

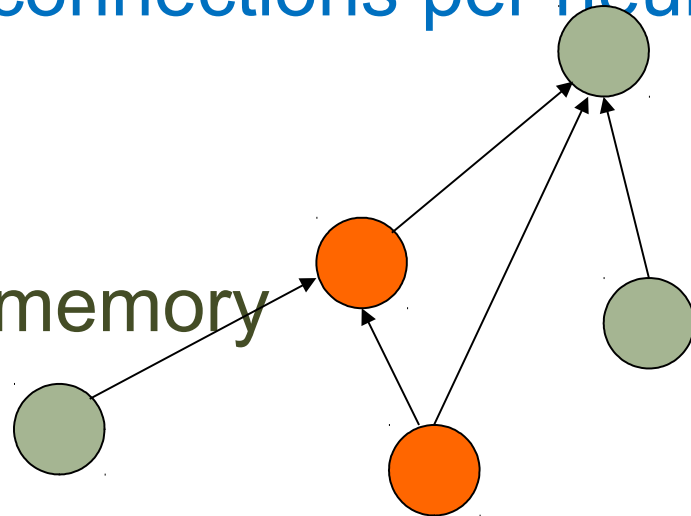
CHAPTER 11:

MULTILAYER PERCEPTRONS

Neural Networks - Inspired by the human brain

2

- Networks of processing units (neurons) with connections (synapses) between them
- Large number of neurons: 10^{10} (other estimates say 10^{11})
- Large connectivity: 10^5 connections per neuron: synapses
- Parallel processing
- Distributed computation/memory
- Robust to noise, failures



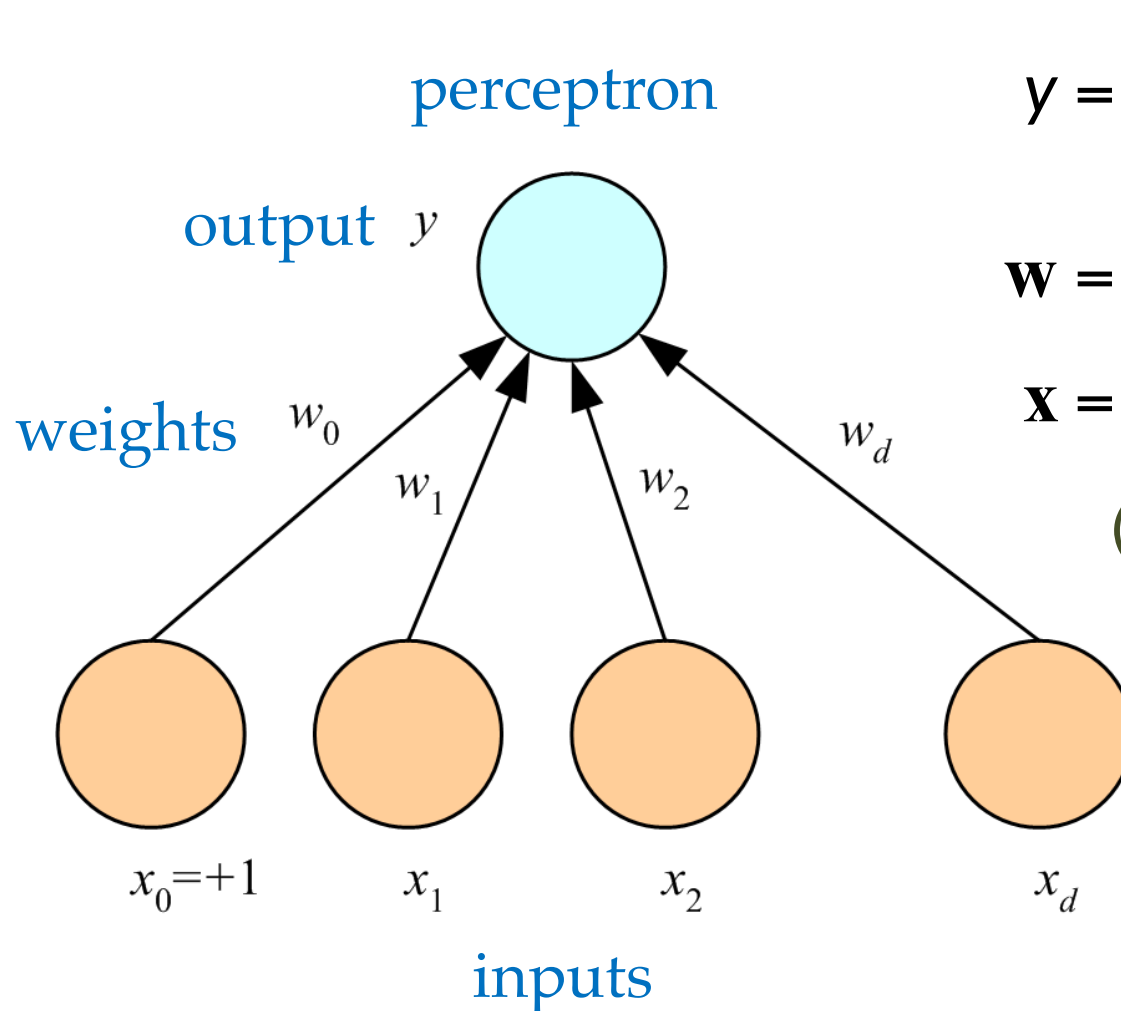
Understanding the Brain

3

- Levels of analysis (Marr, 1982)
 1. Computational theory
 2. Representation and algorithm
 3. Hardware implementation
- Reverse engineering: From hardware to theory
- Parallel processing: SIMD (single instruction, multiple data machines) vs MIMD (multiple instruction, multiple data machines)
Neural net: SIMD with modifiable local memory
Learning: Update by training/experience

Perceptron = one “neuron” or computational unit

4



$$y = \sum_{j=1}^d w_j x_j + w_0 = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

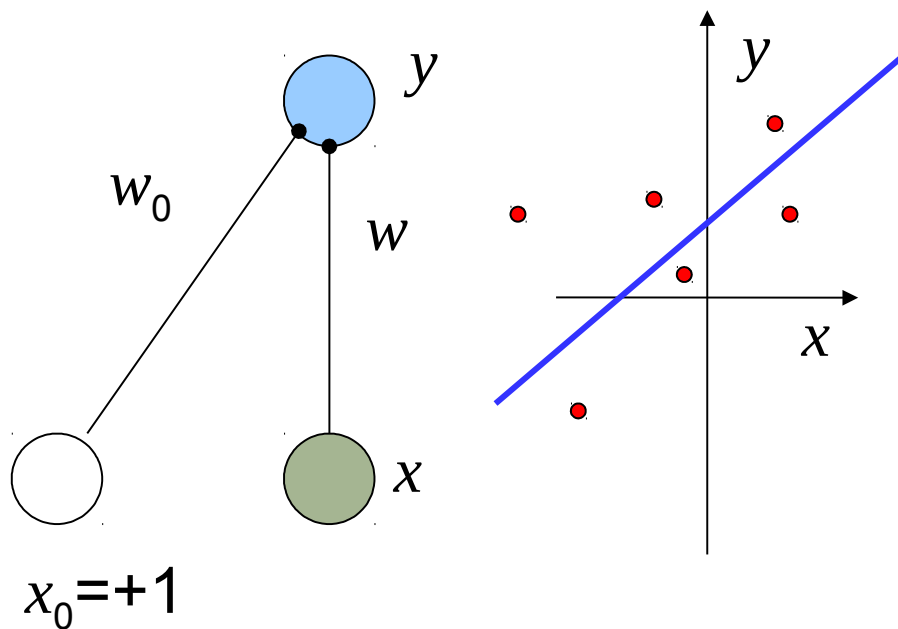
$$\mathbf{x} = [1, x_1, \dots, x_d]^T$$

(Rosenblatt, 1962)

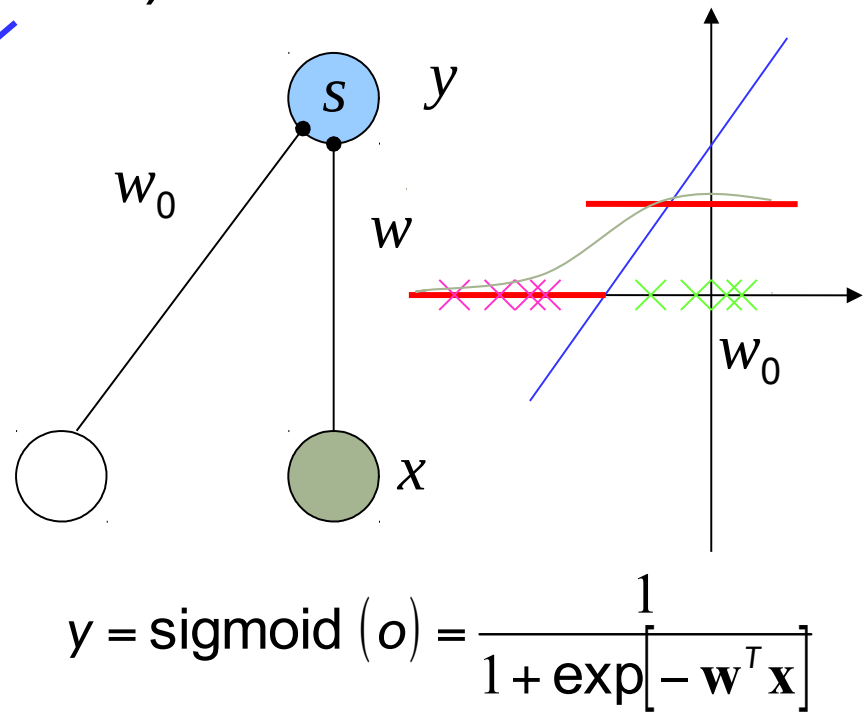
What a Perceptron Does

5

□ Regression: $y = wx + w_0$



□ Classification: $y = 1 (wx + w_0 > 0)$



K Outputs

Regression:

$$y_i = \sum_{j=1}^d w_{ij} x_j + w_{i0} = \mathbf{w}_i^T \mathbf{x}$$

$$\mathbf{y} = \mathbf{W}\mathbf{x}$$

Classification:

