

LINEAR NETWORKS  
FOR DATA-DRIVEN  
CCLASSIFICATION

Prof. FRANCESCO A.N. PALMIERI  
UNIVERSITA' DELLA CAMPANIA  
"LUIGI VANVITELLI"

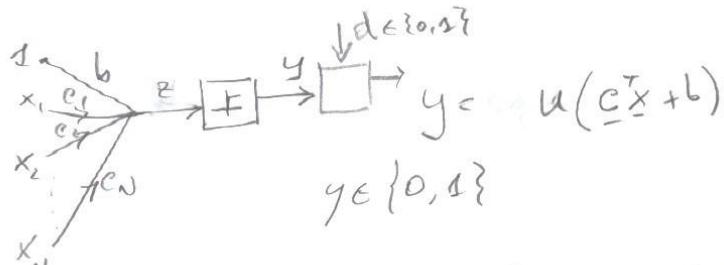
[ CORSO DI SIGNAL PROCESSING  
E DATA FUSION OCT. 2023 ]

(Lec)

## LINEAR NETWORKS FOR CLASSIFICATION

BINARY CLASSIFIER WITH HARD THRESHOLD

The simplest neural network for classification is the linear perceptron shown in the figure.



The perceptron, originally proposed as a simplified model of a biological neuron (thresholding) represents a linear binary classifier. We have found already this structure ~~when we derived~~ when we derived model-based classifiers for Gaussian or exponential-family generative models. The network here is ~~formulated~~ <sup>instead</sup> and must be learned from the training set. More specifically we have a training set of pairs

$$\mathcal{D} = \{(x^{[n]}, d^{[n]}), n=1, \dots, n_r\}, \quad d^{[n]} \in \{0, 1\}$$

the parameters  $c$  and  $b$  fixed

and we seek ~~to minimize~~ the classification error, i.e. the events  $y^{[n]} \neq d^{[n]}$ .

The cost function can be written as

$$E(c, b) = \sum_{n=1}^{n_r} |d^{[n]} - y^{[n]}| = \sum_{n=1}^{n_r} |d^{[n]} - c^T x^{[n]} - b|$$

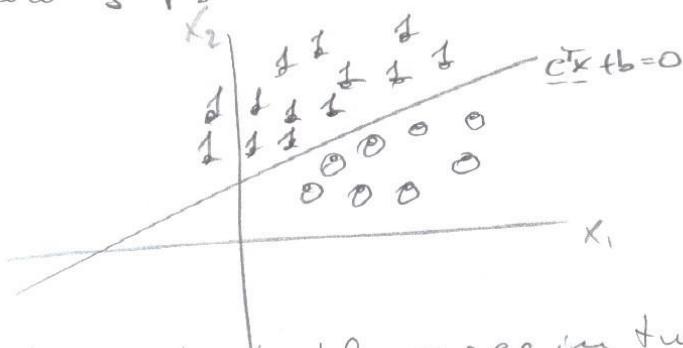
$E(c, b)$  counts the number of errors.

Note that since  $y^{[n]} - d^{[n]} \in \{-1, 0, 1\}$  the cost can be written also as a sum of squares

$$E(c, b) = \sum_{u=1}^n (d(u) - c^T x(u) - b)^2$$

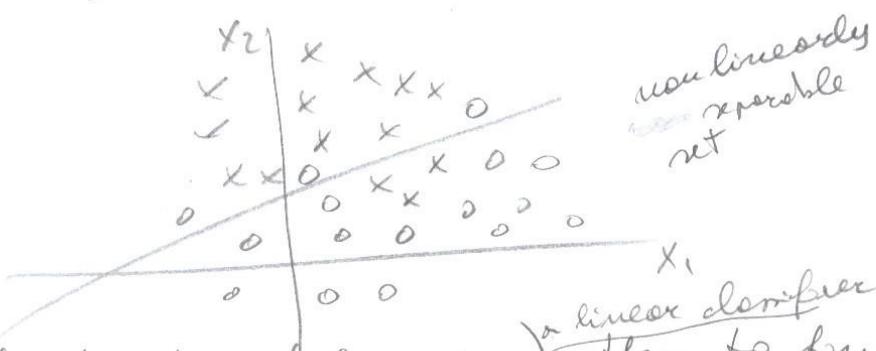
(lec)

The classifier can be visualized in the following space in 2D where the hyperplane



$c^T x + b = 0$  divides the space in two regions. The points above produce output  $y=1$  and the ones below,  $y=0$ .

Clearly the examples in the training set (or in the validation/test set) may not be such that there exist a hyperplane that separates exactly the two classes. (linearly separable)



The objective of learning is then to find the best parameters ( $c, b$ ) that leave the smallest number of examples on the wrong side of the hyperplane.

Unfortunately, there is no closed-form solution to (1.2) for this problem, but we can invoke iterative procedures that adjust progressively the position of the hyperplane until the best configuration is found.

For easier notation we use, from now on, the augmented vectors  $\underline{x}_a = \begin{bmatrix} 1 \\ \underline{x} \end{bmatrix}$  and  $\underline{w} = \begin{bmatrix} b \\ \underline{w} \end{bmatrix}$  so that

$$y[n] = \text{sgn}(\underline{w}^T \underline{x}_a[n]) \quad (u(\varepsilon) \text{ is the indicator function } = 1 \text{ for } \varepsilon > 0 \text{ and } 0 \text{ for } \varepsilon \leq 0)$$

The perceptron rule is a simple iterative algorithm that adjusts  $\underline{w}$  progressively to the best linear classifier. (Rosenblatt, 1962)  
The steps are the following:

- 1. Initialize  $\underline{w}$  with random values
- 2. Present the examples  $(\underline{x}_a[n], d[n])$  of the training set one at a time, i.e. for each example of  $y[n] = d[n]$ ,  $\underline{w}[n] = \underline{w}[n-1]$  (no change)
  - else if  $y[n] = 1$  and  $d[n] = 0$ ,  $\underline{w}[n] = \underline{w}[n-1] + \mu \underline{x}_a[n]$
  - else if  $y[n] = 0$  and  $d[n] = 1$ ,  $\underline{w}[n] = \underline{w}[n-1] + \mu \underline{x}_a[n]$
- 3. Once the whole training set has been presented, take the resulting  $\underline{w}$  and use it as initial condition for a new epoch (perhaps shuffling the training set order).
- 4. Stop at convergence of  $\underline{w}$  and  $E(\underline{w})$ , i.e. when there are no more significant changes.

