**College of**
**Engineering**
**& Applied**
**Science**

UNIVERSITY OF
Cincinnati

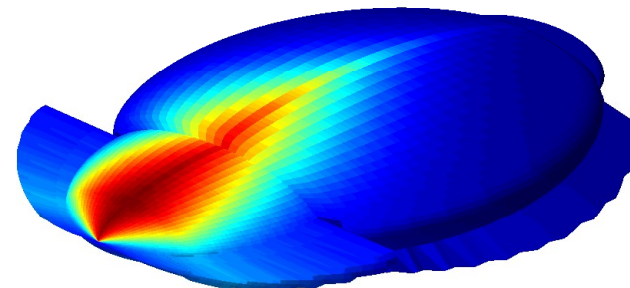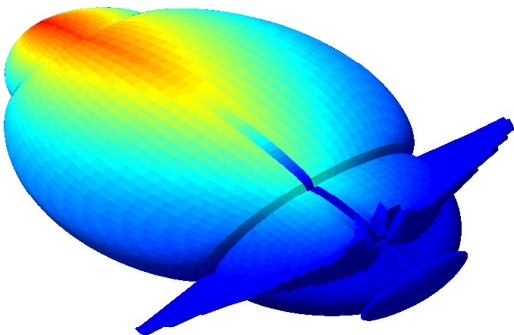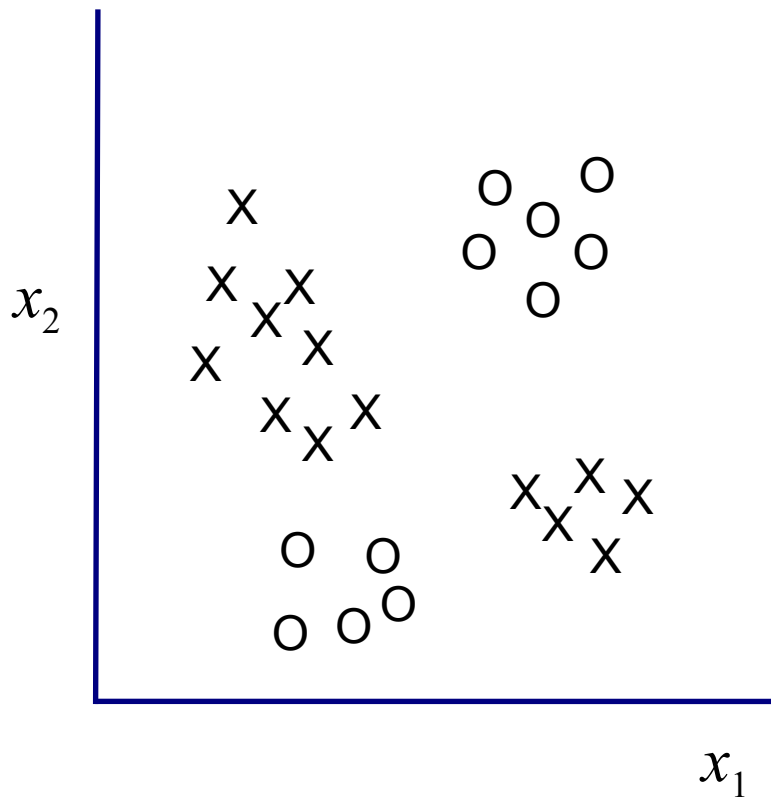Complex Adaptive
Systems Laboratory

# Lecture 7
# Multi-Layer Networks
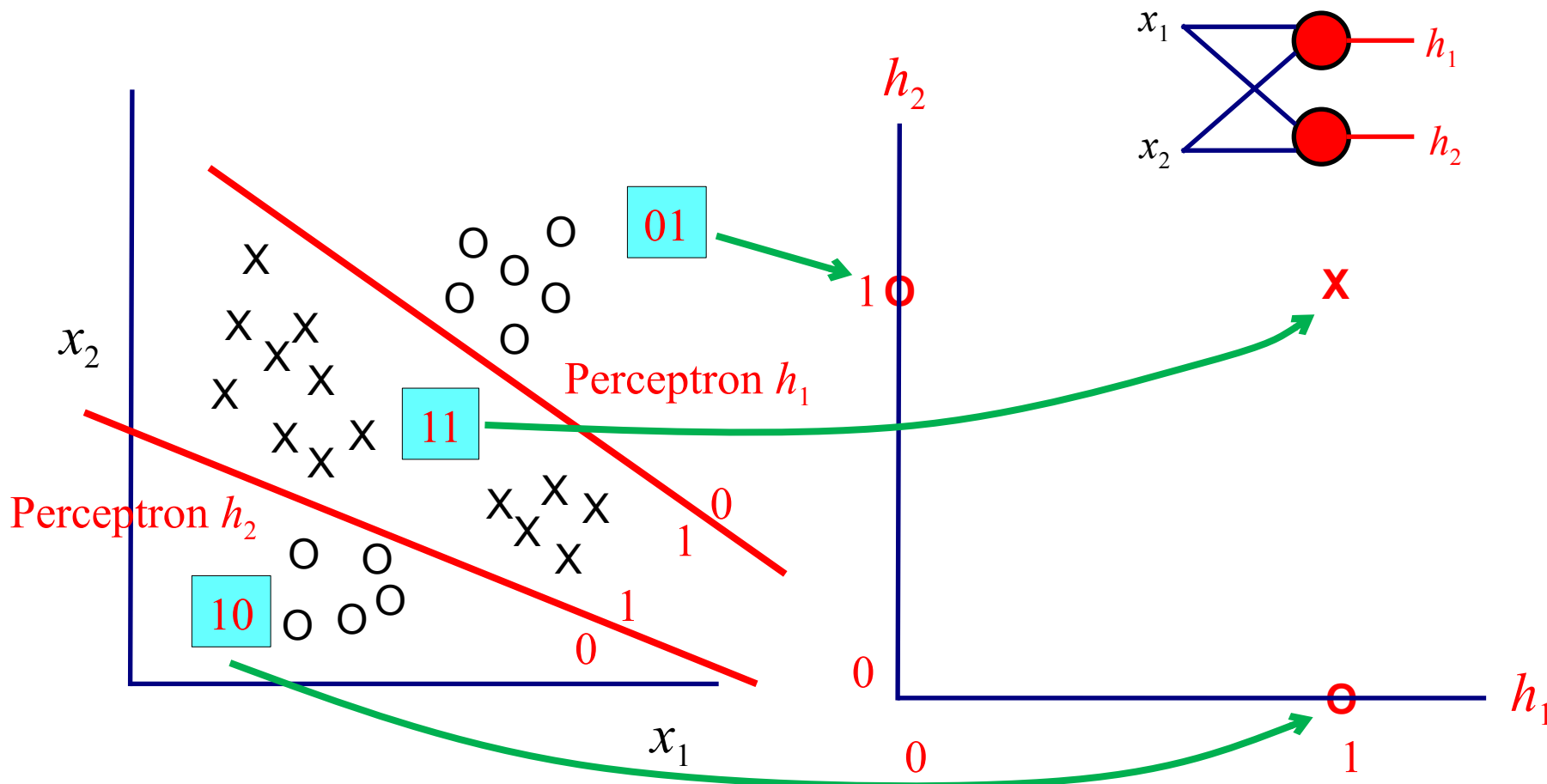# & Backpropagation

# Multilayer Perceptrons

Single perceptrons can only learn linear decision boundaries, so these classes cannot be separated.

