

Neural Networks and Back-Propagation Learning

11/2/05 1

Supervised Learning

- Produce desired outputs for training inputs
- Generalize reasonably & appropriately to other inputs
- Good example: pattern recognition
- Feedforward multilayer networks

11/2/05 2

Feedforward Network

11/2/05 3

Credit Assignment Problem

How do we adjust the weights of the hidden layers?

11/2/05 4

Adaptive System

System Evaluation Function
(Fitness, Figure of Merit)

11/2/05 5

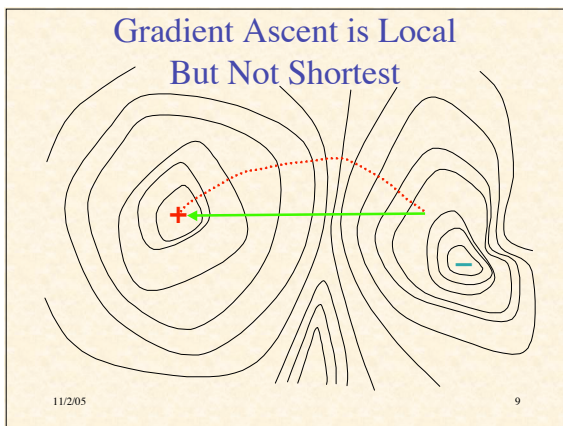
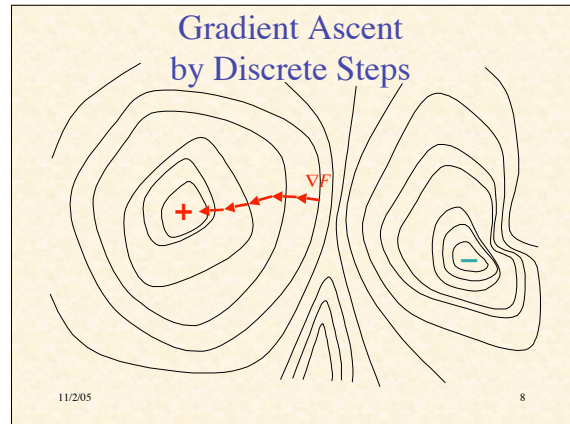
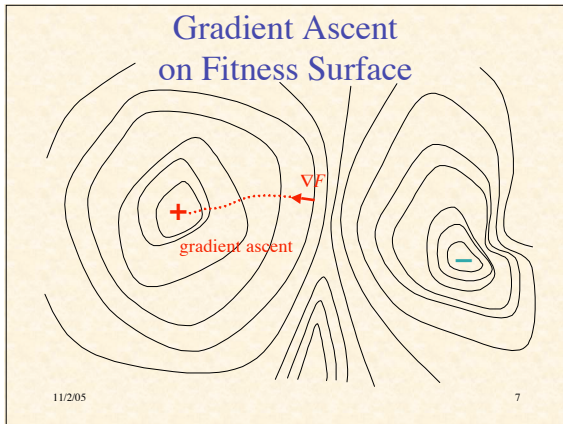
Gradient

$\frac{\partial F}{\partial P_k}$ measures how F is altered by variation of P_k

$$\nabla F = \begin{pmatrix} \frac{\partial F}{\partial P_1} \\ \vdots \\ \frac{\partial F}{\partial P_k} \\ \vdots \\ \frac{\partial F}{\partial P_m} \end{pmatrix}$$

∇F points in direction of maximum increase in F

11/2/05 6



Gradient Ascent Process

$$\dot{\mathbf{P}} = \eta \nabla F(\mathbf{P})$$

Change in fitness:

$$\dot{F} = \frac{dF}{dt} = \sum_{k=1}^m \frac{\partial F}{\partial P_k} \frac{dP_k}{dt} = \sum_{k=1}^m (\nabla F)_k \dot{P}_k$$

$$\dot{F} = \nabla F \cdot \dot{\mathbf{P}}$$

$$\dot{F} = \nabla F \cdot \eta \nabla F = \eta \|\nabla F\|^2 \geq 0$$

Therefore gradient ascent increases fitness (until reaches 0 gradient)