The algorithm can be compactly written as the update

$$\underline{w}[n] = \underline{w}[n-1] + y[n] \underline{x}_a[n] \quad (d[n] - u(\underline{w}[n-1]^T \underline{x}_a[n]))$$

because the error  $d[n] - u(\underline{w}[n-1]^T \underline{x}_a[n])$  can only be -1, 0 or 1.

Note that in the perceptron algorithm we have included a stepsize parameter  $\mu$ . Just as in any gradient algorithm,  $\mu$  must be positive. It could be adaptive  $\mu(t)$ , but in the original algorithm proposed by Rosenblatt it was fixed to  $\mu=1$  even if it is unnecessary to keep it small with  $0 < \mu < 1$ .

The proof of convergence of the algorithm is very interesting and is one of the few rigorous results that can be invoked in the neural network literature. It is derived under the assumption that the examples are linearly separable that  $\mu=1$  and  $w^{(0)}=0$ . We direct the interested student to the specific literature linking ourselves here to a few key observations that relate this algorithm to a gradient search.

Note that when a pattern  $x_a^{(n)}$ , that should be classified with  $d(x_a) = 1$ , and gives  $y_a^{(n)} = 0$ , the update is

$$w^{(n)} = w^{(n-1)} + \gamma x_a^{(n)} [1 - 0]$$

Now if we re-use the same pattern  $x_a^{(n)}$  with the new  $w^{(n)}$  (this is called a staggered update) we have

$$\begin{aligned} w^{(n)} &= w^{(n-1)} + \gamma x_a^{(n)} [1 - u((w^{(n-1)} + \gamma x_a^{(n)})^T x_a^{(n)})] \\ &= w^{(n-1)} + \gamma x_a^{(n)} [1 - u(w^{(n-1)} x_a^{(n)} + \gamma x_a^{(n)} x_a^{(n)})] \end{aligned}$$

Now since  $\mu > 0$  and  $x_a^{(n)} x_a^{(n)} > 0$ , the update tends to push  $u(\cdot) \rightarrow 1$ , i.e. towards a correct classification. The dual situation happens when

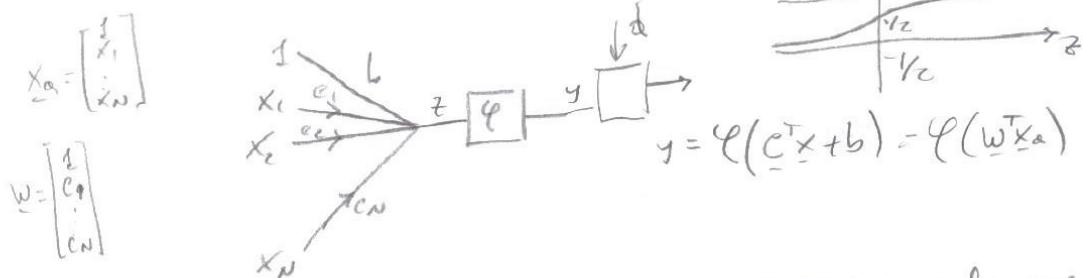
$$d(x_a) = 0 \text{ and } y_a^{(n)} = 1.$$

A standard gradient algorithm can't be derived because  $u(\cdot)$  is not differentiable. We will see in the following how we can find the perceptron algorithm as the limit of a smooth minimization function.

(Lc5)

### BINARY CLASSIFIER WITH A SMOOTH SIGMOIDAL FUNCTION

We consider now the same network structure of the previous section, but with a smooth activation function



We have already encountered this architecture when we considered the model-based solution for Gaussian or exponential family generative models. We found that this network with  $\phi(z) = \frac{1}{1+e^{-z}}$  (logistic function) provides the posterior probability for one of the two classes. Recall that the posterior for the other is the complement to all of the output  $(1-y)$ . In this case we postulate the above architecture and want to learn its parameters from a training set  $\mathcal{D} = \{(x^{[n]}, d^{[n]}), n=1, \dots, m\}$  where the desired output  $d^{[n]}$  is either 1 or 0 for the two classes.

We can adopt a standard cost function, such as the sum of squared errors

$$E(w) = \frac{1}{m} \sum_{n=1}^m \psi(d^{[n]} - \phi(w^T x^{[n]})) = \frac{1}{m} \sum_{n=1}^m (d^{[n]} - \phi(w^T x^{[n]}))^2$$

or more appropriately a binary cross-entropy  
(discussed in a previous section)

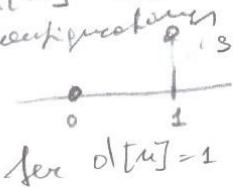
(LC6)

$$E_{ce}(\underline{w}) = \frac{1}{n^2} \sum_{u=1}^{n^2} H(d[u], P(w^T x_a[u]))$$

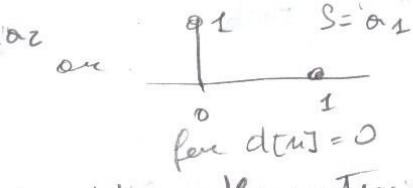
$$= \frac{1}{n^2} \sum_{u=1}^{n^2} \left[ d[u] \log \frac{1}{e^{(w^T x_a[u])}} + (1-d[u]) \log \frac{1}{1-e^{(w^T x_a[u])}} \right]$$

The cross-entropy appears to be a more appropriate choice if we want to attribute to the output of the classifier the meaning of a posterior probability. In fact we can see the desired value

$d[u]$  as the "desired distribution" with the two configurations



for  $d[u]=1$

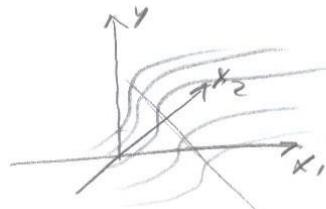
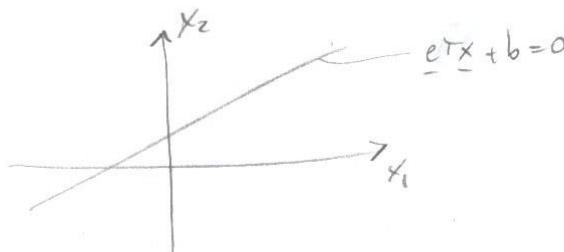


for  $d[u]=0$

and the output being the posterior distribution

$$1-y_{[u]} = P(S=o_2 | x) \quad y_{[u]} = P(S=o_1 | x)$$

The network, just as in the case with the hard threshold, bases its decision on a single hyperplane that divides the input space in two regions. the difference



here is that there is a smooth transition between the two regions (smooth ridge).