But two perceptrons can draw two linear decision boundaries to re-map the data to a new $(p_1, p_2)$ space



$x_1$

$x_2$

$h_1$

$h_2$

$h_2$

01

1

X

Perceptron $h_1$

$x_2$

11

0
1

Perceptron $h_2$

10

1
0

0

$x_1$

$h_1$

0

1

But how can the weights for the network be found?

…. and a third perceptron in the next layer can do the classification.



$x_1$ $h_1$ $y$
$x_2$ $h_2$

$h_2$

01

1

Perceptron $h_1$

$x_2$

11

**X**

1

0

Perceptron $h_2$

10

0
1

1
0

0

Perceptron $y$

0

$x_1$

0

1

$h_1$

But how can the weights for the network be found?

# The Backpropagation Algorithm

- Based on gradient descent.

- Allows training of any feedforward network

- Can be used to train recurrent networks with minor modification.

  - "Discovered" in various forms by:
        Bryson and Ho        1969
        Werbos               1974
        Parker               1985

    The familiar neural network formulation by Rumelhart, Hinton & Williams, 1985
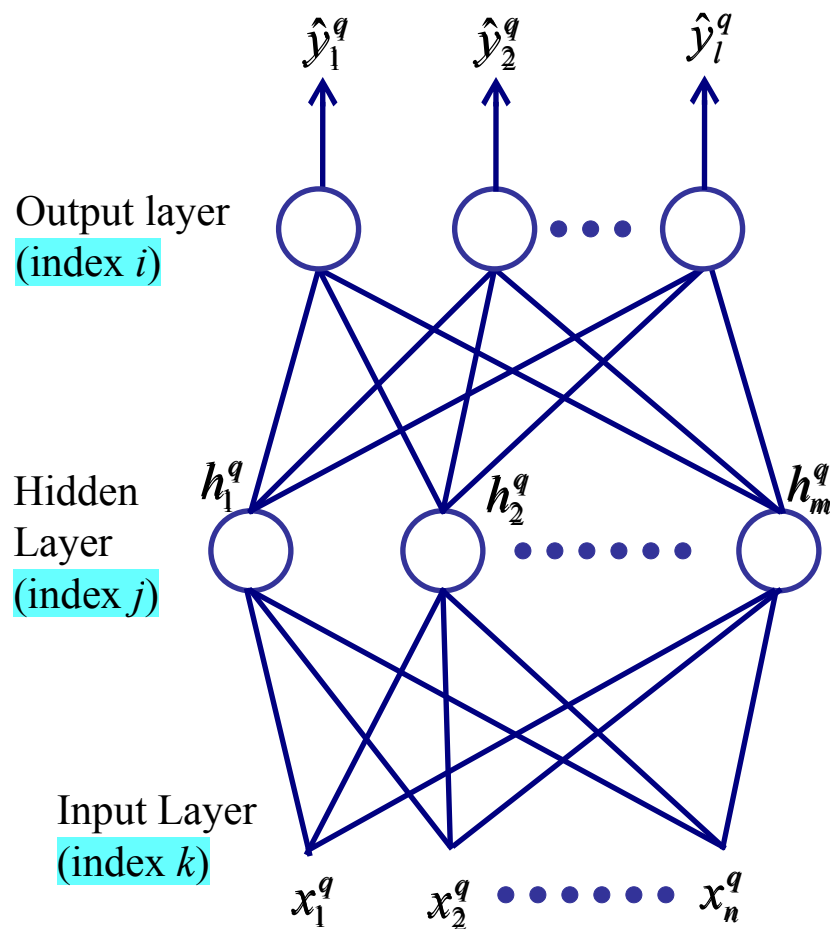
Why is backpropagation important?

Multilayer feedforward networks can classify any data and approximate any function

A general training algorithm for such networks is very useful

# The Backpropagation Algorithm



Output layer (index $i$)

Hidden Layer (index $j$)

Input Layer (index $k$)

$q$ = pattern index $\left[ \text{data-point } q = \left( x^q, y^q \right) \right]$

$n$ = # of inputs (input space dimension)

$m$ = # of hidden neurons (hidden space dim.)

$l$ = # of outputs (output space dimension)

$w_{ij}$ = weight from hidden neuron $j$ to output $i$

$w_{jk}$ = weight from input $k$ to hidden neuron $j$

$$s_j^q = \sum_{k=0}^{n} w_{jk} x_k^q$$

$$h_j^q = f_j\left(s_j^q\right)$$

Hidden Layer

$$s_i^q = \sum_{j=0}^{m} w_{ij} h_j^q$$

$$\hat{y}_i^q = f_i\left(s_i^q\right)$$

Output Layer

$$J^q = \frac{1}{2} \sum_{i=1}^{l} \left(y_i^q - \hat{y}_i^q\right)^2$$

$$= \frac{1}{2} \sum_{i=1}^{l} \left(e_i^q\right)^2$$

Loss function

$$J = \sum_{q=1}^{N} J^q \quad \rightarrow \quad \text{learn by } \Delta\overline{w} = -\eta \nabla J(\overline{w})$$

