

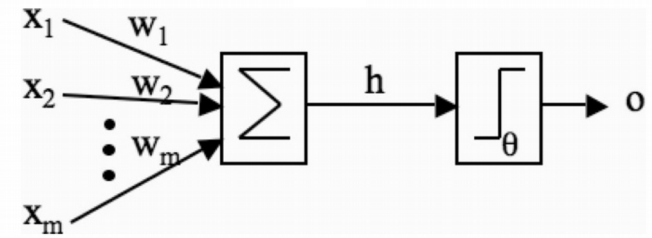
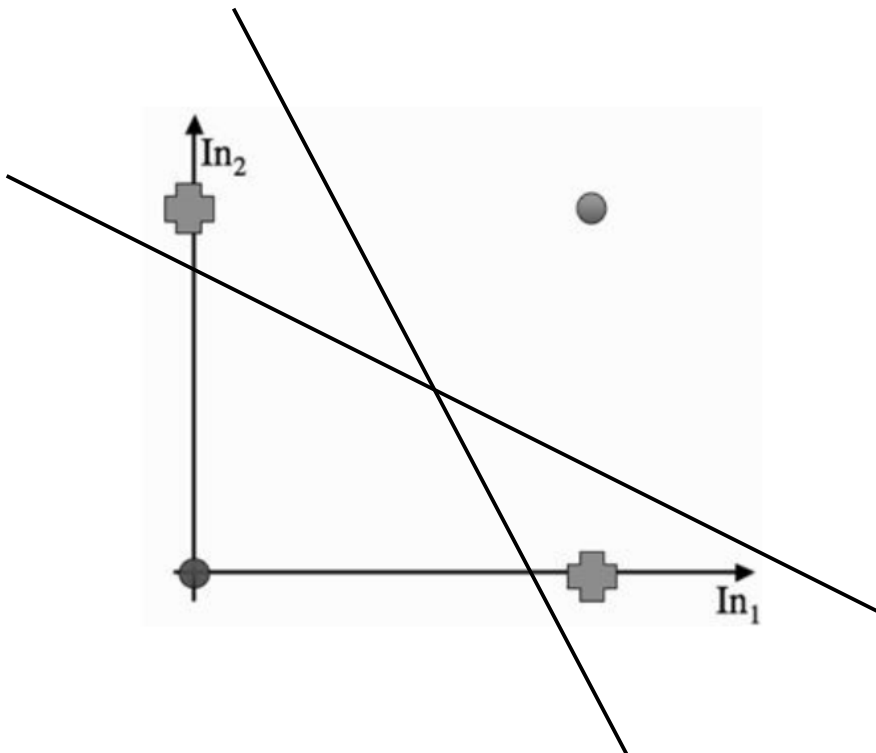


Multi-Layer Perceptrons

- Construct a multi-layer neural network that classifies a given dataset in 2D, overcoming the limitation on learning separability.
- Define an appropriate error to minimise for Feed-forward neural networks.
- Derive the update rule of the weights of the NN, through backpropagation.
- Apply NNs to real-world data sets

Perceptron limitations

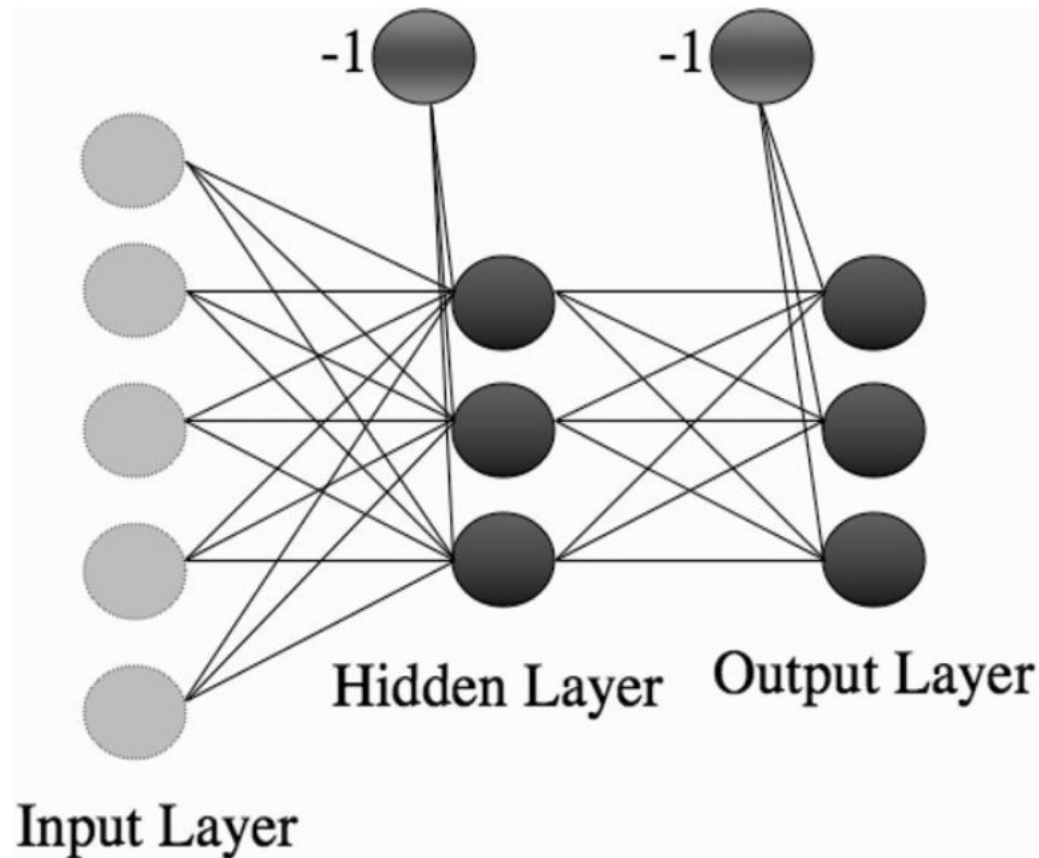
XOR



$$h_w(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 > 0$$

Multi-layer Perceptron

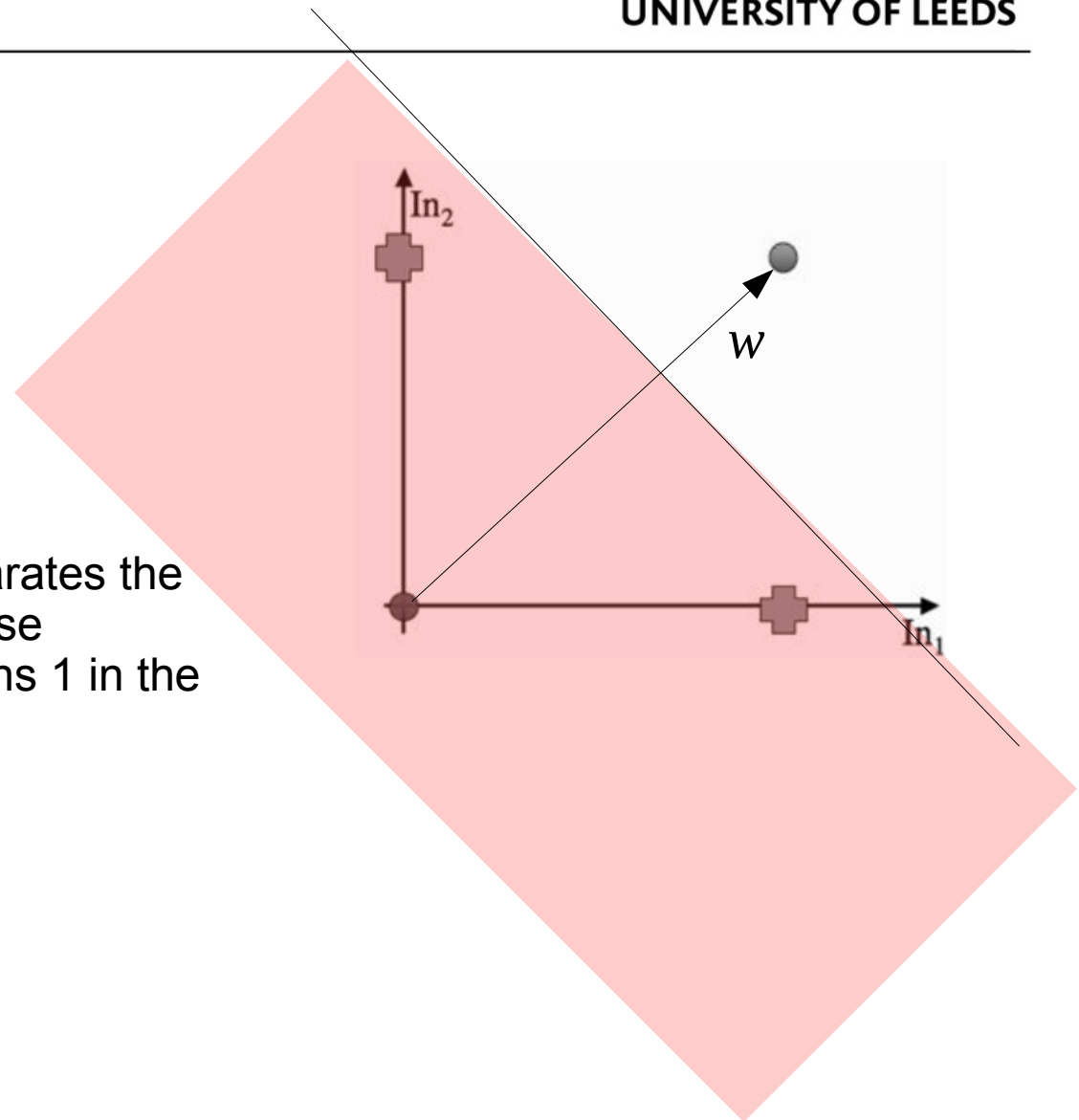
Multi-layer perceptrons are used to obtain non-linear classifiers.