Even if the cross entropy is the most appropriate (CC7) loss function for this case, the squared error produces similar results. Therefore for completeness we consider both.

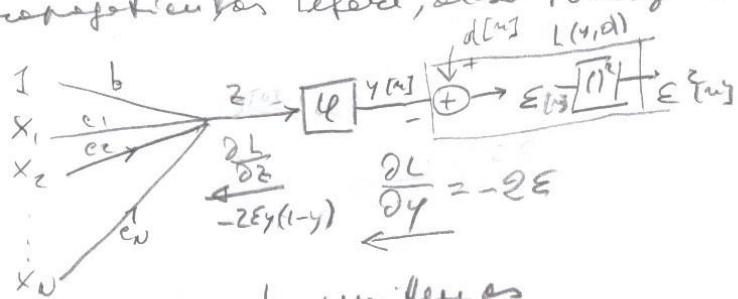
G. GRADIENT WITH SQUARE LOSS  
The gradient with the square loss is

$$\nabla_w \mathcal{E}_{\text{mse}}(\underline{w}) = -2 \sum_{m=1}^{M^2} (d[m] - \Phi(\underline{w}^T \underline{x}_a[m])) \Phi'(\underline{w}^T \underline{x}_a[m]) \underline{x}_a[m]$$

Using the well-known derivative for the log-loss function  
 $\Phi'(s) = \Phi(s)(1 - \Phi(s))$  (see Appendix)

$$\nabla_w \mathcal{E}_{\text{mse}}(\underline{w}) = -2 \sum_{m=1}^{M^2} (d[m] - y[m]) y[m] (1 - y[m]) \underline{x}_a[m]$$

Note the backpropagation flow before, also through the sigmoid block



The batch iteration can be written as  
 $w[k] = w[k-1] + \mu \sum_{m=1}^{M^2} (d[m] - y_{k-1}[m]) y_{k-1}[m] (1 - y_{k-1}[m]) \underline{x}_a[m]$

The notation  $y_{k-1}[m]$  is used to emphasize that we use the "old" values  $w[k-1]$  to compute the output

$$y_k[m] = \underline{w}^T \underline{x}_a[m]$$

We can give <sup>now</sup> more intuitive interpretation to the updates.

We see that  $y = \Phi(w^T x)$  is called the activation. Also

The contribution of the  $n$ th example  
to the previous weight,  $w_i[k-1]$  is made evident as

$$w_i[k] = w_i[k-1] + \mu (d[n] - y_{k-1}[n]) y_{k-1}[n] (1 - y_{k-1}[n]) x_{ai}[n]$$

+ ... (contribution from other examples) ...

Now suppose that the desired value in the  $n$ th example  
is  $d[n] = 1$ , but the output  $y_{k-1}[n] = \delta$ , where  $\delta < 1$   
is a small positive number. We are clearly in an undesirable  
situation because we would like  $y_{k-1}[n]$  to be close to  
 $d[n] = 1$ . Let us see if the update helps to correct the situation.

The contribution to the output is

$$w_i[k] x_{ai}[n] = w_i[k-1] x_{ai}[n] + \mu (d[n] - y_{k-1}[n]) y_{k-1}[n] (1 - y_{k-1}[n]) x_{ai}[n]$$

$\Delta_i + \dots$

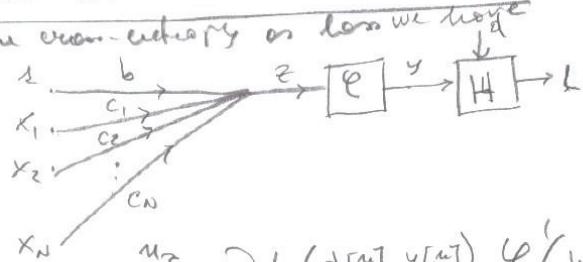
where the error  $d[n] - y_{k-1}[n] = 1 - \delta$  is positive  
 $\Rightarrow y_{k-1}[n] (1 - y_{k-1}[n]) = \delta (1 - \delta) > 0$  and  $x_{ai}[n] > 0$ . Therefore  
the contribution  $\Delta_i$  to the output is a positive and  
it tends to push  $y_{k-1}[n]$  to higher values,  
itowards  $d[n] = 1$ .

A dual situation happens when  $d[n] = 0$  and  
 $y_{k-1}[n] = 1 - \delta$ . The error  $d[n] - y_{k-1}[n] = 0 - 1 + \delta < 0$  is  
negative,  $y_{k-1}[n] (1 - y_{k-1}[n]) = (1 - \delta) \delta > 0$ ,  $x_{ai}[n] > 0$  and the  
update will be negative pushing the output towards  
smaller values.

(Lc9)

### GRADIENT WITH THE CROSS ENTROPY

Using the cross-entropy as loss we have



$$\nabla_w \mathcal{E}_{ce}(\underline{w}) = \frac{1}{m_2} \sum_{u=1}^{m_2} \frac{\partial L(d[u], y[u])}{\partial y[u]} \varphi'(w^T x_{o[u]}) x_{o[u]}$$

$$L(d, y) = f(d, y) = d \log \frac{1}{y} + (1-d) \log \frac{1}{1-y}, \text{ and}$$

$$\frac{\partial L(d, y)}{\partial y} = -\frac{y-d}{y(1-y)} \quad (\text{see Appendix})$$

The gradient is

$$\nabla_w \mathcal{E}_{ce}(\underline{w}) = \frac{1}{m_2} \sum_{u=1}^{m_2} \left[ \frac{y[u] - d[u]}{y[u](1-y[u])} \right] y[u] (1-y[u]) x_{o[u]}.$$

The batch update for cross-entropy loss is

$$\underline{w}[k] = \underline{w}[k-1] + \mu \sum_{u=1}^{m_2} (d[u] - y_{k-1}[u]) x_{o[u]} \quad (\text{Analogous similarity with the RPROP rule})$$

where again  $y_{k-1}[u] = \underline{w}[k-1]^T \underline{x}_o[u]$  is the output computed with the "old values"  $\underline{w}[k-1]$ .

Also here it is quite interesting to give an intuitive interpretation to the update by looking at the contribution of the  $u$ th example to the update of the  $i$ th weight and its impact on the output.