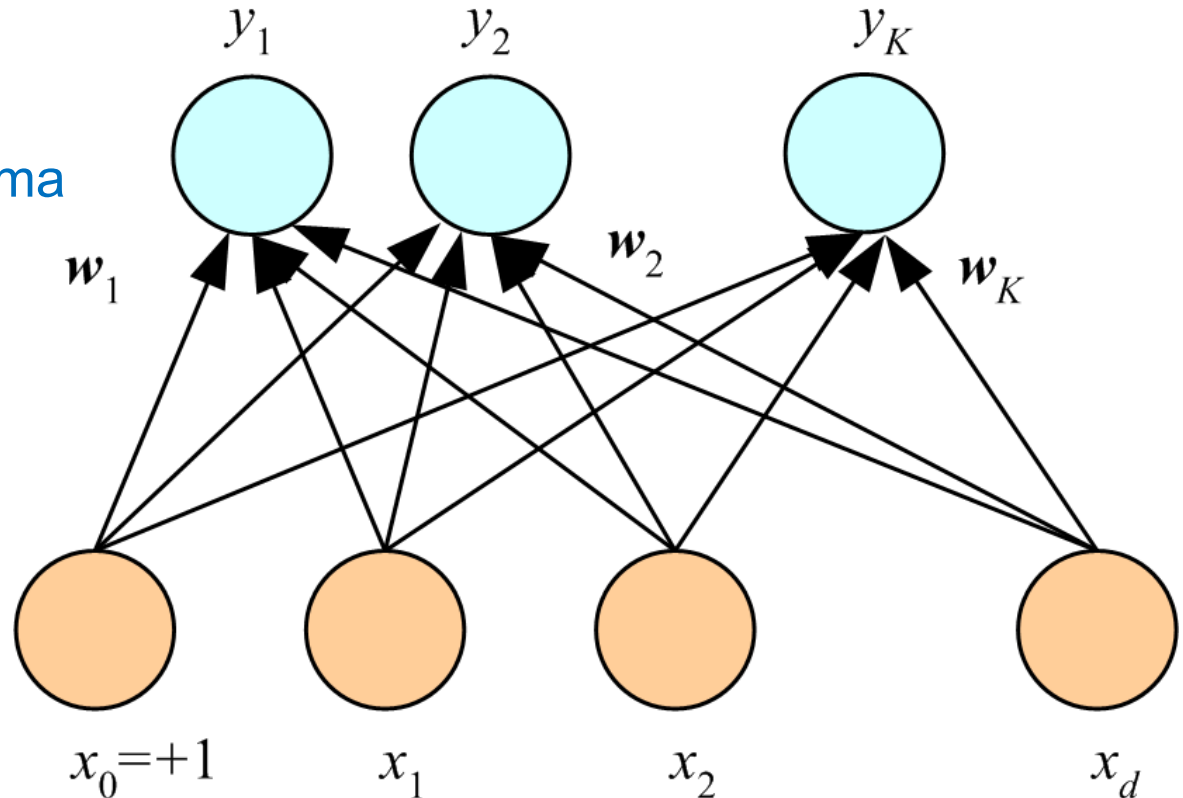
$$o_i = \mathbf{w}_i^T \mathbf{x}$$

$$y_i = \frac{\exp(o_i)}{\sum_k \exp(o_k)} \quad \text{softmax}$$

choose C_i

$$\text{if } y_i = \max_k y_k$$

softmax function or **normalized exponential function** converts a real vector into one with values in the range (0, 1) that add up to 1



Training (i.e., learning the right weights)

7

- Online (instances seen one by one) vs batch (whole sample) learning: **In online (= incremental) learning:**
 - ▣ No need to store the whole sample
 - ▣ Problem may change in time
 - ▣ Wear and degradation in system components
- Stochastic gradient-descent (**= incremental gradient descent**): Update after a single data instance
- Generic update rule (LMS rule):

$$\Delta w_{ij}^t = \eta (r_i^t - y_i^t) x_j^t$$

Update = Learning Factor · (Desired Output – Actual Output) · Input

Training a Perceptron: Regression

- Regression (Linear output):

$$E^t(\mathbf{w} | \mathbf{x}^t, r^t) = \frac{1}{2}(r^t - y^t)^2 = \frac{1}{2}[r^t - (\mathbf{w}^T \mathbf{x}^t)]^2$$

$$\Delta w_j^t = \eta(r^t - y^t)x_j^t$$

Classification

9

□ Single sigmoid output

$$y^t = \text{sigmoid}(\mathbf{w}^T \mathbf{x}^t)$$

$$E^t(\mathbf{w} \mid \mathbf{x}^t, \mathbf{r}^t) = -r^t \log y^t - (1 - r^t) \log (1 - y^t)$$

cross-entropy cost function:
used instead of the quadratic error function to speed up learning

$$\Delta w_j^t = \eta (r^t - y^t) x_j^t$$

□ $K > 2$ softmax outputs

$$y^t = \frac{\exp \mathbf{w}_i^T \mathbf{x}^t}{\sum_k \exp \mathbf{w}_k^T \mathbf{x}^t}$$
$$E^t(\{\mathbf{w}_i\}_i \mid \mathbf{x}^t, \mathbf{r}^t) = -\sum_i r_i^t \log y_i^t$$

cross-entropy cost function

$$\Delta w_{ij}^t = \eta (r_i^t - y_i^t) x_j^t$$

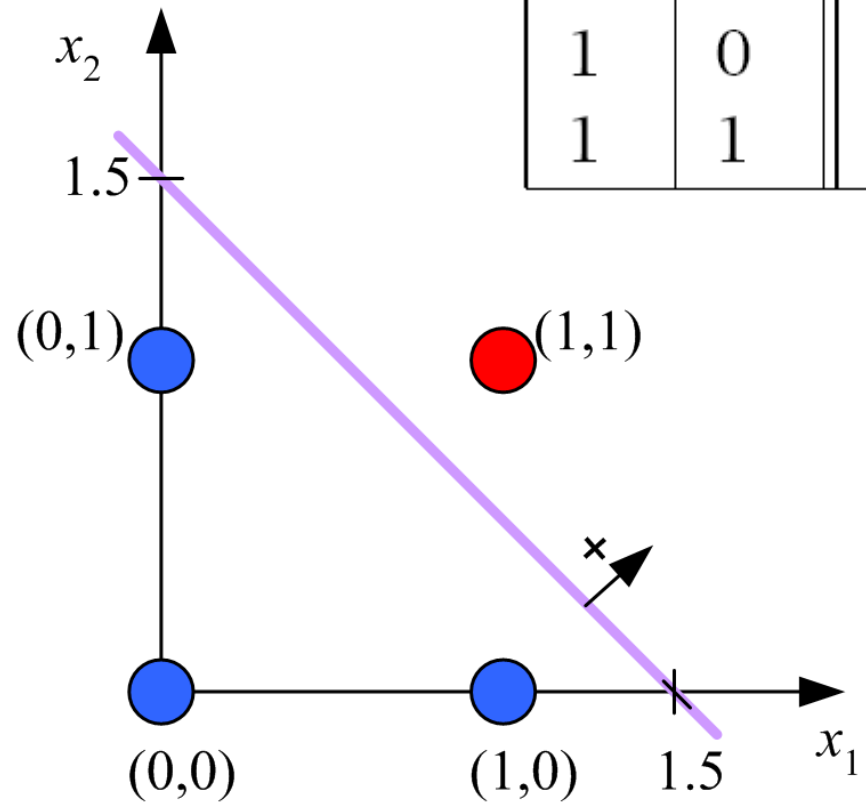
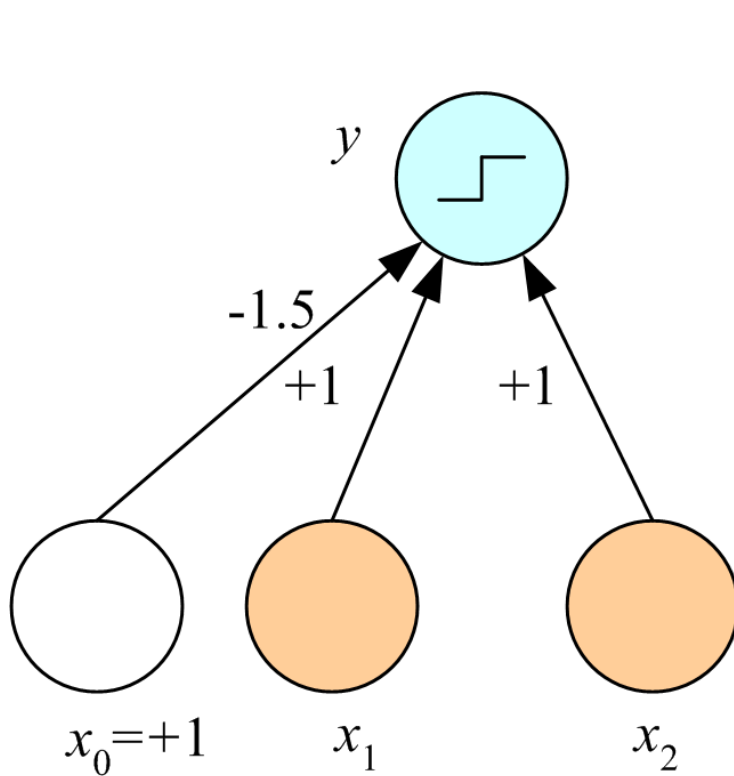
see

[Chapter 3 of Michael Nielsen's "Neural Networks and Deep Learning" online](#)

