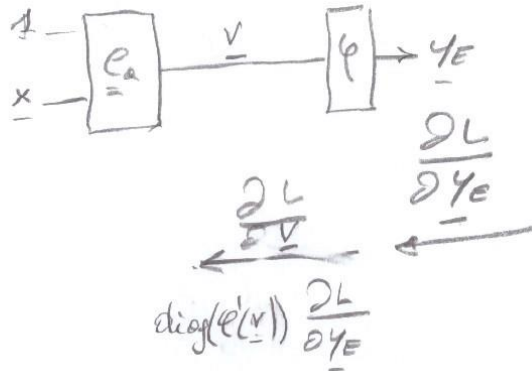
$$\frac{\partial L}{\partial \underline{c}_a^a} = \begin{bmatrix} \frac{\partial L}{\partial c_{10}^a} & \frac{\partial L}{\partial c_{20}^a} & \dots & \frac{\partial L}{\partial c_{NE0}^a} \end{bmatrix} = \begin{bmatrix} \frac{\partial L}{\partial y_{E1}} \frac{\partial y_{E1}}{\partial v_1} x_a, \frac{\partial L}{\partial y_{E2}} \frac{\partial y_{E2}}{\partial v_2} x_a, \dots, \frac{\partial L}{\partial y_{ENE}} \frac{\partial y_{ENE}}{\partial v_{NE}} x_a \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial L}{\partial y_{E1}} \phi'(v_1) x_a, \frac{\partial L}{\partial y_{E2}} \phi'(v_2) x_a, \dots, \frac{\partial L}{\partial y_{ENE}} \phi'(v_{NE}) x_a \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\partial L}{\partial v_1} x_a, \frac{\partial L}{\partial v_2} x_a, \dots, \frac{\partial L}{\partial v_{NE}} x_a \end{bmatrix}$$

$$= \left( \frac{\partial L}{\partial \underline{v}} \right)^T \otimes x_a = \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \text{diag}(\phi'(v)) \otimes x_a$$

Two.47



The rule is then: once reached the linear layer output  $\underline{u}$  with  $\frac{\partial L}{\partial \underline{v}}$ ,

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

In going through a block multiply by the Jacobian transposed.

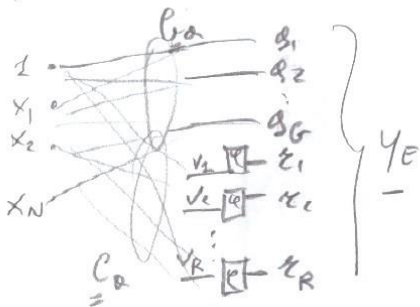
Recall that for ReLU's  $\phi'(v_i) = u(v_i)$  that set to zero some components of  $\frac{\partial L}{\partial y_E}$

For leaky ReLU sigmoid  $\phi'(v_i) = \phi(v_i)(1 - \phi(v_i))$   
 $= y_{Ei}(1 - y_{Ei})$

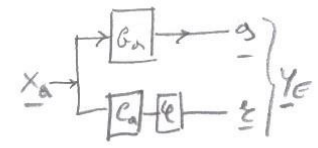
For all other nonlinearities Appendix X gives expressions and derivatives.



SKIP-CONNECTION LAYER



$$\begin{cases} g = \underline{G}_a^T \underline{x}_a \\ z = \varphi(\underline{C}_a^T \underline{x}_a) \end{cases}$$



$$\underline{y}_E = \begin{pmatrix} g \\ z \end{pmatrix}$$

$$\frac{\partial L}{\partial \underline{y}_E} = \begin{pmatrix} \frac{\partial L}{\partial g} \\ \frac{\partial L}{\partial z} \end{pmatrix} \quad \text{The backward gradient on } \underline{y}_E \text{ is split in two parts.}$$

