

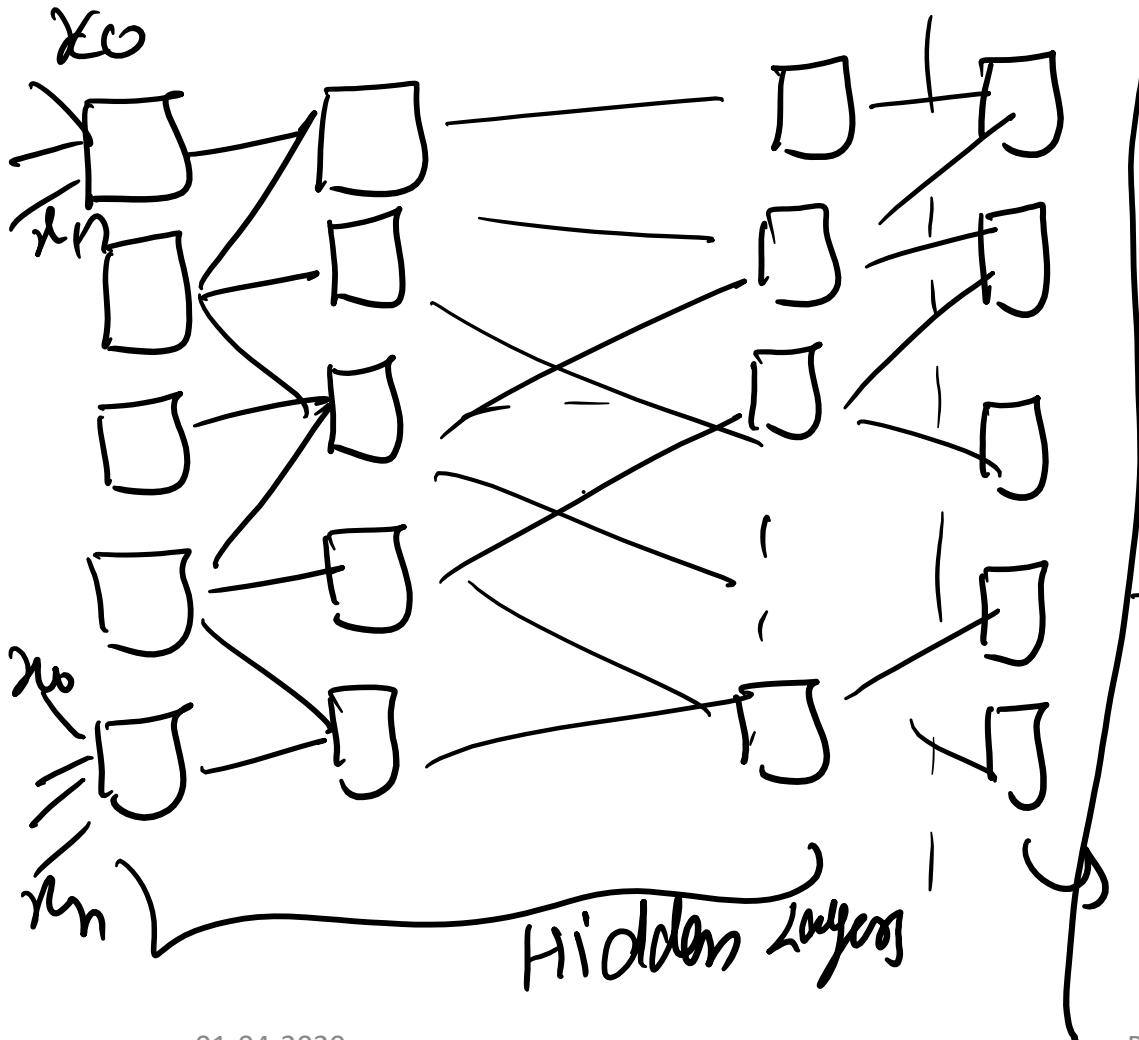
Machine Learning (COL 774)

