

Lecture 28

12/3/07 1

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)

12/3/07 2

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 $|\varphi| < 90^\circ$

12/3/07 3

General Ascent on Fitness Surface

12/3/07 4

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

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

12/3/07 5

Gradient of Fitness

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

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

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

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

12/3/07 6

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

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

12/3/07 7

Derivative of Squared Euclidean Distance

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

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

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

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

12/3/07 8

Gradient of Error on q^{th} Input

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

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

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

$\nabla E^q = 2(\mathbf{J}^q)^T (\mathbf{y}^q - \mathbf{t}^q)$

12/3/07 9

Recap

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

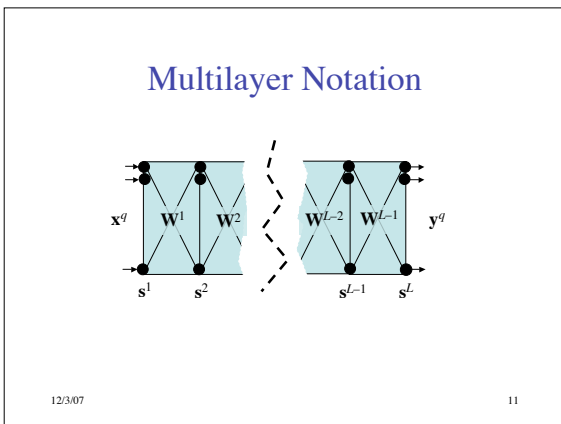
To know how to decrease the differences between actual & desired outputs,

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)

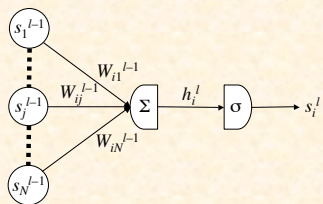
The Jacobian depends on the specific form of the system, in this case, a feedforward neural network

12/3/07 10



- ### 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 ($l = 1, \dots, L-1$)
- 12/3/07 12

Typical Neuron



12/3/07

13

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 ($l = L - 2, \dots, 1$)

12/3/07

14

Delta Values

Convenient to break derivatives by chain rule :

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

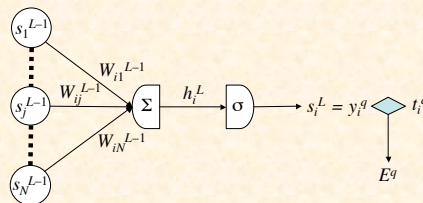
Let $\delta_i^l = \frac{\partial E^q}{\partial h_i^l}$

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

12/3/07

15

Output-Layer Neuron



12/3/07

16

Output-Layer Derivatives (1)

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

12/3/07

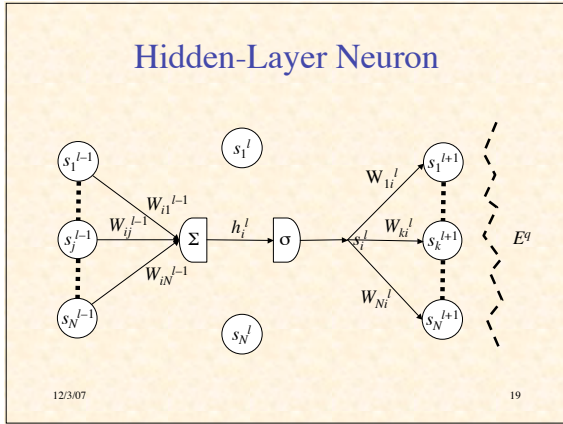
17

Output-Layer Derivatives (2)

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

12/3/07

18



Hidden-Layer Derivatives (1)

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

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

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

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

$12/3/07$ 20

Hidden-Layer Derivatives (2)

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

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

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

$12/3/07$ 21

Derivative of Sigmoid

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

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

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

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

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

$12/3/07$ 22

Summary of Back-Propagation Algorithm

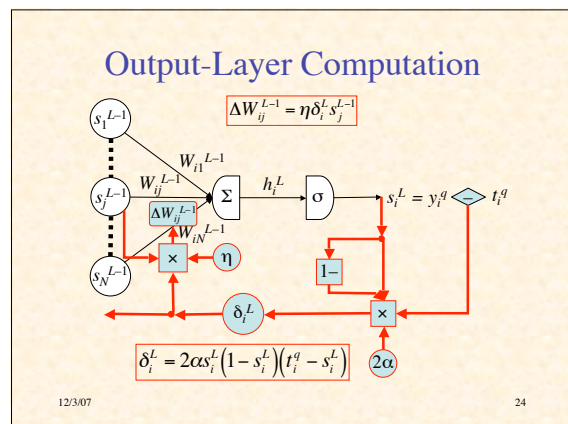
Output layer: $\delta_i^L = 2\alpha s_i^L (1 - s_i^L) (t_i^q - s_i^L)$