MLP and XOR



UNIVERSITY OF LEEDS



Choose a straight line that separates the points as in the figure, and whose corresponding perceptron returns 1 in the highlighted area

Possible solution:

$$-2x_1 - x_2 + 2.5 = 0$$

$$w = \langle -2, -1 \rangle$$

MLP and XOR

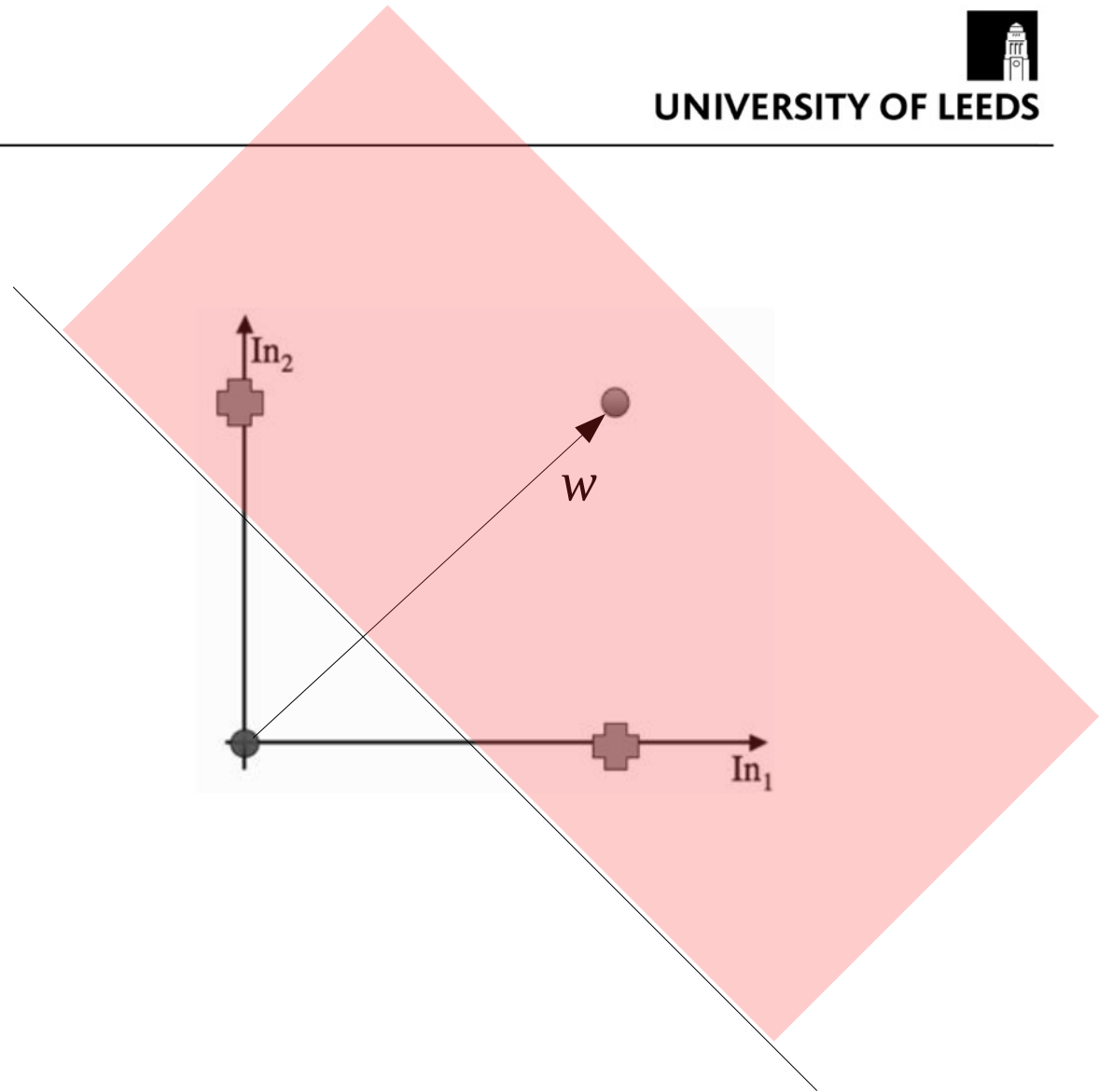


UNIVERSITY OF LEEDS

Now for the other points:

$$x_1 + x_2 - 0.5 = 0$$

$$w = \langle 1, 1 \rangle$$



MLP and XOR

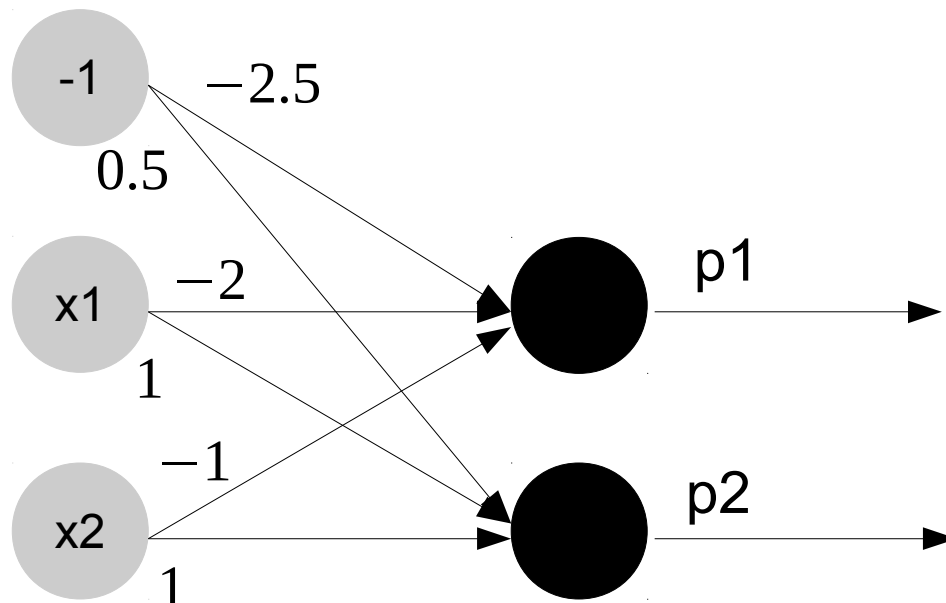
$$-2x_1 - x_2 + 2.5 \geq 0$$

$$x_1 + x_2 - 0.5 \geq 0$$

These are 2
perceptrons with
weights:

