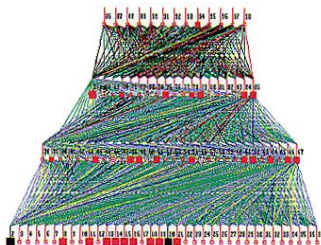# Back-Propagation

COMP4211

THE DEPARTMENT OF
**COMPUTER SCIENCE & ENGINEERING**
計算機科學及工程學系

# Back-Propagation



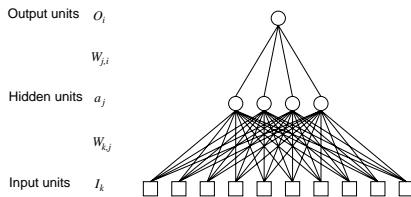Nonlinear activation functions + multi-layer networks

- requires more sophisticated learning algorithms

Back-propagation

## Idea: Gradient descent

- start with initial value for w
- repeat until convergence
  - compute the gradient vector of the error function for current w
  - move in the opposite direction
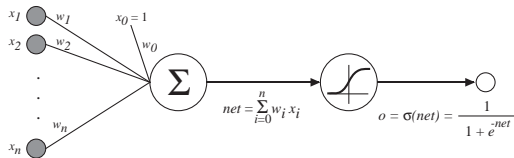
# How to Compute the Gradient?



- define an error function which is a differentiable function of the network outputs
- network with differentiable activation functions
  $\rightarrow$ outputs are differentiable functions of input and of the weights (and biases)
- $\rightarrow$ error is a differentiable function of the weights

## chain rule

$$x \rightarrow u \rightarrow y, \qquad \frac{dy}{dx} = \frac{dy}{du}\frac{du}{dx}$$
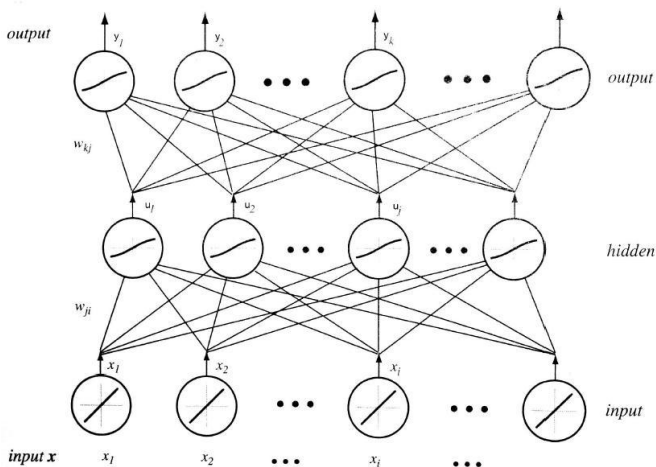
# Network with One Sigmoid Unit

we first derive gradient decent rules to train one sigmoid unit



$$\frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_i}\frac{1}{2}(t-o)^2 = \frac{1}{2}2(t-o)\frac{\partial}{\partial w_i}(t-o)$$

$$= (t-o)\left(-\frac{\partial o}{\partial w_i}\right) = -(t-o)\frac{\partial o}{\partial net}\frac{\partial net}{\partial w_i}$$
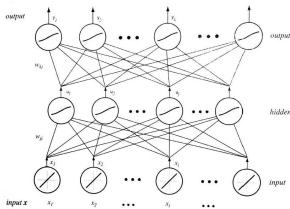
$$\frac{\partial o}{\partial net} = \frac{\partial \sigma(net)}{\partial net} = o(1-o), \qquad \frac{\partial net}{\partial w_i} = \frac{\partial(w'x)}{\partial w_i} = x_i$$

$$\frac{\partial E}{\partial w_i} = -(t-o)o(1-o)x_i$$

- $N_o$ output units
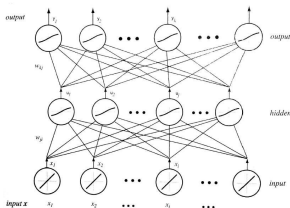- error for one training example: $E_d = \sum_{k=1}^{N_o}(t_{d,k} - o_{d,k})^2$