Neural Networks: Backpropagation

Apr 1, 2020

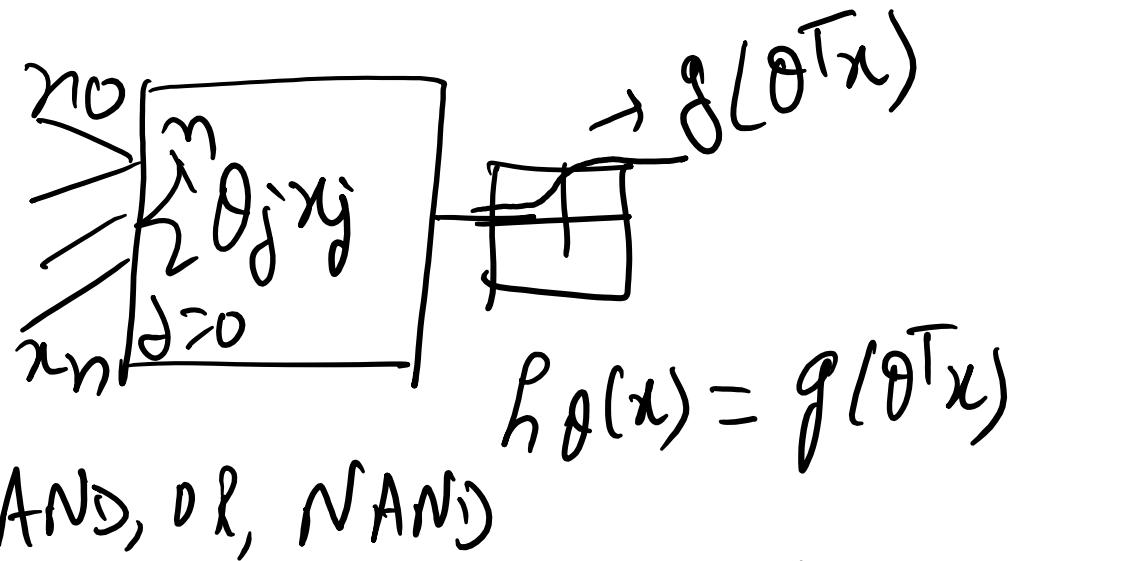
① Neural Networks :-

Last Class :- Perception



01-04-2020

Parag Singla @ IIT Delhi



↳ AND, OR, NAND

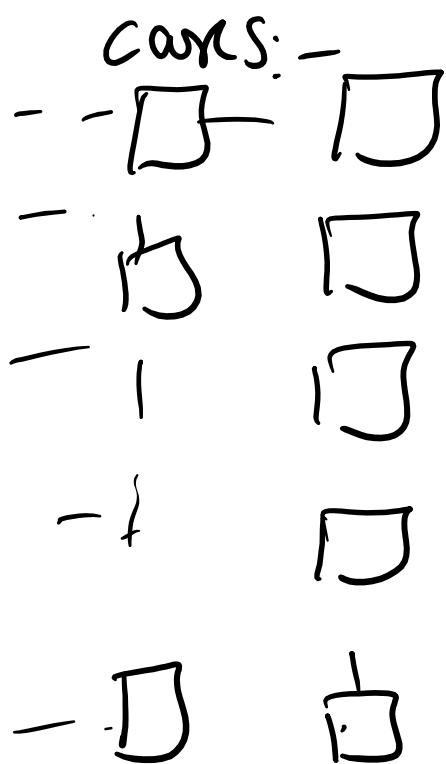
↳ XOR → Multiple such perceptions (arranged in a layered manner)

$$J(\theta) = \frac{1}{2m} \sum_{l=1}^L \sum_{i=1}^{m_l} (y_i^{(l)} - \theta_l)^2$$

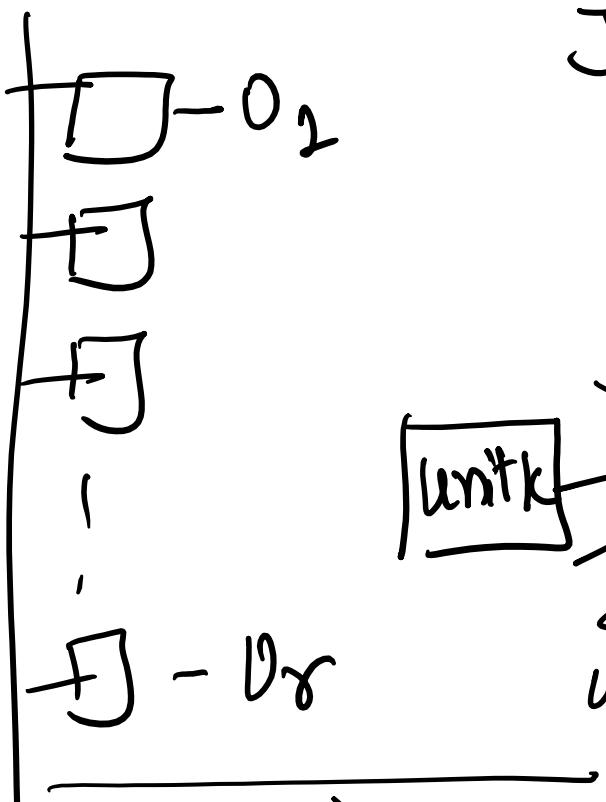
output layers

$\frac{\partial J(\theta)}{\partial \theta_{jk}}$ → θ_{jk} denotes weight of an interconnection

To compute the gradient, we will divide this problem into two cases:-



Hidden layers



$$\frac{\partial J(\theta)}{\partial \theta_{jk}}$$

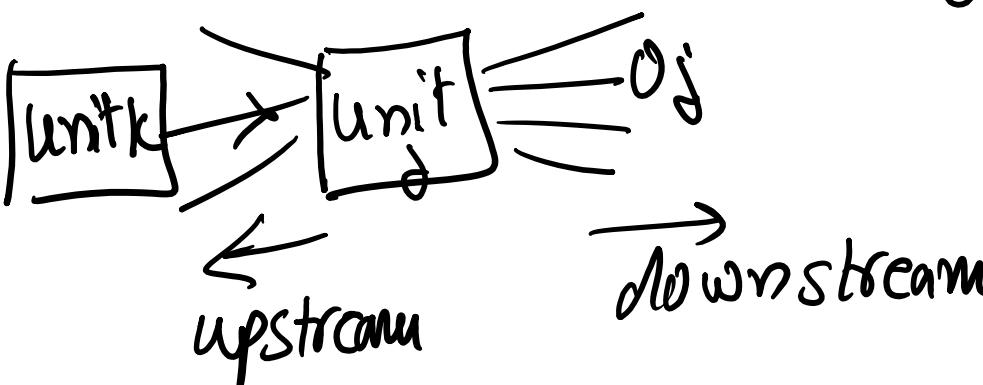
$$\frac{\partial J(\theta)}{\partial net_j}$$