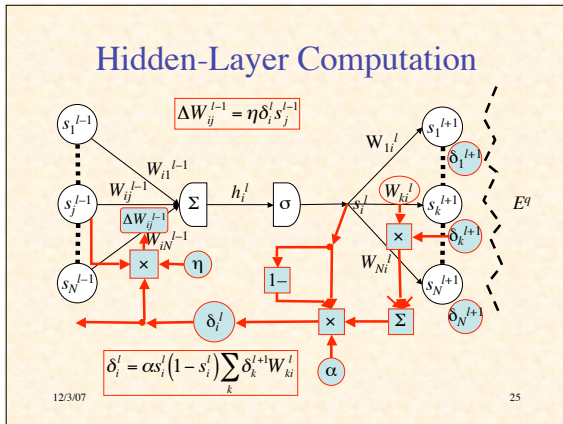
$$\frac{\partial E^q}{\partial W_{ij}^{L-1}} = \delta_i^L s_j^{L-1}$$

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

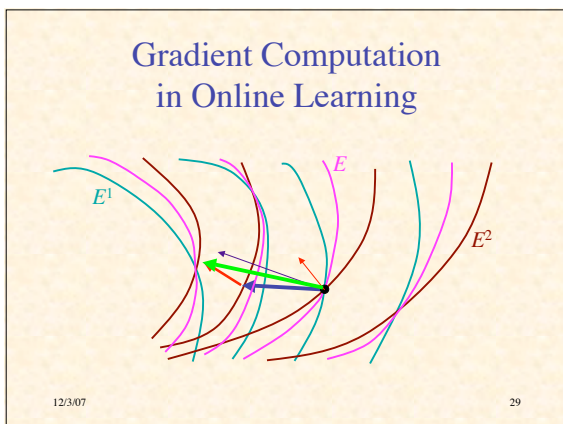
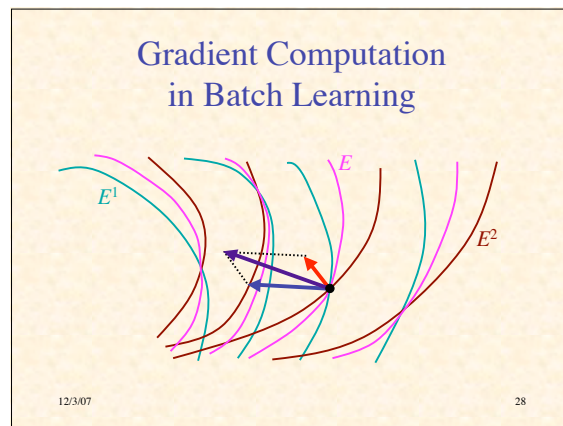
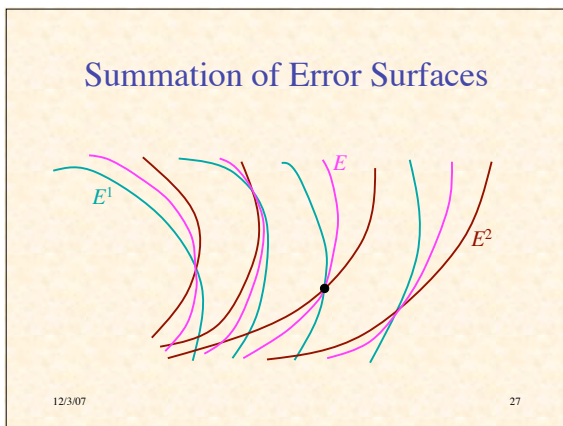
$$\frac{\partial E^q}{\partial W_{ij}^{l-1}} = \delta_i^l s_j^{l-1}$$

$12/3/07$ 23





- ### 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
- 12/3/07 26



The Golden Rule of Neural Nets

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

12/3/07 30

VIII. Review of Key Concepts

12/3/07

31

Complex Systems

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

12/3/07

32

Many Interacting Elements

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

12/3/07

33

Complementary Interactions

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

12/3/07

34

Emergence & Self-Organization

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

12/3/07

35

Pattern Formation

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

12/3/07

36

Stigmergy

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

12/3/07

37

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

12/3/07

38

Attractors

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

12/3/07

39

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

12/3/07

40

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

12/3/07

41

Biased Randomness

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

12/3/07

42

Natural Computation

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

12/3/07

43

Student Course Evaluation!
(We will do it in class this time)

12/3/07

44