$\overline{w} \equiv$ vector of all weights

$$\Delta w_{ij} = -\eta_o \frac{\partial J^q}{\partial w_{ij}}$$

$$\Delta w_{jk} = -\eta_h \frac{\partial J^q}{\partial w_{jk}}$$

$$s_j^q = \sum_{k=0}^{n} w_{jk} x_k^q$$

$$h_j^q = f_j\left(s_j^q\right)$$

Hidden Layer

$$s_i^q = \sum_{j=0}^{m} w_{ij} h_j^q$$

$$\hat{y}_i^q = f_i\left(s_i^q\right)$$

Output Layer

$$J^q = \frac{1}{2}\sum_{i=1}^{l}\left(y_i^q - \hat{y}_i^q\right)^2$$

$$= \frac{1}{2}\sum_{i=1}^{l}\left(e_i^q\right)^2$$

Loss function

$$J = \sum_{q=1}^{N} J^q \qquad \rightarrow \qquad \text{learn by } \Delta\overline{w} = -\eta\nabla J(\overline{w})$$
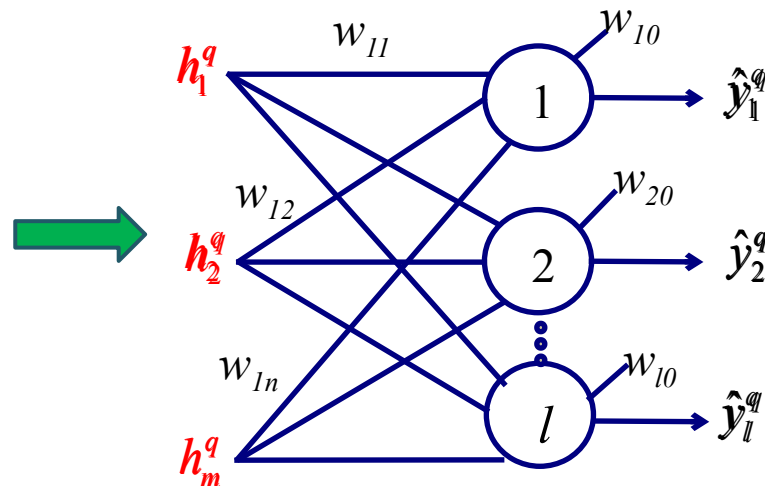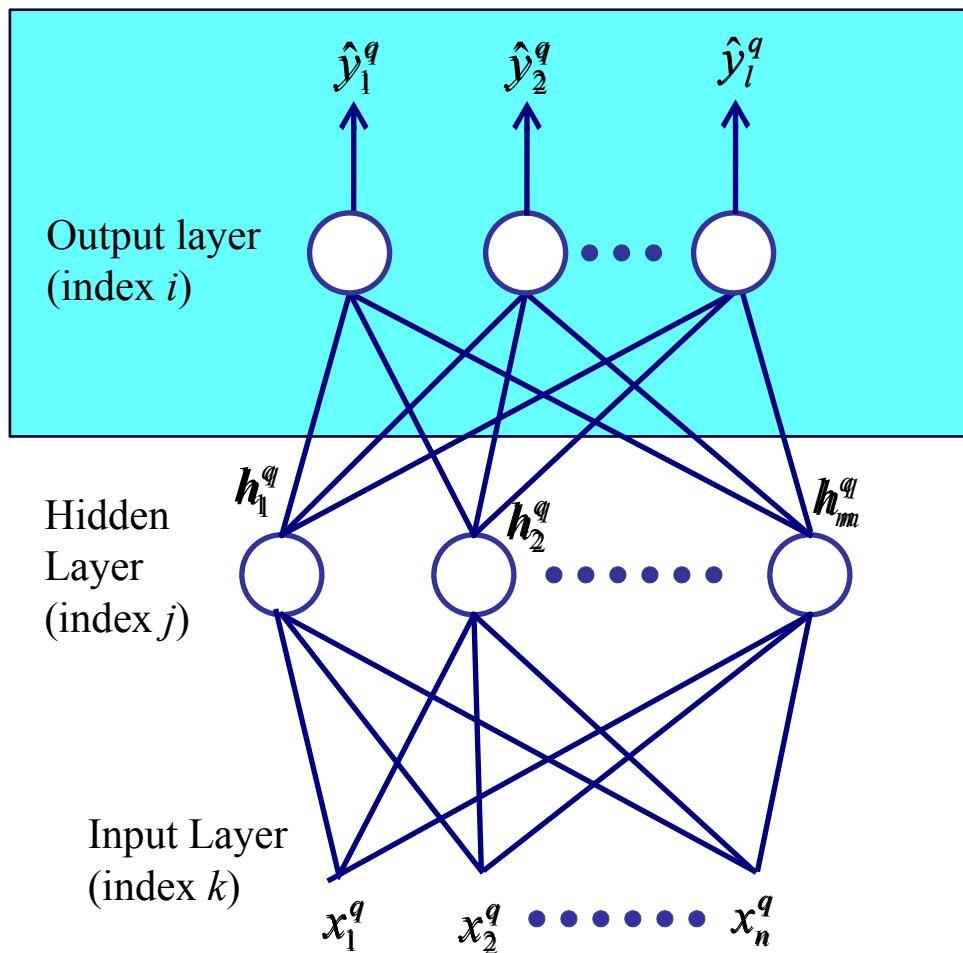
$\overline{w} \equiv$ vector of all weights

$$\Delta w_{ij} = -\eta_o \frac{\partial J^q}{\partial w_{ij}}$$ ——— This is easy to calculate using LMS

$$\Delta w_{jk} = -\eta_h \frac{\partial J^q}{\partial w_{jk}}$$ ——— But how to calculate this?

College of
Engineering
& Applied
Science

Department of
Electrical&
Computer
Engineering

# Learning for the Output Layer

Output layer
(index $i$)

$\hat{y}_1^q$ $\hat{y}_2^q$ $\hat{y}_l^q$

Hidden
Layer
(index $j$)

$h_1^q$ $h_2^q$ $h_m^q$

Input Layer
(index $k$)

$x_1^q$ $x_2^q$ $x_n^q$

$h_1^q$ $w_{11}$ $w_{10}$ $\hat{y}_1^q$
1

$h_2^q$ $w_{12}$ $w_{20}$ $\hat{y}_2^q$
2

$h_m^q$ $w_{1n}$ $w_{l0}$ $\hat{y}_l^q$
$l$

$$\frac{\partial J^q}{\partial w_{ij}} = - e_i^q f'\!\left(s_i^q\right) h_j^q$$

$$= - \left(y_i^q - \hat{y}_i^q\right) f'\!\left(s_i^q\right) h_j^q$$

**LMS**

$$\Delta w_{ij} = -\eta \frac{\partial J^q}{\partial w_{ij}} = \eta \left(y_i^q - \hat{y}_i^q\right) f'\!\left(s_i^q\right) h_j^q$$

To see this repeat the LMS calculation for $\dfrac{\partial J^q}{\partial w_{ij}}$ :

$$\frac{\partial J^q}{\partial w_{ij}} = \frac{\partial J^q}{\partial e_i^q} \cdot \frac{\partial e_i^q}{\partial \hat{y}_i^q} \cdot \frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{ij}}$$

$$= (e_i^q) \cdot (-1) \cdot (f_i'(s_i^q)) \cdot (h_j^q)$$

$$\equiv - e_i^q \cdot f_i'(s_i^q) \cdot h_j^q$$

where $f'(u) \equiv \dfrac{df}{du}\Big|_u$

If $f(u) = \dfrac{1}{1+e^{-u}}$

$f'(u) = \dfrac{1}{1+e^{-u}} \dfrac{e^{-u}}{1+e^{-u}} \equiv f(u)(1-f(u))$

$$\frac{\partial J^q}{\partial w_{ij}} = -\left(y_i^q - \hat{y}_i^q\right) f_i'(s_i^q) h_j^q \quad \text{①}$$

$$\equiv \delta_i^q$$

$$= -\delta_i^q h_j^q$$

$$\therefore \quad w_{ij} = w_{ij} + \eta_o \delta_i^q h_j^q$$

or $\quad \Delta w_{ij} = \eta_o \delta_i^q h_j^q \quad \text{②}$

# Learning for the Hidden Layer

**Calculating** $\dfrac{\partial J^q}{\partial w_{jk}}$ :

Note that

$$J^q = \frac{1}{2}\sum_{i=1}^{l}\left[y_i^q - \hat{y}_i^q\right]^2 = \frac{1}{2}\sum_{i=1}^{l}\left(e_i^q\right)^2 \quad \text{(3)}$$