Parag Singla @ IIT Delhi

$$net_j =$$

$$\sum_{k \in \text{UpNbrs}(j)} \theta_{jk} x_{jk}$$

θ_{jk} : weight of connection from unit k to j
 x_{jk} : input going from unit k to j

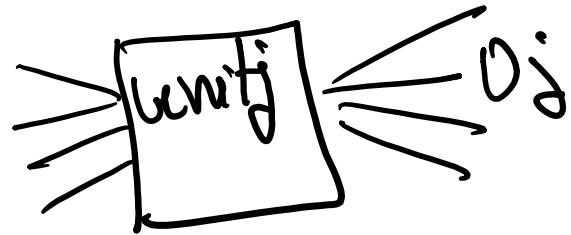


$$J(\theta) = \sum_{l=1}^L \left(y_l - \hat{y}_l \right)^2$$

$$m=1$$

$$o_j = g(\text{net}_j)$$

↳ sigmoid



$$\frac{\partial J(\theta)}{\partial \theta_{jk}} = \frac{\partial J(\theta)}{\partial \text{net}_j} \cdot \left[\frac{\partial \text{net}_j}{\partial \theta_{jk}} \right] = \frac{\partial J(\theta)}{\partial \text{net}_j} \cdot x_{jk}$$

$$\text{net}_j = \sum_{k \in \text{upNbrs}(j)} \theta_{jk} x_k$$

$$o_j = g(\text{net}_j)$$

Backpropagation

$$\frac{\partial o_j}{\partial \text{net}_j} = \frac{\partial g(\text{net}_j)}{\partial \text{net}_j} = g(\text{net}_j) (1 - g(\text{net}_j)) \\ = o_j (1 - o_j)$$

Our objective:-

$\frac{\partial J(\theta)}{\partial \text{net}_j}$ + units j in the network

(A) $j \in$ Output Layer

(B) $j \in$ Hidden Layer



(A) $\frac{\partial J(\theta)}{\partial \text{net}_j} = \frac{\partial}{\partial \text{net}_j} \left[\frac{1}{2} \sum_{c=1}^C (y_c - o_c)^2 \right] = \frac{\cancel{1/2} * \cancel{\sum} (y_j - o_j)}{\cancel{(c=1)}} (-1) * \frac{\partial o_j}{\partial \text{net}_j} = -(y_j - o_j) o_j(1-o_j)$

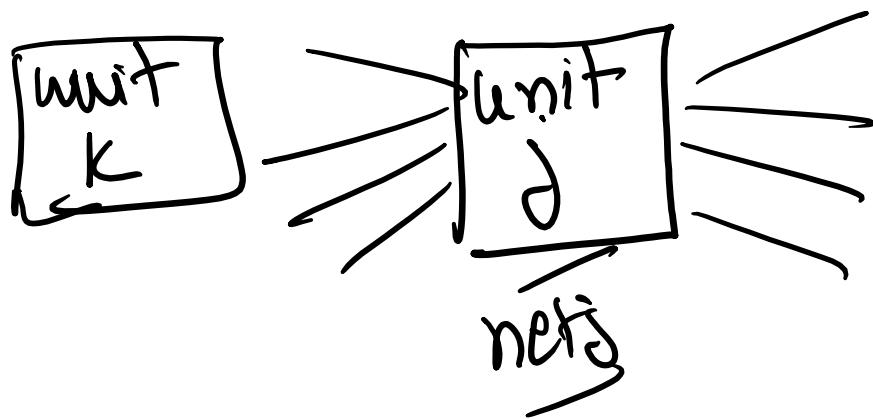
\in Output layers

Define:-
$$\delta_j = - \frac{\partial J(\theta)}{\partial \text{net}_j}$$

For output layers:-
$$\boxed{\delta_j = (y_j - o_j) o_j(1-o_j)}$$

(B) $\frac{\partial J(\theta)}{\partial \text{net}_j} \rightarrow \frac{\partial \epsilon}{\partial \text{hidden layer}}$





$$J(\theta)$$

unit
e

Aside:-

$$f(y_1 - \dots - y_K)$$

y_k is a function $x_1 - x_n$

$f(y_1(x_1 - x_n), y_2(x_1 - x_n))$
Objective :- compute

Express $\frac{\partial J(\theta)}{\partial \text{net}_j}$ in terms of
where $\ell \in \text{down Nbr}(j)$
 $\frac{\partial J(\theta)}{\partial \text{net}_e}$ has already been
computed

$$\frac{\partial f}{\partial x_j} = y_k(x_e - x_n) + x_j$$



$$\frac{\partial f(y_1, \dots, y_k)}{\partial x_j} = \sum_{k=1}^K \frac{\partial f(y_1, \dots, y_k)}{\partial y_k} \cdot \frac{\partial y_k}{\partial x_j}$$