$$\langle -2, -1, -2.5 \rangle$$

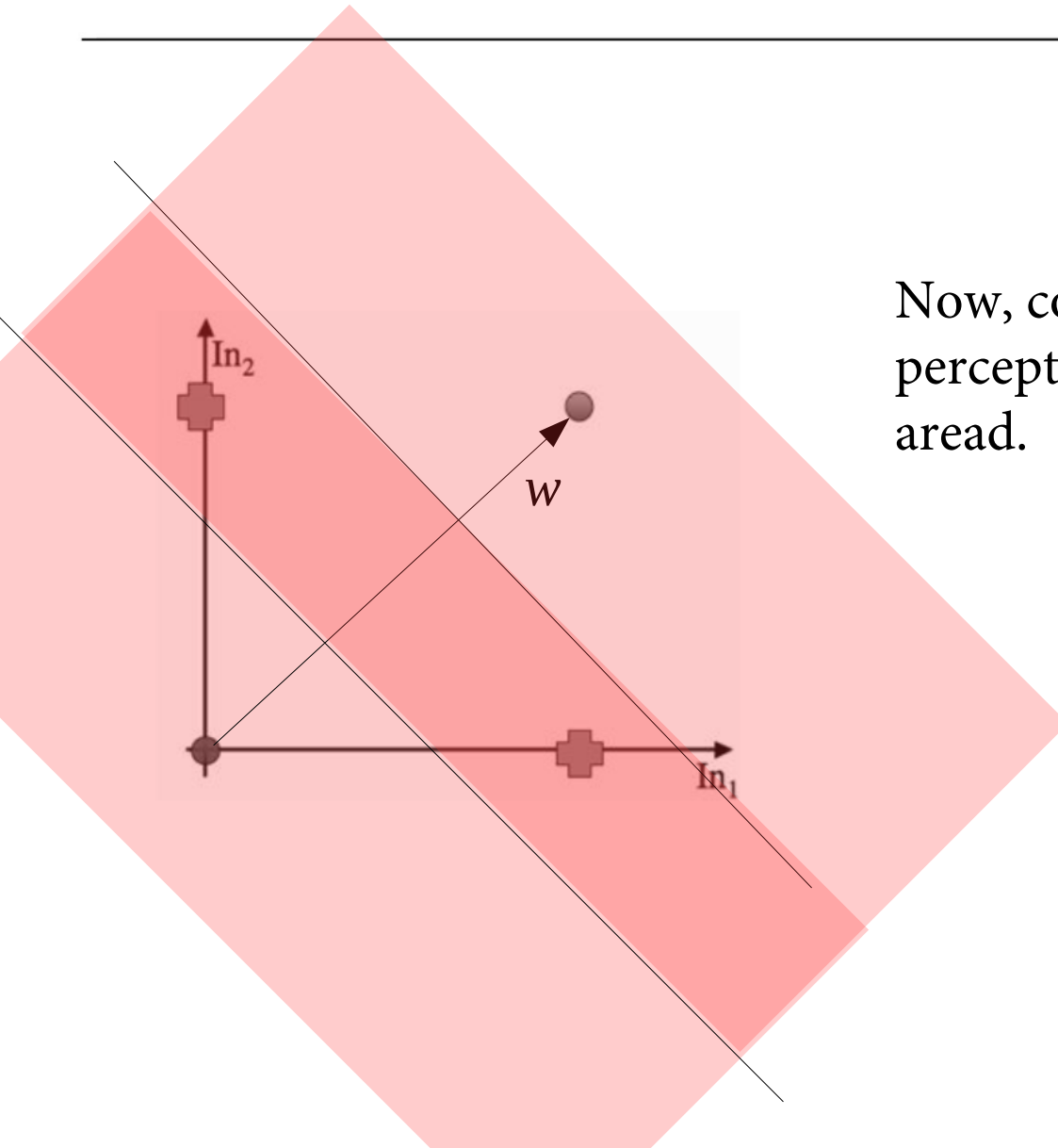
$$\langle 1, 1, 0.5 \rangle$$



MLP and XOR



UNIVERSITY OF LEEDS



Now, combine the outputs of those perceptrons by intersecting the shaded aread.

MLP and XOR



UNIVERSITY OF LEEDS

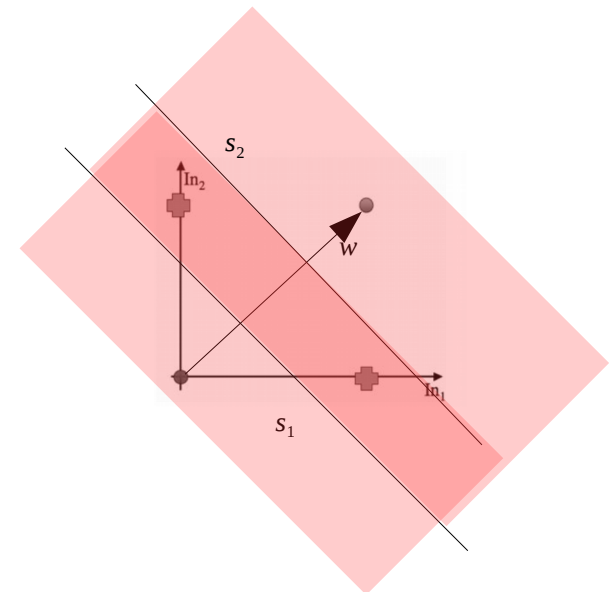
Their outputs are:

x1	x2	p1	p2	o
0	0	1	0	0
0	1	1	1	1
1	0	1	1	1
1	1	0	1	0

What we want

$$p_1 = -2x_1 - x_2 + 2.5 \geq 0$$

$$p_2 = x_1 + x_2 - 0.5 \geq 0$$

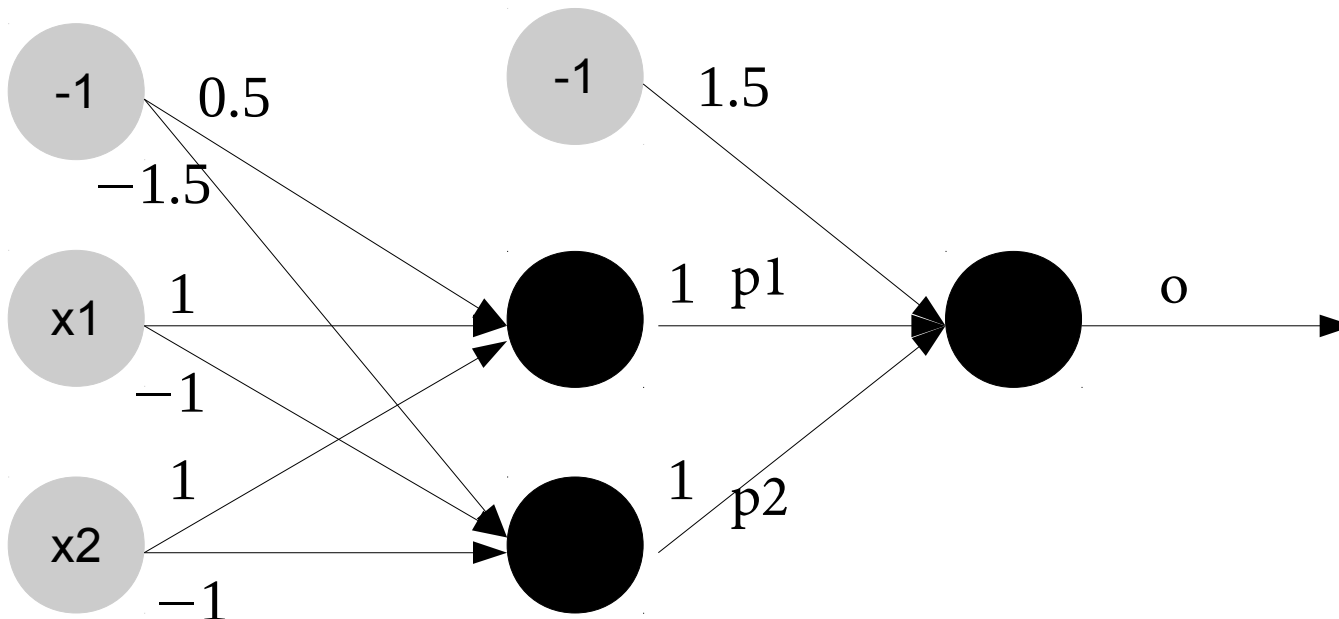


MLP and XOR

$$p_1 \equiv x_1 + x_2 - 0.5 \geq 0$$

$$p_2 \equiv -x_1 - x_2 + 1.5 \geq 0$$

$$o \equiv p_1 + p_2 - 1.5 \geq 0$$



MLP: A Universal Approximator

$$g(x) = \sum_j^N w_j \sigma(y_j^T x + \theta_j) \quad \text{given} \quad f(x) \quad \epsilon > 0$$
$$|g(x) - f(x)| < \epsilon$$

MLP is a universal function approximator, that is, it can represent any function. From a theoretical point of view, this can be done with a single hidden layer.

In practice, that hidden layer would have to be incredibly large. Recent development in neural networks (often referred to as “deep” learning), favour “deep” structures, where many layers can be used to create an hierarchical classification.