The contribution to the  $i$ -th weight could be evidenced as  
 $w_i[k] = w_i[k-1] + \mu [d[u] - y_{k-1}[u]] x_{ai} + \dots$  contributions from other examples  
 over its contribution to the output  $z$

(Lec 10)

$$w_i[k] x_{ai}[u] = w_i[k-1] x_{ai}[u] + \mu [d[u] - y_{k-1}[u]] x_{ai}^2 \\ + \dots \text{contribution from other examples} \dots$$

Now if  $d[u]=1$  and  $y_{k-1}[u]=\delta$  with  $\delta > 0$  small value, we are in an undesirable situation because we would like to move  $y_{k-1}[u]$  close to one.

The correction term is

$$\mu(d[u] - y_{k-1}[u]) x_{ai}[u]^2 = \mu(1-\delta)x_{ai}^2 > 0$$

It has pushed the output  $z$  and then  $y$  to higher values.

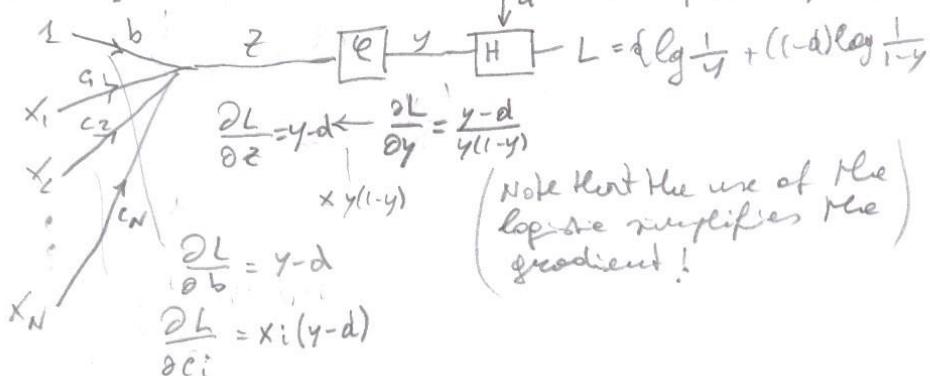
The dual situation is when  $d[u]=0$  and  $y_{k-1}[u]=1-\delta$ . Again undesirable because we would like to move  $y_{k-1}[u]$  close to zero. The correction term is

$$\mu(d[u] - y_{k-1}[u]) x_{ai}[u]^2 = \mu(0-1+\delta)x_{ai}^2 < 0$$

and the output  $z$  and then  $y$  is pulled towards lower values.

The backpropagation flow is also interesting.

On a generic slice  $(x^{[u]}, d^{[u]})$  (drop  $u$  for simplicity)

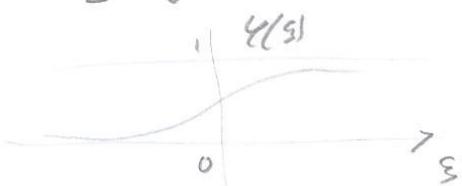


10.11

## BINARY SMOOTH CLASSIFIERS WITH OTHER SIGMOIDAL FUNCTIONS

We have argued in the previous sections that the smooth sigmoid placed on the output of the classifier is a logistic function. This choice was the consequence of the results on model-based classifiers where we had knowledge of the likelihoods (i.e. for which are exponential family). Certainly special densities may be hidden in the data, but we generally do not know much in data-driven methods.

Therefore any sigmoid-shaped function may be used



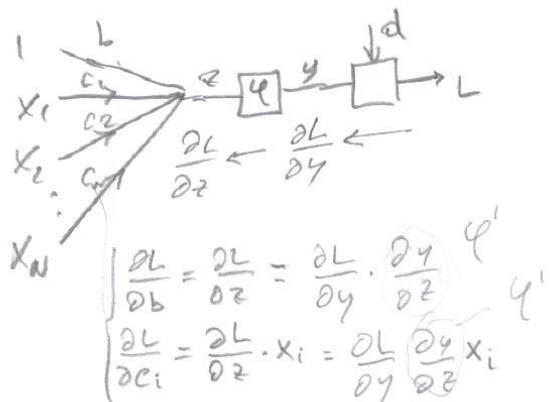
It must be positive, non-decreasing and such that  $y(s) \xrightarrow[s \rightarrow \infty]{} 1$  and  $y(s) \xrightarrow[s \rightarrow -\infty]{} 0$ .

Recall that this is the behaviour of any cumulative distribution function (cdf) where its derivative is the probability density function (pdf).

Therefore we could pick any (pdf, cdf) pair we know, such as (uniform, piecewise linear), (gaussian, cumulative gaussian), (laplacian, laplacian cdf)

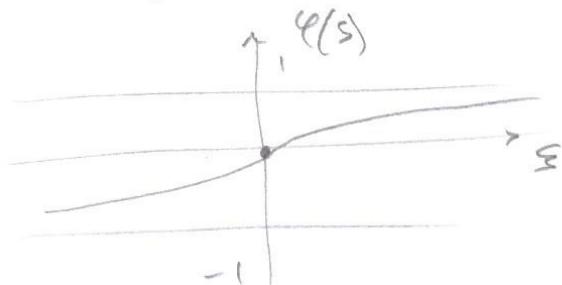
Various choices are detailed in the Appendix on Activation functions.

In all cases we backpropagate the loss in LC.12  
the sigmoidal function  $\varphi$



In some cases the binary classifier can be slightly changed with target values ~~for the two classes~~ ~~at the output~~  
 $d \in \{-1, 1\}$ , instead of  $d \in \{0, 1\}$ .

The sigmoid function must be scaled and shifted



and be such that  $\lim_{s \rightarrow \infty} \varphi(s) \rightarrow 1$  and  $\lim_{s \rightarrow -\infty} \varphi(s) \rightarrow -1$ .

Examples of these sigmoid are the hyperbolic tangent and the trigonometric tangent described in the Appendix are activation functions.