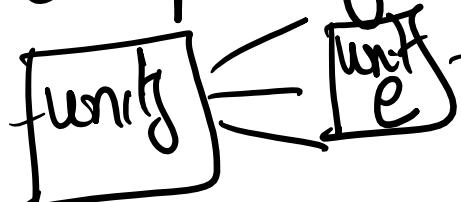
Each y_k is a fn. of (x_1, \dots, x_n)

$$\frac{f(y)}{\partial x} = \frac{\partial f(y)}{\partial y} \cdot \frac{\partial y}{\partial x} \rightarrow \text{standard chain rule}$$

y is a fn of x

extension
when defen-
dence is
through
multiple
variables

→ Back to computing

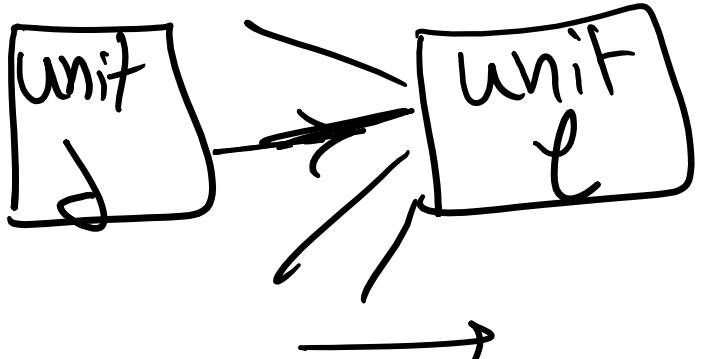
- 

$$\frac{\partial J(\theta)}{\partial \text{net}_j} = \sum_j \left[\frac{\partial J(\theta)}{\partial \text{net}_j} \right] \left[\frac{\partial \text{net}_j}{\partial \text{net}_j} \right]$$

EE down $Nb(x(j))$ → already computed



$$\frac{\partial \text{net}_e}{\partial \text{net}_j} = \frac{\partial \sum_{j \in \text{upNbr}(e)} x_{ej} \theta_{ej}}{\partial \text{net}_j}$$



$$= \frac{\partial \sum_{j \in \text{upNbr}(e)} \theta_j \cdot \theta_{ej}}{\partial \text{net}_j} = \frac{\partial \theta_j \cdot \theta_{ej}}{\partial \text{net}_j} = \theta_{ej} \cdot \frac{\partial \theta_j}{\partial \text{net}_j}$$

$$= \theta_{ej} \cdot (\theta_j (1 - \theta_j))$$

(B) $\frac{\partial J(\theta)}{\partial \text{net}_j} = \sum_{e \in \text{downNbr}(j)} \left[\frac{\partial J(\theta)}{\partial \text{net}_e} \right] \cdot \left[\frac{\partial \text{net}_e}{\partial \text{net}_j} \right]$ $\delta e j \circ_j (1 - \circ_j)$