Learning Boolean AND

10

x_1	x_2	r
0	0	0
0	1	0
1	0	0
1	1	1



XOR

x_1	x_2	r
0	0	0
0	1	1
1	0	1
1	1	0

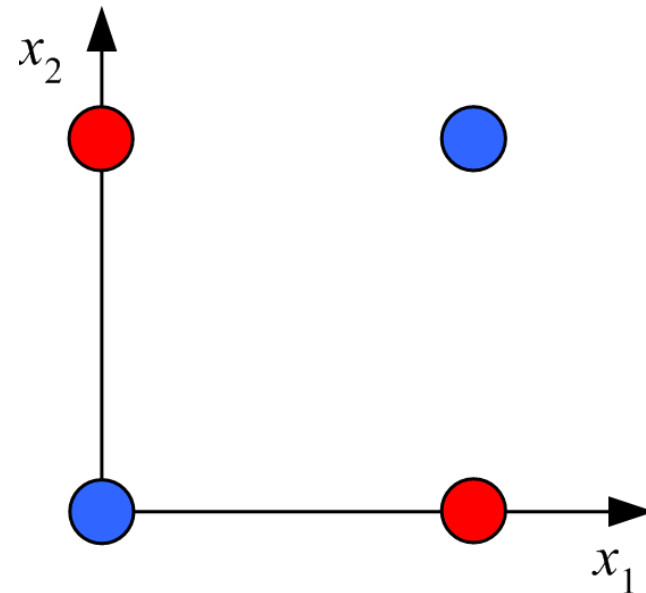
□ No w_0, w_1, w_2 satisfy:

$$w_0 \leq 0$$

$$w_2 + w_0 > 0$$

$$w_1 + w_0 > 0$$

$$w_1 + w_2 + w_0 \leq 0$$

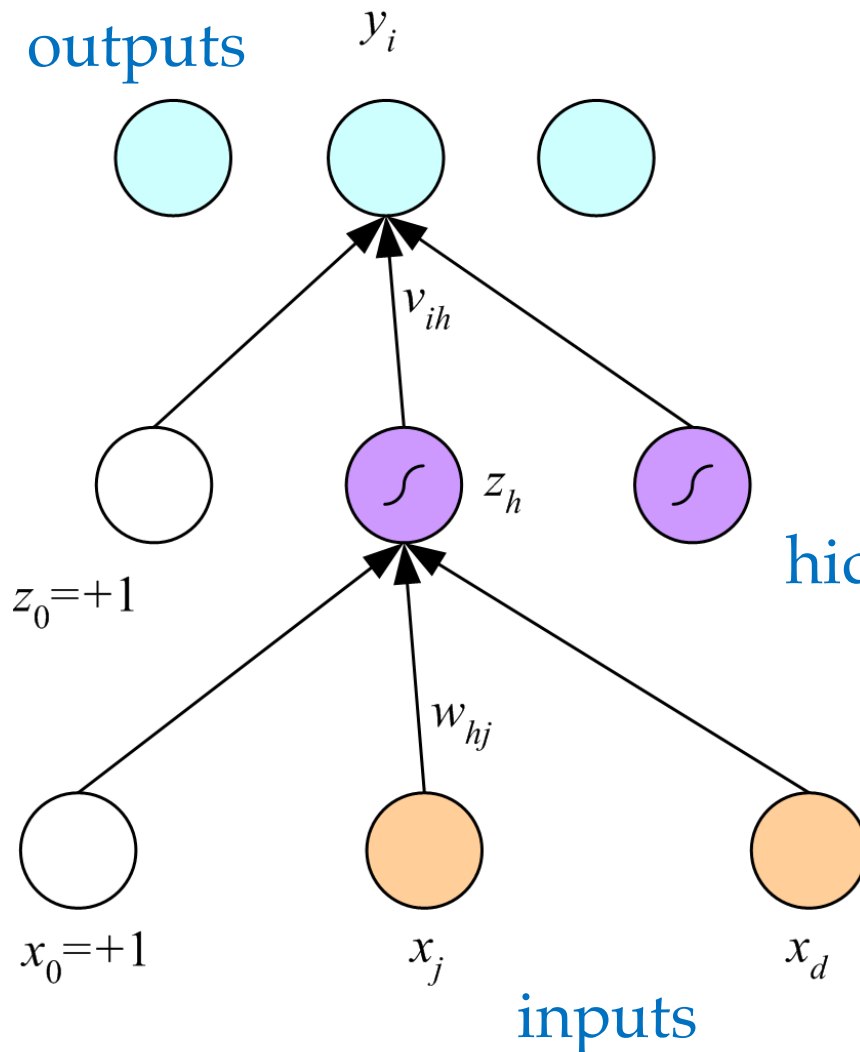


(Minsky and Papert, 1969)

Multilayer Perceptrons

12

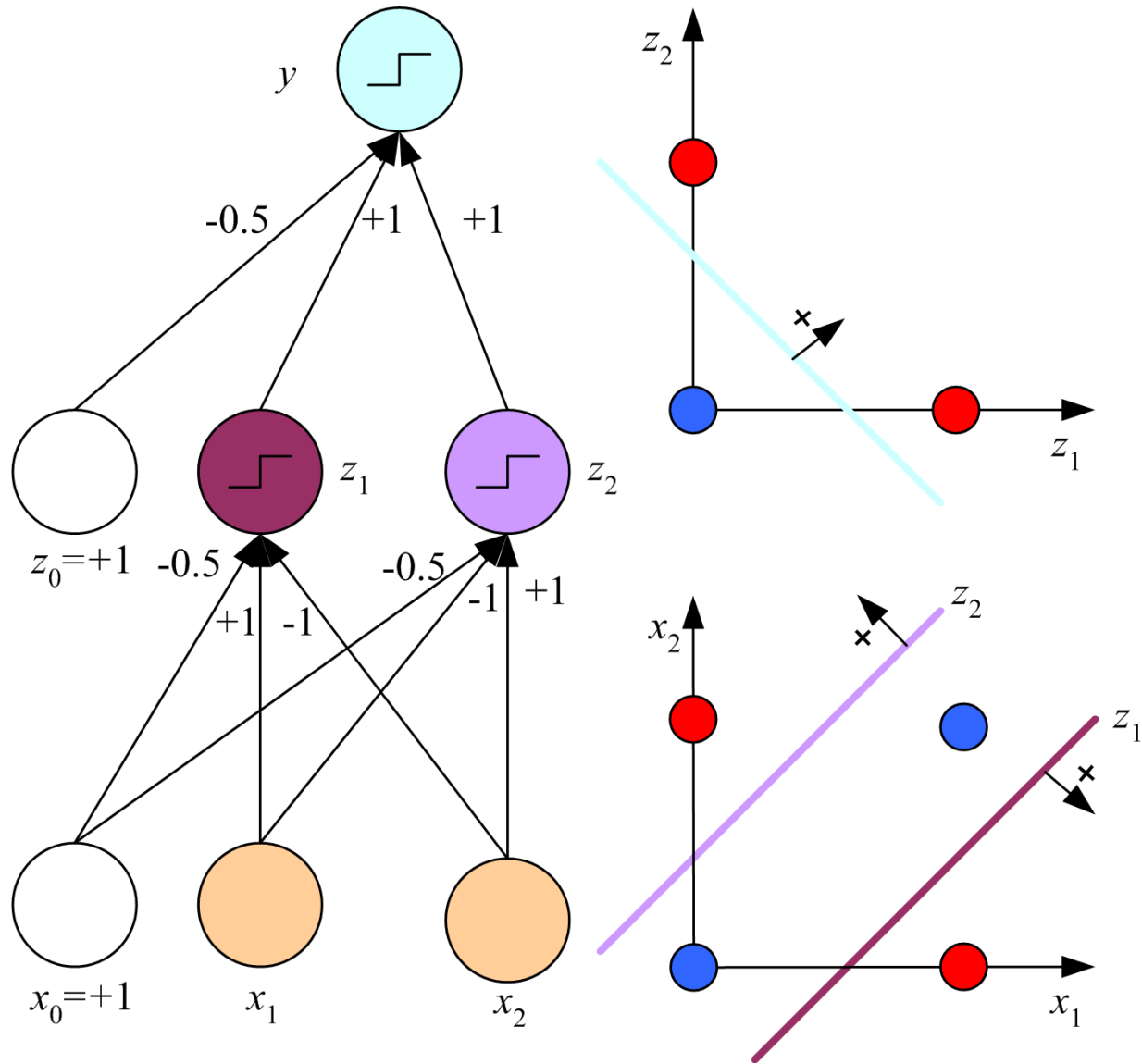
outputs



$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$
$$= \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^d w_{hj} x_j + w_{h0}\right)\right]}$$

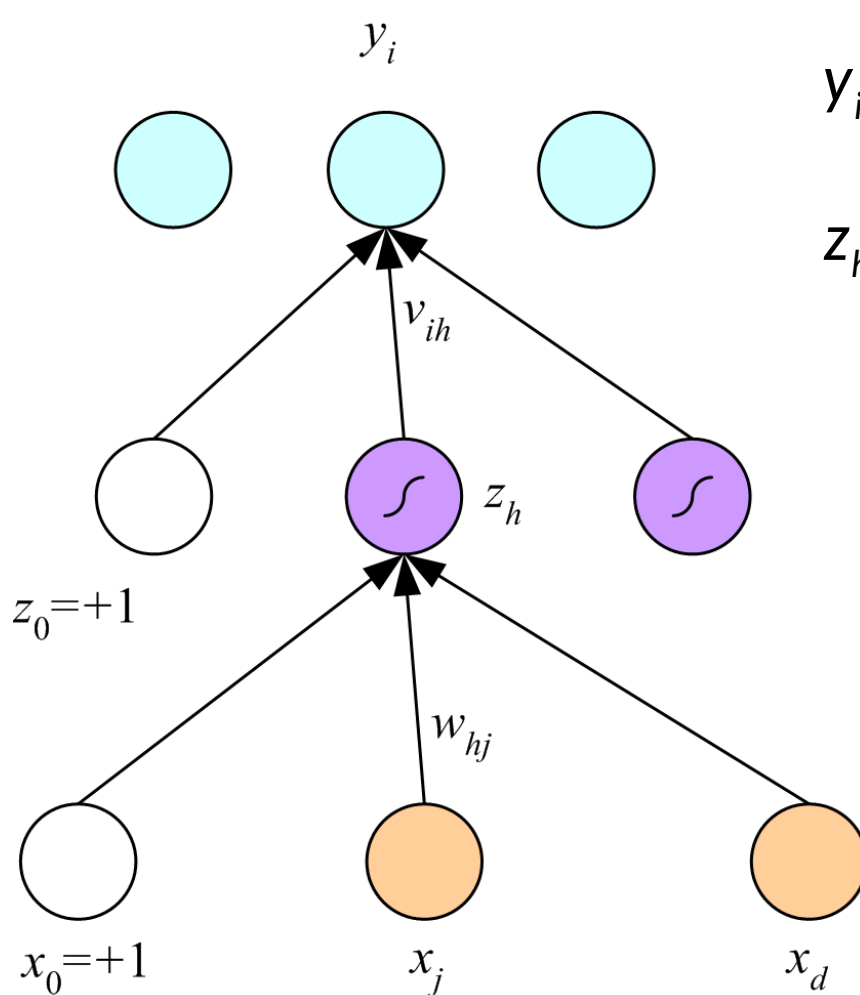
(Rumelhart, Hinton, and
Williams, 1986)



$$x_1 \text{ XOR } x_2 = (x_1 \text{ AND } \sim x_2) \text{ OR } (\sim x_1 \text{ AND } x_2)$$

