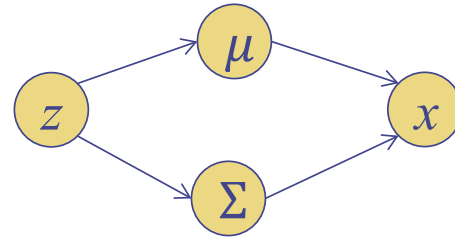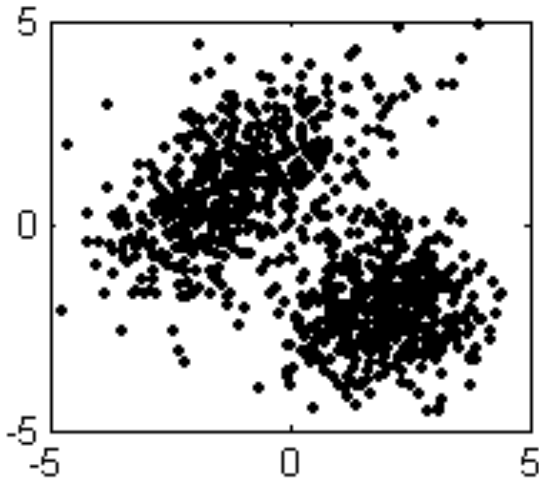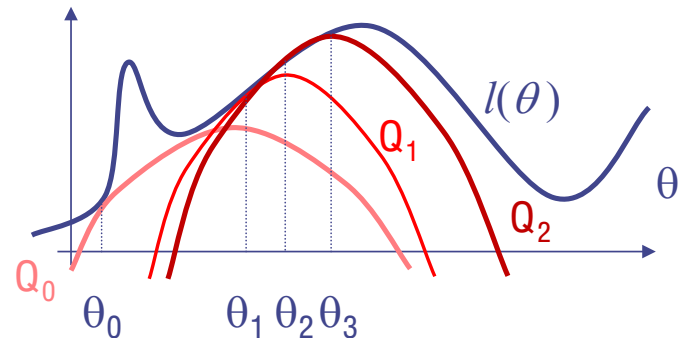# **Machine Learning**

## 4771

Instructor: Itsik Pe'er

# Reminder: EM for Gauss. Mix.



Expectation-Maximization:
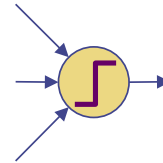Iteratively improve
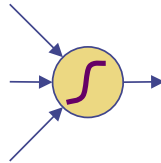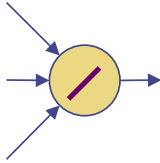Expected-log-likelihood

# Class 16

- Multi-Layer Neural Networks

- Back-Propagation

# Multi-Layer Neural Networks

- Perceptron/linear/logistic/threshold neurons

# Multi-Layer Neural Networks

•Perceptron/linear/logistic/threshold



- Different functions of the linear combination of inputs

$$f(x) = g(\theta^T x)$$

- Different loss functions

- Different strategies for minimizing empirical risk

# Multi-Layer Neural Networks

- Need to introduce non-linearities between layers



- Neural network can adjust the basis functions themselves...



$$f(\boldsymbol{x}) = g\left( \sum_{i=1}^{P} \theta_i g\left( \widetilde{\theta}_i^{\boldsymbol{T}} \boldsymbol{x} \right) \right)$$

# Multi-Layer Neural Networks

- Multi-Layer Network can handle more complex decisions
- 1-layer: is linear, can't handle XOR
- Each layer adds more flexibility (but more parameters!)
- Each node splits its input space with linear hyperplane
- 2-layer: if last layer is AND operation, get convex hull
- 3-layer: can do almost anything multi-layer can
        by fanning out the inputs at $2^{nd}$ layer



- Note: Without loss of generality, we can omit the 1 and $\theta_0$

# Parameterizing Neural Networks



$x_i$

$x_1$

$x_2$

$x_3$

$z_i =$

$w_{ij}$

$z_j$

$a_j = \sum_i x_i w_{ij}$

$z_j = g(a_j)$

$w_{jk}$

$a_k = \sum_j z_j w_{jk}$

$z_k = g(a_k)$

$w_{kl}$

$a_l = \sum_k z_k w_{kl}$

$z_l = g(a_l)$

# Parameterizing Neural Networks

- Parameters are *weights*   $\theta = \{w_{ij}, w_{jk}, w_{kl}\}$
- Weights define linear combinations of inputs...