11/2/05 10

General Ascent in Fitness

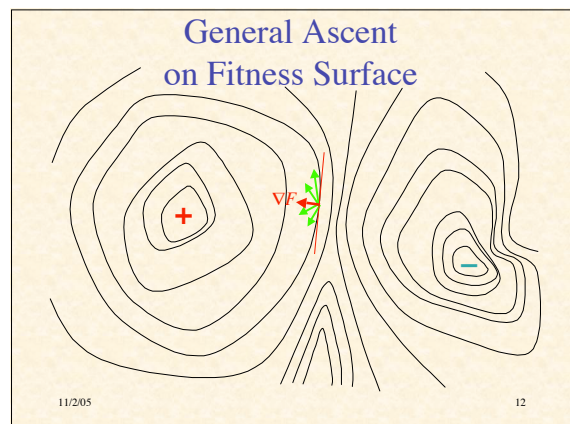
Note that any adaptive process $\mathbf{P}(t)$ will increase fitness provided:

$$0 < \dot{F} = \nabla F \cdot \dot{\mathbf{P}} = \|\nabla F\| \|\dot{\mathbf{P}}\| \cos \varphi$$

where φ is angle between ∇F and $\dot{\mathbf{P}}$

Hence we need $\cos \varphi > 0$
or $|\varphi| < 90^\circ$

11/2/05 11



Fitness as Minimum Error

Suppose for Q different inputs we have target outputs $\mathbf{t}^1, \dots, \mathbf{t}^Q$

Suppose for parameters \mathbf{P} the corresponding actual outputs are $\mathbf{y}^1, \dots, \mathbf{y}^Q$

Suppose $D(\mathbf{t}, \mathbf{y}) \in [0, \infty)$ measures difference between target & actual outputs

Let $E^q = D(\mathbf{t}^q, \mathbf{y}^q)$ be error on q th sample

Let $F(\mathbf{P}) = -\sum_{q=1}^Q E^q(\mathbf{P}) = -\sum_{q=1}^Q D[\mathbf{t}^q, \mathbf{y}^q(\mathbf{P})]$

11/2/05 13

Gradient of Fitness

$$\nabla F = \nabla \left(-\sum_q E^q \right) = -\sum_q \nabla E^q$$

$$\frac{\partial E^q}{\partial P_k} = \frac{\partial}{\partial P_k} D(\mathbf{t}^q, \mathbf{y}^q) = \sum_j \frac{\partial D(\mathbf{t}^q, \mathbf{y}^q)}{\partial y_j^q} \frac{\partial y_j^q}{\partial P_k}$$

$$= \frac{dD(\mathbf{t}^q, \mathbf{y}^q)}{d\mathbf{y}^q} \cdot \frac{\partial \mathbf{y}^q}{\partial P_k}$$

$$= \nabla_{\mathbf{y}^q} D(\mathbf{t}^q, \mathbf{y}^q) \cdot \frac{\partial \mathbf{y}^q}{\partial P_k}$$

11/2/05 14

Jacobian Matrix

Define Jacobian matrix $\mathbf{J}^q = \begin{pmatrix} \frac{\partial y_1^q}{\partial P_1} & \dots & \frac{\partial y_1^q}{\partial P_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_n^q}{\partial P_1} & \dots & \frac{\partial y_n^q}{\partial P_m} \end{pmatrix}$

Note $\mathbf{J}^q \in \mathbb{R}^{n \times m}$ and $\nabla D(\mathbf{t}^q, \mathbf{y}^q) \in \mathbb{R}^{n \times 1}$