Backpropagation (Rumelhart, Hinton & Williams, 1986)

14



$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$

$$= \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^d w_{hj} x_j + w_{h0}\right)\right]}$$

$$\frac{\partial E}{\partial w_{hj}} = \frac{\partial E}{\partial y_i} \frac{\partial y_i}{\partial z_h} \frac{\partial z_h}{\partial w_{hj}}$$

gradient of the error function

Regression

$$E(\mathbf{W}, \mathbf{v} | \mathbf{X}) = \frac{1}{2} \sum_t (r^t - y^t)^2$$

$$y^t = \sum_{h=1}^H v_h z_h^t + v_0$$

Forward

$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x})$$

\mathbf{x}

$$\Delta v_h = \sum_t (r^t - y^t) z_h^t$$

Backward

$$\Delta w_{hj} = -\eta \frac{\partial E}{\partial w_{hj}}$$

$$= -\eta \sum_t \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial w_{hj}}$$

$$= -\eta \sum_t - (r^t - y^t) v_h \boxed{z_h^t (1 - z_h^t)} x_j^t$$

$$= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t$$

Regression with Multiple Outputs

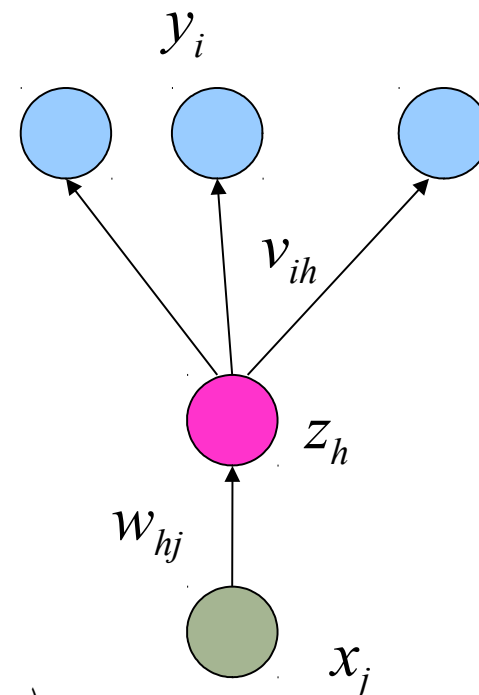
16

$$E(\mathbf{W}, \mathbf{V} | \mathbf{X}) = \frac{1}{2} \sum_t \sum_i (r_i^t - y_i^t)^2$$

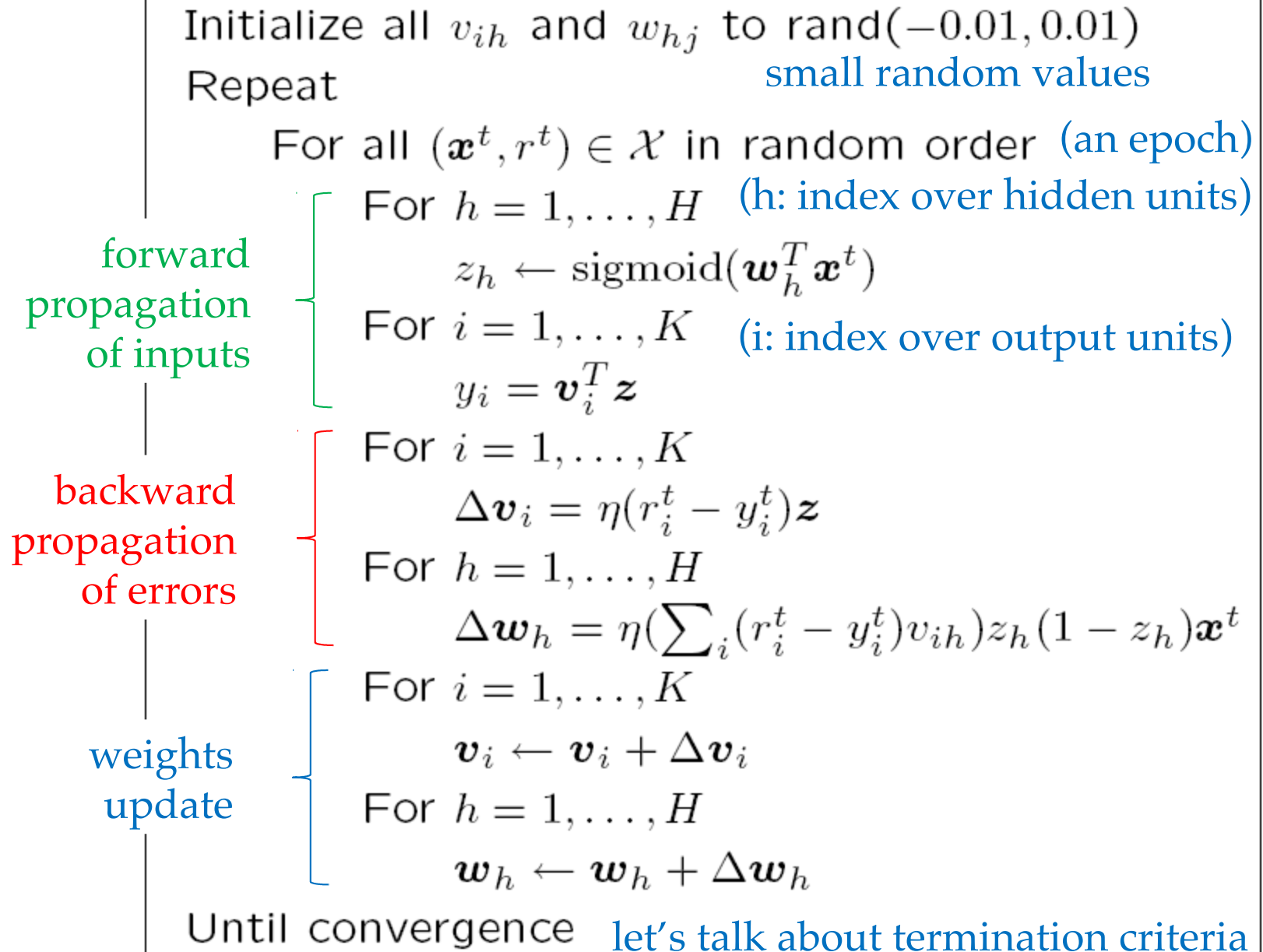
$$y_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0}$$

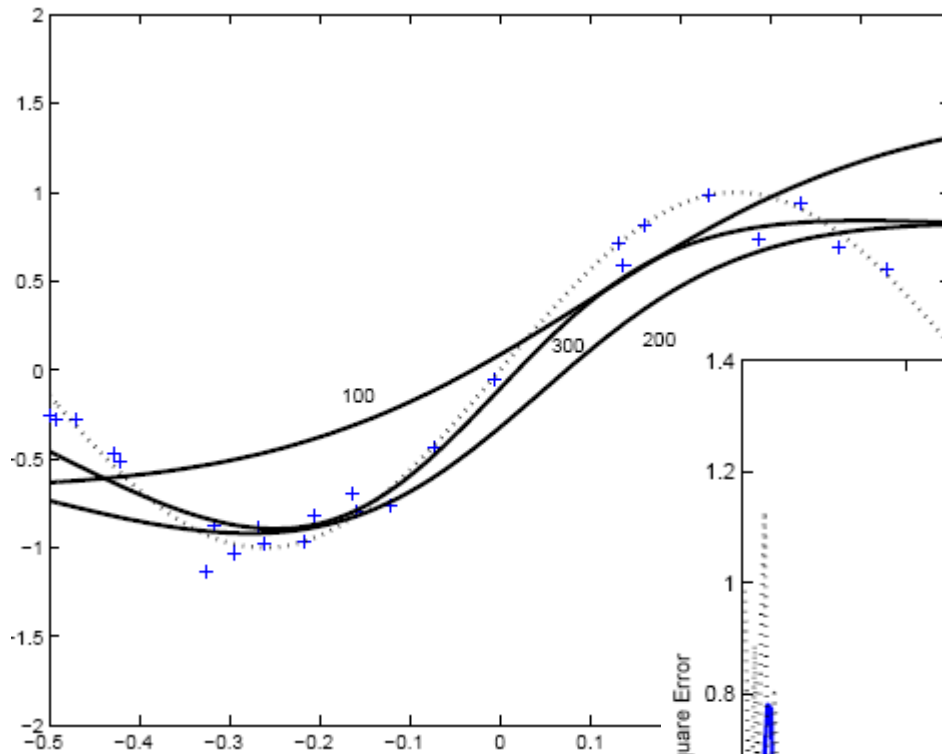
$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t$$