Error definition



UNIVERSITY OF LEEDS

$$E(\mathbf{X}) = \sum_{\mathbf{x}_n \in \mathbf{X}} |y_n - t_n|$$

Number of errors on the training set

$$E_p(\mathbf{X}) = \sum_{\mathbf{x}_n \in \mathbf{X}} \mathbf{w}^T \mathbf{x}_n (y_n - t_n)$$

The Perceptron error

$$E_m(\mathbf{X}) = \frac{1}{2} \sum_{\mathbf{x}_n \in \mathbf{X}} (y_n - t_n)^2$$

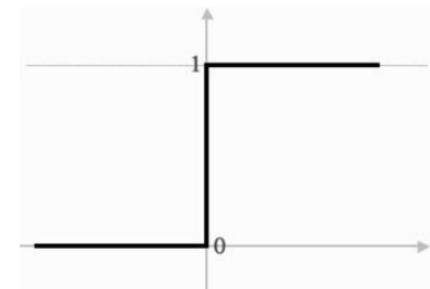
Squared error function (differentiable!)
Usually known as the Mean Squared Error (MSE)

$$y = h\left(\sum_{i=1}^M w_i x_i\right)$$

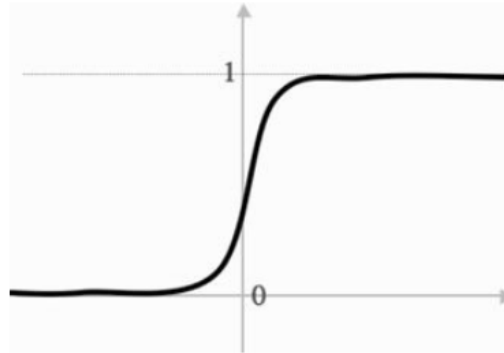
Output is differentiable if h is

$h =$

Not good



A different activation function



The sigmoid function: $h(x) = \frac{1}{1 + e^{-\beta x}} \equiv \sigma_{\beta}$

$$\sigma_{\beta}'(x) = ?$$

The derivative of the sigmoid



UNIVERSITY OF LEEDS

The sigmoid function: $h(x) = \frac{1}{1 + e^{-\beta x}} \equiv \sigma_{\beta}$

$$\sigma_{\beta}'(x) = ?$$

Two useful properties of derivatives:

$$f(x) = e^x \quad f'(x) = e^x$$

Example: $(e^{x^2})' = e^{x^2} \cdot 2x$

Chain rule: $(f \circ g)'(x) = f'(g(x)) \cdot g'(x)$

Hint: $\frac{1}{1 + e^{-\beta x}} = (1 + e^{-\beta x})^{-1}$

The derivative of the sigmoid



UNIVERSITY OF LEEDS

The sigmoid function: $h(x) = \frac{1}{1+e^{-\beta x}} \equiv \sigma_\beta$

We derive the most external function first

$$\sigma_\beta'(x) = \left((1+e^{-\beta x})^{-1} \right)' = -1(1+e^{-\beta x})^{-2} \cdot (1+e^{-\beta x})'$$

Then this

$$= -1(1+e^{-\beta x})^{-2} \cdot e^{-\beta x} \cdot (-\beta x)' = -1(1+e^{-\beta x})^{-2} \cdot e^{-\beta x} \cdot (-\beta)$$

and finally this one

$$\sigma_\beta'(x) = -1(1+e^{-\beta x})^{-2} \cdot e^{-\beta x} \cdot (-\beta) = \frac{\beta e^{-\beta x}}{(1+e^{-\beta x})^2}$$

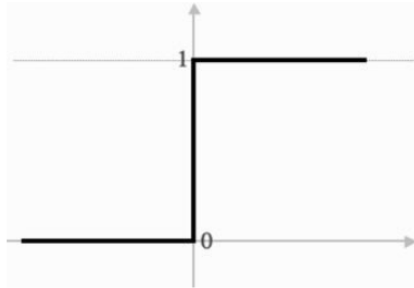
Let's note that:

$$1 - \sigma_\beta = 1 - \frac{1}{1+e^{-\beta x}} = \frac{1+e^{-\beta x} - 1}{1+e^{-\beta x}} = \frac{e^{-\beta x}}{1+e^{-\beta x}} \Rightarrow \sigma_\beta' = \beta \sigma_\beta (1 - \sigma_\beta)$$

A different activation function

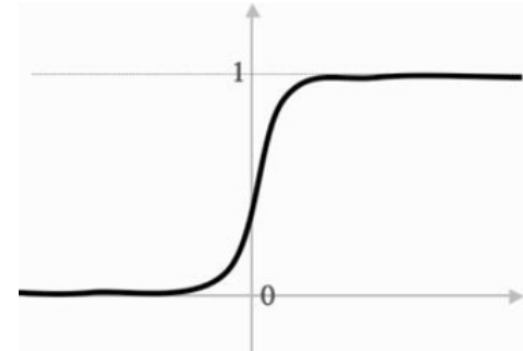


UNIVERSITY OF LEEDS



before

$$h(\mathbf{w}^T \mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} > 0 \\ 0 & \text{if } \mathbf{w}^T \mathbf{x} \leq 0 \end{cases}$$

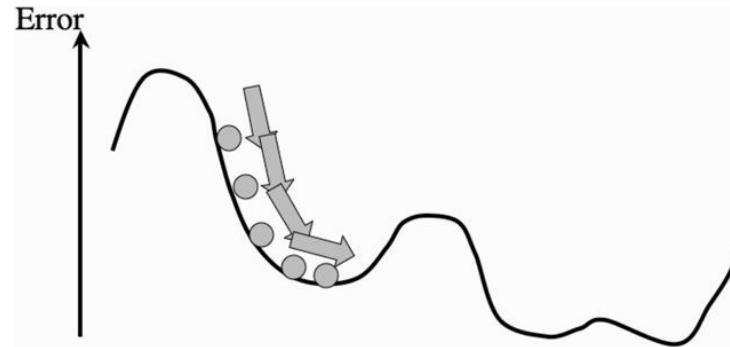


after

$$h(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-\beta \mathbf{w}^T \mathbf{x}}}$$

$$\sigma_{\beta}'(x) = \beta \frac{e^{-\beta x}}{(1 + e^{-\beta x})^2} = \beta \sigma_{\beta}(x)(1 - \sigma_{\beta}(x))$$

Gradient descent (again)



$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla E(\mathbf{x})$$

Perceptron

$$E_p(\mathbf{X}) = \sum_{\mathbf{x}_n \in \mathbf{X}} \mathbf{w}^t \mathbf{x}_n (y_n - t_n)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta (y - t) \mathbf{x}$$

Multi-Layer P

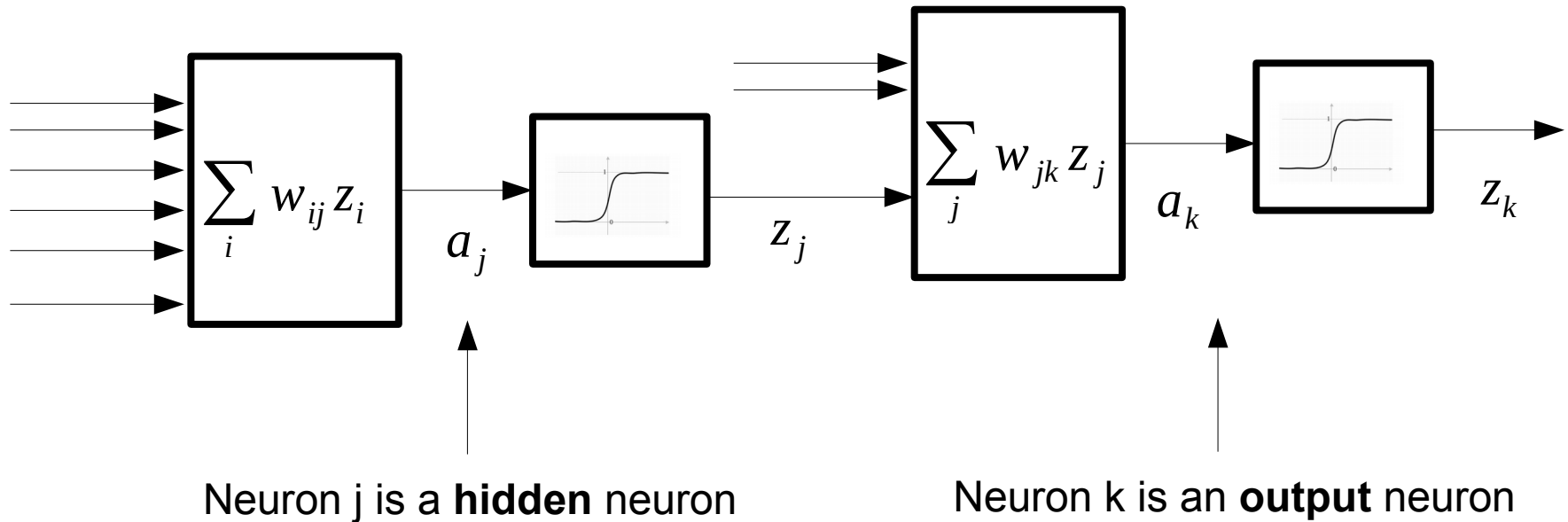
$$E_m(\mathbf{X}) = \frac{1}{2} \sum_{\mathbf{x}_n \in \mathbf{X}} (y_n - t_n)^2$$

?