To interpret the output of the sigmoid as a posterior we must scale and shift to get

$$P(\text{class } | x) = \frac{y+1}{2}$$

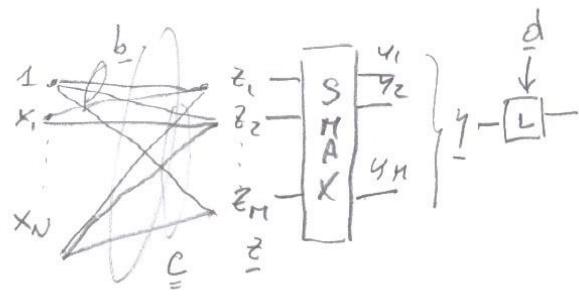
or the inverse

$$y = 2P(\text{class } | x) - 1$$

26.13

## MULTI-CLASS LINEAR CLASSIFIER

The typical multi-class linear classifier has the following architecture



$$z = \underline{C}^T \underline{x} + \underline{b} \quad y = S_{\max}(z)$$

Using the augmented vector  $\underline{x}_a = \begin{bmatrix} 1 \\ \underline{x}_1 \\ \vdots \\ \underline{x}_N \end{bmatrix} = \begin{bmatrix} 1 \\ \underline{x} \end{bmatrix}$  and the augmented matrix  $\underline{W} = \begin{bmatrix} \underline{b}^T \\ \underline{C}^T \end{bmatrix}$

$$\text{we have } \underline{z} = \underline{W}^T \underline{x}_a, \quad y = S_{\max}(\underline{z})$$

[In note(\*) we explode the augmented vector and vector for better visualization of indices.]

$S_{\max}(z)$  is the usual softmax function already encountered in model-based architectures for M-class classifiers. The linear part comes out from the model if we assume gaussian (an exponential family) classes with the same covariance matrix.

Hence, instead, we have no knowledge of the generative model. The structure is postulated and we have to learn it from the training set

$$\mathcal{D} = \{(\underline{x}_{\text{train}}^u, \underline{d}_{\text{train}}^u), u=1, \dots, n\}$$

NOTE  
(\*) For better visualization of the indices in the matrices, let us explode them.

LC14

$\underline{C}$  and  $\underline{b}$  are defined as

$$\underline{C} = \begin{bmatrix} C_{11} & C_{12} & \dots & C_{1M} \\ C_{21} & C_{22} & \dots & C_{2M} \\ \vdots & & & \\ C_{N1} & C_{N2} & \dots & C_{NM} \end{bmatrix} \quad \underline{b} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_N \end{bmatrix}$$

note(s)

The output  $\underline{z}$  is computed as

$$\underline{z} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_M \end{bmatrix} = \underline{C}^T \underline{x} + \underline{b} = \begin{bmatrix} C_{11} & C_{21} & \dots & C_{N1} \\ C_{12} & C_{22} & \dots & C_{N2} \\ \vdots & & & \\ C_{1M} & C_{2M} & \dots & C_{NM} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_M \end{bmatrix}$$

Therefore  $c_{ij}$  connects input  $x_i$  with output  $z_j$



In the augmented representation

$$\underline{z} = \begin{bmatrix} b_1 & C_{11} & C_{21} & \dots & C_{N1} \\ b_2 & C_{12} & C_{22} & \dots & C_{N2} \\ \vdots & & & & \\ b_M & C_{1M} & C_{2M} & \dots & C_{NM} \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = \begin{bmatrix} \underline{b} & \underline{C}^T \end{bmatrix} \underline{x}_a = \underline{W}^T \underline{x}_a$$

$$\underline{b} = \begin{bmatrix} \vdots \\ b \\ \vdots \\ \vdots \end{bmatrix}_{N+1}$$

$$\underline{W} = \begin{bmatrix} \underline{b}^T \\ \underline{C}^T \end{bmatrix} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \\ \vdots \end{bmatrix}_M^T$$

For convenience, define a  $\underline{w}$  with a row index that starts from zero  $b^+$

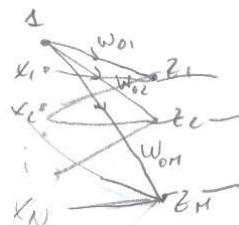
(CL.15)

$$\underline{W} = \begin{bmatrix} w_{01} & w_{02} & \dots & w_{0M} \\ w_{11} & w_{12} & \dots & w_{1M} \\ w_{21} & w_{22} & \dots & w_{2M} \\ \vdots & \ddots & \ddots & \vdots \\ w_{N1} & w_{N2} & \dots & w_{NM} \end{bmatrix} \quad N+1$$

$$\underline{z} = \underline{W}^T \underline{x}_a = \begin{bmatrix} w_{01} & w_{11} & w_{21} & \dots & w_{N1} \\ w_{02} & w_{12} & w_{22} & \dots & w_{N2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ w_{0M} & w_{1M} & w_{2M} & \dots & w_{NM} \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

note ( $x$ )

Therefore  $w_{0j}$  connects  $s$  to  $z_j$   
 $w_{ij}$  connects  $x_i$  to  $z_j$



end of note ( $x$ )

As explained in one of the previous sections,  $\underline{d}^{[n]}$  are distribution vectors, typically one-hot

Lc.16

$$\underline{d} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ for } a_i \quad \underline{d} = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \text{ for } a_2 \dots \dots \quad \underline{d} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \text{ for } a_m$$

We already pointed out that, more generally,  $\underline{d}^{[n]}$  may be a distribution expressing uncertainty about the class in the training set. For example when we are not sure about the label and we have two concuring classes with the same probability, we could use a 2-hot distribution with  $y_2$  in the 2 corresponding positions. More generally, the training set may not be "sharp" and provide only a distribution on some samples.

The problem is formulated as

$$\underline{w}^* = \underset{\underline{w}}{\operatorname{arg\min}} E(\underline{w}) = \underset{\underline{w}}{\operatorname{arg\min}} \frac{1}{M} \sum_{m=1}^{M_2} L(\underline{d}^{[m]}, \underline{y}^{[m]})$$

As usual we seek the solution with an iterative update that requires the computation of the gradient

$$\nabla_{\underline{w}} E(\underline{w}) = \left[ \frac{\partial E}{\partial w_{ij}} \right]_{\substack{i=0, \dots, N \\ j=1, \dots, M}}$$

Note that the gradient is a  $(N+1) \times M$  matrix (gradient flow)  
More specifically

$$\frac{\partial E}{\partial w_{ij}} = \sum_{e=1}^M \frac{\partial E_e}{\partial z_e} \frac{\partial z_e}{\partial w_{ij}} = \frac{\partial E_e}{\partial z_j} \frac{\partial z_j}{\partial w_{ij}} .$$

because only  $z_j$  depends on  $w_{ij}$ . Therefore

$$\frac{\partial E_e}{\partial w_{ij}} = \frac{\partial E_e}{\partial z_j} \times_{ai}$$

(Recall that  $x_{00} = 1$  and corresponds to the connection that takes care of the bias)

(Lec 17)

Now to compute  $\frac{\partial E}{\partial z_j}$ ,  $j = 1, \dots, M$ , we need to define a loss function that measures the discrepancy between  $y$  and  $\underline{d}$ . Any loss function could be used such as those used for regression. The mean square being obviously the sum of squares  $L(\underline{d}, y) = \sum_{j=1}^M (d_j - y_j)^2$ . However since both  $y$  and  $\underline{d}$  are distributions, the most appropriate choice would be the cross-entropy  $L(\underline{d}, y) = \sum_{j=1}^M d_j \log \frac{1}{y_j}$ . For completeness, just as we did for the binary classifier, we consider both the sum of squares and the cross-entropy.