Since $(\nabla E^q)_k = \frac{\partial E^q}{\partial P_k} = \sum_j \frac{\partial y_j^q}{\partial P_k} \frac{\partial D(\mathbf{t}^q, \mathbf{y}^q)}{\partial y_j^q}$,

$\therefore \nabla E^q = (\mathbf{J}^q)^T \nabla D(\mathbf{t}^q, \mathbf{y}^q)$

11/2/05 15

Derivative of Squared Euclidean Distance

Suppose $D(\mathbf{t}, \mathbf{y}) = \|\mathbf{t} - \mathbf{y}\|^2 = \sum_i (t_i - y_i)^2$

$$\frac{\partial D(\mathbf{t} - \mathbf{y})}{\partial y_j} = \frac{\partial}{\partial y_j} \sum_i (t_i - y_i)^2 = \sum_i \frac{\partial (t_i - y_i)^2}{\partial y_j}$$

$$= \frac{d(t_j - y_j)^2}{d y_j} = -2(t_j - y_j)$$

$\therefore \frac{dD(\mathbf{t}, \mathbf{y})}{d\mathbf{y}} = 2(\mathbf{y} - \mathbf{t})$

11/2/05 16

Gradient of Error on q th Input

$$\frac{\partial E^q}{\partial P_k} = \frac{dD(\mathbf{t}^q, \mathbf{y}^q)}{d\mathbf{y}^q} \cdot \frac{\partial \mathbf{y}^q}{\partial P_k}$$

$$= 2(\mathbf{y}^q - \mathbf{t}^q) \cdot \frac{\partial \mathbf{y}^q}{\partial P_k}$$

$$= 2 \sum_j (y_j^q - t_j^q) \frac{\partial y_j^q}{\partial P_k}$$

11/2/05 17

Recap

$$\dot{\mathbf{P}} = \eta \sum_q (\mathbf{J}^q)^T (\mathbf{t}^q - \mathbf{y}^q)$$

To know how to decrease the differences between actual & desired outputs,

we need to know elements of Jacobian, $\frac{\partial y_j^q}{\partial P_k}$,

which says how j th output varies with k th parameter (given the q th input)

The Jacobian depends on the specific form of the system, in this case, a feedforward neural network

11/2/05 18

Equations

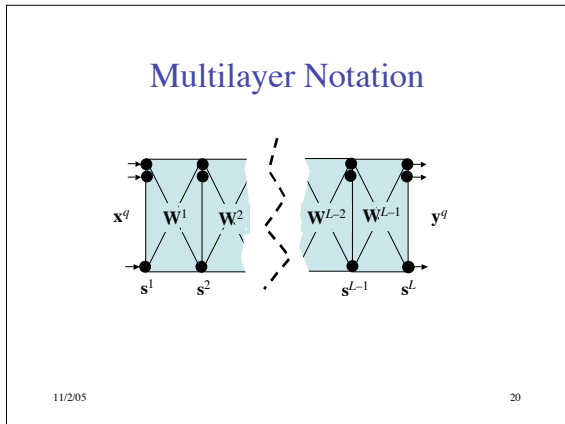
Net input:
$$h_i = \left(\sum_{j=1}^n w_{ij} s_j \right) - \theta$$