Hidden layer

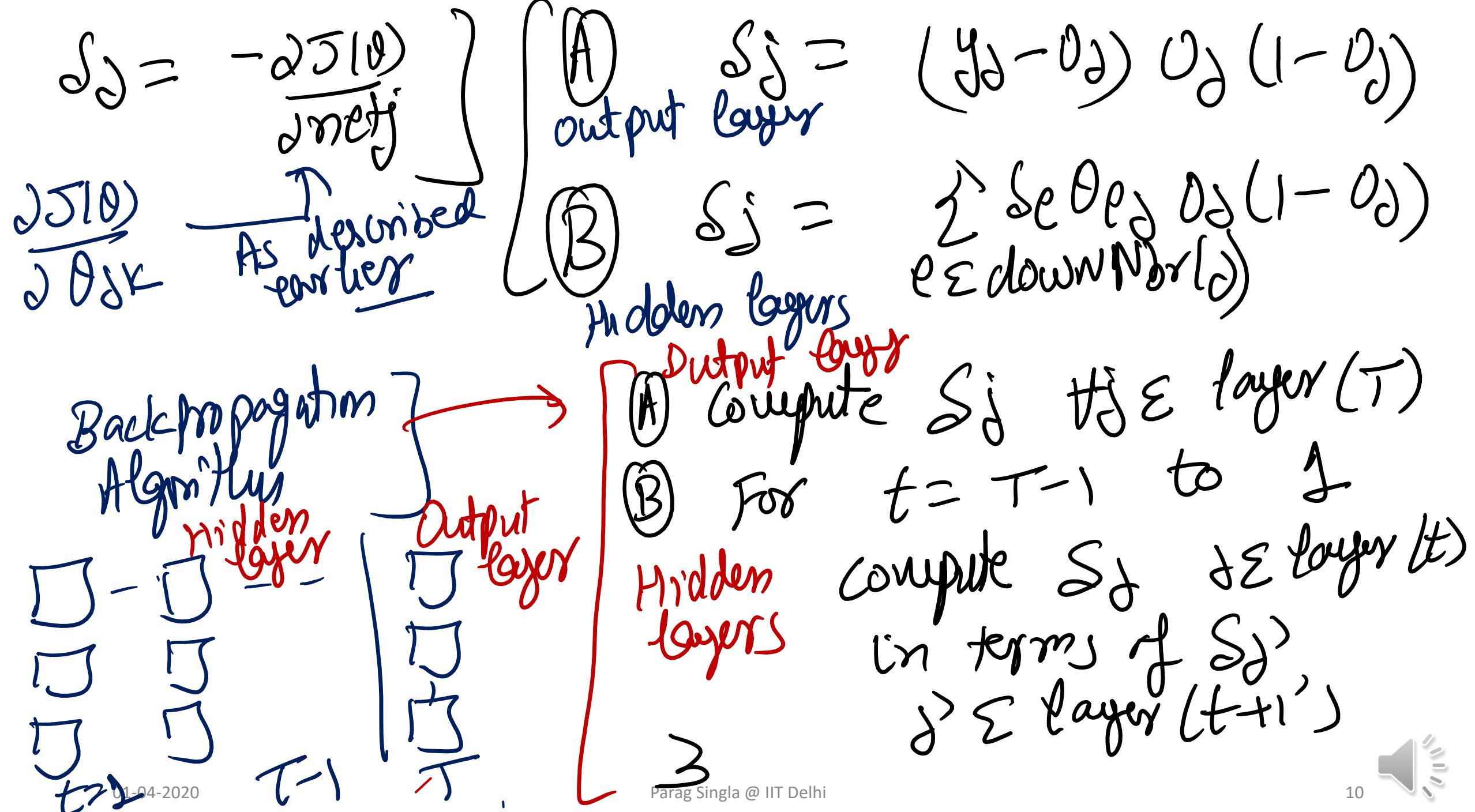
$$= \sum_{e \in \text{downNbr}(j)} -\delta e \cdot \circ_j \overline{\circ_e} \circ_j \circ_d (1 - \circ_j)$$

(A) $\frac{\partial J(\theta)}{\partial \text{net}_j} = \frac{(y_j - \circ_j)(-\circ_j)}{\circ_j(1 - \circ_j)}$

$\delta e = -\frac{\partial J(\theta)}{\partial \text{net}_e}$

$\frac{\partial J(\theta)}{\partial \text{net}_j} \quad \begin{cases} (A) (y_j - \circ_j)(-\circ_j) \\ (B) \sum_{e \in \text{downNbr}(j)} -\delta e \circ_j (1 - \circ_j) \end{cases}$

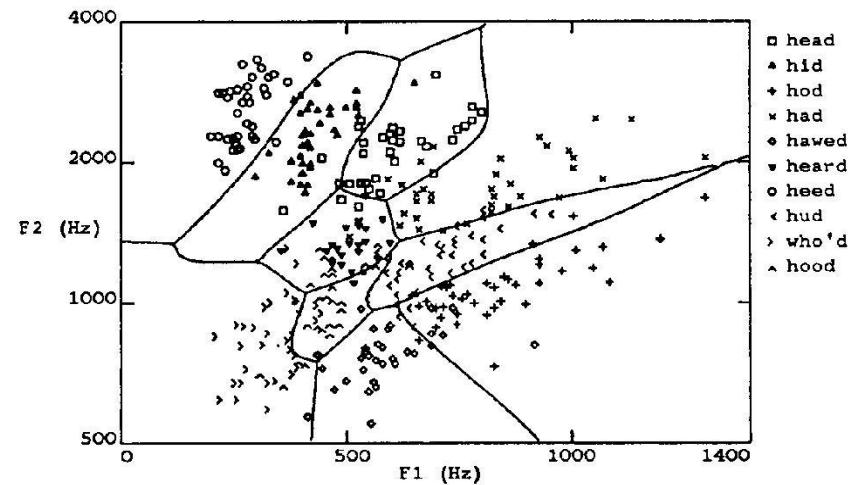
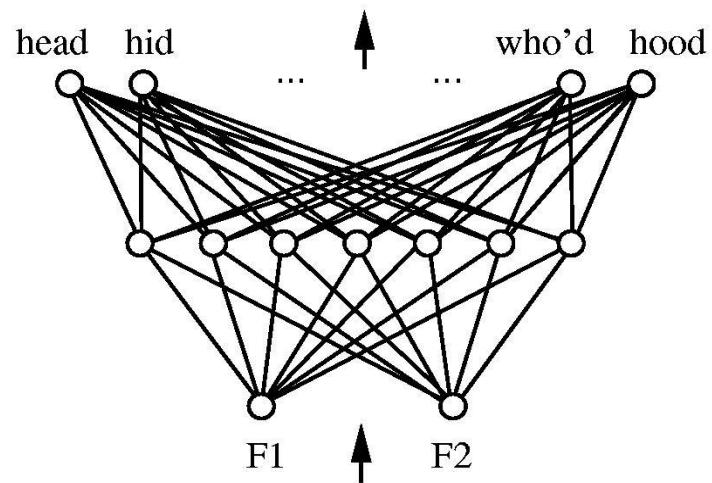
01-04-2020 Parag Singla @ IIT Delhi 9 