$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t$$

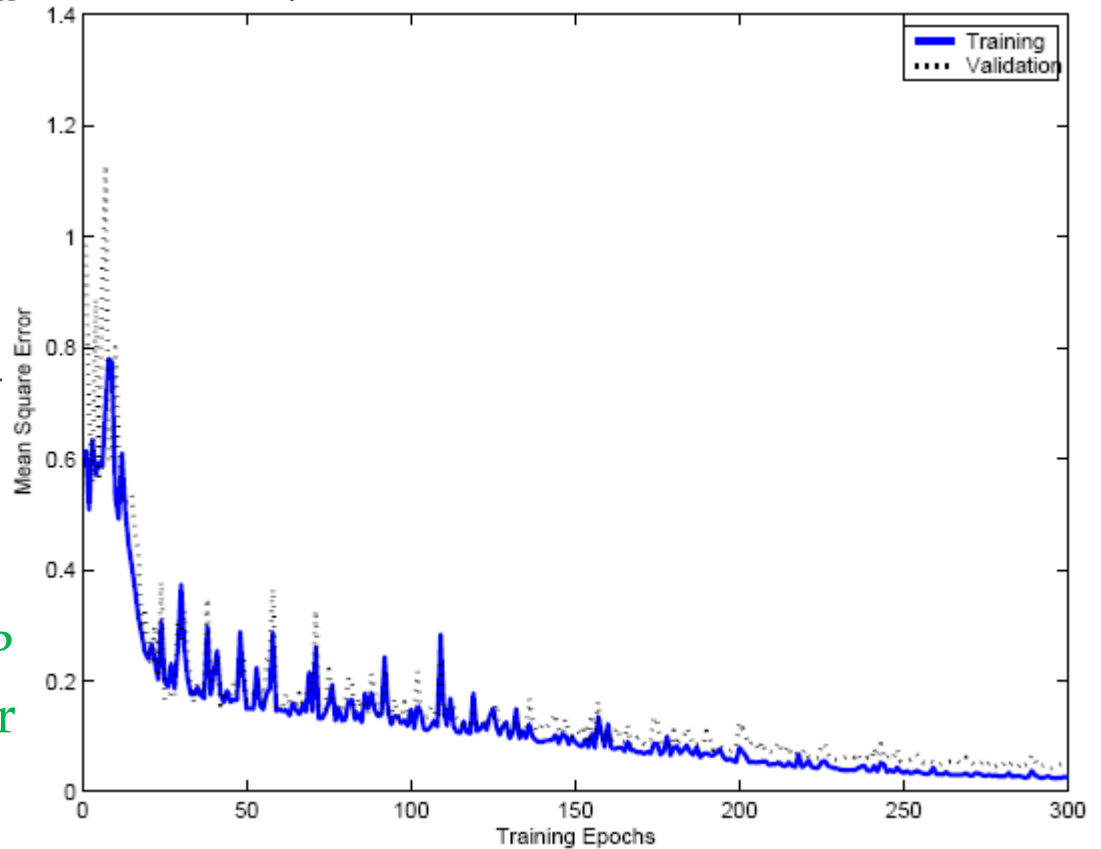


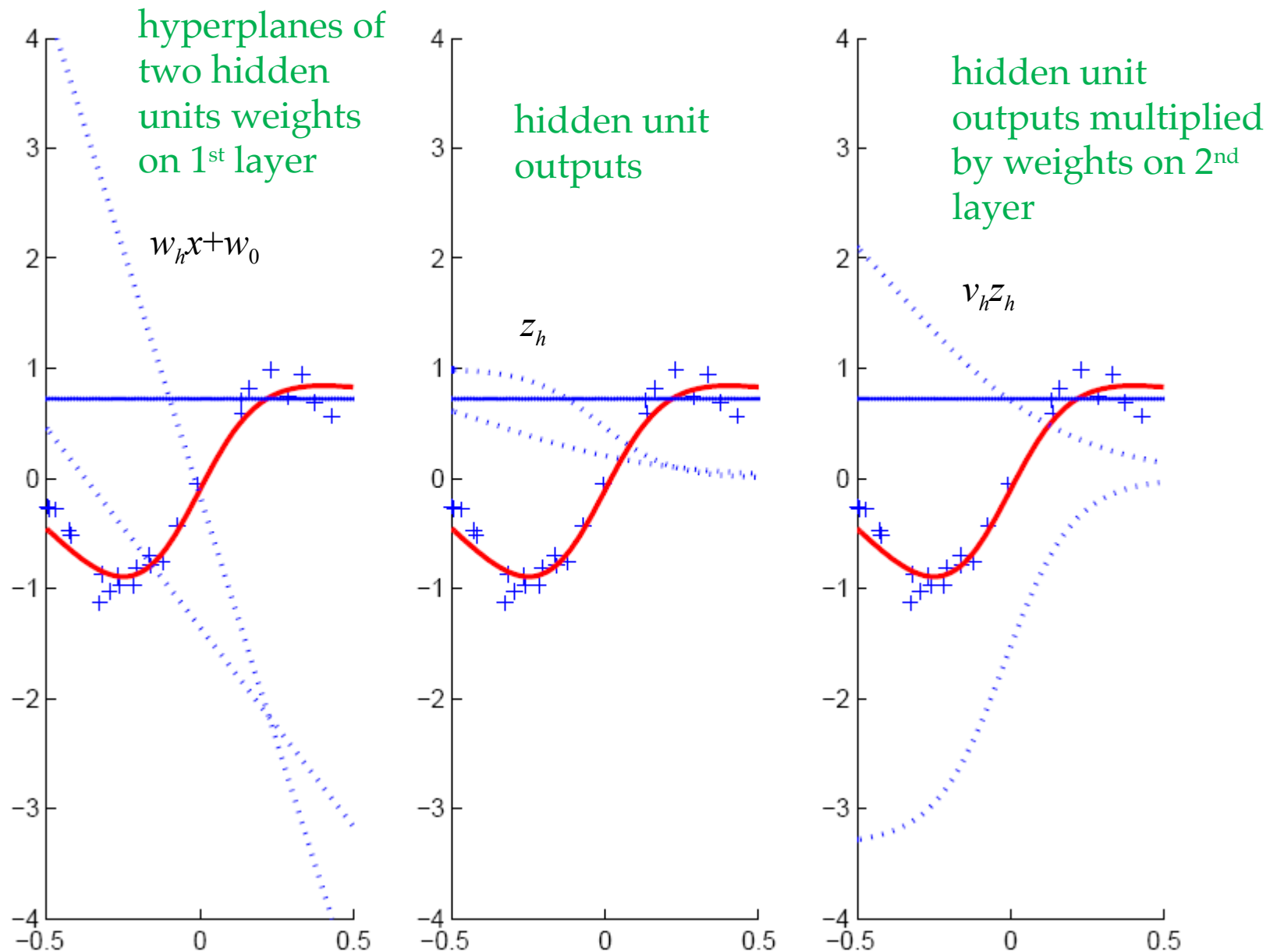
Error Back Propagation Algorithm using Stochastic Gradient Descent





- +: training instances
- ...: $f(x) = \sin(6x)$
- 100,200, 300: Fit of an MLP with two hidden units after 100, 200, and 300 epochs





MLP fit is formed as the sum of the outputs of the hidden units

Two-Class Discrimination

20

- One sigmoid output y^t for $P(C_1|\mathbf{x}^t)$ and $P(C_2|\mathbf{x}^t) \equiv 1 - y^t$

$$y^t = \text{sigmoid} \left(\sum_{h=1}^H v_h z_h^t + v_0 \right)$$

$$E(\mathbf{W}, \mathbf{v} | \mathbf{X}) = - \sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

$$\Delta v_h = \eta \sum_t (r^t - y^t) z_h^t$$

$$\Delta w_{hj} = \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t$$

$K > 2$ Classes

21

softmax

\times

$$o_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad y_i^t = \frac{\exp o_i^t}{\sum_k \exp o_k^t} \equiv P(C_i | \mathbf{x}^t)$$

$$E(\mathbf{W}, \mathbf{v} | \mathbf{X}) = - \sum_t \sum_i r_i^t \log y_i^t$$

$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t$$

$$\Delta w_{hj} = \eta \sum_t \left[\sum_i (r_i^t - y_i^t) v_{ih} \right] z_h^t (1 - z_h^t) x_j^t$$

Multiple Hidden Layers

22

- MLP (multi-layer perceptrons) with one hidden layer is a universal approximator (Hornik et al., 1989), but using multiple layers may lead to simpler networks

$$z_{1h} = \text{sigmoid} \left(\mathbf{w}_{1h}^T \mathbf{x} \right) = \text{sigmoid} \left(\sum_{j=1}^d w_{1hj} x_j + w_{1h0} \right), h = 1, \dots, H_1$$

$$z_{2l} = \text{sigmoid} \left(\mathbf{w}_{2l}^T \mathbf{z}_1 \right) = \text{sigmoid} \left(\sum_{h=1}^{H_1} w_{2lh} z_{1h} + w_{2l0} \right), l = 1, \dots, H_2$$

$$y = \mathbf{v}^T \mathbf{z}_2 = \sum_{l=1}^{H_2} v_l z_{2l} + v_0$$

Improving Convergence

23

□ Momentum

$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \underbrace{\alpha \Delta w_i^{t-1}}$$

where t is time (= epoch) and α is a constant, say between 0.5 and 1

□ Adaptive learning rate

$$\Delta \eta = \begin{cases} +a & \text{if } E^{t+\tau} < E^t \\ -b\eta & \text{otherwise} \end{cases}$$