$x_1$

$x_2$

$x_3$

$w_{ij}$      $w_{jk}$      $w_{kl}$

$z_i$     $a_j$    $z_j$     $a_k$   $z_k$      $a_l$

# Parameterizing Neural Networks

- Parameters are *weights* $\quad \theta = \{w_{ij}, w_{jk}, w_{kl}\}$
- Weights define linear combinations of inputs...

$$x_1 \quad x_2 \quad x_3$$

$$w_{ij} \qquad w_{jk} \qquad w_{kl}$$

$$z_i \qquad a_j \qquad z_j \qquad a_k \qquad z_k \qquad a_l$$

$$a_j = \sum_i w_{ij} z_i \qquad a_k = \sum_j w_{jk} z_j \qquad a_l = \sum_k w_{kl} z_k$$

$$z_j = g(a_j) \qquad z_k = g(a_k) \qquad z_l = g(a_l)$$

...that activate neurons...
...that linearly combine...
...to activate neurons...
...that linearly combine to produce output

# Back-Propagation

- Gradient descent on squared loss is done layer by layer

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

$$z_j = g(a_j)$$

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

$$z_k = g(a_k)$$

$$a_l = \sum_k w_{kl} z_k$$

$$z_l = g(a_l)$$

- Back-Propagation: Splits layer into its inputs & outputs
- Get gradient on output...back-track chain rule until input

# Back-Propagation

- Cost function: $R(\theta) = \frac{1}{N} \sum_{n=1}^{N} L(y^n - f(x^n))$

$$= \frac{1}{N} \sum_{n}^{N} \frac{1}{2} \left( y^n - g\left( \sum_{k} w_{kL}\, g\left( \sum_{j} w_{jk}\, g\left( \sum_{i} w_{ij} x_i \right) \right) \right) \right)^2$$



$x_1$    $x_2$    $x_3$

$z_i$   $W_{ij}$   $a_j$   $z_j$   $W_{jk}$   $a_k$   $z_k$   $W_{kl}$   $a_l$   $z_l$

$f(x)$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} L(y^n - f(\boldsymbol{x}^n))$$

- Cost function:

$$= \frac{1}{N} \sum_{n=1}^{N} \frac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$



$x_1$

$x_2$

$x_3$

$z_i \quad W_{ij} \quad a_j \quad z_j \quad W_{jk} \quad a_k \quad z_k \quad W_{kl} \quad a_l \quad z_l \quad \mathfrak{Z}_l^n$

- First compute output layer derivative:

$$L^n \overset{\text{def}}{=} \frac{1}{2}(y^n - f(\boldsymbol{x}^n))^2$$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \frac{\partial L^n}{\partial w_{kl}} = \frac{1}{N} \sum_n (y^n - f(x^n)) \cdot g'(a_l) \frac{\partial a_l}{\partial w_{kl}}$$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} L(y^n - f(x^n))$$

- Cost function:

$$= \frac{1}{N} \sum_{n=1}^{N} \frac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$



$$x_1 \quad x_2 \quad x_3$$

$$z_i \quad w_{ij} \quad a_j \quad z_j \quad w_{jk} \quad a_k \quad z_k \quad w_{kl} \quad a_l \quad z_l$$

- First compute output layer derivative:

$$L^n \stackrel{\text{def}}{=} \frac{1}{2}(y^n - f(x^n))^2$$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right)$$