### LINEAR MULTICLASS WITH THE SQUARE LOSS

The cost function to minimize using the square loss is

$$E(W) = \frac{1}{M^2} \sum_{m=1}^{M^2} \sum_{l=1}^M (d_{e[m]} - y_{e[m]})^2$$

As shown in the previous section, we need

$$\frac{\partial E_{\text{sq}}}{\partial z_j}, \quad j = 1, \dots, M \quad \text{which can be written as}$$

$$\frac{\partial E_{\text{sq}}}{\partial z_j} = -\frac{1}{M^2} \sum_{m=1}^{M^2} \sum_{l=1}^M \delta(d_{e[m]} - y_{e[m]}) \frac{\partial y_e}{\partial z_j}$$

The derivative  $\frac{\partial y_e}{\partial z_j}$  is computed across the softmax function and is (see Appendix)

$$\frac{\partial y_e}{\partial z_j} = \delta_{je} y_e - y_e y_j$$

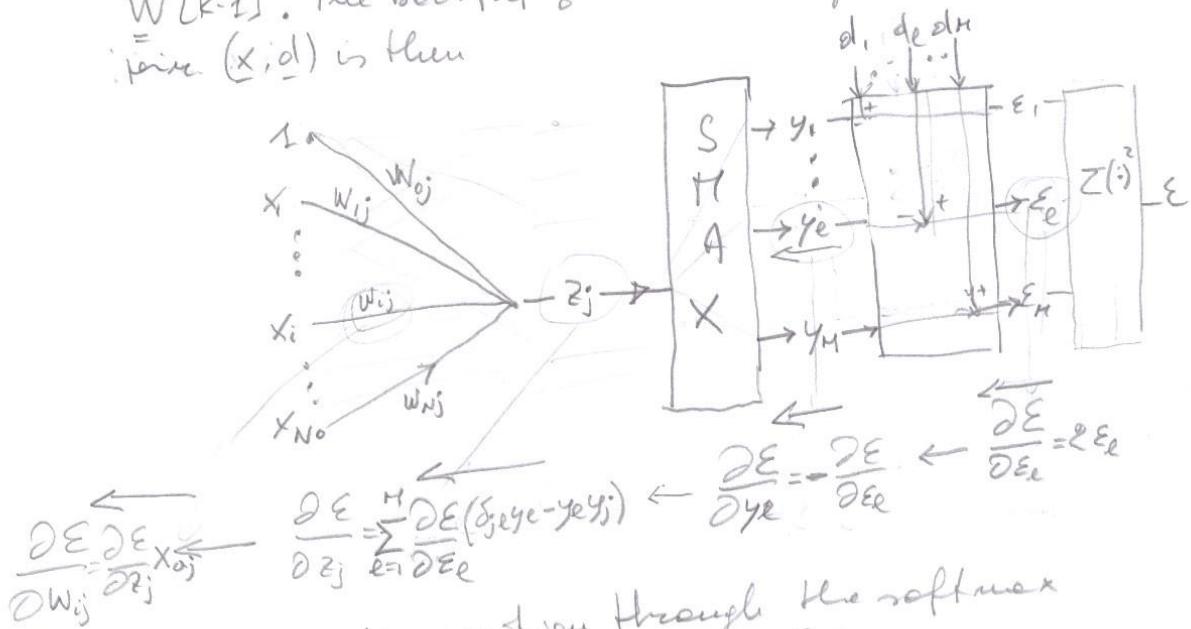
Then the update for the generic weight is

LC.18

$$W_{ij}[k] = W_{ij}[k-1] + \mu \frac{1}{m} \sum_{n=1}^m \sum_{e=1}^E (\delta_{je}^{[n]} - y_e^{[n]} y_j^{[n]}) (S_{je} y_e^{[n]} - y_e^{[n]} y_j^{[n]})$$

where as usual the outputs  $y_i^{[n]}$  are computed using  $W[k-1]$ . The backpropagation on a single example:

pair  $(x, d)$  is then



Note that the backpropagation through the softmax requires on each  $z_j$  all the derivatives  $\frac{\partial E}{\partial y_l}$   $l=1, \dots, M$ .

(Lec. 13)

## LINEAR MULTICLASS WITH CROSS-ENTROPY

The cost function to minimize using the cross-entropy loss is

$$E_{ce}(w) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^n d_{ai}^{[i]} \log \frac{1}{y_{ai}^{[i]}}$$

We need as before

$$\frac{\partial E_{ce}}{\partial z_j}, j=1, \dots, n$$

$$\frac{\partial E_{ce}}{\partial z_j} = \sum_{e=1}^n \frac{\partial E_{ce}}{\partial y_e} \frac{\partial y_e}{\partial z_j}$$

$$\frac{\partial E_{ce}}{\partial y_e} = \frac{\partial}{\partial y_e} \sum_{m=1}^n d_m \log \frac{1}{y_{ai}^{[m]}} = -\frac{\partial}{\partial y_e} \text{d} \log y_e = -\frac{de}{y_e}$$

Using the derivative across the softmax, we

$$\frac{\partial y_e}{\partial z_j} = \delta_{ej} y_e - y_e y_j$$

$$\frac{\partial E_{ce}}{\partial z_j} = \sum_{e=1}^n -\frac{de}{y_e} \cdot (\delta_{ej} y_e - y_e y_j) = -\sum_{e=1}^n de \delta_{ej}$$

$$= -\sum_{e=1}^n de \delta_{ej} + \left[ \sum_{e=1}^n de y_j \right]$$

$= -d_j + y_j$   
the derivative with respect to the generic weight  $w_{ij}$  is