NOTE: Following slides have been borrowed from Pedro Domingos' Machine Learning course offered at the University of Washington



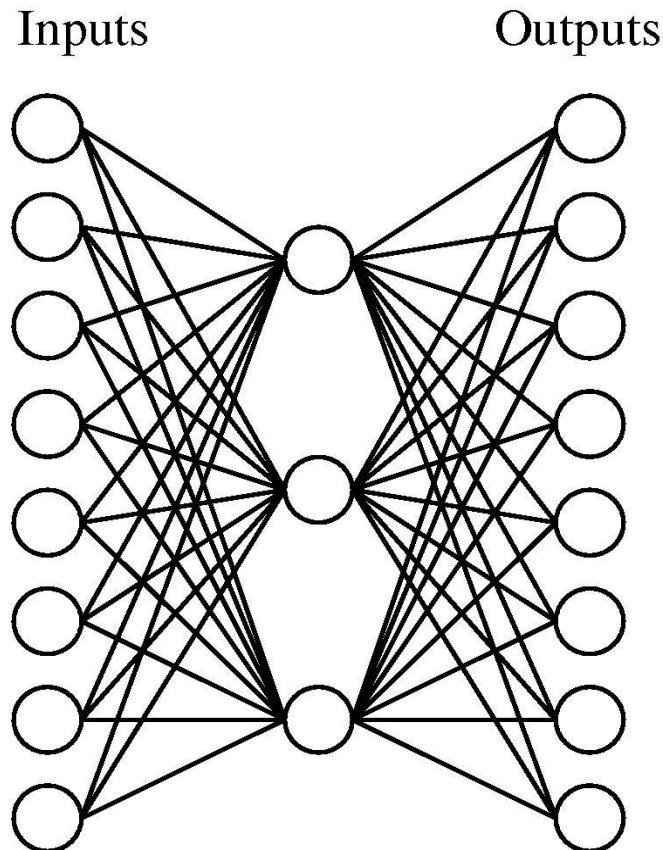
Multilayer Networks of Sigmoid Units



Source: Pedro Domingos' course offered at Univ. of Washington



Learning Hidden Layer Representations



Source: Pedro Domingos' course offered at Univ. of Washington



A target function:

Input	Output
10000000	→ 10000000
01000000	→ 01000000
00100000	→ 00100000
00010000	→ 00010000
00001000	→ 00001000
00000100	→ 00000100
00000010	→ 00000010
00000001	→ 00000001

Can this be learned?