$$= \frac{1}{2}\sum_{i=1}^{l}\left\{y_i^q - f_i\left[\sum_{j}w_{ij}f_j\left(\underbrace{\sum_{k}w_{jk}x_k^q}_{h_j^q}\right)\right]\right\}^2$$

$$\frac{\partial J^q}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}}\left[\frac{1}{2}\sum_{i}\left(e_i^q\right)^2\right]$$

$$= \frac{1}{2}\sum_{i=1}^{l}\frac{\partial}{\partial w_{jk}}\left(e_i^q\right)^2$$

$$= \frac{1}{2}\sum_{i}\frac{\partial\left(e_i^q\right)^2}{\partial e_i^q}\cdot\frac{\partial e_i^q}{\partial w_{jk}}$$

$$\frac{\partial J^q}{\partial w_{jk}} \equiv \sum_{i}e_i^q\cdot\frac{\partial e_i^q}{\partial w_{jk}} \quad \text{(4)}$$

$$\frac{\partial e_i^q}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}}\left(y_i^q - \hat{y}_i^q\right)$$

$$= -\frac{\partial \hat{y}_i^q}{\partial w_{jk}} \quad \text{(since } y_i^q \text{ is fixed)}$$

$$= -\frac{\partial \hat{y}_i^q}{\partial s_i^q} \cdot \frac{\partial s_i^q}{\partial w_{jk}} = -f_i''\left(s_i^q\right)\frac{\partial s_i^q}{\partial w_{jk}}$$

$$\boxed{\frac{\partial e_i^q}{\partial w_{jk}} \equiv -f_i''\left(s_i^q\right)\frac{\partial s_i^q}{\partial w_{jk}}} \qquad \text{(5)}$$

$$\frac{\partial s_i^q}{\partial w_{jk}} = \frac{\partial}{\partial w_{jk}}\sum_{j'=0}^{m} w_{ij'} h_{j'}^q$$

$$= \sum_{j'}\frac{\partial}{\partial w_{jk}} w_{ij'} h_{j'}^q$$

$$\boxed{\frac{\partial s_i^q}{\partial w_{jk}} = \sum_{j'} w_{ij'}\frac{\partial h_{j'}^q}{\partial w_{jk}}} \qquad \text{(6)}$$

$$j' \in \text{Hidden Layer}$$

$$\text{Note}: \text{we use } j' \text{ because } j$$
$$\text{is already used in } w_{jk}$$

$$\frac{\partial s_i^q}{\partial w_{jk}} = \sum_{j'} w_{ij'} \frac{\partial h_{j'}^q}{\partial w_{jk}}$$

⑥

⑦ → ⑥

$$\frac{\partial h_{j'}^q}{\partial w_{jk}} = \frac{\partial h_{j'}^q}{\partial s_{j'}^q} \frac{\partial s_{j'}^q}{\partial w_{jk}}$$

Since $w_{jk}$ has no effect on $s_{j'}^q$ if $j \neq j'$

$$\frac{\partial s_i^q}{\partial w_{jk}} = w_{ij} f_j'\left(s_j^q\right) x_k^q$$

⑧

$$= \begin{cases} f_j'\left(s_j^q\right) x_k^q & \text{if } j' = j \\ 0 & \text{if } j' \neq j \end{cases}$$

⑧ → ⑤

$$\therefore \frac{\partial h_{j'}^q}{\partial w_{jk}} = \begin{cases} f_j'\left(s_j^q\right) x_k^q & \text{if } j' = j \\ 0 & \text{if } j' \neq j \end{cases}$$

⑦

$$\frac{\partial e_i^q}{\partial w_{jk}} = - f_i'\left(s_i^q\right) w_{ij} f_j'\left(s_j^q\right) x_k^q$$

⑨

$9 \rightarrow 4$

$$\frac{\partial J^q}{\partial w_{jk}} = - \sum_i w_{ij} \overbrace{e_i^q f_i'\left(s_i^q\right)}^{\delta_i^q} f_j'\left(s_j^q\right) x_k^q$$

$$= - \sum_i w_{ij} \delta_i^q f_j'\left(s_j^q\right) x_k^q \qquad \textcircled{10}$$

$$\frac{\partial J^q}{\partial w_{jk}} = - \underbrace{\left[ f_j'\left(s_j^q\right) \sum_i w_{ij} \delta_i^q \right]}_{\delta_j^q} x_k^q$$

$$\therefore \quad \boxed{\frac{\partial J^q}{\partial w_{jk}} = - \delta_j^q x_k^q}$$

where $\delta_j^q = f_j'\left(s_j^q\right) \sum_i w_{ij} \delta_i^q$

and

$$w_{jk} = w_{jk} + \eta_h \delta_j^q x_k^q$$

$$\boxed{\Delta w_{jk} = \eta_h \delta_j^q x_k^q} \qquad \textcircled{11}$$

# Gradient Calculation Process

$$\boxed{1} \quad \frac{\partial J^q}{\partial w_{ij}} = - f_i'\left(s_i^q\right)\left[ y_i^q - \hat{y}_i^q \right] h_j^q \quad \longleftarrow \quad \text{same as LMS}$$

$$\boxed{10} \quad \frac{\partial J^q}{\partial w_{jk}} = - f_j'\left(s_j^q\right) \sum_i w_{ij} \delta_i^q \, x_k^q \quad \longleftarrow \quad \text{same form as 1}$$

College of
Engineering
& Applied
Science

Department of
Electrical&
Computer
Engineering

$$\textcolor{red}{\delta_i^q}$$

$$\text{(1)} \quad \frac{\partial J^q}{\partial w_{ij}} = -\boxed{f_i'\left(s_i^q\right)\left[y_i^q - \hat{y}_i^q\right]} h_j^q$$

$$\text{(10)} \quad \frac{\partial J^q}{\partial w_{jk}} = -\boxed{f_j'\left(s_j^q\right)\sum_i w_{ij}\delta_i^q} x_k^q$$

$$\textcolor{red}{\delta_j^q}$$

-Signals flow forward
-Errors (δ's) propagate backwards

Note how $\delta_j^q$ for $j \in$ hidden layer
is calculated recursively

$$\delta_1^q \qquad \delta_2^q \qquad \delta_l^q$$

$$\delta_j^q = f_j'\left(s_j^q\right)\sum_{i=1}^{l} w_{ij}\delta_i^q$$

can be calculated
at $j$ if $\delta_i^q$ are known

# Form of the Learning Rule

Compare ② and ⑪ . Both have the form

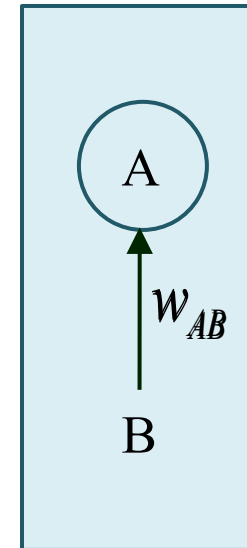$$\Delta w_{AB} = \eta \cdot \delta_A^q \cdot (\text{input on line B})$$