$\mathbf{h} = \mathbf{W}\mathbf{s} - \theta$

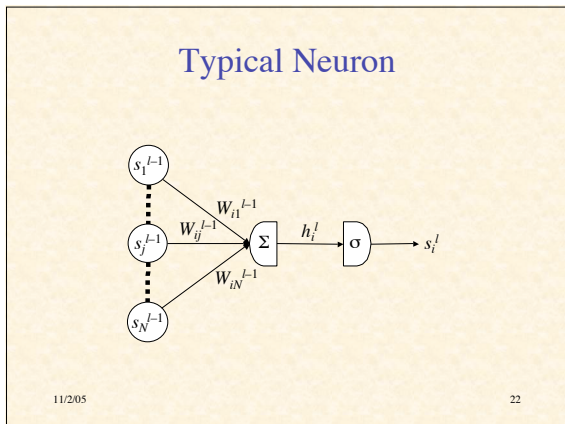
Neuron output:
$$s'_i = \sigma(h_i)$$

$\mathbf{s}' = \sigma(\mathbf{h})$

11/2/05 19



- ### Notation
- L layers of neurons labeled $1, \dots, L$
 - N_l neurons in layer l
 - \mathbf{s}^l = vector of outputs from neurons in layer l
 - input layer $\mathbf{s}^1 = \mathbf{x}^q$ (the input pattern)
 - output layer $\mathbf{s}^L = \mathbf{y}^q$ (the actual output)
 - \mathbf{W}^l = weights between layers l and $l+1$
 - Problem: find how outputs y_i^q vary with weights W_{jk}^l ($l = 1, \dots, L-1$)
- 11/2/05 21



Error Back-Propagation

We will compute $\frac{\partial E^q}{\partial W_{ij}^l}$ starting with last layer ($l = L - 1$) and working back to earlier layers ($l = L - 2, \dots, 1$)

11/2/05 23

Delta Values

Convenient to break derivatives by chain rule:

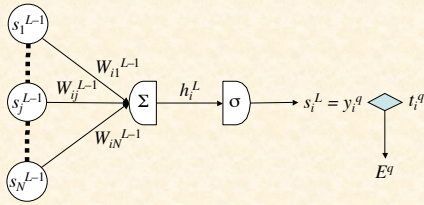
$$\frac{\partial E^q}{\partial W_{ij}^{l-1}} = \frac{\partial E^q}{\partial h_i^l} \frac{\partial h_i^l}{\partial W_{ij}^{l-1}}$$

Let $\delta_i^l = \frac{\partial E^q}{\partial h_i^l}$

So
$$\frac{\partial E^q}{\partial W_{ij}^{l-1}} = \delta_i^l \frac{\partial h_i^l}{\partial W_{ij}^{l-1}}$$

11/2/05 24

Output-Layer Neuron



11/2/05

25

Output-Layer Derivatives (1)

$$\begin{aligned} \delta_i^L &= \frac{\partial E^q}{\partial h_i^L} = \frac{\partial}{\partial h_i^L} \sum_k (s_k^L - t_k^q)^2 \\ &= \frac{d(s_i^L - t_i^q)^2}{dh_i^L} = 2(s_i^L - t_i^q) \frac{ds_i^L}{dh_i^L} \\ &= 2(s_i^L - t_i^q) \sigma'(h_i^L) \end{aligned}$$

11/2/05

26

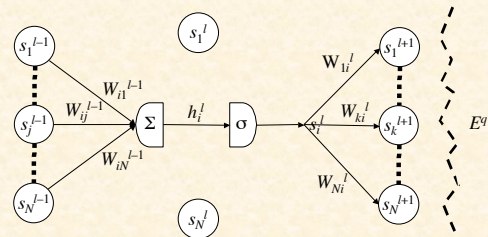
Output-Layer Derivatives (2)

$$\begin{aligned} \frac{\partial h_i^L}{\partial W_{ij}^{L-1}} &= \frac{\partial}{\partial W_{ij}^{L-1}} \sum_k W_{ik}^{L-1} s_k^{L-1} = s_j^{L-1} \\ \therefore \frac{\partial E^q}{\partial W_{ij}^{L-1}} &= \delta_i^L s_j^{L-1} \\ \text{where } \delta_i^L &= 2(s_i^L - t_i^q) \sigma'(h_i^L) \end{aligned}$$