Backpropagation of errors, notation



UNIVERSITY OF LEEDS

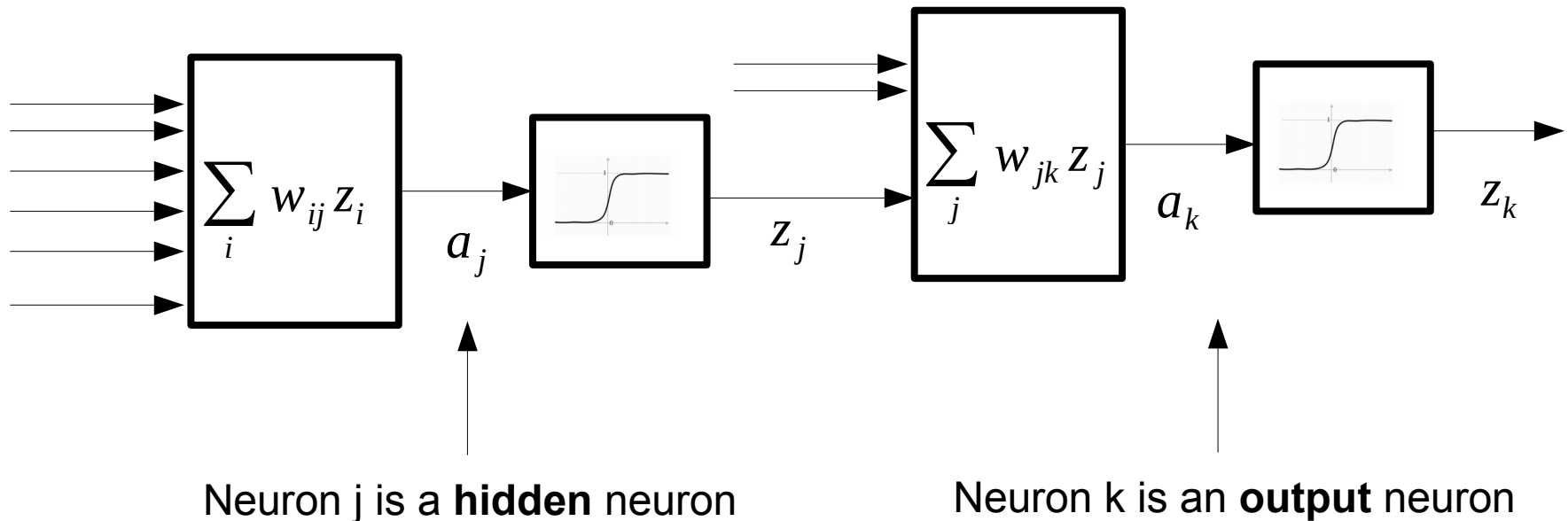


$$\dots \quad a_j = \sum_{i=1}^N w_{ij} z_i \quad z_j = h(a_j) \quad a_k = \sum_{j=1}^M w_{jk} z_j \quad z_k = h(a_k)$$

Forward pass



UNIVERSITY OF LEEDS



$$\dots \quad a_j = \sum_{i=1}^N w_{ij} z_i \quad z_j = h(a_j) \quad a_k = \sum_{j=1}^M w_{jk} z_j \quad z_k = h(a_k)$$

Forward pass: compute all the z



Backward pass, output neuron



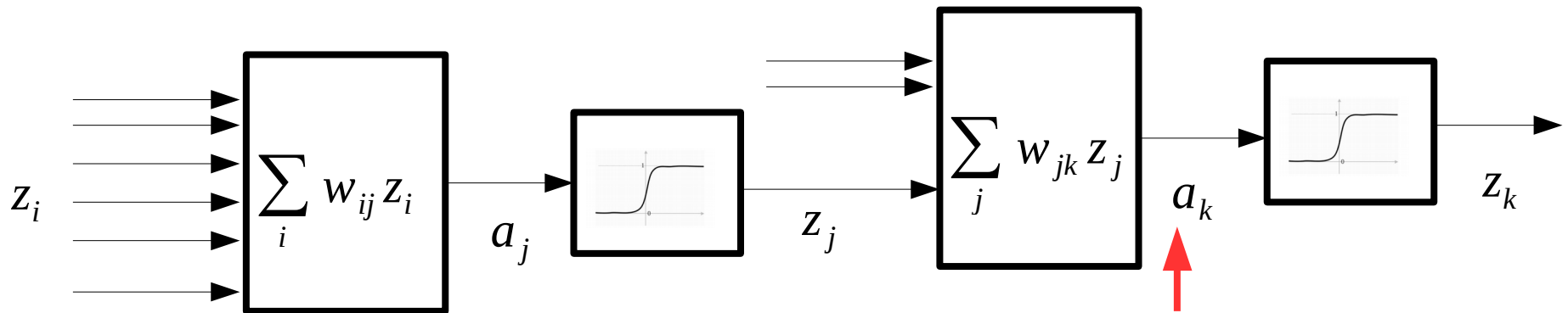
UNIVERSITY OF LEEDS

How does a_k affect the error?

$$E(\mathbf{x}) = \frac{1}{2} (y - t)^2 = \frac{1}{2} \sum_k (z_k - t_k)^2$$

$$\frac{\partial E}{\partial a_k} = \frac{\partial}{\partial a_k} \frac{1}{2} (z_k - t_k)^2 = \frac{\partial}{\partial a_k} \frac{1}{2} (\sigma(a_k) - t_k)^2 = (\sigma(a_k) - t_k) \sigma(a_k) (1 - \sigma(a_k))$$

Now this useful, because it cancels out the exponent in the derivation



We call this derivative δ : $\delta_k = (z_k - t_k) z_k (1 - z_k)$

Backward pass



Backward pass, output neuron



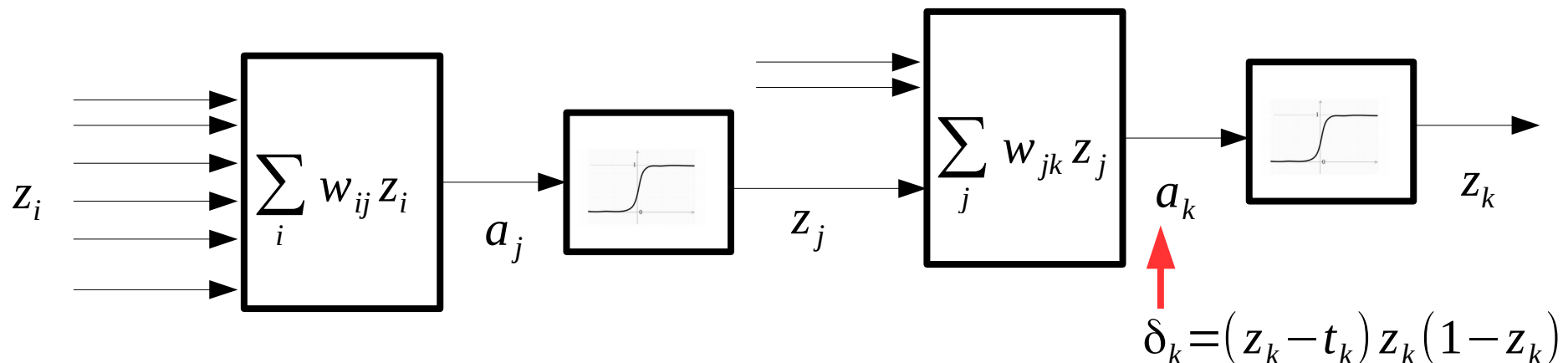
UNIVERSITY OF LEEDS

One step backward, inside the box: how does w_{jk} affect the error?

$$a_k = \sum_j w_{jk} z_j$$

We apply the chain rule again: $\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial w_{jk}} = \delta_k \cdot ?$

$$\frac{\partial a_k}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} w_{0k} z_0 + w_{1k} z_1 + w_{2k} z_2 + \dots + w_{jk} z_j = ?$$