**Chain Rule**

$$= \frac{1}{N} \sum_n \left[ \frac{\partial \frac{1}{2}(y^n - g(a_l^n))^2}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right)$$
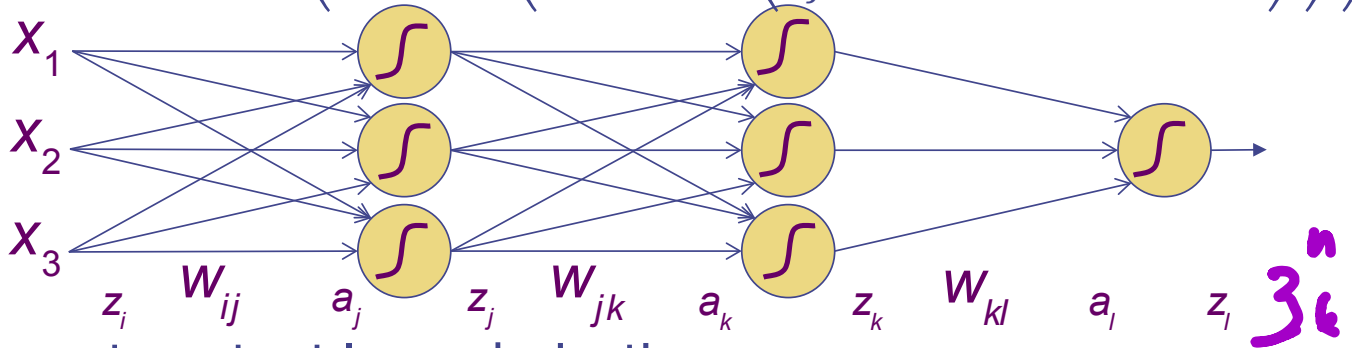
# Back-Propagation

$$R(\theta) = \frac{1}{N}\sum_{n=1}^{N} L(y^n - f(x^n))$$

- Cost function:

$$= \frac{1}{N}\sum_{n=1}^{N} \frac{1}{2}\left(y^n - g\left(\sum_k w_{kl} g\left(\sum_j w_{jk} g\left(\sum_i w_{ij} x_i^n\right)\right)\right)\right)^2$$

$x_1$

$x_2$

$x_3$

$z_i$   $W_{ij}$   $a_j$   $z_j$   $W_{jk}$   $a_k$   $z_k$   $W_{kl}$   $a_l$   $z_l$

- First compute output layer derivative:

$$L^n \stackrel{\text{def}}{=} \frac{1}{2}(y^n - f(x^n))^2$$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N}\sum_n \left[\frac{\partial L^n}{\partial a_l^n}\right]\left(\frac{\partial a_l^n}{\partial w_{kl}}\right)$$
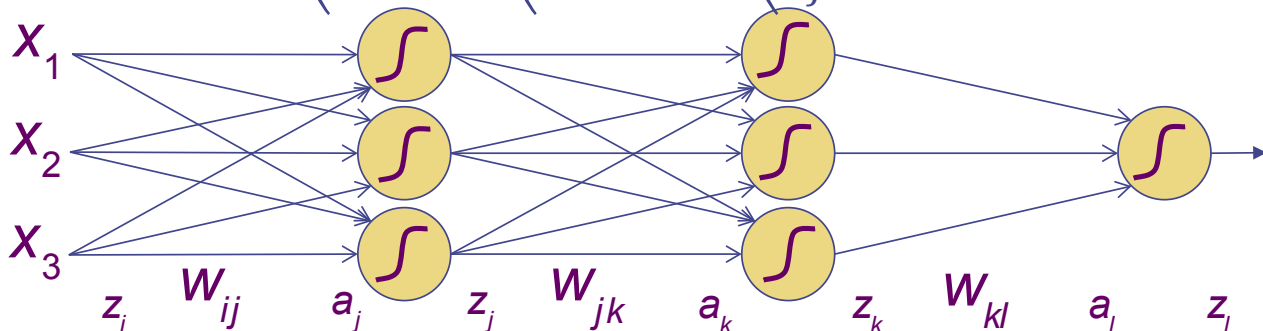
**Chain Rule**

$$= \frac{1}{N}\sum_n \left[\frac{\partial \frac{1}{2}(y^n - g(a_l^n))^2}{\partial a_l^n}\right]\left(\frac{\partial a_l^n}{\partial w_{kl}}\right) = \frac{1}{N}\sum_n \underbrace{[-(y^n - z_l^n)g'(a_l^n)]}_{\delta_l^n}(z_k^n)$$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} L(y^n - f(x^n))$$

- Cost function:

$$= \frac{1}{N} \sum_{n=1}^{N} \frac{1}{2} \left( y^n - g\left( \sum_k w_{kl} g\left( \sum_j w_{jk} g\left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$