if $A \in$ output layer

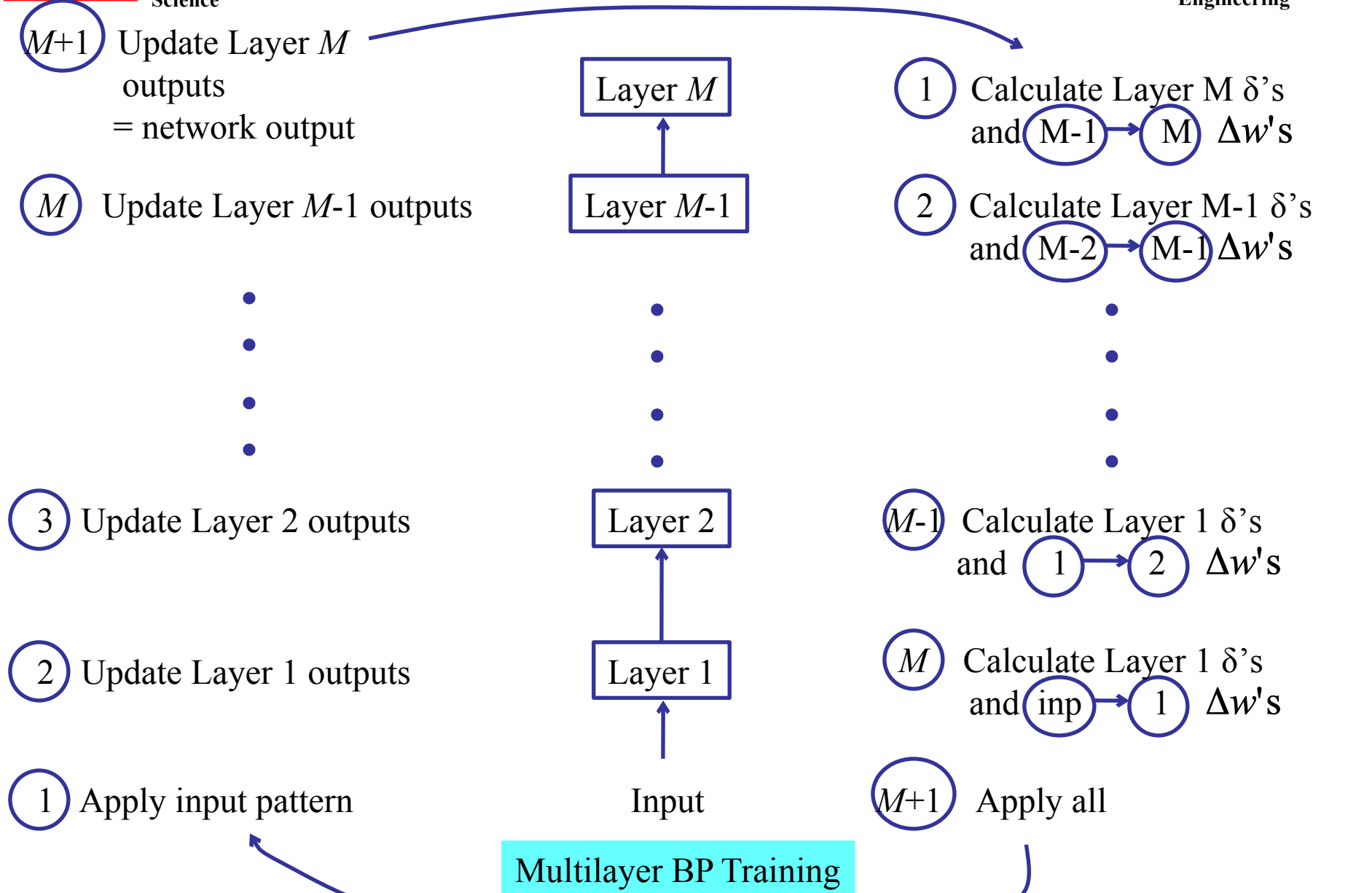$$\delta_A^q = f_A'\left(s_A^q\right)\left[y_A^q - \hat{y}_A^q\right]$$

if $A \in$ hidden layer

$$\delta_A^q = f_A'\left(s_A^q\right)\sum_r w_{rA}\delta_r^q$$

where $r$ indexes the neurons **_to_** which A connects.

$M$+1  Update Layer $M$ outputs
 = network output

$M$  Update Layer $M$-1 outputs

• • •

3  Update Layer 2 outputs

2  Update Layer 1 outputs

1  Apply input pattern

Layer $M$

Layer $M$-1

• • •

Layer 2

Layer 1

Input

1  Calculate Layer M δ's
 and M-1 → M  $\Delta w$'s

2  Calculate Layer M-1 δ's
 and M-2 → M-1 $\Delta w$'s

• • •

$M$-1  Calculate Layer 1 δ's
 and 1 → 2  $\Delta w$'s

$M$  Calculate Layer 1 δ's
 and inp → 1  $\Delta w$'s

$M$+1  Apply all

Multilayer BP Training

# A General Formulation

Networks may use:

- Different loss functions.
- Different activation functions (even different between layers or within a layer.)

We can define a general formulation to handle this under the assumption that:

$$J^q = \sum_i J_i^q$$  i.e., total loss is added over the output neurons.

**For output layer weights:**

$$\frac{\partial J^q}{\partial w_{ij}} = \frac{\partial J^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q} \frac{\partial s_i^q}{\partial w_{ij}} = \boxed{\frac{\partial J_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q}} \boxed{\frac{\partial s_i^q}{\partial w_{ij}}}$$

$$\boxed{-\delta_i^q} \qquad \boxed{h_j^q}$$

**Mean-squared loss function**
**Linear composition function**

**For the hidden layer weights:**

$$\frac{\partial J^q}{\partial w_{jk}} = \sum_i \frac{\partial J_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q} \frac{\partial s_i^q}{\partial h_j^q} \frac{\partial h_j^q}{\partial s_j^q} \frac{\partial s_j^q}{\partial w_{jk}}$$

$$= \left[ \sum_i \frac{\partial J_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q} \frac{\partial s_i^q}{\partial h_j^q} \right] \left[ \frac{\partial h_j^q}{\partial s_j^q} \right] \left[ \frac{\partial s_j^q}{\partial w_{jk}} \right]$$

## For the hidden layer weights:

$$\frac{\partial J^q}{\partial w_{jk}} = \sum_i \frac{\partial J_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q} \frac{\partial s_i^q}{\partial h_j^q} \frac{\partial h_j^q}{\partial s_j^q} \frac{\partial s_j^q}{\partial w_{jk}}$$

$$= \left[ \sum_i \frac{\partial J_i^q}{\partial \hat{y}_i^q} \frac{\partial \hat{y}_i^q}{\partial s_i^q} \frac{\partial s_i^q}{\partial h_j^q} \right] \left[ \frac{\partial h_j^q}{\partial s_j^q} \right] \left[ \frac{\partial s_j^q}{\partial w_{jk}} \right]$$

$$-\delta_i^q \qquad w_{ij} \qquad f_j'(s_j^q) \qquad x_k^q$$

$$-(y_i^q - \hat{y}_i^q) f_i'(s_i^q)$$

$$\frac{\partial J^q}{\partial w_{jk}} = - f_j'(s_j^q) \sum_i w_{ij} \delta_i^q x_k^q$$

**Mean-squared loss function**
**Linear composition function**

# Practical Choices

# Using Backpropagation for Classification Problems

Backpropagation requires that $\dfrac{\partial J^q}{\partial w_{jk}}$ exists $\forall w_{ij}$

$\Rightarrow f_i'(s_i^q)$ must exist

$\Rightarrow f(u)$ (activation function) must be differentiable

$$\rightarrow \text{use } f(u) = \frac{1}{1 + e^{-u}}$$

$$\text{or } f(u) = \tanh(u)$$

instead of hard threshold.