Backward pass, output neuron

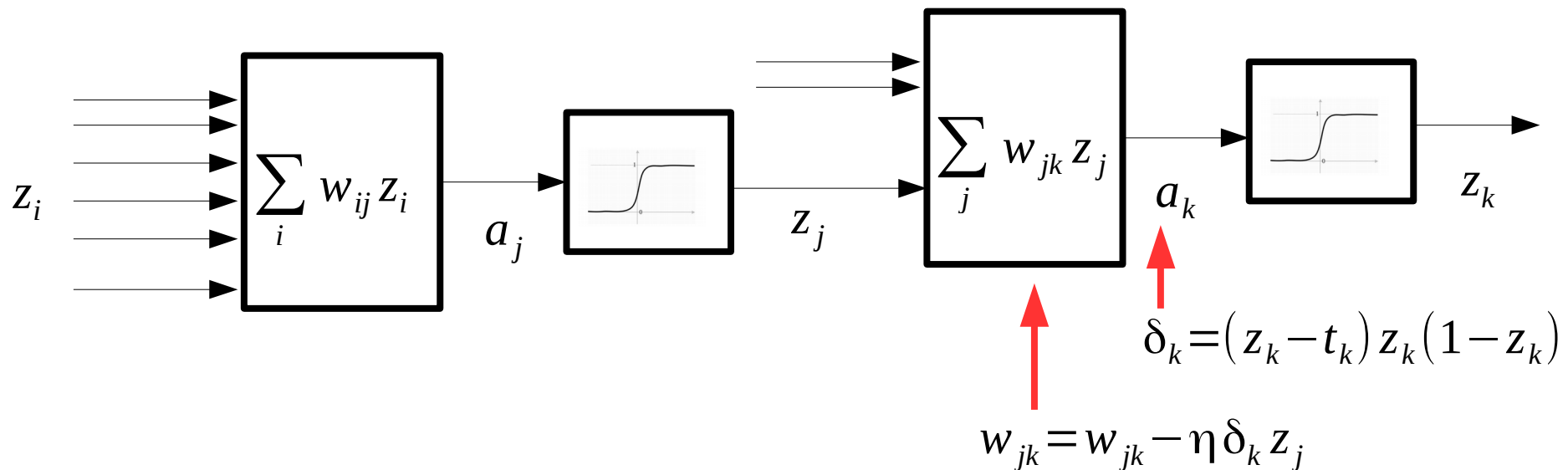


UNIVERSITY OF LEEDS

One step backward, inside the box: how does w_{jk} affect the error?

We apply the chain rule again:
$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial w_{jk}} = \delta_k z_j$$

$$\frac{\partial a_k}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}} w_{0k} z_0 + w_{1k} z_1 + w_{2k} z_2 + \dots + w_{jk} z_j = z_j$$



Backward pass, hidden neuron

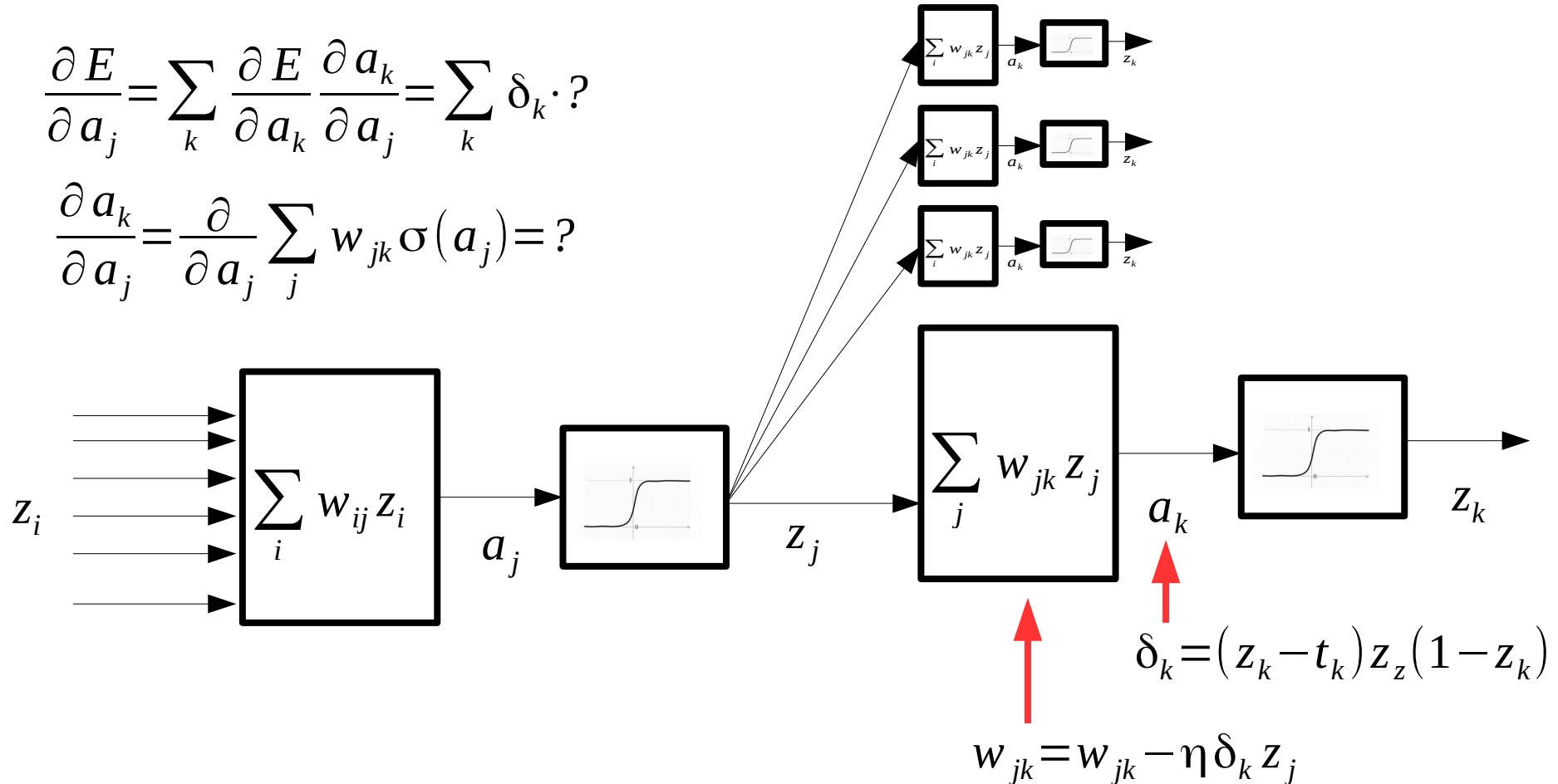


UNIVERSITY OF LEEDS

One step backward: how does a_j affect the error?

$$\frac{\partial E}{\partial a_j} = \sum_k \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_k \delta_k \cdot ?$$

$$\frac{\partial a_k}{\partial a_j} = \frac{\partial}{\partial a_j} \sum_j w_{jk} \sigma(a_j) = ?$$



Backward pass, hidden neuron

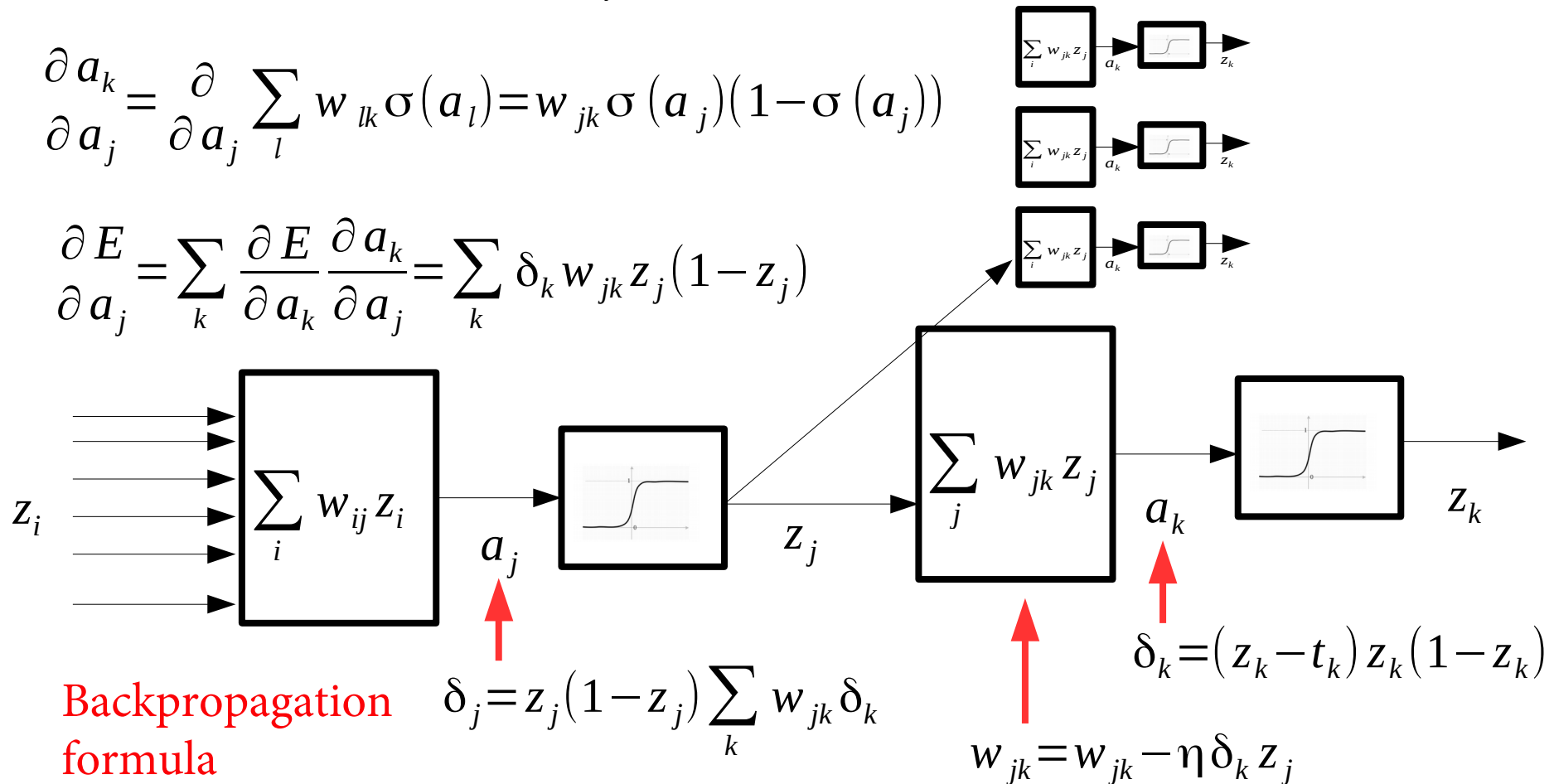


UNIVERSITY OF LEEDS

One step backward: how does a_j affect the error?

$$\frac{\partial a_k}{\partial a_j} = \frac{\partial}{\partial a_j} \sum_l w_{lk} \sigma(a_l) = w_{jk} \sigma(a_j) (1 - \sigma(a_j))$$

$$\frac{\partial E}{\partial a_j} = \sum_k \frac{\partial E}{\partial a_k} \frac{\partial a_k}{\partial a_j} = \sum_k \delta_k w_{jk} z_j (1 - z_j)$$



Computing delta

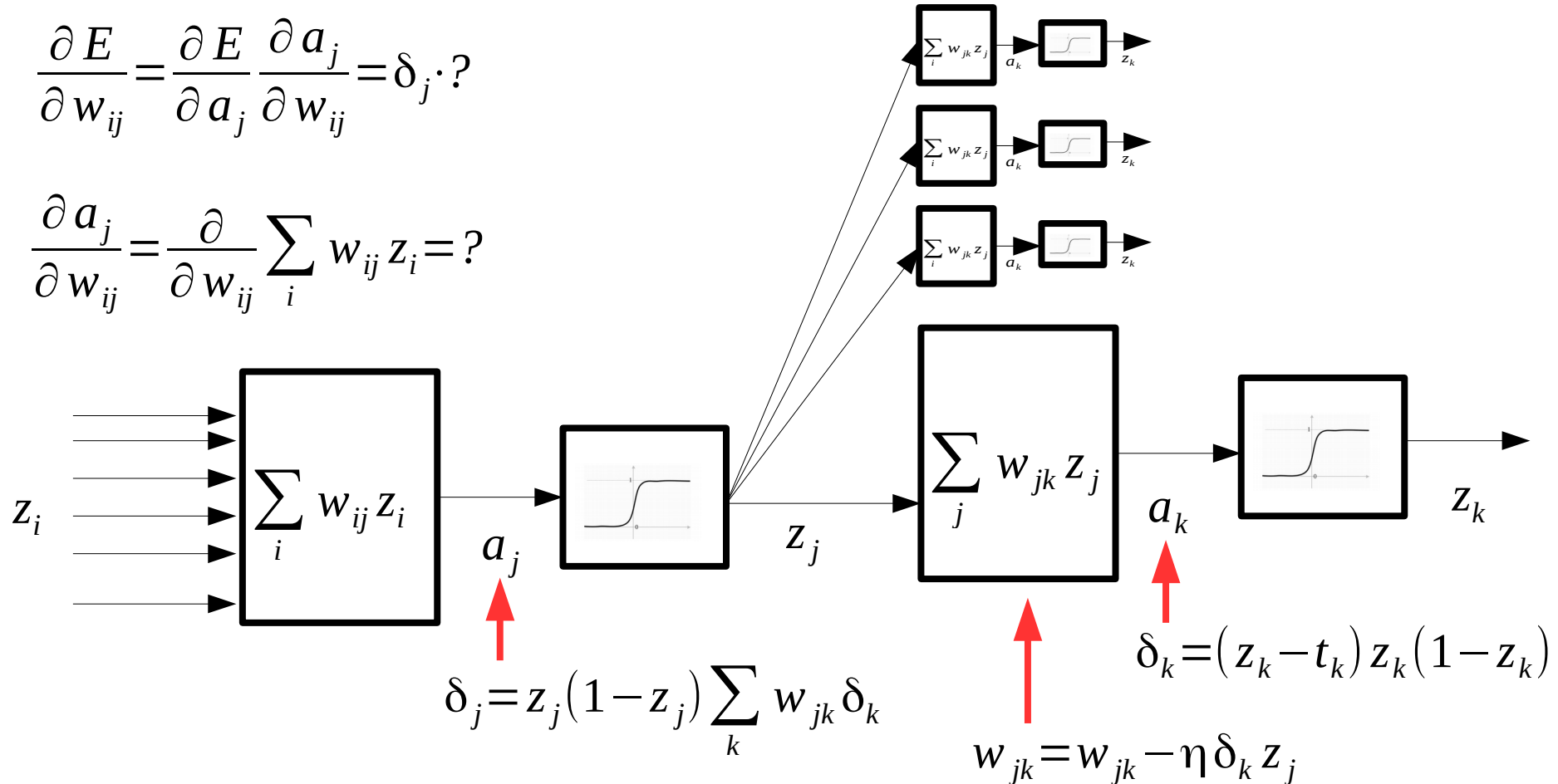


UNIVERSITY OF LEEDS

One step backward, inside the box: how does w_{ij} affect the error?

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}} = \delta_j \cdot ?$$

$$\frac{\partial a_j}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_i w_{ij} z_i = ?$$



Computing delta

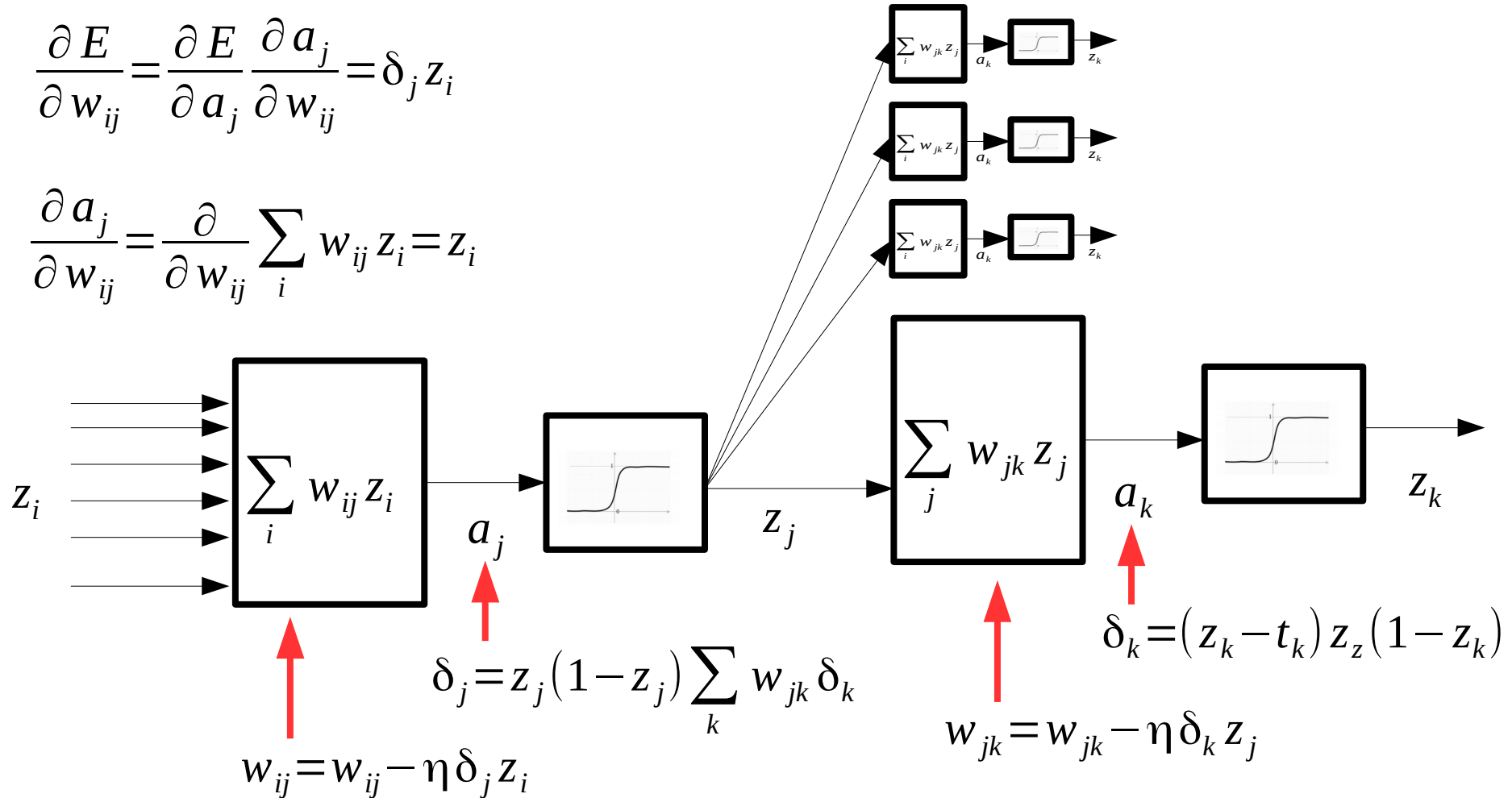


UNIVERSITY OF LEEDS

One step backward, inside the box: how does w_{ij} affect the error?

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial a_j} \frac{\partial a_j}{\partial w_{ij}} = \delta_j z_i$$

$$\frac{\partial a_j}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \sum_i w_{ij} z_i = z_i$$



$$\delta_k = (z_k - t_k) z_k (1 - z_k)$$

$$w_{jk} = w_{jk} - \eta \delta_k z_j$$

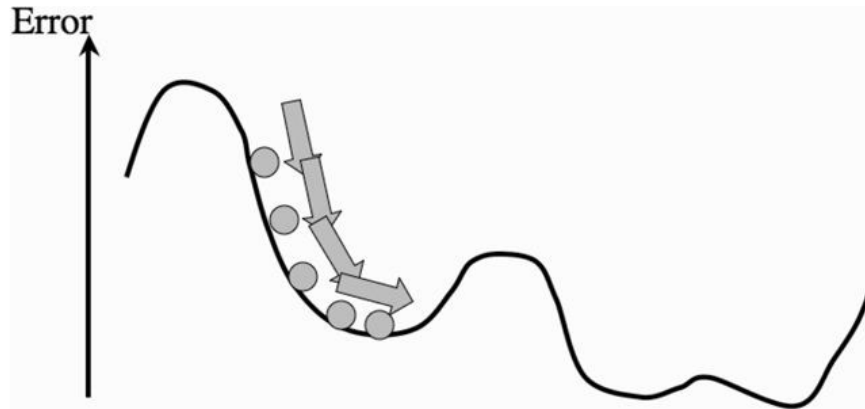
$$\delta_j = z_j (1 - z_j) \sum_k w_{jk} \delta_k$$

$$w_{ij} = w_{ij} - \eta \delta_j z_i$$

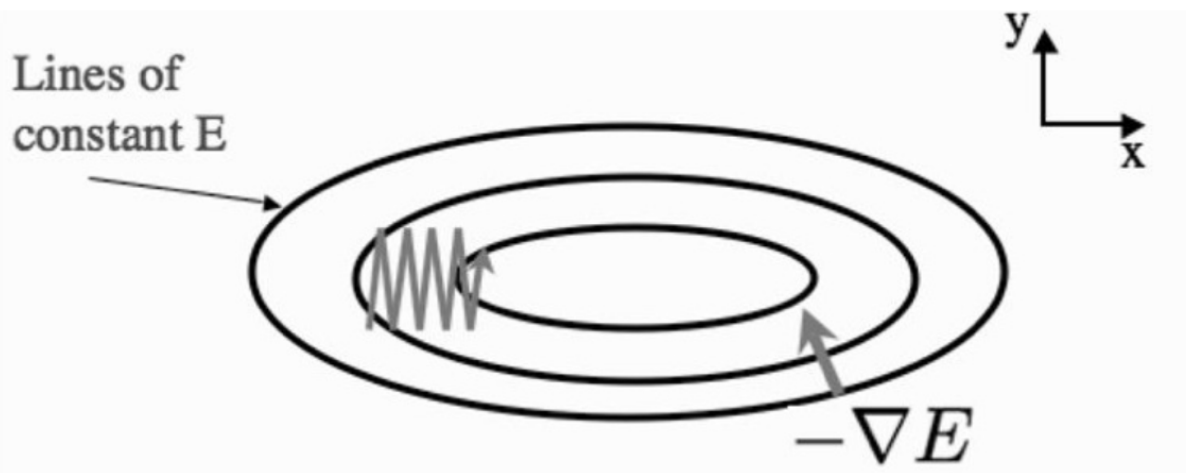
Local Minima



UNIVERSITY OF LEEDS



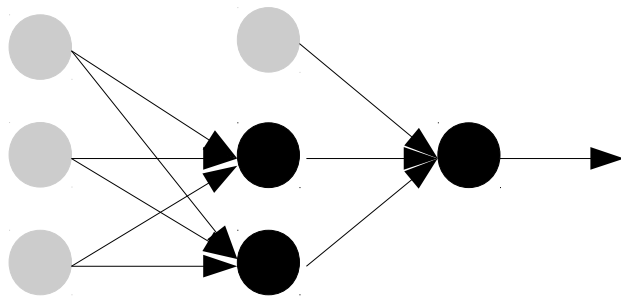
Start with weights close to 0: where the decision is actually made



Multiple random restarts

Using MLPs

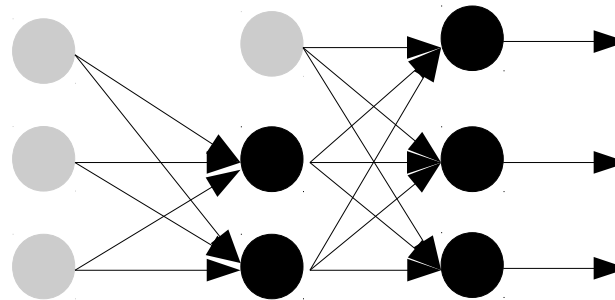
Regression



$$h(a) = a$$

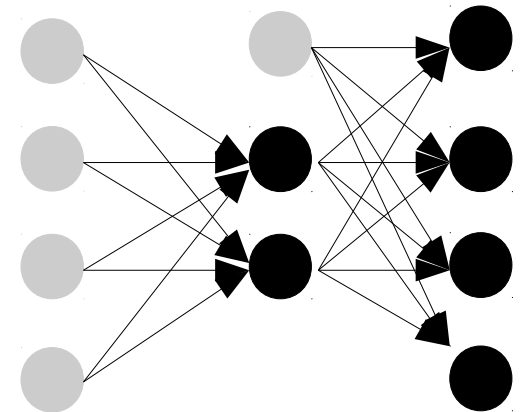
Last neuron
linear

Classification



One output
per class, pick
highest

Compression



Middle
“bottleneck”
layer

Training “recipe”



UNIVERSITY OF LEEDS

Choose features

Normalize
(rescale) data:

$$x' = \frac{x - \bar{x}}{\sigma}$$

or

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

Create training,
validation, and
test sets

Zero
mean, unit
variance

in [0,1]

Decide whether you need hidden layers and how big.
Try several ones.

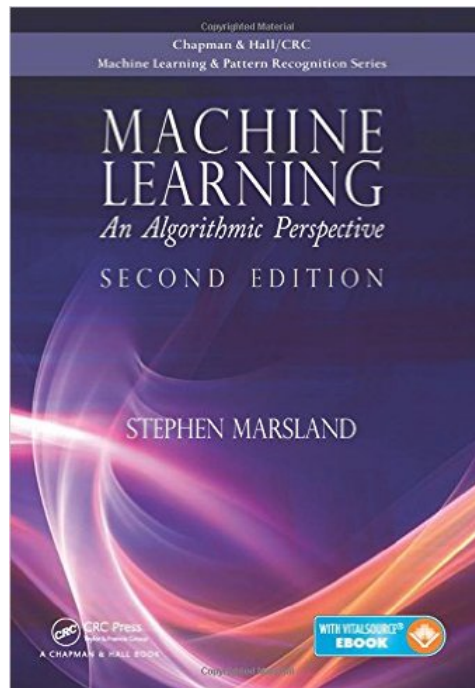
Train

Test



Conclusion

- Construct a multi-layer neural network that classifies a given dataset in 2D, overcoming the limitation on learning separability.
- Define an appropriate error to minimise for Feed-forward neural networks.
- Derive the update rule of the weights of the NN, through backpropagation.
- Apply NNs to real-world data sets



Chapter 4