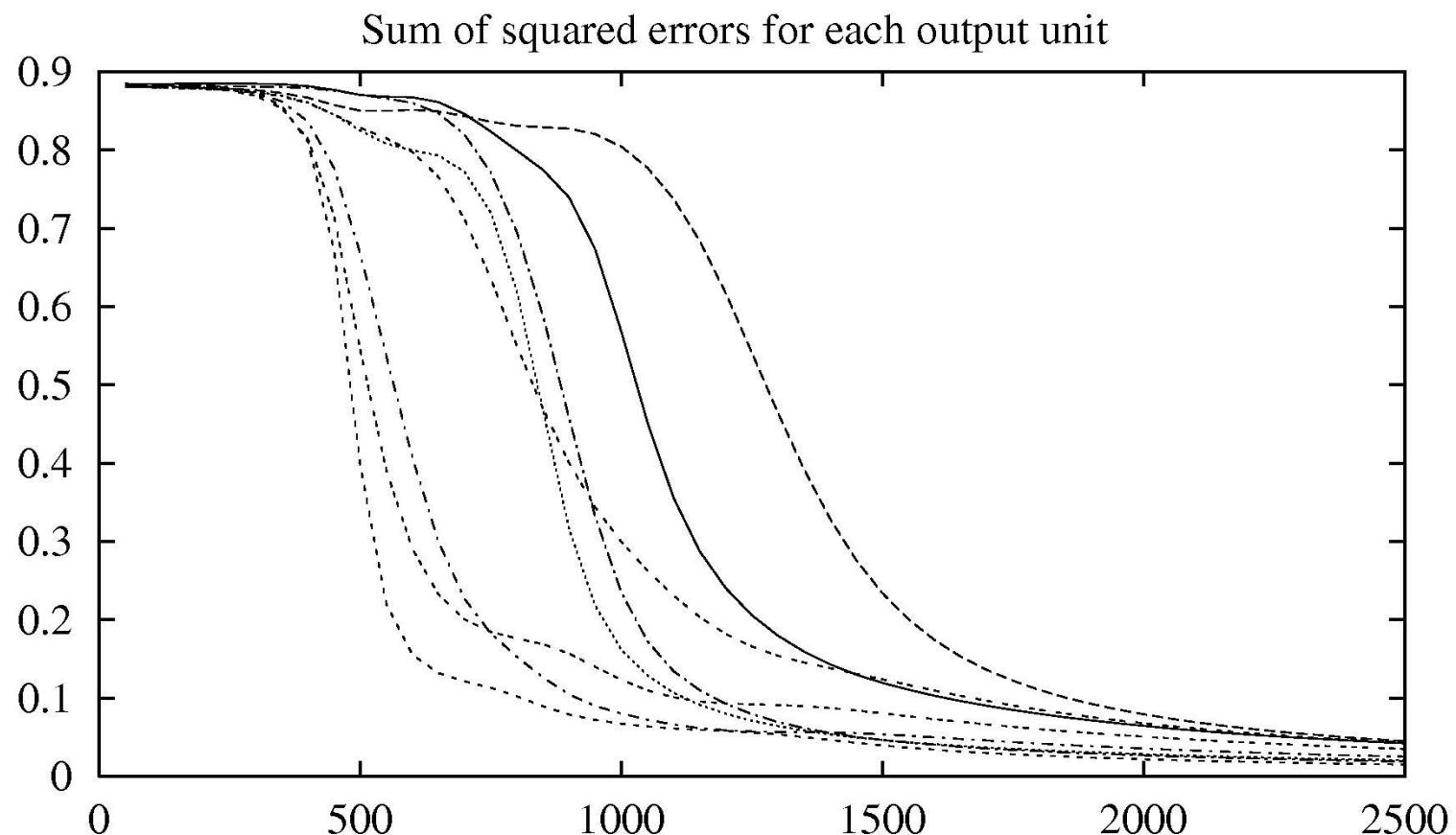


Learned hidden layer representation:

Input	Hidden			Output	
	Values				
10000000	→	.89	.04	.08	→ 10000000
01000000	→	.01	.11	.88	→ 01000000
00100000	→	.01	.97	.27	→ 00100000
00010000	→	.99	.97	.71	→ 00010000
00001000	→	.03	.05	.02	→ 00001000
00000100	→	.22	.99	.99	→ 00000100
00000010	→	.80	.01	.98	→ 00000010
00000001	→	.60	.94	.01	→ 00000001