> But classification requires that output neurons have 0/1 or +1/-1 output

## ?????????

College of
Engineering
& Applied
Science

Department of
Electrical&
Computer
Engineering

# Possible Solutions

- Use 0/1 or +1/-1 as target values $y_i^q$

- Set operating parameters $H, L, L<H$ such that:

$$\hat{y}_i \geq H \qquad \text{is considered a match for } y_i = +1$$

$$\hat{y}_i \leq L \qquad \text{is considered a match for } y_i = 0 \text{ or } -1$$

- During training, use 0/1 or +1/-1 as targets but if $\hat{y}_i$ matches the operating targets, make no weight change

  Typical values for $H, L$ are:

  $H$=0.75     $L$=0.25    for 0-1 sigmoid
  $H$=0.75     $L$=-0.75    for +/- sigmoid

# Multi-Class Classification

If there are more than 2 classes, how should the output be represented?

## Solution:

- If there are $M$ classes, have $M$ neurons in the output layer.

- During training, use one-hot codes as targets for each class, e.g. if $M = 3$,
  Class 1 = [1 0 0 ]    Class 2 = [0 1 0]      Class 3 = [0 0 1]

- After training, for a test input $x$, choose the class as the output neuron with the highest output, e.g., Output = [0.1 0.6 0.2 ] → [0 1 0]

- If "no choice" is allowed, require that an output must be higher than some threshold $\theta$ to be considered 1.

A better solution: softmax classifiers

# Data Scaling

Given input vectors $\quad \bar{x} = \begin{bmatrix} x_1 & x_2 & \cdot & \cdot & \cdot & x_n \end{bmatrix}^T$

Scale each $x_k$ such that it is 0-mean and has approximately the same range as the other components. Scaling between -1 and +1 is usually a good idea.

Note:  if $x_k \geq 0$  in all cases

$$\Delta w_{jk} = \eta \delta_j x_k$$

$\Rightarrow$    All weights for hidden neuron $j$ either increase together or decrease together

$\rightarrow$      less discrimination in each training step

$\rightarrow$      slower training

# Data Scaling

In binary input situations, use +/- rather than 0/1:

Recall that (inp) ⟶ (hidden) weights are modified as:

$$\Delta w_{jk} = \eta \delta_j^q \, x_k^q$$

$$x_k^q = 0 \quad \longrightarrow \quad \Delta w_{jk} = 0$$

Using $\pm$ allows an active weight change at every steps rather than only at the $x = +1$ steps.

⇒ 0 inputs cause no weight change even if $\delta_j^q$ is large

⇒ Learning on these inputs happens only passively (by omission).

Potentially faster learning

# Choice of Activation Function

Use $f(u) = a \tanh(bu)$

Same reason as scaling inputs 0-mean

$$f(u) = \frac{1}{1 + e^{-u}} \quad \Rightarrow \quad 0 < f(u) < 1$$

$\Rightarrow$     all weights to a target neuron increase or decrease together

$$f(u) = a \tanh(bu) \Rightarrow -a < f(u) < a$$

which allows greater discrimination

Le Cun (1989, 1993) suggests

$a = 1.7159 \quad b = 2/3$

This gives $f(1) = 1 \quad f(-1) = -1$
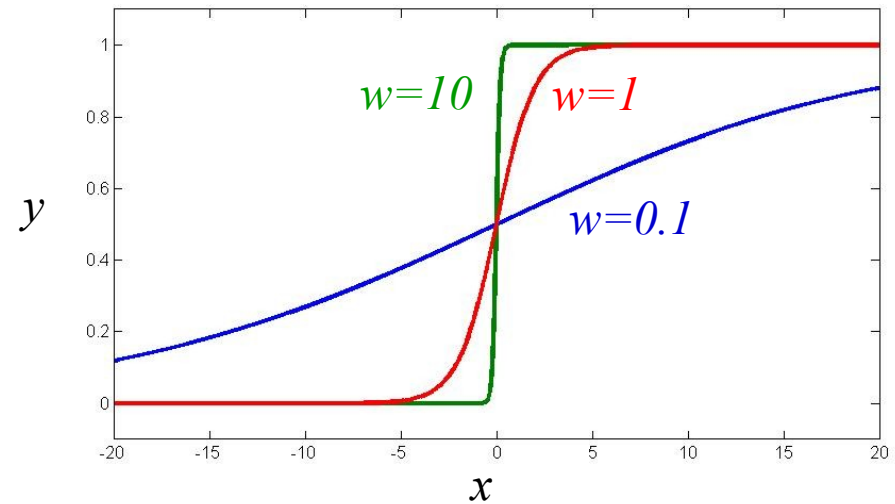
In some cases, a rectified linear unit (ReLU) may work well.

# Weight Initialization

Consider

$$y = \frac{1}{1 + e^{-wx}}$$



- If $w$ is too small, $f'(wx)$ is a small value for most $x \rightarrow$ slow learning.

- If $w$ is too large, $wx$ is large for most $x$ (except near $x = 0$) and $f'(wx) \approx 0$ (saturation) $\rightarrow$ slow learning.

- For $x \in \mathfrak{R}^n$, the same argument applies w.r.t. values of $\|w\|$

$\Rightarrow$ The more input lines a neuron has, the smaller its range of weights should be to avoid saturating the sigmoid and causing the gradient to vanish.

Gaussian with mean = 0, standard deviation σ ~

Uniform between (*-a, +a*) where *a* ~

Where $n$ = number of inputs to the neuron.