11/2/05

27

Hidden-Layer Neuron



11/2/05

28

Hidden-Layer Derivatives (1)

$$\begin{aligned} \text{Recall } \frac{\partial E^q}{\partial W_{ij}^{l-1}} &= \delta_i^l \frac{\partial h_i^l}{\partial W_{ij}^{l-1}} \\ \delta_i^l &= \frac{\partial E^q}{\partial h_i^l} = \sum_k \frac{\partial E^q}{\partial h_k^{l+1}} \frac{\partial h_k^{l+1}}{\partial h_i^l} \\ \frac{\partial h_k^{l+1}}{\partial h_i^l} &= \frac{\partial \sum_m W_{km}^l s_m^l}{\partial h_i^l} = \frac{\partial W_{ki}^l s_i^l}{\partial h_i^l} = W_{ki}^l \frac{d\sigma(h_i^l)}{dh_i^l} = W_{ki}^l \sigma'(h_i^l) \\ \therefore \delta_i^l &= \sum_k \delta_k^{l+1} W_{ki}^l \sigma'(h_i^l) = \sigma'(h_i^l) \sum_k \delta_k^{l+1} W_{ki}^l \end{aligned}$$

11/2/05

29

Hidden-Layer Derivatives (2)

$$\begin{aligned} \frac{\partial h_i^l}{\partial W_{ij}^{l-1}} &= \frac{\partial}{\partial W_{ij}^{l-1}} \sum_k W_{ik}^{l-1} s_k^{l-1} = \frac{dW_{ij}^{l-1} s_j^{l-1}}{dW_{ij}^{l-1}} = s_j^{l-1} \\ \therefore \frac{\partial E^q}{\partial W_{ij}^{l-1}} &= \delta_i^l s_j^{l-1} \\ \text{where } \delta_i^l &= \sigma'(h_i^l) \sum_k \delta_k^{l+1} W_{ki}^l \end{aligned}$$

11/2/05

30

Derivative of Sigmoid

Suppose $s = \sigma(h) = \frac{1}{1 + \exp(-ah)}$ (logistic sigmoid)

$$D_h s = D_h [1 + \exp(-ah)]^{-1} = -[1 + \exp(-ah)]^{-2} D_h (1 + e^{-ah})$$

$$= -(1 + e^{-ah})^{-2} (-\alpha e^{-ah}) = \alpha \frac{e^{-ah}}{(1 + e^{-ah})^2}$$

$$= \alpha \frac{1}{1 + e^{-ah}} \frac{e^{-ah}}{1 + e^{-ah}} = \alpha s \left(\frac{1 + e^{-ah}}{1 + e^{-ah}} - \frac{1}{1 + e^{-ah}} \right)$$

$$= \alpha s(1 - s)$$

31

Summary of Back-Propagation Algorithm

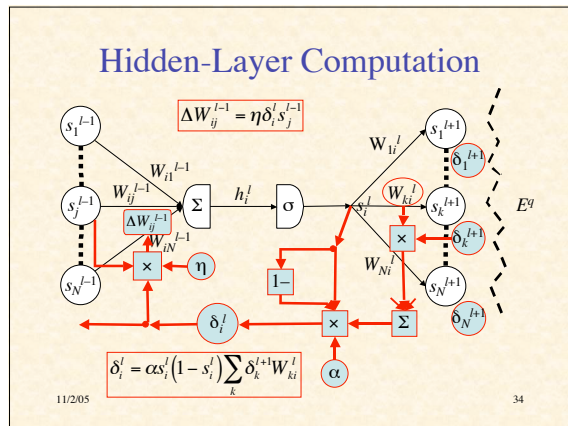
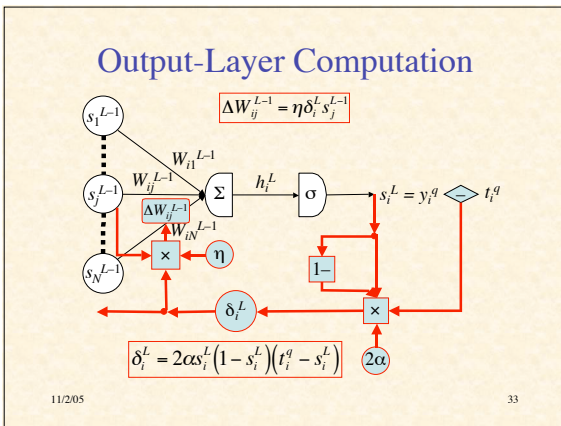
Output layer: $\delta_i^L = 2\alpha s_i^L (1 - s_i^L) (s_i^L - t_i^q)$