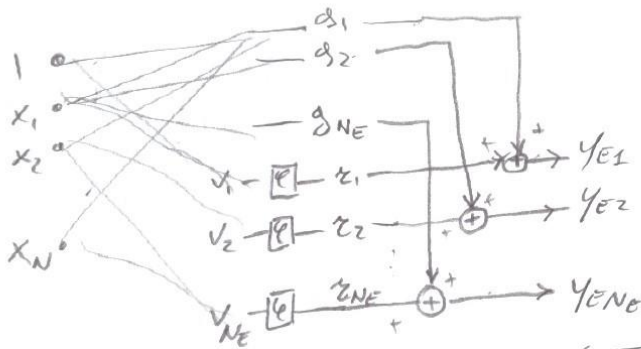
The gradient for the update of  $\underline{C}_a$ , similarly to the above cases becomes

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

The gradient for the update of matrix  $\underline{C}_a$  is

$$\frac{\partial L}{\partial \underline{C}_a} = \left( \left( \frac{\partial L}{\partial z} \right)^T \text{diag}(\varphi'(z)) \right) \otimes \underline{x}_a$$

## RESIDUAL LAYER



$$\underline{y}_E = \underline{C}_a^T \underline{x}_a + \mathcal{L}(\underline{C}_a^T \underline{x}_a) \quad \leftarrow \frac{\partial L}{\partial \underline{y}_E}$$

The gradient for the update of  $\underline{C}_a$  is

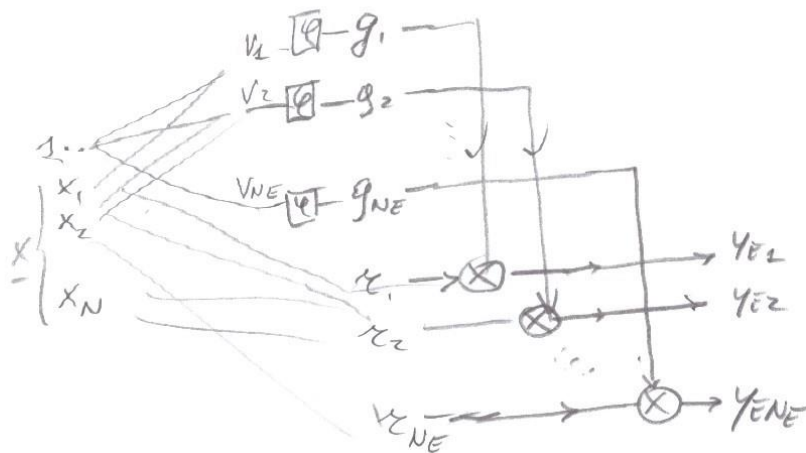
$$\frac{\partial L}{\partial \underline{C}_a} = \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \otimes \underline{x}_a$$

and the one for the update of  $\underline{C}_a$

$$\frac{\partial L}{\partial \underline{C}_a} = \left[ \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \text{diag}(\mathcal{L}'(\underline{v})) \right] \otimes \underline{x}_a$$

# GATED LAYER WITH SIGMOIDS

TWO.50.



$$\underline{y}_E = \underline{g} \odot \underline{z} = \varphi \left( \underset{\substack{\uparrow \\ N_E}}{\underline{G}_a^T \underline{x}_a} \right) \odot \left( \underset{\substack{\uparrow \\ N_E}}{\underline{C}_a^T \underline{x}_a} \right)$$

To go through the multiplier we can consider

$$\frac{\partial L}{\partial \underline{z}} \quad \text{and} \quad \frac{\partial L}{\partial \underline{g}} \quad \text{that give}$$

$$\frac{\partial L}{\partial \underline{z}} = \frac{\partial L}{\partial \underline{y}_E} \odot \underline{g}$$