Xavier initialization (Glorot & Bengio, 2010):
Uniform between (*-a, +a*) where *a* ~

$N_s$ = number of neurons in source layer
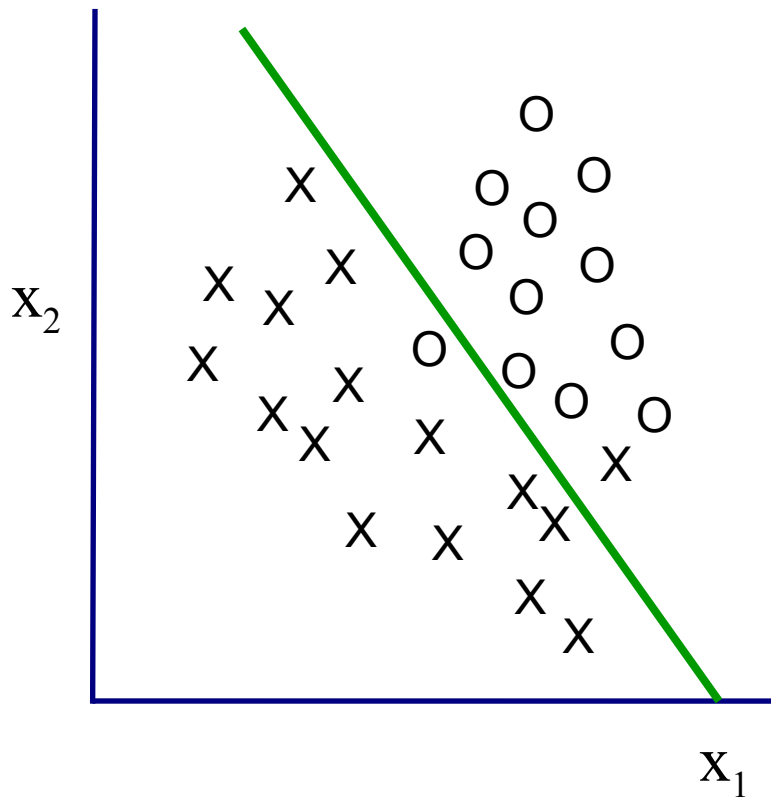$N_t$ = number of neurons in target layer

Xavier Glorot, Yoshua Bengio, Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, PMLR 9:249-256, 2010.

- Don't initialize all weights to identical values. This may cause all hidden neurons to learn the same thing and preclude overall learning.

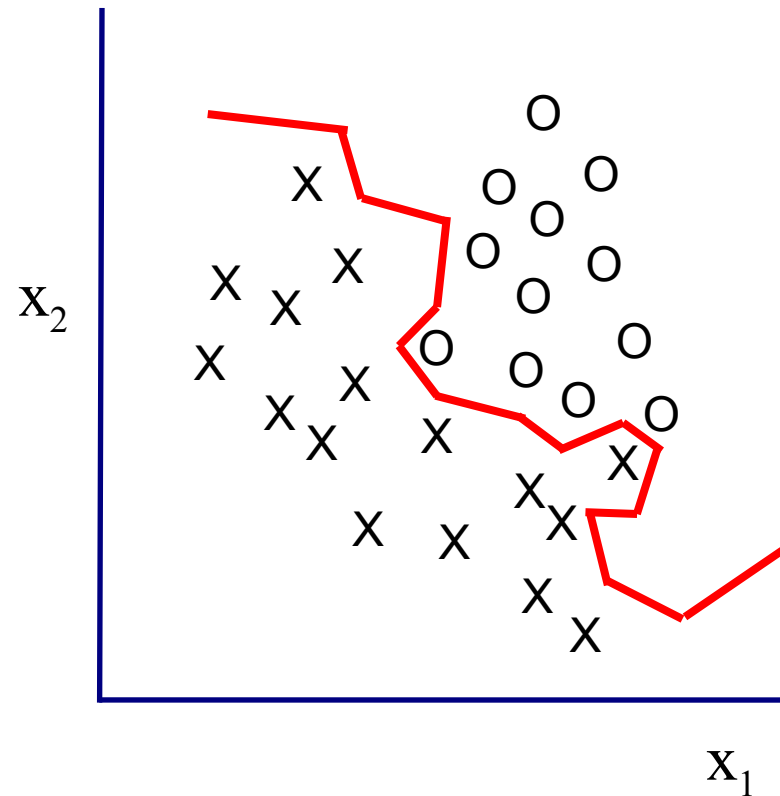- Initialize weights to be both positive and negative instead of just positive.

# Generalization and Validation
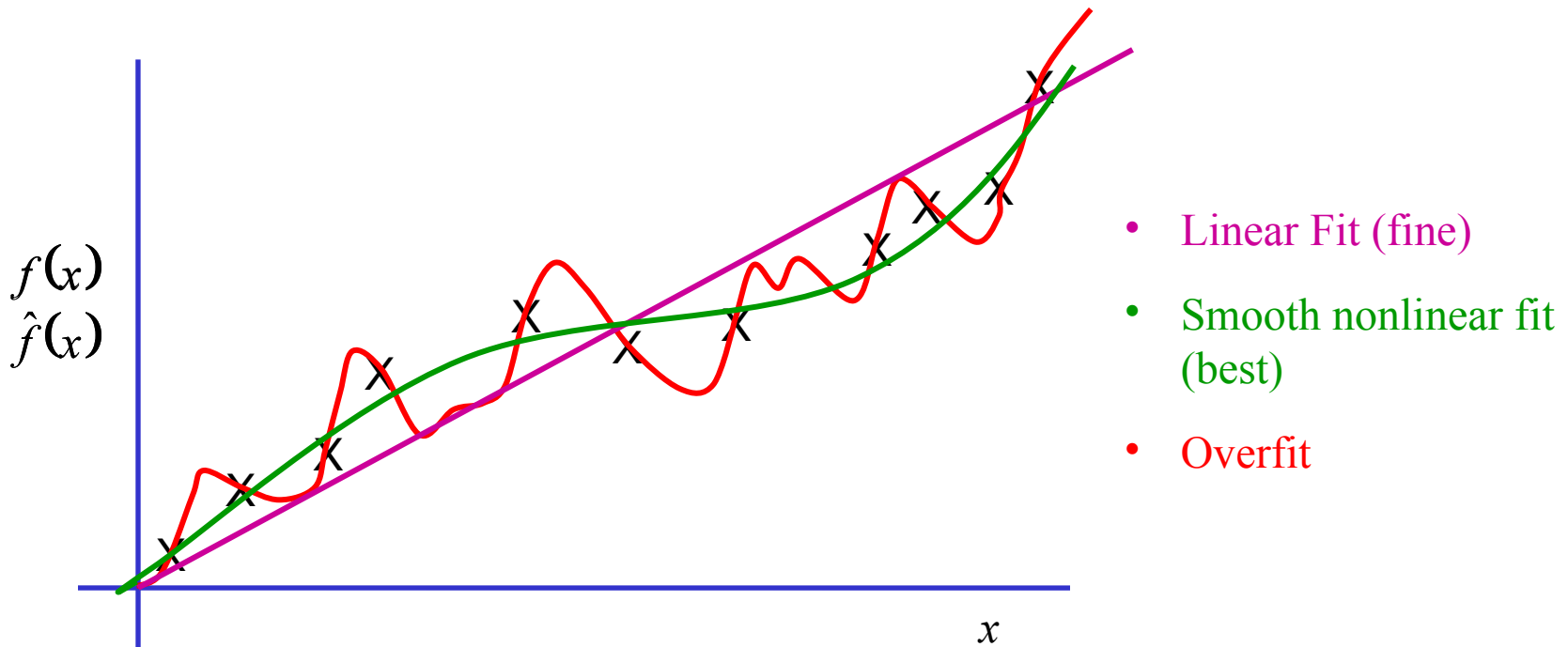
# Fitting vs. Overfitting (Classification)



Good Fit

Overfit

College of
Engineering
& Applied
Science

Department of
Electrical&
Computer
Engineering

# Fitting vs. Overfitting
# (Approximation)

$$f(x)$$
$$\hat{f}(x)$$

$x$

- Linear Fit (fine)

- Smooth nonlinear fit (best)

- Overfit

# Fitting vs. Overfitting

The fine variations seen in overfitting are fitting the "noise" rather than the "signal"

$\rightarrow \rightarrow$ poor generalization

What causes overfitting?