$x_1$

$x_2$

$x_3$

$z_i \quad W_{ij} \quad a_j \quad z_j \quad W_{jk} \quad a_k \quad z_k \quad W_{kl} \quad a_l \quad z_l$

- First compute output layer derivative:

$$L^n \stackrel{\text{def}}{=} \frac{1}{2}(y^n - f(x^n))^2$$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right)$$
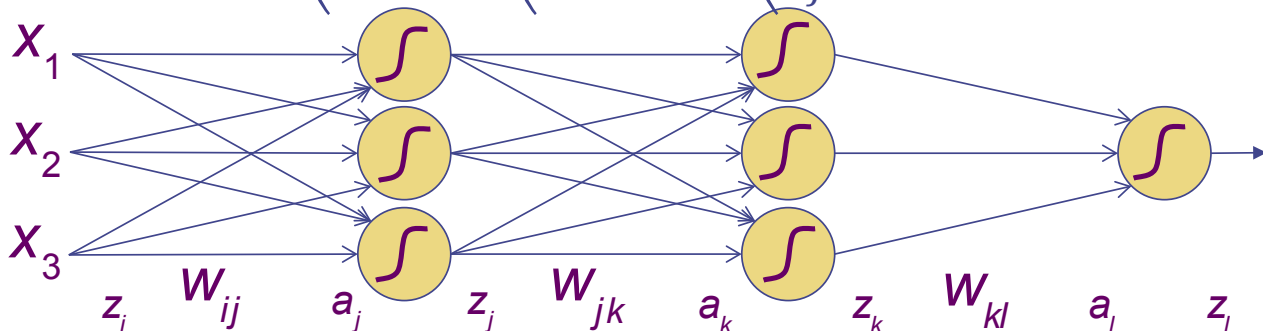
**Chain Rule**

$$= \frac{1}{N} \sum_n \left[ \frac{\partial \frac{1}{2}(y^n - g(a_l^n))^2}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right) = \frac{1}{N} \sum_n \left[ -(y^n - z_l^n) g'(a_l^n) \right](z_k^n) = \frac{\sum_n \delta_l^n z_k^n}{N}$$

**Define as** $\delta$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} \tfrac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$



$x_1$ $x_2$ $x_3$

$z_i$ $\quad W_{ij}$ $\quad a_j$ $\quad z_j$ $\quad W_{jk}$ $\quad a_k$ $\quad z_k$ $\quad W_{kl}$ $\quad a_l$ $\quad z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\par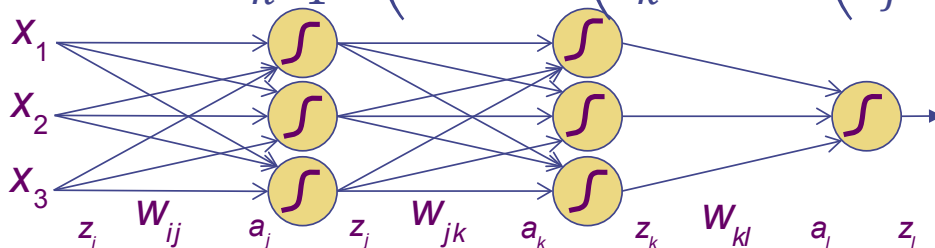tial a_l^n}{\partial w_{kl}} \right) = \frac{1}{N} \sum_n [-(y^n - z_l^n) g'(a_l^n)] z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$

• Next, hidden layer derivative:

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_n \frac{\partial L^n}{\partial w_{jk}} = \frac{1}{N} \sum_{n,l} \frac{\partial L^n}{\partial a_l} \frac{\partial a_l}{\partial W_{jk}} = \sum_{n,l} \frac{\partial L^n}{\partial a_l} \frac{\partial a_l}{\partial a_k} \frac{\partial a_k}{\partial W_{jk}}$$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} \frac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$