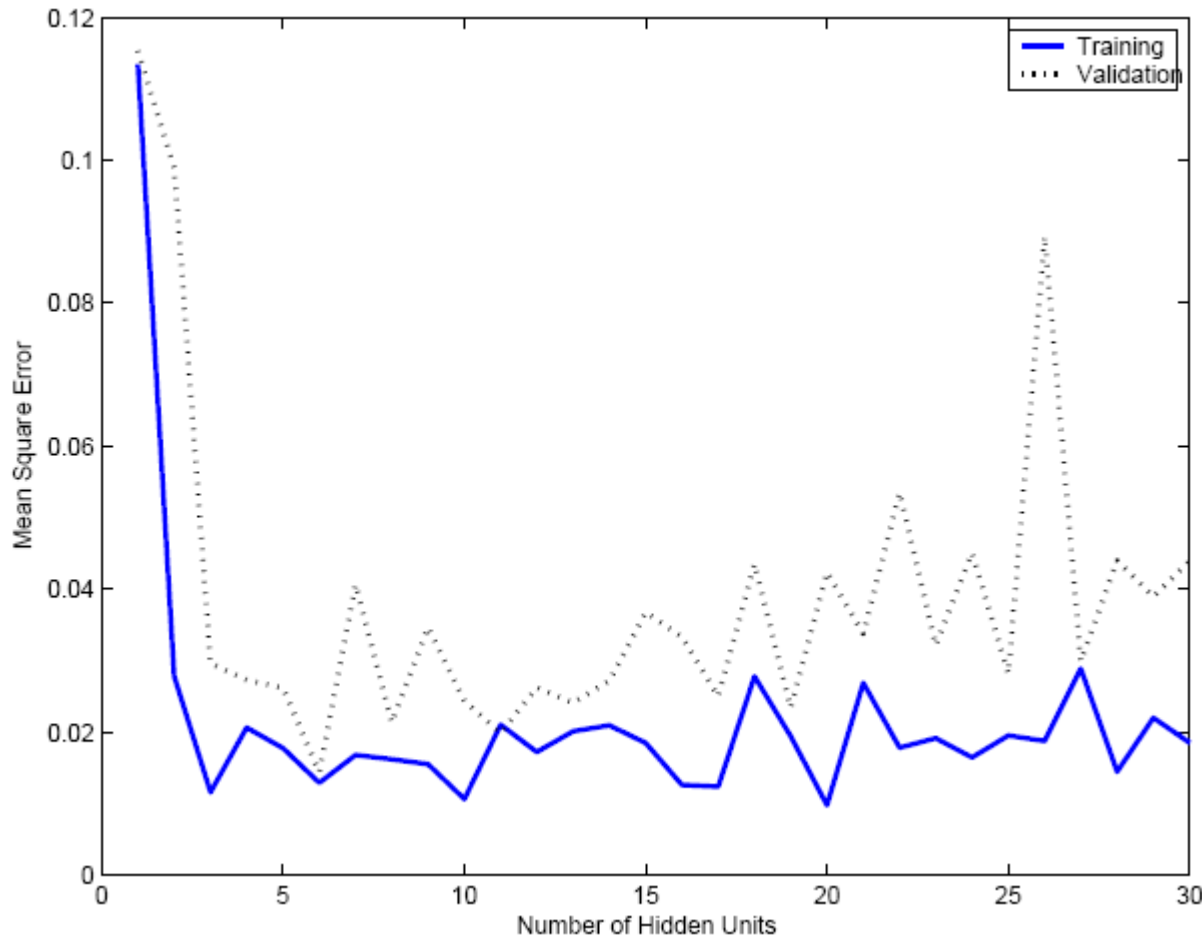
if error rate is decreasing, increase learning rate by a constant amount; otherwise, decrease learning rate geometrically

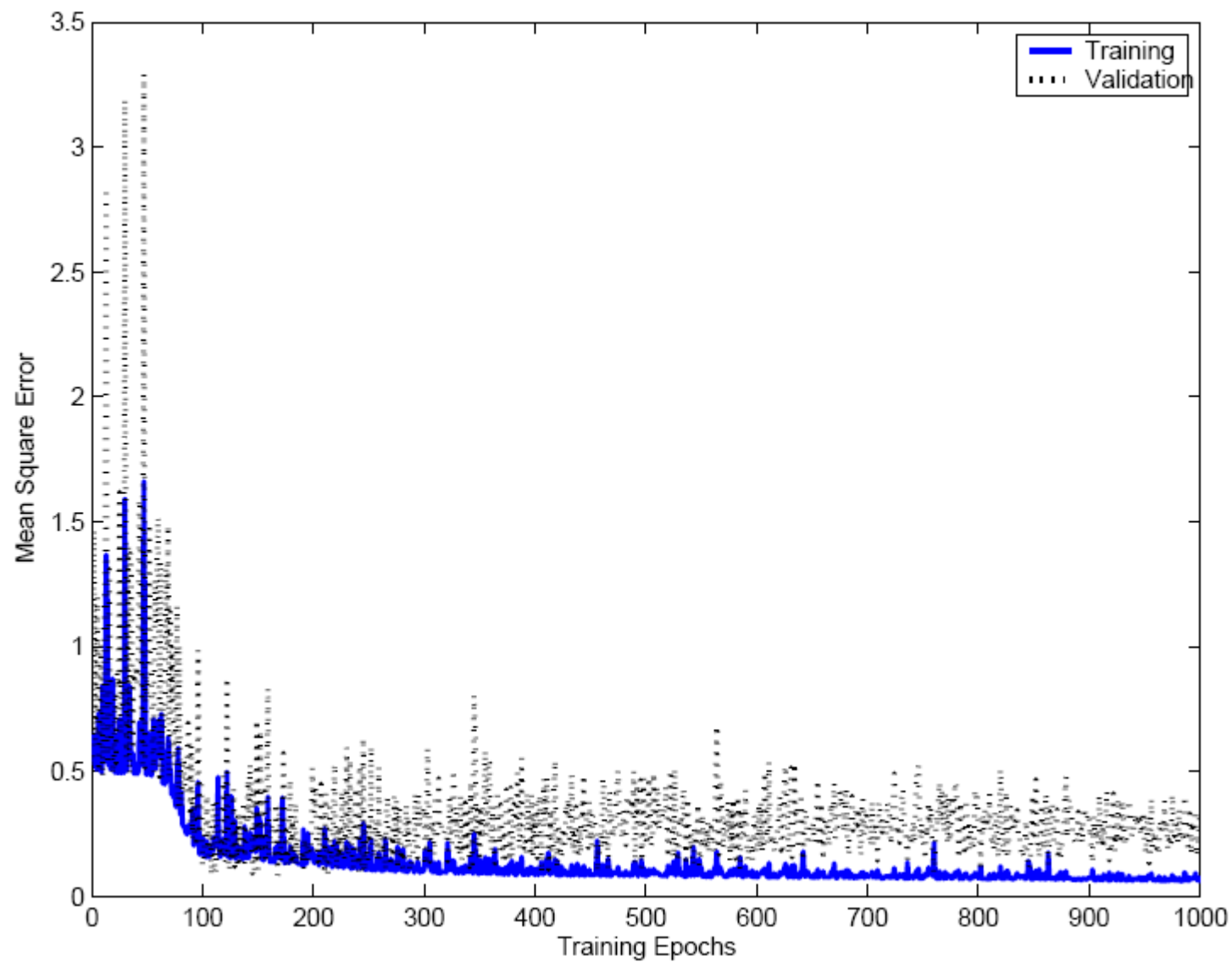
Overfitting/Overtraining

24

Number of weights: $H(d+1) + (H+1)K$

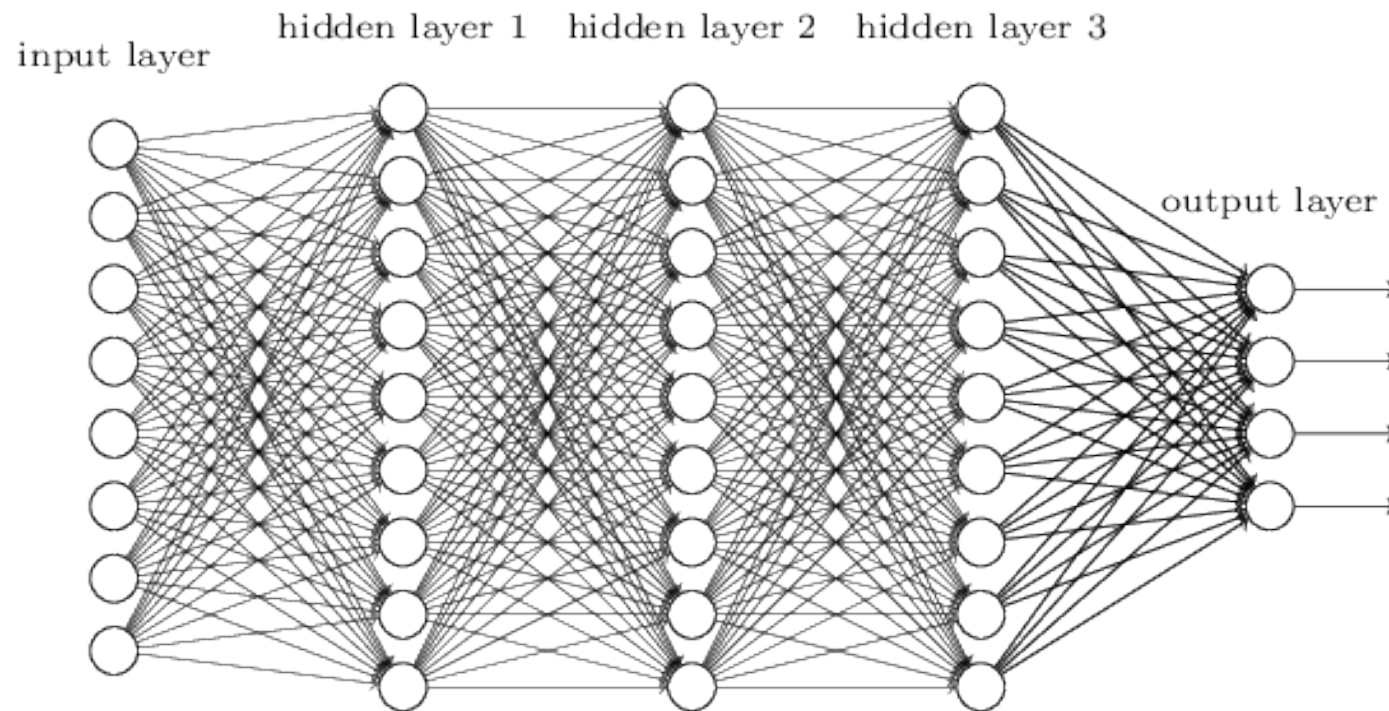
MLP with d inputs, H hidden units, and K outputs: has $H \times (d+1)$ weights in 1st layer and $K \times (H+1)$ in 2nd layer