**Too many "degrees of freedom" = ways of variation.**

In feed-forward networks:

Degrees of freedom $\leftrightarrow$ Number of weights and hidden neurons

More weights and neurons arrow more complicated functions and decision boundaries $\rightarrow$ Overfitting $\rightarrow$ Poor generalization

use fewer hidden neurons and/or fewer weights/neuron.

# Regularization

Modifications to the learning process, network behavior, or network structure that force better generalization.
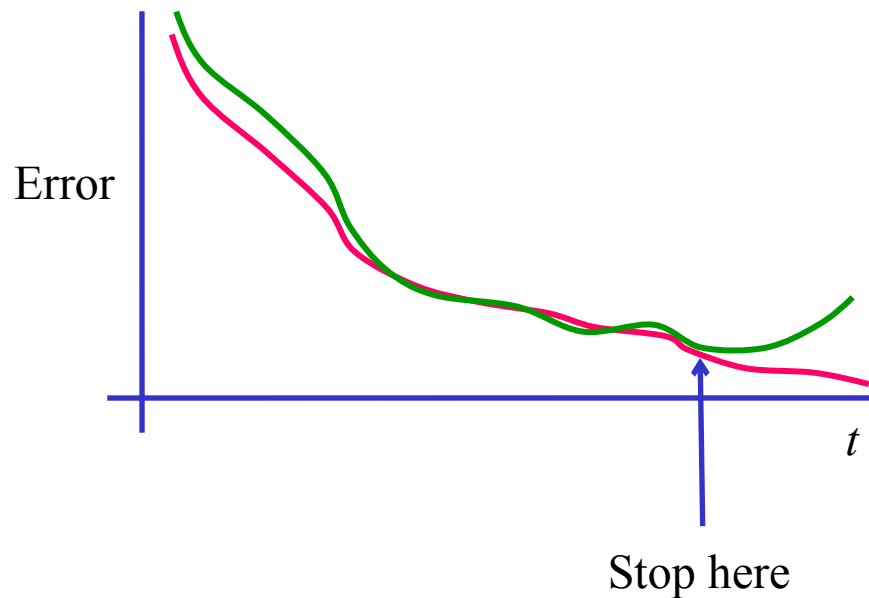
Some common regularization methods:

- Early stopping.

- Weight decay.

- Pruning.

- Dropout.

- Sparsity constraints.

# Early Stopping

- During training, periodically evaluate network performance on a ___validation set___.

- If training and testing error are both going down, keep training.

- If validation error starts to go up, stop training <u>even if</u> training error is still going down.

— Training Error

— Validation Error

Error

*t*

Stop here

- Assumes [1] "Signal" is learned before "noise".

[2] Once validation error starts going up, it will never reverse course

⟶ Both are heuristics

# Weight Decay

Force the total weight of the network to be as low as possible.

For example, add a weight penalty term to the loss function:

$$J^q = \frac{1}{2} \sum_{i \in Outputs} \left(y_i^q - \hat{y}_i^q\right)^2 + \boxed{\lambda \sum_{Layers\ L} \sum_{j \in Layer\ L} \sum_{i \in Layer\ L+1} \left(w_{ij}\right)^2} \longleftarrow \text{Weight Penalty}$$

Why does this work?

It forces the network to use weights more efficiently in one or both of two ways:

- Make all weights small → smoother function
- Use fewer non-zero weights → fewer degrees-of-freedom

Both → better generalization.

# **Pruning and Construction**

## **Pruning:**

- During training, encourage the network to use as few weights as possible

$\rightarrow$ some weights become $\approx 0$

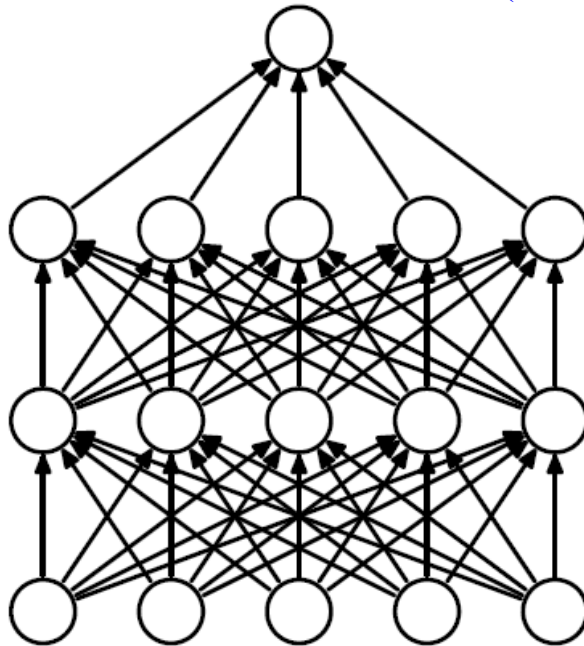- Remove weights of very small size

Notes: | 1 | Some algorithms do only the second step

| 2 | Pruning algorithms can also remove
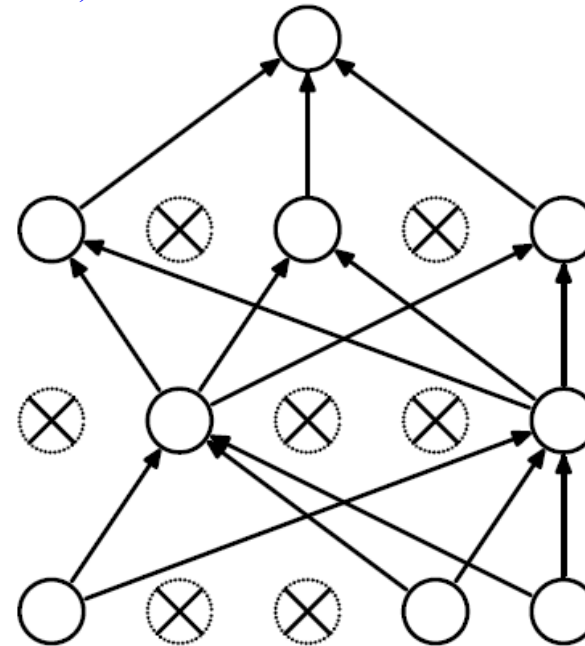neurons to become useless

## **Construction:**

Type I:   Start with a small network, train as far as possible, add a neuron
and keep training.  Repeat this process till satisfied.

Type II:  Start with a large network, train, prune, train, add, ….

# Dropout

(Srivastava et al., 2014)



(a) Standard Neural Net

(b) After applying dropout.

Randomly leave out some inputs and hidden neurons during each iteration of the learning update, to prevent them from co-adapting – forcing them to learn somewhat different "views" of the data.

Srivastava et al. (2014) Dropout: A Simple Way to Prevent Neural Networks from Overfitting, *Journal of Machine Learning Research* 15: 1929-1958