$x_1$   $x_2$   $x_3$

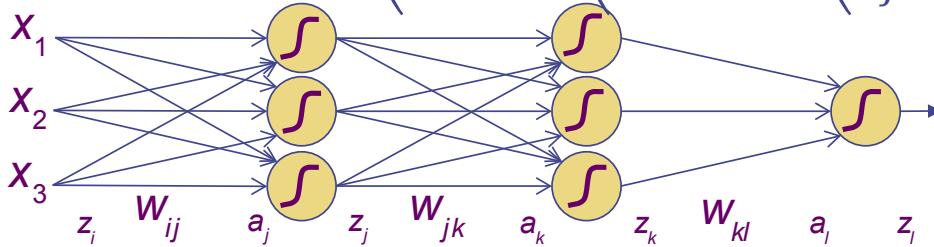$z_i$   $w_{ij}$   $a_j$   $z_j$   $w_{jk}$   $a_k$   $z_k$   $w_{kl}$   $a_l$   $z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right) = \frac{1}{N} \sum_n [-(y^n - z_l^n) g'(a_l^n)] z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$

• Next, hidden layer derivative:

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_k^n} \right] \left( \frac{\partial a_k^n}{\partial w_{jk}} \right)$$

# Back-Propagation

$$R(\theta) = \frac{1}{N}\sum_{n=1}^{N}\tfrac{1}{2}\left(y^n - g\left(\sum_k w_{kl}\,g\left(\sum_j w_{jk}\,g\left(\sum_i w_{ij}x_i^n\right)\right)\right)\right)^2$$



$x_1$
$x_2$
$x_3$

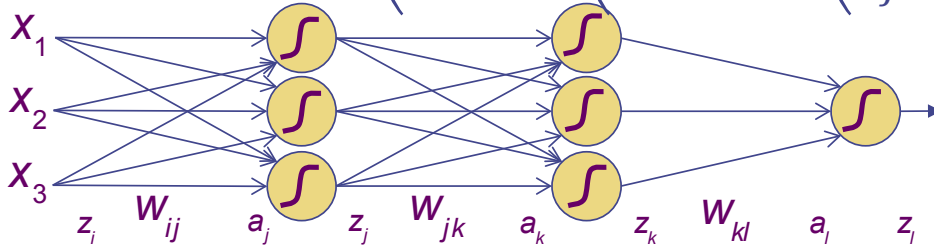$z_i$  $w_{ij}$  $a_j$  $z_j$  $w_{jk}$  $a_k$  $z_k$  $w_{kl}$  $a_l$  $z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_l^n}\right]\left(\frac{\partial a_l^n}{\partial w_{kl}}\right) = \frac{1}{N}\sum_n[-(y^n - z_l^n)g'(a_l^n)]z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$

- Next, hidden layer derivative:

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right) = \frac{1}{N}\sum_n\left[\sum_l\frac{\partial L^n}{\partial a_l^n}\frac{\partial a_l^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right)$$

**Multivariate Chain Rule**

# Back-Propagation

$$R(\theta) = \frac{1}{N}\sum_{n=1}^{N}\tfrac{1}{2}\left(y^n - g\left(\sum_k w_{kl}g\left(\sum_j w_{jk}g\left(\sum_i w_{ij}x_i^n\right)\right)\right)\right)^2$$



$x_1$
$x_2$
$x_3$

$z_i$  $W_{ij}$  $a_j$  $z_j$  $W_{jk}$  $a_k$  $z_k$  $W_{kl}$  $a_l$  $z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_l^n}\right]\left(\frac{\partial a_l^n}{\partial w_{kl}}\right) = \frac{1}{N}\sum_n[-(y^n - z_l^n)g'(a_l^n)]z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$
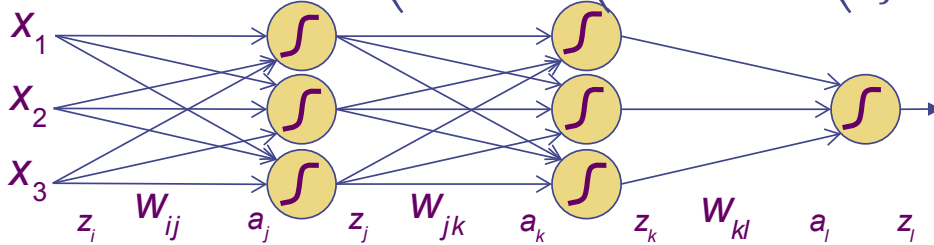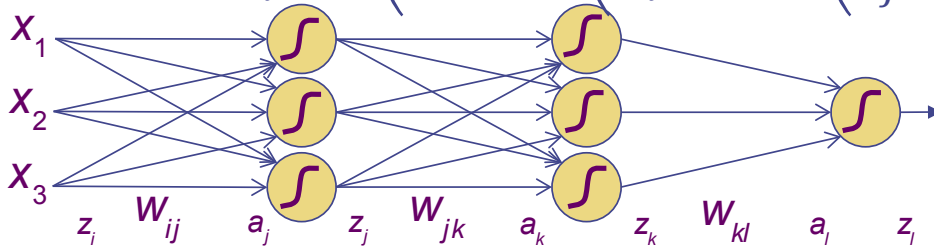
• Next, hidden layer derivative:

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right) = \frac{1}{N}\sum_n\left[\sum_l\frac{\partial L^n}{\partial a_l^n}\frac{\partial a_l^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right)$$

$$= \frac{1}{N}\sum_n\left[\sum_l \delta_l^n \boxed{\frac{\partial a_l^n}{\partial a_k^n}}\right](z_j^n)$$

$$\sum_k w_{kl}\, g'(a_k^n)$$

recall $a_l = \sum_k w_{kl}g(a_k)$

# Back-Propagation

$$R(\theta) = \frac{1}{N}\sum_{n=1}^{N}\tfrac{1}{2}\left(y^n - g\left(\sum_k w_{kl}g\left(\sum_j w_{jk}g\left(\sum_i w_{ij}x_i^n\right)\right)\right)\right)^2$$

$x_1$

$x_2$

$x_3$

$z_i$  $W_{ij}$  $a_j$  $z_j$  $W_{jk}$  $a_k$  $z_k$  $W_{kl}$  $a_l$  $z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_l^n}\right]\left(\frac{\partial a_l^n}{\partial w_{kl}}\right) = \frac{1}{N}\sum_n[-(y^n - z_l^n)g'(a_l^n)]z_k^n = \frac{\sum_n\delta_l^n z_k^n}{N}$$

- Next, hidden layer derivative:

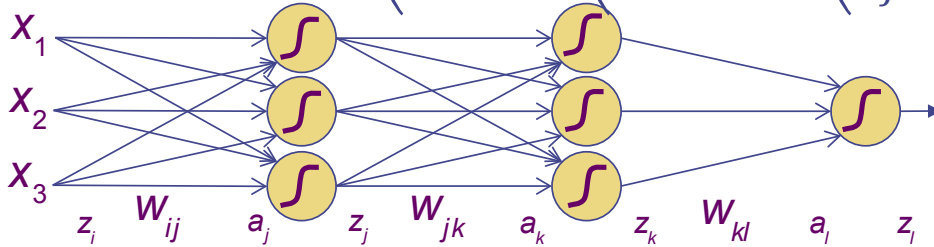$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N}\sum_n\left[\frac{\partial L^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right) = \frac{1}{N}\sum_n\left[\sum_l\frac{\partial L^n}{\partial a_l^n}\frac{\partial a_l^n}{\partial a_k^n}\right]\left(\frac{\partial a_k^n}{\partial w_{jk}}\right)$$

$$= \frac{1}{N}\sum_n\left[\sum_l\delta_l^n\frac{\partial a_l^n}{\partial a_k^n}\right](z_j^n) = \frac{1}{N}\sum_n\left[\sum_l\delta_l^n w_{kl}g'(a_k^n)\right]z_j^n$$

recall $a_l = \sum_k w_{kl}g(a_k)$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} \tfrac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$

$x_1$
$x_2$
$x_3$

$z_i \quad W_{ij} \quad a_j \quad z_j \quad W_{jk} \quad a_k \quad z_k \quad W_{kl} \quad a_l \quad z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right) = \frac{1}{N} \sum_n [-(y^n - z_l^n) g'(a_l^n)] z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$
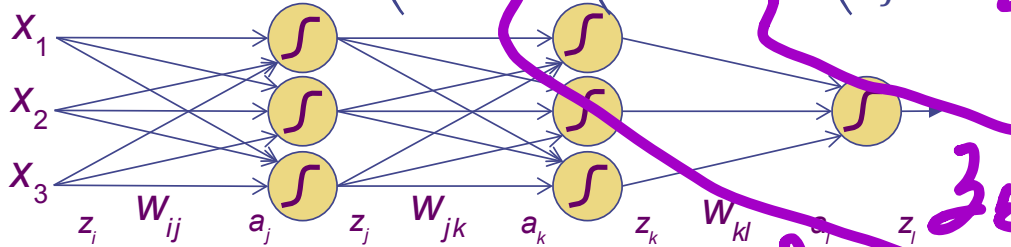
- Next, hidden layer derivative:

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_k^n} \right] \left( \frac{\partial a_k^n}{\partial w_{jk}} \right) = \frac{1}{N} \sum_n \left[ \sum_l \frac{\partial L^n}{\partial a_l^n} \frac{\partial a_l^n}{\partial a_k^n} \right] \left( \frac{\partial a_k^n}{\partial w_{jk}} \right)$$

$$= \frac{1}{N} \sum_n \left[ \sum_l \delta_l^n \frac{\partial a_l^n}{\partial a_k^n} \right] (z_j^n) = \frac{1}{N} \sum_n \left[ \sum_l \delta_l^n w_{kl} g'(a_k^n) \right] z_j^n = \frac{1}{N} \sum_n \delta_k^n z_j^n$$

recall $a_l = \sum_k w_{kl} g(a_k)$    **Define as** $\delta$

# Back-Propagation

$$R(\theta) = \frac{1}{N} \sum_{n=1}^{N} \tfrac{1}{2} \left( y^n - g \left( \sum_k w_{kl} g \left( \sum_j w_{jk} g \left( \sum_i w_{ij} x_i^n \right) \right) \right) \right)^2$$

$x_1$
$x_2$
$x_3$

$z_i$  $w_{ij}$  $a_j$  $z_j$  $w_{jk}$  $a_k$  $z_k$  $w_{kl}$  $a_l$  $z_l$

$$\frac{\partial R}{\partial w_{kl}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_l^n} \right] \left( \frac{\partial a_l^n}{\partial w_{kl}} \right) = \frac{1}{N} \sum_n [-(y^n - z_l^n)g'] z_k^n = \frac{\sum_n \delta_l^n z_k^n}{N}$$

$$\frac{\partial R}{\partial w_{jk}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_k^n} \right] \left( \frac{\partial a_k^n}{\partial w_{jk}} \right) = \frac{1}{N} \sum_n \left[ \sum_l \delta_l^n w_{kl} g'(a_k^n) \right] z_j^n \frac{\sum_n \delta_k^n z_j^n}{N}$$

What is this last z?

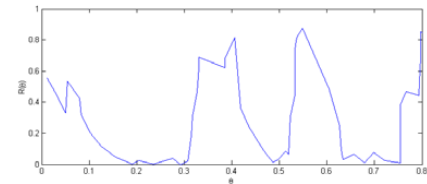• Any previous (input) layer derivative: repeat the formula!

$$\frac{\partial R}{\partial w_{ij}} = \frac{1}{N} \sum_n \left[ \frac{\partial L^n}{\partial a_j^n} \right] \left( \frac{\partial a_j^n}{\partial w_{ij}} \right) = \frac{1}{N} \sum_n \left[ \sum_k \frac{\partial L^n}{\partial a_k^n} \frac{\partial a_k^n}{\partial a_j^n} \right] \left( \frac{\partial a_j^n}{\partial w_{ij}} \right) = \frac{\sum_n \delta_j^n z_i^n}{N}$$