# Error Gradient of Weights to Output Unit $k$



$w_{kj}$: weight for the link from unit $j$ to (output) unit $k$

$$
\begin{aligned}
\frac{\partial E_d}{\partial w_{kj}} &= \frac{1}{2} \frac{\partial}{\partial w_{kj}} \sum_{m=1}^{N_o} (t_{d,m} - o_{d,m})^2 \quad \text{(dropping } d \text{ for simplicity)} \\
&= \frac{1}{2} \frac{\partial}{\partial w_{kj}} (t_k - o_k)^2 = (t_k - o_k) \frac{\partial(-o_k)}{\partial w_{kj}} \\
&= -(t_k - o_k) \frac{\partial o_k}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{kj}} \\
&= -\underbrace{(t_k - o_k) o_k (1 - o_k)}_{\delta_k = -\frac{\partial E_d}{\partial net_k}} \cdot u_j \quad (u_j \text{ is the output of unit } j)
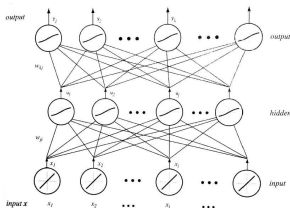\end{aligned}
$$

$w_{ji}$: weight for the link from unit $i$ to (hidden) unit $j$

$$\frac{\partial E_d}{\partial w_{ji}} = \frac{1}{2}\frac{\partial}{\partial w_{ji}}\sum_{k=1}^{N_o}(t_{d,k}-o_{d,k})^2 \quad \text{(dropping } d \text{ for simplicity)}$$

$$= \sum_{k=1}^{N_o}(t_k-o_k)\frac{\partial(-o_k)}{\partial w_{ji}} = -\sum_{k=1}^{N_o}(t_k-o_k)\frac{\partial o_k}{\partial net_k}\cdot\frac{\partial net_k}{\partial w_{ji}}$$

$$= -\sum_{k=1}^{N_o}(t_k-o_k)\frac{\partial o_k}{\partial net_k}\cdot\frac{\partial net_k}{\partial u_j}\cdot\frac{\partial u_j}{\partial w_{ji}}$$

$$= -\sum_{k=1}^{N_o}(t_k-o_k)o_k(1-o_k)\cdot w_{kj}\cdot\frac{\partial u_j}{\partial net_j}\cdot\frac{\partial net_j}{\partial w_{ji}}$$

# Error Gradient of Weights to Hidden Unit $j$...



$w_{ji}$: weight for the link from unit $i$ to (hidden) unit $j$

$$
\begin{aligned}
\frac{\partial E_d}{\partial w_{ji}} &= -\sum_{k=1}^{N_o}(t_k - o_k)o_k(1 - o_k) \cdot w_{kj} \cdot \frac{\partial u_j}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ji}} \\
&= -\sum_{k=1}^{N_o}(t_k - o_k)o_k(1 - o_k)w_{kj} \cdot u_j(1 - u_j) \cdot u_i \\
&= -\underbrace{\left[\sum_{k=1}^{N_o}\delta_k w_{kj}\right] \cdot u_j(1 - u_j)}_{\delta_j = -\frac{\partial E_d}{\partial net_j}} \cdot u_i
\end{aligned}
$$

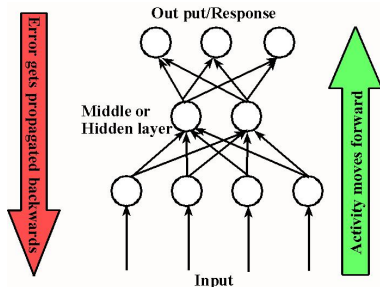- note that $i$ may be an input unit. In that case, $u_i$ is just $x_i$