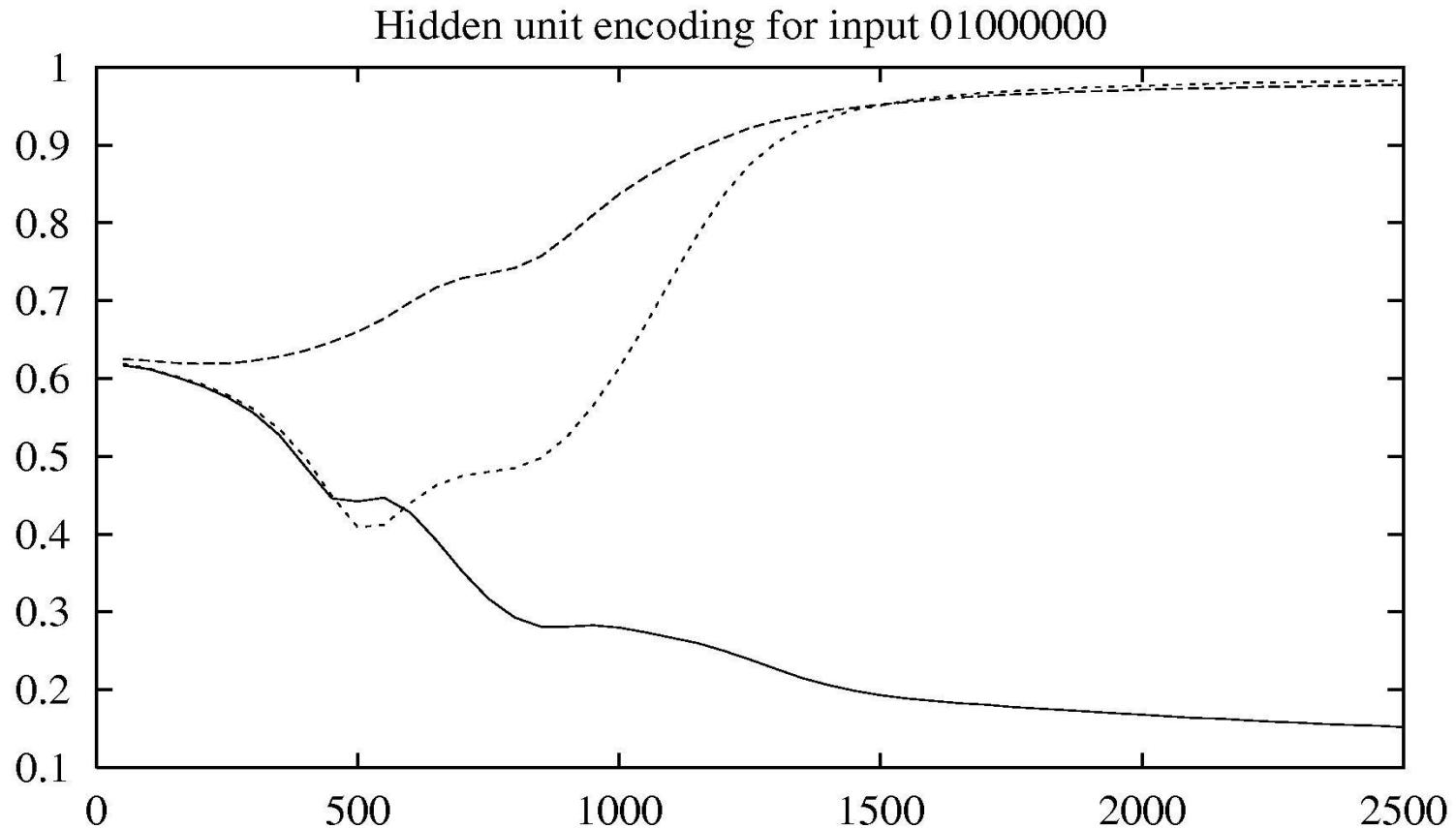
Training



Source: Pedro Domingos' course offered at Univ. of Washington



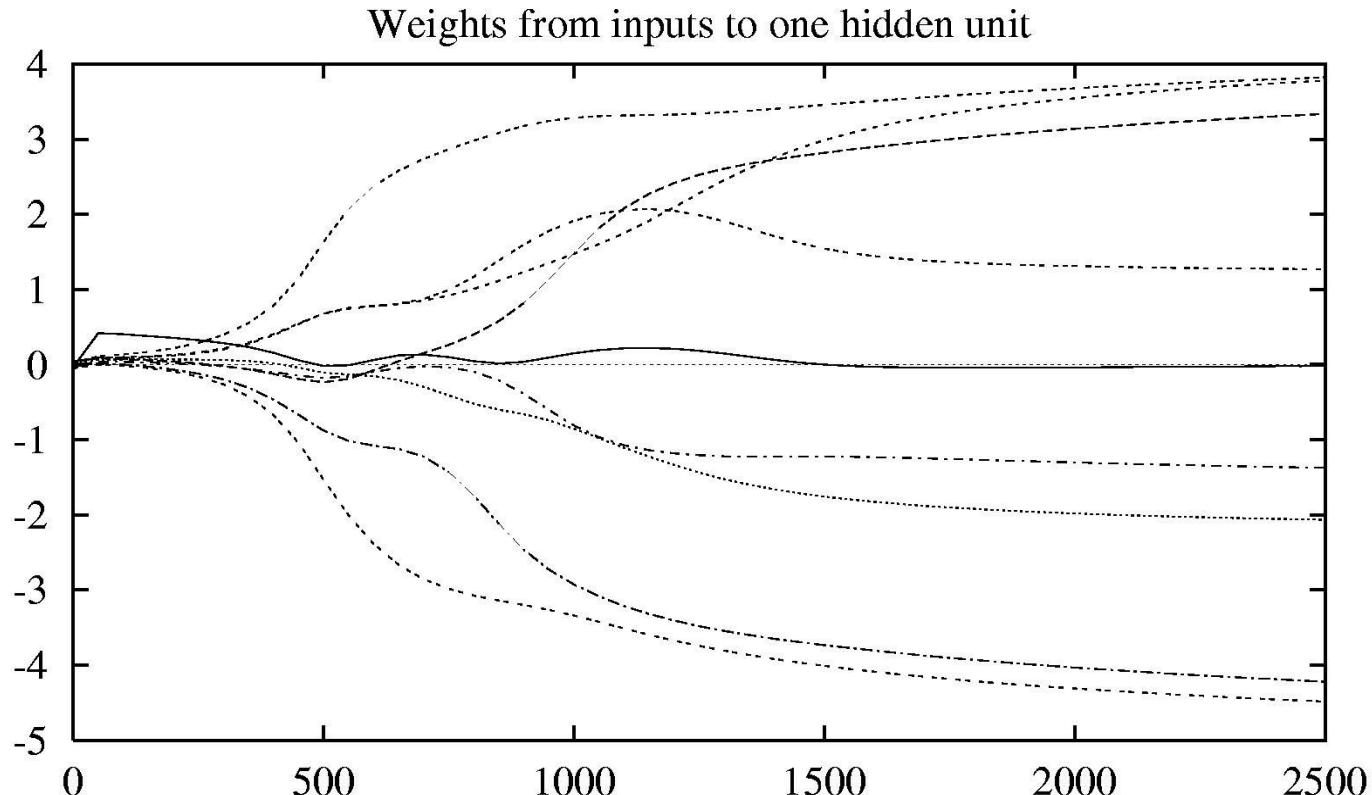
Training



Source: Pedro Domingos' course offered at Univ. of Washington



Training



Convergence of Backpropagation

Gradient descent to some local minimum

- Perhaps not global minimum...
- Add momentum
- Stochastic gradient descent
- Train multiple nets with different initial weights

Nature of convergence

- Initialize weights near zero
- Therefore, initial networks near-linear
- Increasingly non-linear functions possible as training progresses



Expressiveness of Neural Nets

Boolean functions:

- Every Boolean function can be represented by network with single hidden layer
- But might require exponential (in number of inputs) hidden units

Continuous functions:

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers