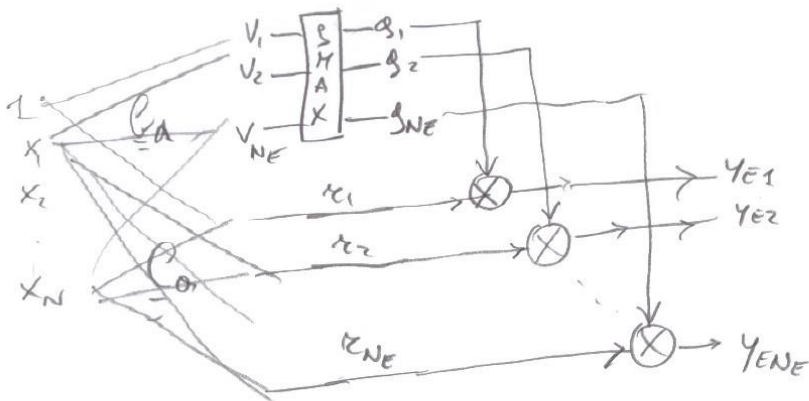
$$\frac{\partial L}{\partial \underline{g}} = \frac{\partial L}{\partial \underline{y}_E} \odot \underline{z}$$

$$\left. \begin{aligned} \frac{\partial L}{\partial \underline{C}_a} &= \left( \frac{\partial L}{\partial \underline{z}} \right)^T \otimes \underline{x}_a = \left( \underline{g}^T \odot \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \right) \otimes \underline{x}_a \end{aligned} \right\}$$

$$\left. \begin{aligned} \frac{\partial L}{\partial \underline{C}_a} &= \left( \frac{\partial L}{\partial \underline{g}} \right)^T \text{diag}(\varphi'(v_i)) \otimes \underline{x}_a = \left[ \underline{z}^T \odot \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \right] \text{diag}(\varphi'(v_i)) \otimes \underline{x}_a \end{aligned} \right\}$$

# GATED LAYER WITH SOFTMAX

TWO.51



$$\underline{y}_E = \underline{g} \odot \underline{e} = \sum_{\text{softmax}} \left( \underline{C}_a^T \underline{x}_a \right) \odot \left( \underline{C}_a^T \underline{x}_a \right)$$

For the lower branch, we have just as for the gated layer with sigmoidals

$$\frac{\partial L}{\partial \underline{C}_a} = \left( \underline{g}^T \odot \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \right) \otimes \underline{x}_a$$

For the upper branch we need the Jacobian across the softmax that we recall is

$$\frac{\partial \underline{g}}{\partial \underline{v}} = \text{diag}(\underline{g}) - \underline{g}\underline{g}^T$$

$$\frac{\partial L}{\partial \underline{C}_a} = \left( \frac{\partial L}{\partial \underline{v}} \right)^T \otimes \underline{x}_a = \left( \left( \frac{\partial \underline{g}}{\partial \underline{v}} \right)^T \frac{\partial L}{\partial \underline{g}} \right)^T \otimes \underline{x}_a$$

$$= \left[ \left( \frac{\partial L}{\partial \underline{g}} \right)^T \left( \text{diag}(\underline{g}) - \underline{g}\underline{g}^T \right) \right]^T \otimes \underline{x}_a$$

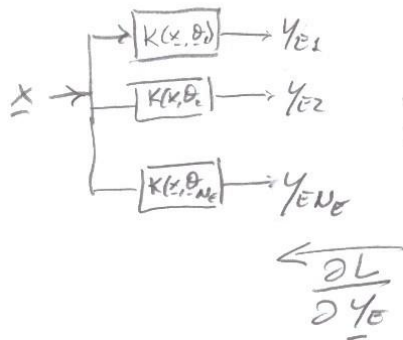
$$= \left[ \left( \frac{\partial L}{\partial \underline{y}_E} \right)^T \odot \underline{e} \left( \text{diag}(\underline{g}) - \underline{g}\underline{g}^T \right) \right]^T \otimes \underline{x}_a$$

TWO, 52

# GRADIENT FOR <sup>SOFT</sup> KERNEL-BASED ARCHITECTURES

Kernel-based architectures, such as the radial-basis functions, are usually trained in an unsupervised way, i.e. only on the input sequence  $\{\underline{x}_i\}$  of the training set.

However, their parameters, at least in principle, can be adjusted according to a backward gradient flow coming from the output.



For each kernel function  $K(\underline{x}, \underline{\theta}_i)$  we need the gradient for  $\underline{\theta}_i$ :

$$\frac{\partial L}{\partial \underline{\theta}_i} = \frac{\partial L}{\partial y_{Ei}} \frac{\partial y_{Ei}}{\partial \underline{\theta}_i}$$

$$\frac{\partial y_{Ei}}{\partial \underline{\theta}_i} = \frac{\partial K(\underline{x}, \underline{\theta}_i)}{\partial \underline{\theta}_i}$$

For radial functions

$$K(\underline{x}, \underline{\theta}_i) = K\left(\frac{\|\underline{x} - \underline{\mu}_i\|^2}{h_i}\right), \quad \underline{\theta}_i = \begin{pmatrix} \underline{\mu}_i \\ h_i \end{pmatrix}$$

where  $\underline{\mu}_i$  is the cluster centers and  $h_i$  controls the size of the sphere around  $\underline{\mu}_i$ .

For gaussian radial functions

$$K(x, \underline{\theta}_i) = e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}}$$

$$\begin{aligned} \frac{\partial K}{\partial \mu_i} &= -e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}} \frac{\partial}{\partial \mu_i} \left[ \frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2} \right] \\ &= -e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}} \frac{\partial}{\partial \mu_i} \left[ \frac{x^T x - 2\mu_i^T x + \mu_i^T \mu_i}{2h_i^2} \right] \\ &= e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}} \frac{\partial}{\partial \mu_i} \left[ \frac{x - \mu_i}{h_i^2} \right] = \gamma_{Ei} \left( \frac{x - \mu_i}{h_i^2} \right) \end{aligned}$$

$$\begin{aligned} \frac{\partial K}{\partial h_i} &= -e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}} \frac{\partial}{\partial h_i} \left( \frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2} \right) \\ &= -e^{-\frac{(x-\mu_i)^T(x-\mu_i)}{2h_i^2}} \frac{1}{2} \frac{(x-\mu_i)^T(x-\mu_i)}{h_i^3} \left( -\frac{2}{h_i^3} \right) \\ &= \gamma_{Ei} \frac{\|x - \mu_i\|^2}{h_i^3} \end{aligned}$$

$$\begin{cases} \frac{\partial K}{\partial \mu_i} = \gamma_{Ei} \frac{x - \mu_i}{h_i^2} \\ \frac{\partial K}{\partial h_i} = \gamma_{Ei} \frac{\|x - \mu_i\|^2}{h_i^3} \end{cases}$$

FOR MULTIQADRATIC FUNCTIONS

TWO.54

$$K(\underline{x}, \theta_i) = \sqrt{\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2}$$

$$\frac{\partial K}{\partial \underline{\mu}_i} = \frac{-(\underline{x} - \underline{\mu}_i)}{\sqrt{\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2}} = -y_{Ei} (\underline{x} - \underline{\mu}_i)$$

$$\frac{\partial K}{\partial \beta_i} = \frac{\beta_i}{\sqrt{\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2}} = y_{Ei} \beta_i$$

$$\begin{aligned} \frac{\partial \|\underline{x} - \underline{\mu}_i\|^2}{\partial \underline{\mu}_i} \frac{\partial (x^T x - 2\underline{\mu}_i^T x + \underline{\mu}_i^T \underline{\mu}_i)}{\partial \underline{\mu}_i} \\ = -2\underline{x} + 2\underline{\mu}_i = \\ = 2(\underline{\mu}_i - \underline{x}) \end{aligned}$$

FOR THE INVERSE MULTIQADRATIC

$$K(\underline{x}, \theta_i) = \frac{1}{\sqrt{\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2}}$$

$$\frac{\partial K}{\partial \underline{\mu}_i} = +\frac{1}{\cancel{2}} \frac{1}{(\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2)^{3/2}} \cancel{2} (\underline{x} - \underline{\mu}_i)$$

$$= y_{Ei}^3 (\underline{x} - \underline{\mu}_i)$$

$$\frac{\partial K}{\partial \beta_i} = -\frac{1}{2} \frac{2\beta_i}{(\|\underline{x} - \underline{\mu}_i\|^2 + \beta_i^2)^{3/2}} = -y_{Ei}^3 \beta_i$$