## Multilayer Notation



$\mathbf{x}^q$  $\mathbf{W}^1$  $\mathbf{W}^2$   $\mathbf{W}^{L-2}$  $\mathbf{W}^{L-1}$  $\mathbf{y}^q$

$\mathbf{s}^1$  $\mathbf{s}^2$   $\mathbf{s}^{L-1}$  $\mathbf{s}^L$

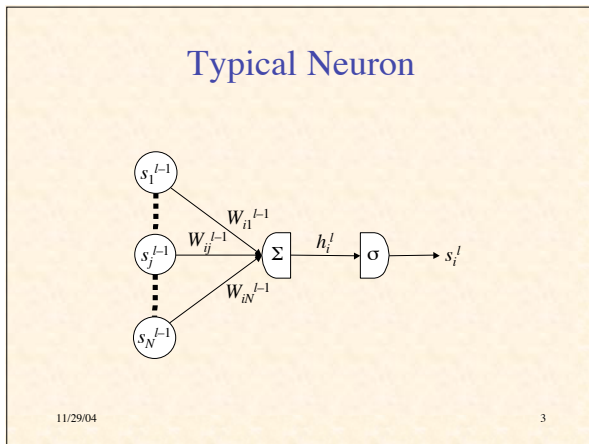11/29/04                                                                1

## Notation

- $L$ layers of neurons labeled 1, …, $L$
- $N_l$ 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 $W_{jk}^l$ ($l = 1, …, L-1$)

11/29/04                                                                2

## Typical Neuron



$s_1^{l-1}$

$W_{i1}^{l-1}$

$s_j^{l-1}$  $W_{ij}^{l-1}$  $\Sigma$  $h_i^l$  $\sigma$  $s_i^l$

$W_{iN}^{l-1}$

$s_N^{l-1}$

11/29/04                                                                3

## Error Back-Propagation

We will compute $\dfrac{\partial E^q}{\partial W_{ij}^l}$ starting with last layer ($l = L-1$)

and working back to earlier layers ($l = L-2,…,1$)

11/29/04                                                                4

## Delta Values

Convenient to break derivatives by chain rule :

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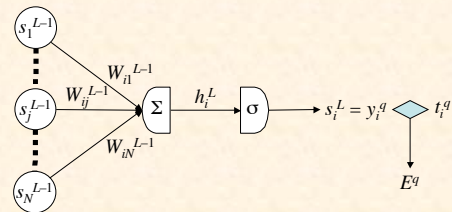
Let $\delta_i^l = \frac{\partial E^q}{\partial h_i^l}$

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

11/29/04                                                                                     5

## Output-Layer Neuron



11/29/04                                                                                     6

## Output-Layer Derivatives (1)

$$\delta_i^L = \frac{\partial E^q}{\partial h_i^L} = \frac{\partial}{\partial h_i^L} \sum_k \left( s_k^L - t_k^q \right)^2$$

$$= \frac{d\left( s_i^L - t_i^q \right)^2}{d h_i^L} = 2\left( s_i^L - t_i^q \right) \frac{d s_i^L}{d h_i^L}$$

$$= 2\left( s_i^L - t_i^q \right) \sigma'\left( h_i^L \right)$$

11/29/04                                                                                     7

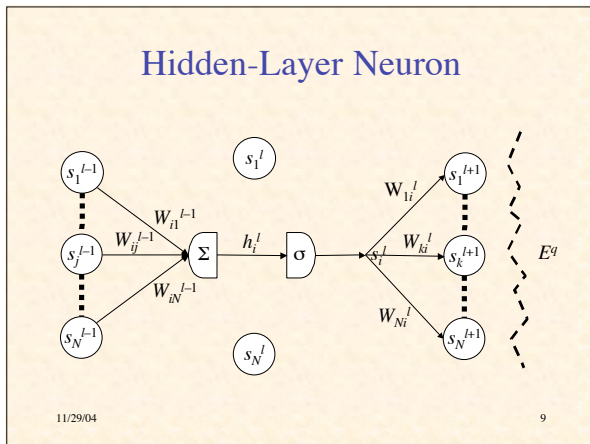## Output-Layer Derivatives (2)

$$\frac{\partial h_i^L}{\partial W_{ij}^{L-1}} = \frac{\partial}{\partial W_{ij}^{L-1}} \sum_k W_{ik}^{L-1} s_k^{L-1} = s_j^{L-1}$$