Fully Connected, Layered MLP Architecture

26



Taken from
[Chapter 6 of Michael Nielsen's "Neural Networks and Deep Learning" online textbook](#)

Expressive Capabilities of ANNs

27

(Taken from Tom Mitchell's "Machine Learning" Textbook)

- Boolean Functions:

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

- Continuous Functions:

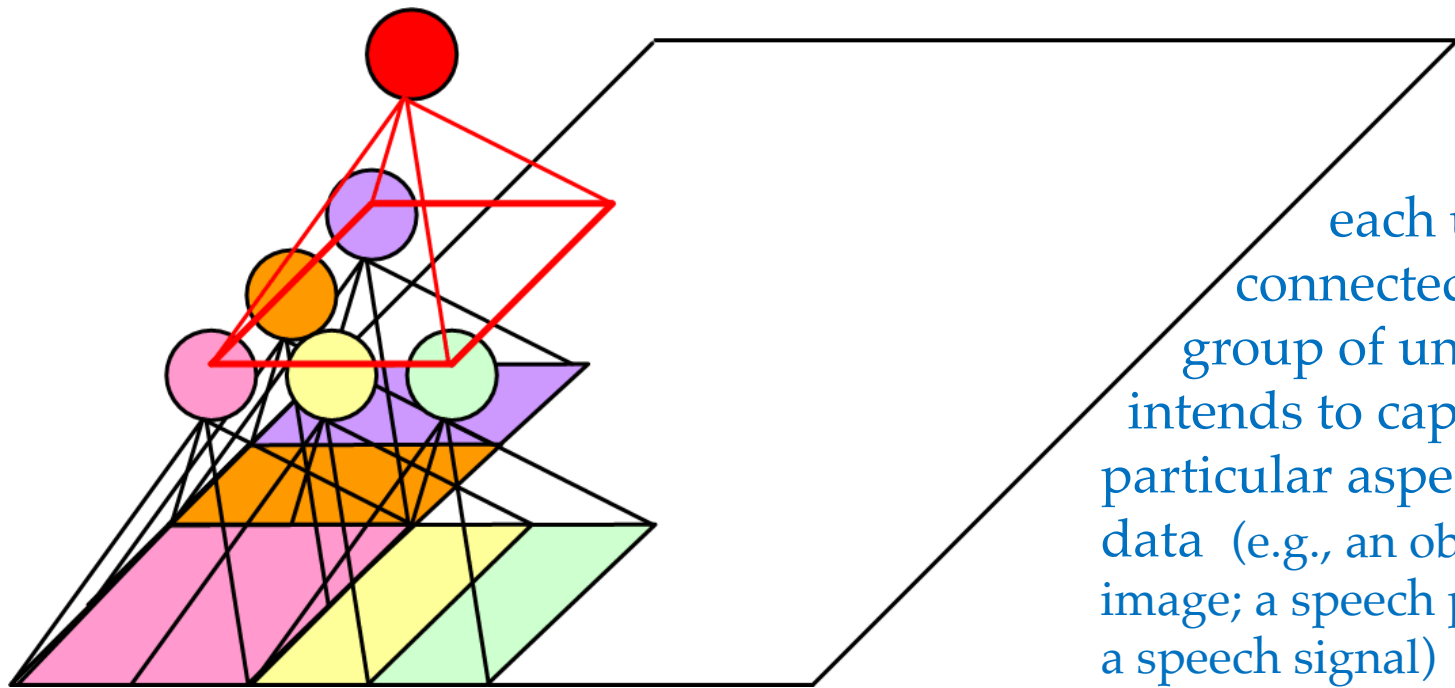
- Every bounded continuous function can be approximated with arbitrary small error by a network with one hidden layer [Cybenko 1989; Hornik et al. 1989]
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].

Other ANN Architectures: Structured MLP

28

□ Convolutional networks (Deep learning)

Useful on data that has “local structure” (e.g., images; sound)



(Le Cun et al, 1989)

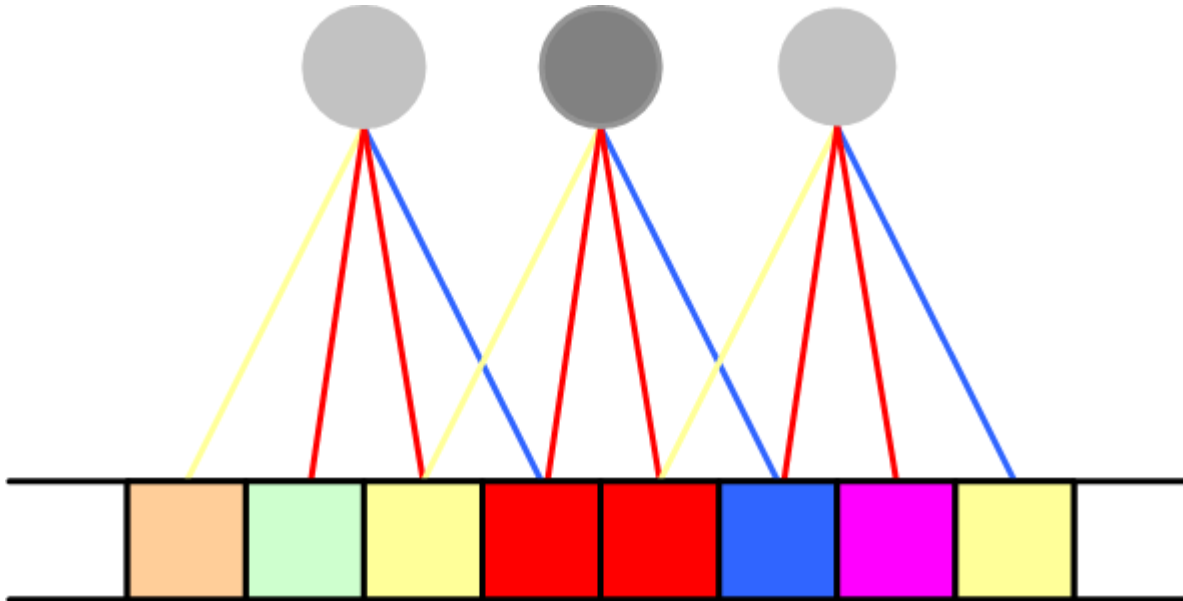
see

[Chapter 6 of Michael Nielsen's "Neural Networks and Deep Learning" online](#)

Weight Sharing

29

three translations over the data of the same “filter”



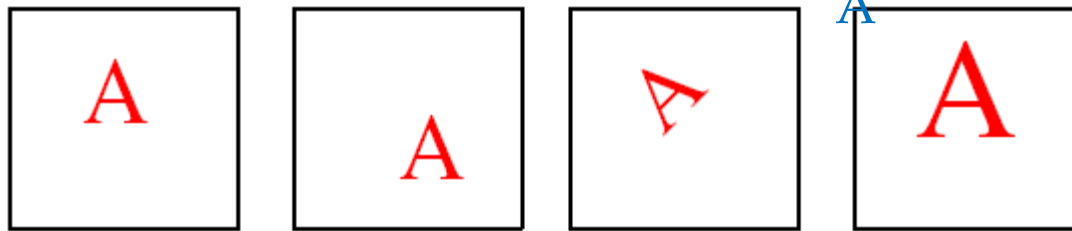
Each “filter” is applied on overlapping segments of the data, and uses the same weights each time that is applied during the same epoch (“translation invariance”).

These weights are updated after each epoch.

Hints

30

- Invariance to translation, rotation, size the “A” remains an “A”



- Virtual examples (Abu-Mostafa, 1995)
- Augmented error: $E' = E + \lambda_h E_h$

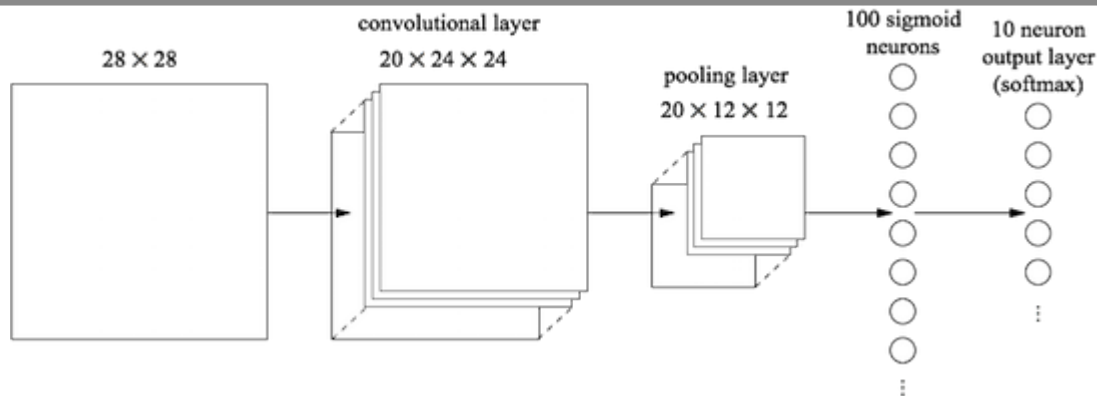
If x' and x are the “same”: $E_h = [g(x|\theta) - g(x'|\theta)]^2$

Approximation hint:

$$E_h = \begin{cases} 0 & \text{if } g(x|\theta) \in [a_x, b_x] \\ (g(x|\theta) - a_x)^2 & \text{if } g(x|\theta) < a_x \\ (g(x|\theta) - b_x)^2 & \text{if } g(x|\theta) > b_x \end{cases}$$