# Weight Update Rule

- output weight: $\Delta w_{kj} = \eta \delta_k u_j$

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

- hidden weight: $\Delta w_{ji} = \eta \delta_j u_i$

$$\delta_j = u_j(1 - u_j) \left[ \sum_{k=1}^{N_o} \delta_k w_{kj} \right]$$



- we need to "propagate error back" when computing the gradient ector

## Backpropagation Algorithm (Stochastic Version)

**begin**
    initialize all weights to small random numbers;
    **repeat**
        **for** *each training example* **do**
            /* propagate input forward                    */
            input the example and compute the network outputs;
            /* propagate errors backward               */
            **for** *each output unit k* **do** $\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$;
            **for** *each hidden unit j* **do** $\delta_j \leftarrow o_j(1 - o_j)\sum_{k=1}^{N_o} w_{kj}\, \delta_k$;
            /* update weights                            */
            **for** *each network weight $w_{ji}$ (weight from i to j)* **do**
                $\Delta w_{ji} = \eta \delta_j u_i$ ;
                $w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$;
            **end**
        **end**
    **until** *convergence*;
**end**

# Backpropagation Algorithm (Batch Version)

**begin**
    initialize all weights to small random numbers;
    **repeat**
        **for** *each $(i, j)$* **do** initialize each $\Delta w_{ji}$ to zero ;
        **for** *each training example* **do**
            /* propagate input forward                */
            input the example and compute the network outputs;
            /* propagate errors backward            */
            **for** *each output unit k* **do** $\delta_k \leftarrow o_k(1 - o_k)(t_k - o_k)$;
            **for** *each hidden unit j* **do** $\delta_j \leftarrow o_j(1 - o_j) \sum_{k=1}^{N_o} w_{kj} \, \delta_k$;
            **for** *each $(i, j)$* **do** $\Delta w_{ji} \leftarrow \Delta w_{ji} + \eta \delta_j u_i$;
        **end**
        /* update weights                         */
        **for** *each network weight $w_{ji}$ (weight from i to j)* **do**
            $w_{ji} \leftarrow w_{ji} + \Delta w_{ji}$;
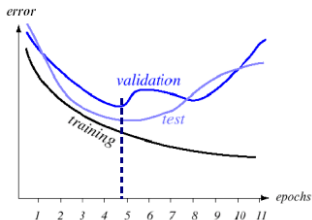        **end**
    **until** *convergence*;
**end**

how to initialize the weight values?

- initialize to some small random values

when to stop training?

1. after a fixed number of iterations through the loop
2. once the training error falls below some threshold
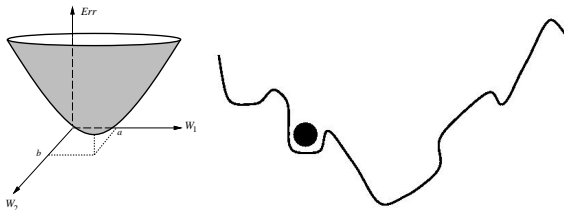3. stop at a minimum of the error on the validation set

# Speed

- testing is fast
- training can be very slow in networks with multiple hidden layers
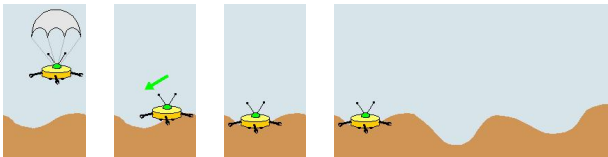
How to speed up BP training?

1. use of momentum term
   - give each weight some inertia or momentum

   $$\Delta w_{ji}(t+1) = -\eta \frac{\partial E}{\partial w_{ji}} + \alpha \Delta w_{ji}(t)$$

   - $0 < \alpha < 1$: momentum parameter (e.g., $\alpha = 0.9$)
2. dynamic adapt $\eta$
3. higher-order information of error surface
4. more sophisticated optimization algorithms

The error surface can have multiple local minima
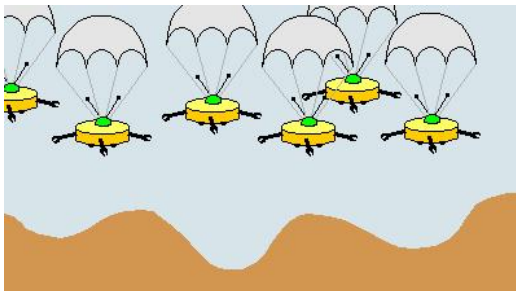


Gradient descent is only guaranteed to converge toward some local minimum, and **not** necessarily to the global minimum

how to escape from locally optimal solutions?



- train multiple networks using the same data, but initialize each network with different random weights