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

where $\delta_i^L = 2\left( s_i^L - t_i^q \right) \sigma'\left( h_i^L \right)$

11/29/04                                                                                     8

2

## Hidden-Layer Neuron



9

## Hidden-Layer Derivatives (1)

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

$$\delta_i^l = \frac{\partial E^q}{\partial h_i^l} = \sum_k \frac{\partial E^q}{\partial h_k^{l+1}} \frac{\partial h_k^{l+1}}{\partial h_i^l} = \sum_k \delta_k^{l+1} \frac{\partial h_k^{l+1}}{\partial h_i^l}$$

$$\frac{\partial h_k^{l+1}}{\partial h_i^l} = \frac{\partial \sum_m W_{km}^l s_m^l}{\partial h_i^l} = \frac{\partial W_{ki}^l s_i^l}{\partial h_i^l} = W_{ki}^l \frac{d\sigma(h_i^l)}{d h_i^l} = W_{ki}^l \sigma'(h_i^l)$$

$$\therefore \delta_i^l = \sum_k \delta_k^{l+1} W_{ki}^l \sigma'(h_i^l) = \sigma'(h_i^l) \sum_k \delta_k^{l+1} W_{ki}^l$$

10

## Hidden-Layer Derivatives (2)

$$\frac{\partial h_i^l}{\partial W_{ij}^{l-1}} = \frac{\partial}{\partial W_{ij}^{l-1}} \sum_k W_{ik}^{l-1} s_k^{l-1} = \frac{d W_{ij}^{l-1} s_j^{l-1}}{d W_{ij}^{l-1}} = s_j^{l-1}$$

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

where $\delta_i^l = \sigma'(h_i^l) \sum_k \delta_k^{l+1} W_{ki}^l$

11

## Derivative of Sigmoid

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

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

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

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

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

12

3

## Summary of Back-Propagation Algorithm

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

$$\frac{\partial E^q}{\partial W_{ij}^{L-1}} = \delta_i^L s_j^{L-1}$$
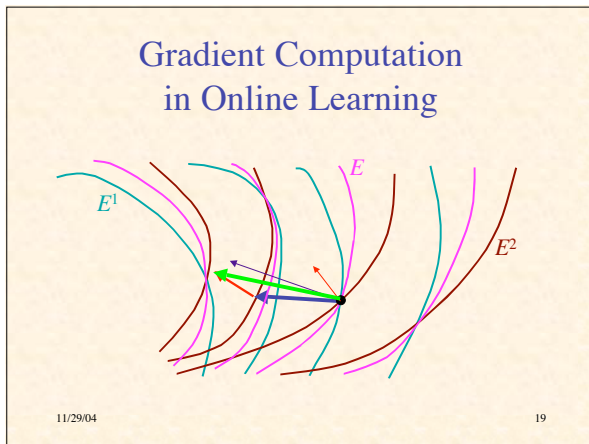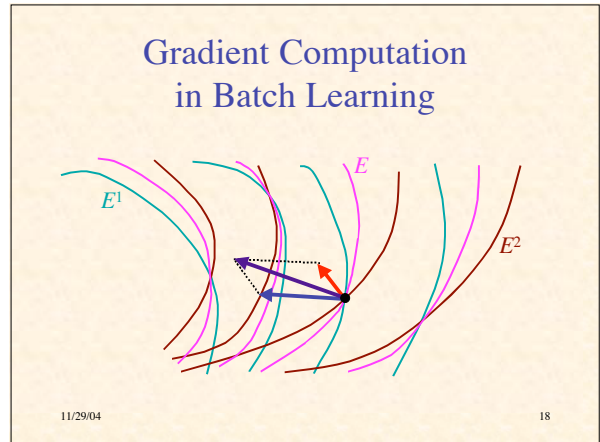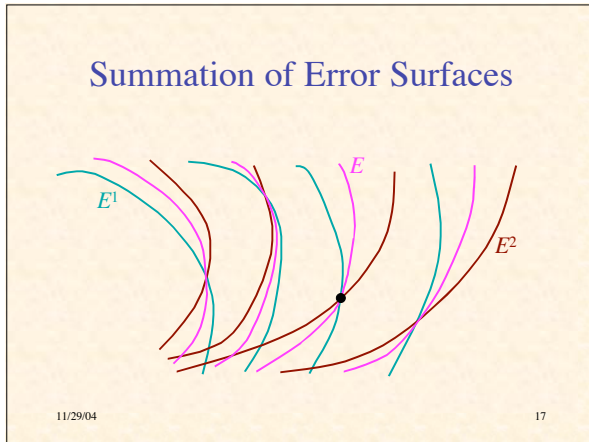
Hidden layers : $\delta_i^l = \alpha s_i^l\left(1 - s_i^l\right)\sum_k \delta_k^{l+1} W_{ki}^l$