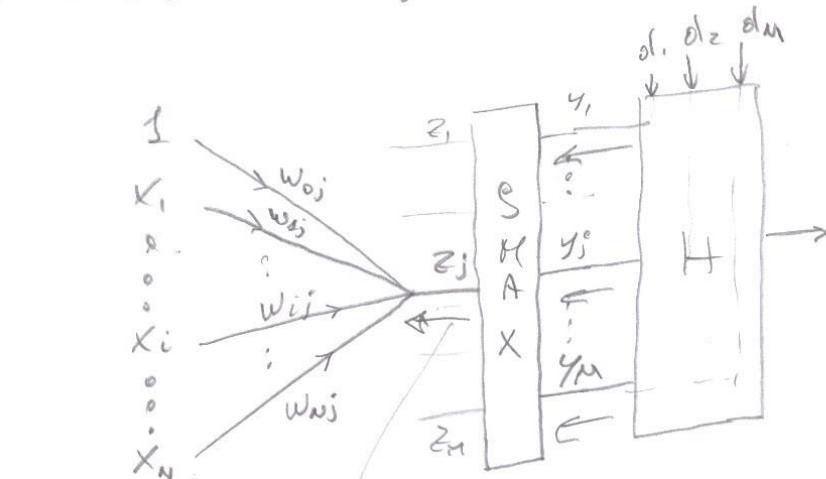
$$\frac{\partial E_{ce}}{\partial w_{ij}} = \frac{\partial E}{\partial z_j} x_{ai} = (y_j - d_j) x_{ai}$$

The gradient update is then

(c.20)

$$w_{ij}[k] = w_{ij}[k-1] + \mu \sum_{u=1}^{n^2} (d_j[u] - y_j[u]) x_{ai}[u] \quad \text{for } i, j$$

where as usual  $y_j[u]$  is computed using  $\hat{w}[k-1]$ .



$$\frac{\partial E_{ac}}{\partial z_j} = y_j - d_j$$

Note that despite the fact that each outputs of the softmax  $y_j$  depend on all the inputs  $x_1, \dots, x_M$ , the derivative using the cross-entropy simplifies greatly the back propagation through both blocks.

To dig deeper into the updating equation, observe as usual that only one of the  $d_j$  is equal to one (we pointed out above that this is generalizable to any distribution on  $d$ ). Say that  $d_9=1$  and the others are zero.

The weights converging on  $z_j$  will be updated according to the following simple rules

$$\left\{ \begin{array}{l} \Delta w_{ij} = \mu (1-y_9)x_{ai} \text{ for } j=9 \\ \Delta w_{ij} = -\mu y_9 x_{ai} \text{ for } j \neq 9 \end{array} \right.$$

(Lc.21)

The variation on  $\theta_j$  caused by this update at a given time step is

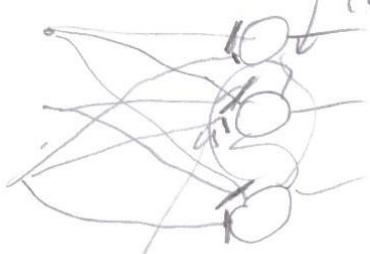
$$\Delta w_{ij} x_i = \begin{cases} \mu(1-y_q)x_i^2 & \text{for } j=q \\ -\mu y_q x_i^2 & \text{for } j \neq q \end{cases}$$

The gradient algorithm pushes up the desired winner  $y_q$  because  $(1-y_q)x_i^2 > 0$ . Conversely the other softmax outputs are pushed down as  $-y_q x_i^2 < 0$  for  $j \neq q$ .

### OBSERVATION (Hebbian learning)

In neurobiology one of the fundamental assumptions about neural adaptation is that synapses become stronger if they succeed in bringing the target neuron to fire. Conversely when they do not succeed they tend to retreat (Hebbian hypothesis). Clearly there is no biological neuron here, but an analogy can be drawn because the competition that happens to the softmax block does resemble lateral inhibition with a winning firing neuron.

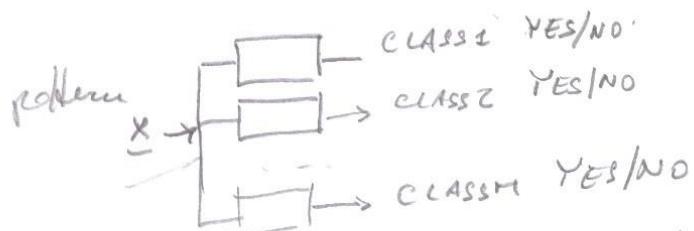
lateral inhibition  
(typically inhibitory)



Hebbian  
excitatory synapses

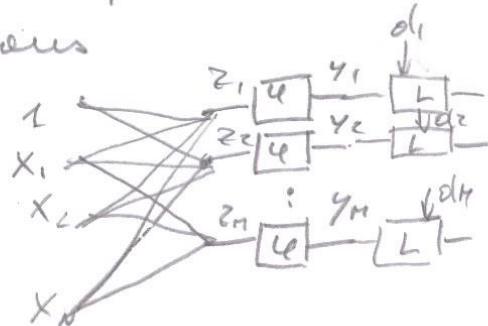
## MULTIPLE BINARY CLASSIFIERS

In discriminating among  $M$  classes instead of using a single multi-class decision maker, we could think of  $M$  binary classifiers that tell us if a pattern belongs to a class or not.



In this architecture there is clearly no "competition" among the classifiers that simply perform a binary decision, or more generally provide a posterior probability (one of the two responses (the other is the complement to one))

A neural architecture with  $M$  binary linear classifier is a parallel of  $M$  perceptrons



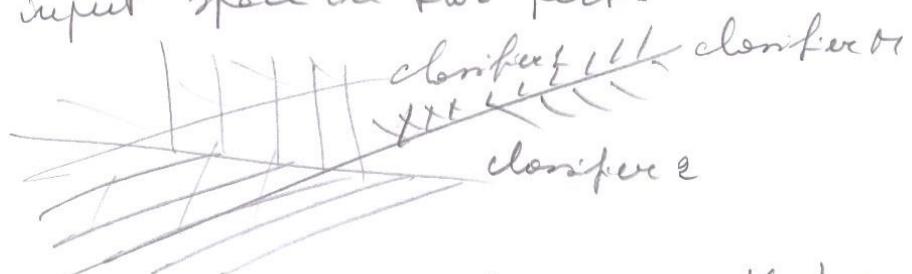
We have already discussed the linear binary classifier where  $\phi$  may be the logistic function and the loss  $L$  may be one of many. The desired outputs can be as usual organized in one-hot vectors

(cc.23)

$$\underline{d} = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \underline{d} = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \quad \dots \quad \underline{d} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix}$$

one use different configurations.