# Back-Propagation

•Again, take small step in direction opposite to gradient

$$w_{ij}^{t+1} = w_{ij}^t - \eta \frac{\partial R}{\partial w_{ij}} \qquad w_{jk}^{t+1} = w_{jk}^t - \eta \frac{\partial R}{\partial w_{jk}} \qquad w_{kl}^{t+1} = w_{kl}^t - \eta \frac{\partial R}{\partial w_{kl}}$$
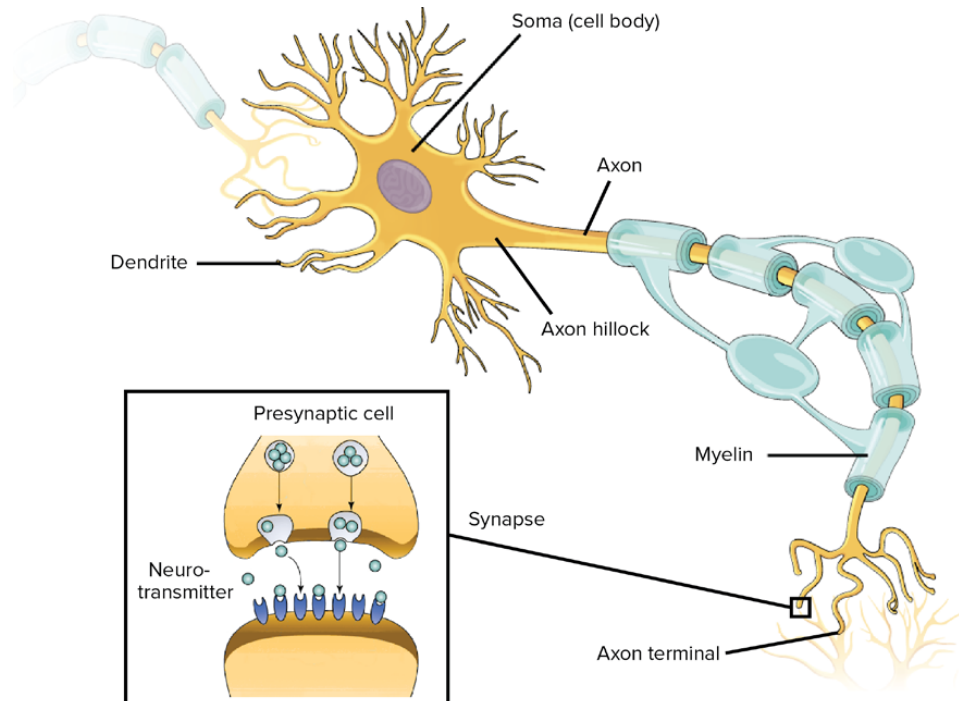
•Problems with back-prop
     is that MLP over-fits...

•Other problems: hard to interpret, black-box
•What are the hidden inner layers doing?

•Other main problem:       minimum training error
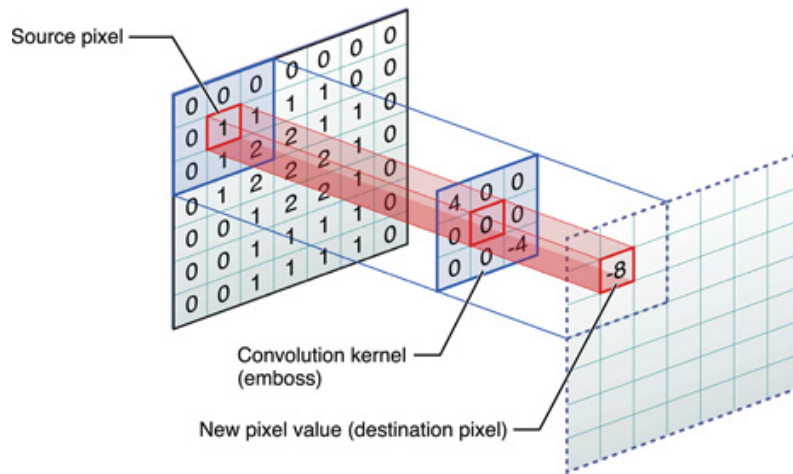                           not minimum testing error...

# Neural Networks - Upside

- Live neurons inspiration

# Neural Networks - Upside

- Live neurons inspiration

- Flexibility, parameter efficiency, modularity



Source pixel

Convolution kernel
(emboss)

New pixel value (destination pixel)

# Neural Networks - Upside

◆ Live neurons inspiration

◆ Flexibility, parameter efficiency, modularity

◆ Success across data-rich domains, tasks
  ▪ Vision, robotics, security, language, genomics...
  ▪ Classification, dimensionality reduction...