$$\frac{\partial E^q}{\partial W_{ij}^{L-1}} = \delta_i^L s_j^{L-1}$$

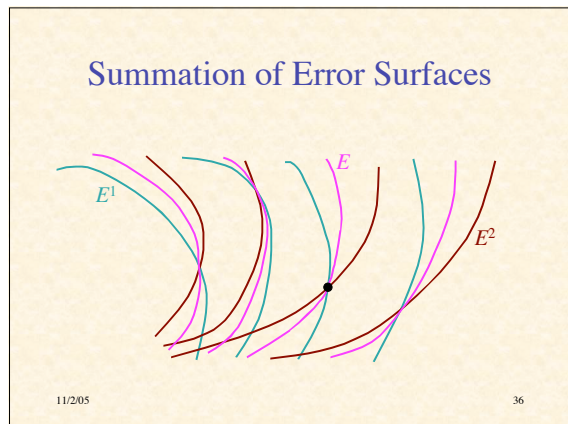
Hidden layers: $\delta_i^l = \alpha s_i^l (1 - s_i^l) \sum_k \delta_k^{l+1} W_{ki}^l$

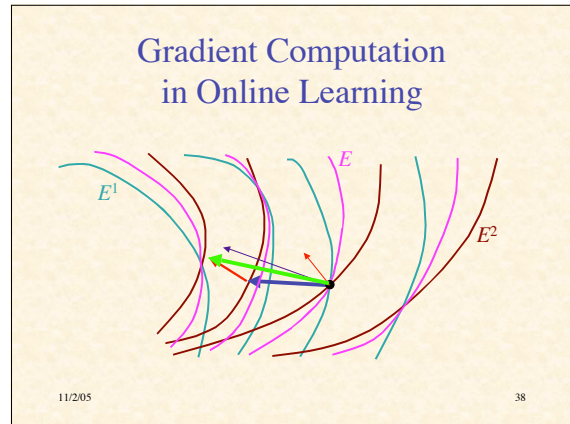
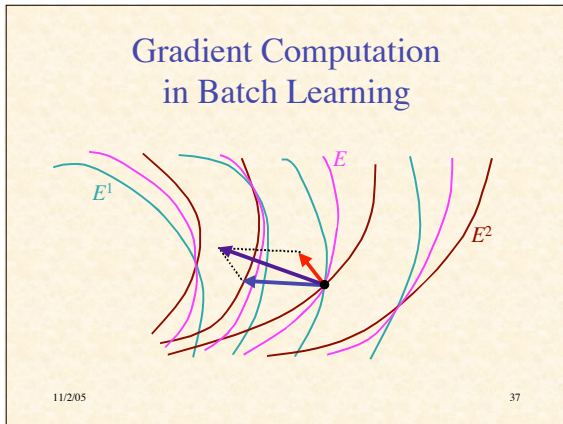
$$\frac{\partial E^q}{\partial W_{ij}^{l-1}} = \delta_i^l s_j^{l-1}$$

32



- ### Training Procedures
- **Batch Learning**
 - on each *epoch* (pass through all the training pairs),
 - weight changes for all patterns accumulated
 - weight matrices updated at end of epoch
 - accurate computation of gradient
 - **Online Learning**
 - weight are updated after back-prop of each training pair
 - usually randomize order for each epoch
 - approximation of gradient
 - Doesn't make much difference
- 35





The Golden Rule of Neural Nets

Neural Networks are the *second-best* way to do *everything*!

11/2/05 39