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

11/29/04                                                                 13

## Output-Layer Computation



$$\Delta W_{ij}^{L-1} = \eta\delta_i^L s_j^{L-1}$$

$$\delta_i^L = 2\alpha s_i^L\left(1 - s_i^L\right)\left(t_i^q - s_i^L\right)$$

11/29/04                                                                 14

## Hidden-Layer Computation



$$\Delta W_{ij}^{l-1} = \eta\delta_i^l s_j^{l-1}$$

$$\delta_i^l = \alpha s_i^l\left(1 - s_i^l\right)\sum_k \delta_k^{l+1} W_{ki}^l$$

11/29/04                                                                 15

## 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

11/29/04                                                                 16

## Summation of Error Surfaces

$E^1$   $E$   $E^2$

11/29/04                                                                                          17

## Gradient Computation in Batch Learning

$E^1$   $E$   $E^2$

11/29/04                                                                                          18

## Gradient Computation in Online Learning

$E^1$   $E$   $E^2$

11/29/04                                                                                          19

## The Golden Rule of Neural Nets

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

11/29/04                                                                                          20

## VIII. Review of Key Concepts

## Complex Systems

- Many interacting elements
- Local vs. global order: entropy
- Scale (space, time)
- Phase space
- Difficult to understand
- Open systems

## Many Interacting Elements

- Massively parallel
- Distributed information storage & processing
- Diversity
  - avoids premature convergence
  - avoids inflexibility

## Complementary Interactions

- Positive feedback / negative feedback
- Amplification / stabilization
- Activation / inhibition
- Cooperation / competition
- Positive / negative correlation

## Emergence & Self-Organization

- Microdecisions lead to macrobehavior
- Circular causality (macro / micro feedback)
- Coevolution
  - predator/prey, Red Queen effect
  - gene/culture, niche construction, Baldwin effect

## Pattern Formation

- Excitable media
- Amplification of random fluctuations
- Symmetry breaking
- Specific difference vs. generic identity
- Automatically adaptive

## Stigmergy

- Continuous (quantitative)
- Discrete (qualitative)
- Coordinated algorithm
  - non-conflicting
  - sequentially linked

## Emergent Control

- Stigmergy
- Entrainment (distributed synchronization)
- Coordinated movement
  - through attraction, repulsion, local alignment
  - in concrete or abstract space
- Cooperative strategies
  - nice & forgiving, but reciprocal
  - evolutionarily stable strategy

## Attractors

- Classes
  - point attractor
  - cyclic attractor
  - chaotic attractor
- Basin of attraction
- Imprinted patterns as attractors
  - pattern restoration, completion, generalization, association

11/29/04                                                                29

## Wolfram's Classes

- Class I: point
- Class II: cyclic
- Class III: chaotic
- Class IV: complex (edge of chaos)
  - persistent state maintenance
  - bounded cyclic activity
  - global coordination of control & information
  - order for free

11/29/04                                                                30

## Energy / Fitness Surface

- Descent on energy surface / ascent on fitness surface
- Lyapunov theorem to prove asymptotic stability / convergence
- Soft constraint satisfaction / relaxation
- Gradient (steepest) ascent / descent
- Adaptation & credit assignment

11/29/04                                                                31

## Biased Randomness

- Exploration vs. exploitation
- Blind variation & selective retention
- Innovation vs. incremental improvement
- Pseudo-temperature
- Diffusion
- Mixed strategies

11/29/04                                                                32

## Natural Computation

- Tolerance to noise, error, faults, damage
- Generality of response
- Flexible response to novelty
- Adaptability
- Real-time response
- Optimality is secondary

11/29/04                                                                      33

## Student Course Evaluation!
(Do it online)

11/29/04                                                                      34