Note that linear classifiers can only divide the input space in two parts



and cannot provide the complex regions that are obtained using the softmax. We will provide an explanation in the following.

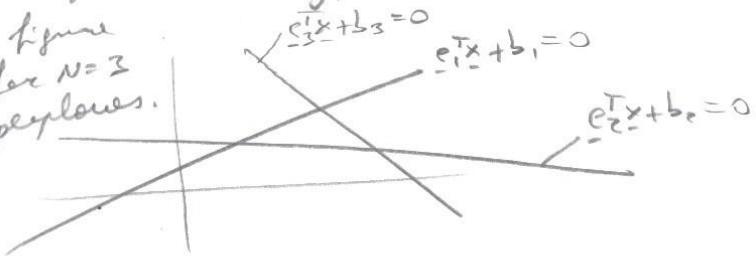
(Le. 24)

## HOW THE LINEAR CLASSIFIERS WORK

In both (x) and (xx) we have M affine transformations

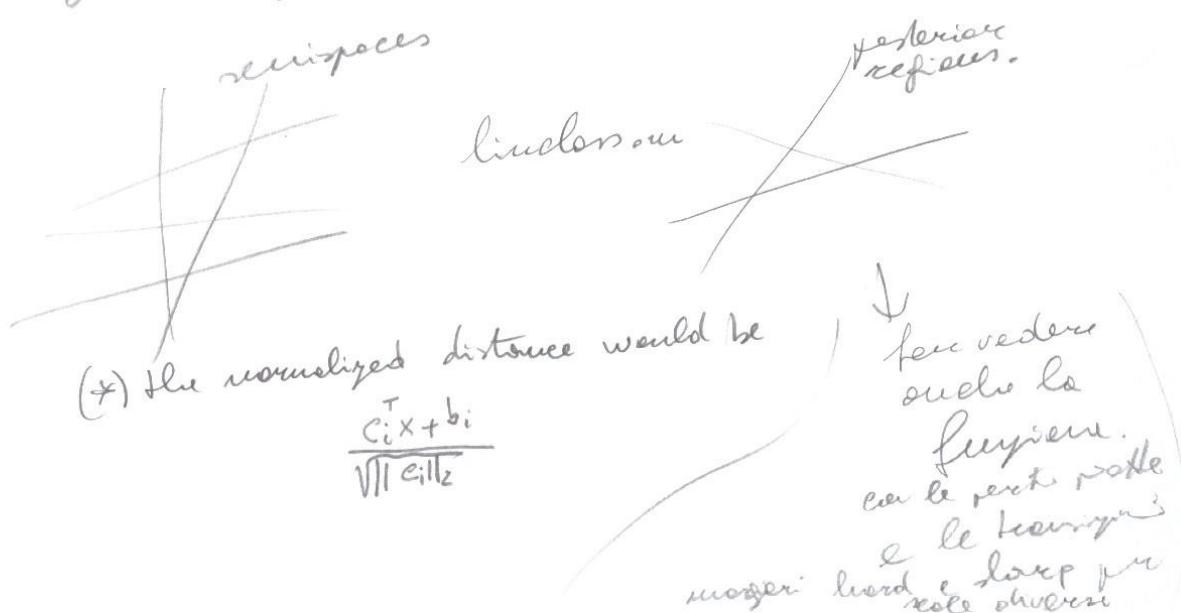
$$z_i = c_i^T x + b_i \quad i = 1, \dots, M$$

that correspond to hyperplanes in the space of  $x$ . The figure shows for  $N=3$  and  $M=3$  hyperplanes.



Now each  $z_i$  is the unnormalized distance from hyperplane  $i$ , with a positive value on one side and negative on the other.

In the set of linear classifiers, the M hyperplanes just divide the space in two semi-spaces, while in the architecture with the softmax they "couple" and can give rise to more complex regions.

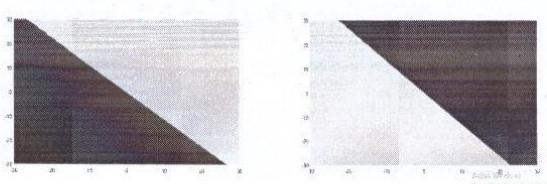
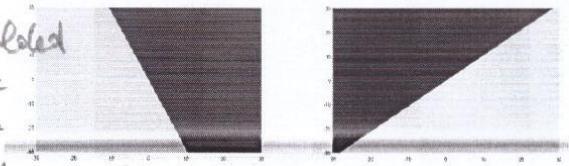


In this following picture we display the regions for  $\text{sigmoid}(\Sigma_i, b_i)$   $i=1..4$  for 4 binary classifiers and for the softmax.

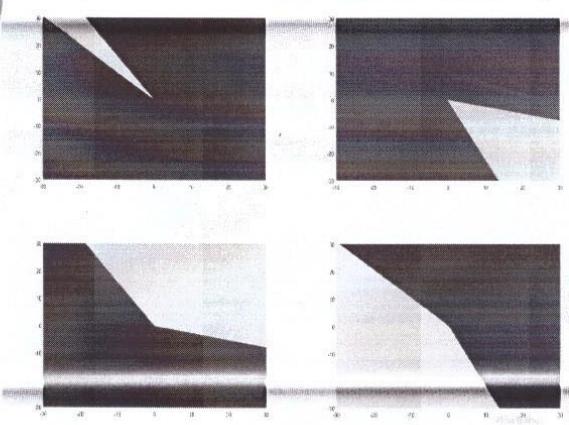
LC.25

(we have thresholded the values for visualization.)

Recall that the smoothness of the transitions is governed by the magnitudes of  $\Sigma_i$ . Large  $\Sigma_i$ 's correspond to sharp transitions (in both cases)



4 binary classifiers



Decision regions from a softmax  
 $M=4$

Note that the boundaries of the softmax regions are not necessarily the hyperplanes seen in the binary decisions. The competition scores from the unnormalized distances from the hyperplanes move the boundaries in different positions.

Therefore the architecture with the softmax function is much more powerful than the parallel of perceptrons in shaping the decision regions. The structure with multiple binary perceptrons is usually inserted in multi-layer architectures.