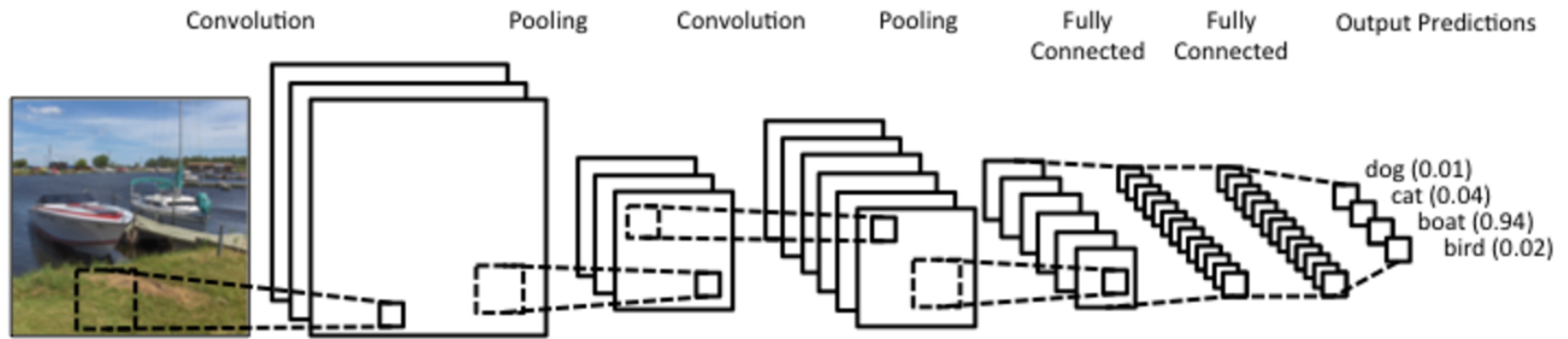
Convolutional Neural Networks in Practice

31



Taken from

[Chapter 6 of Michael Nielsen's "Neural Networks and Deep Learning" online textbook](#)



Taken from ["Understanding Convolutional Neural Networks for NLP"](#)

Tuning the Network Size

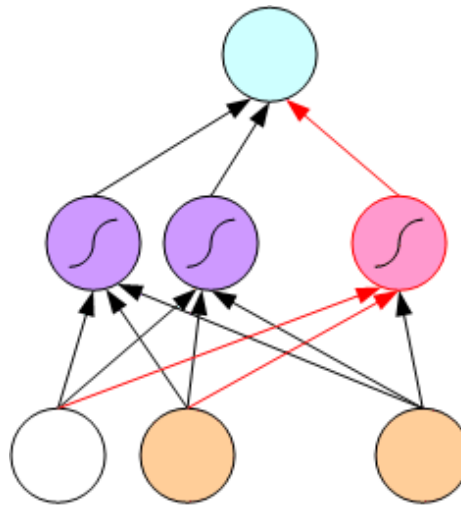
32

□ Destructive

Weight decay:

$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i$$

$$E' = E + \frac{\lambda}{2} \sum_i w_i^2$$

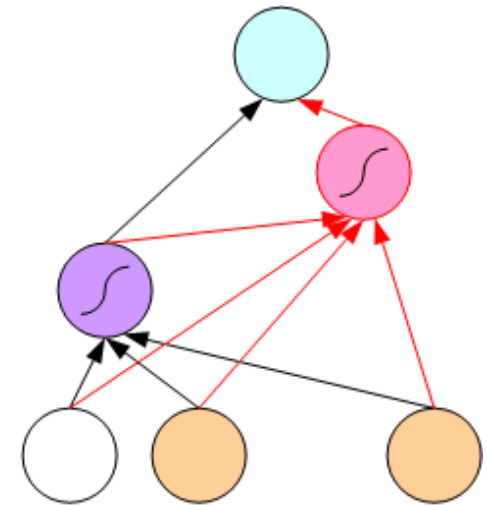


Dynamic Node Creation

(Ash, 1989)

□ Constructive

Growing networks



Cascade Correlation

(Fahlman and Lebiere, 1989)

Bayesian Learning

33

- Consider weights w_i as random vars, prior $p(w_i)$

$$p(\mathbf{w} | \mathbf{X}) = \frac{p(\mathbf{X} | \mathbf{w})p(\mathbf{w})}{p(\mathbf{X})} \quad \hat{\mathbf{w}}_{MAP} = \arg \max_{\mathbf{w}} \log p(\mathbf{w} | \mathbf{X})$$

$$\log p(\mathbf{w} | \mathbf{X}) = \log p(\mathbf{X} | \mathbf{w}) + \log p(\mathbf{w}) + C$$

$$p(\mathbf{w}) = \prod_i p(w_i) \quad \text{where} \quad p(w_i) = c \cdot \exp\left[-\frac{w_i^2}{2(1/2\lambda)}\right]$$

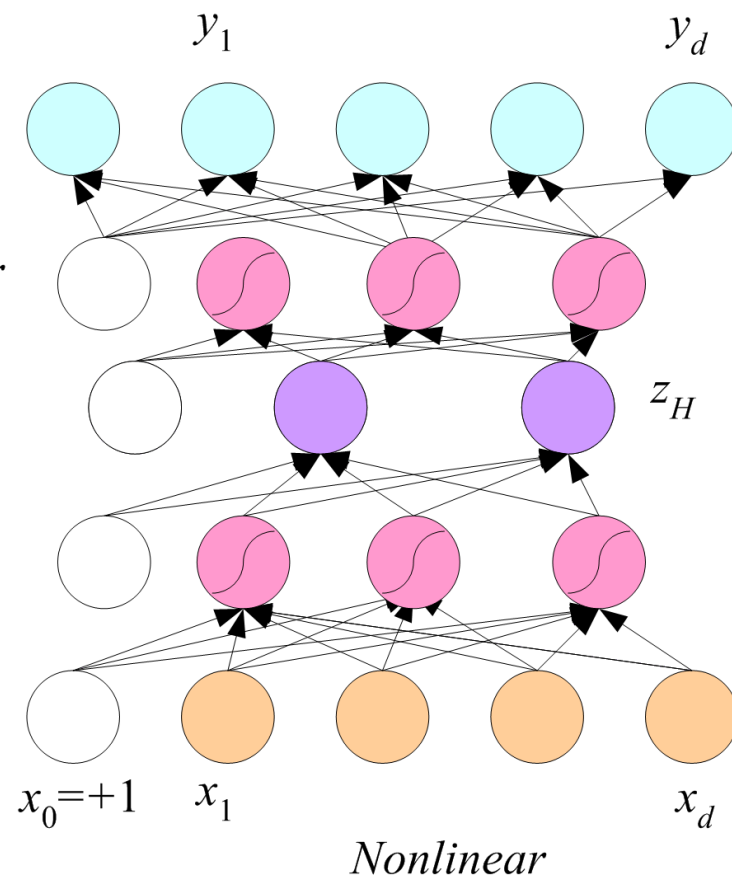
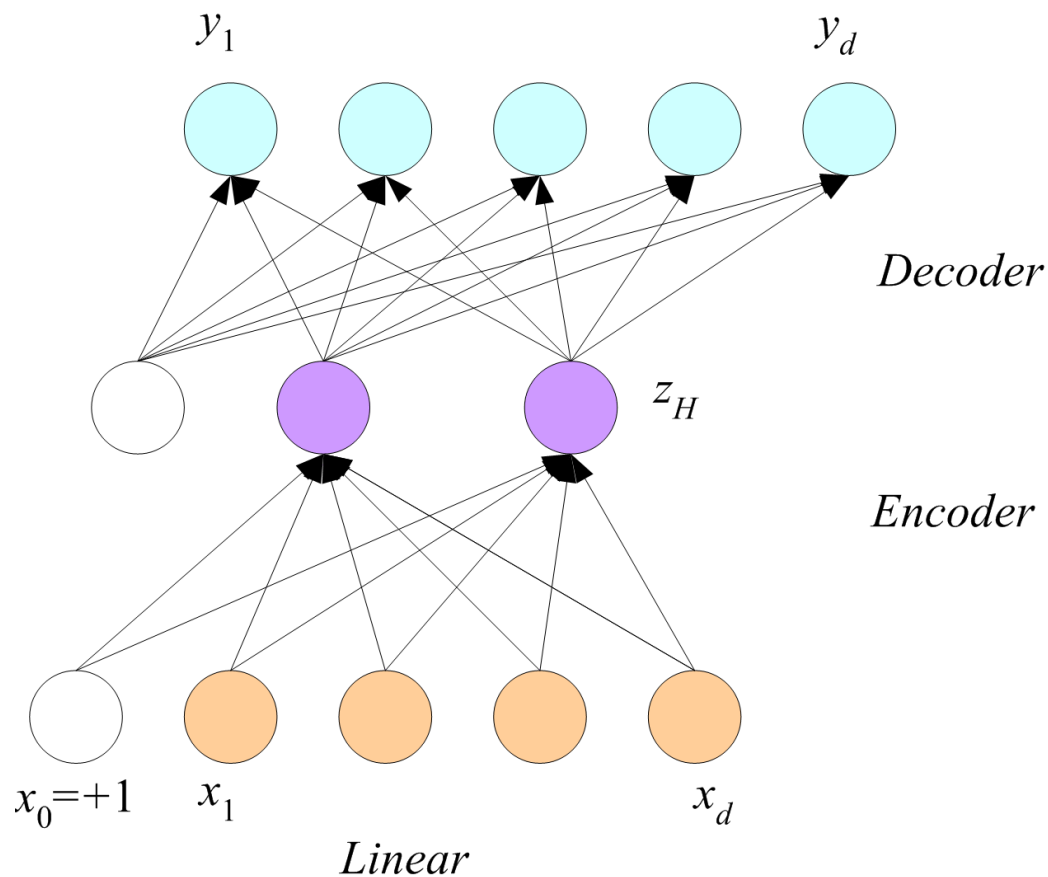
$$E' = E + \lambda \|\mathbf{w}\|^2$$

- Weight decay, ridge regression, regularization
cost=data-misfit + λ complexity

More about Bayesian methods in chapter 14

Dimensionality Reduction

34



Autoencoder networks

Autoencoder: An Example

35

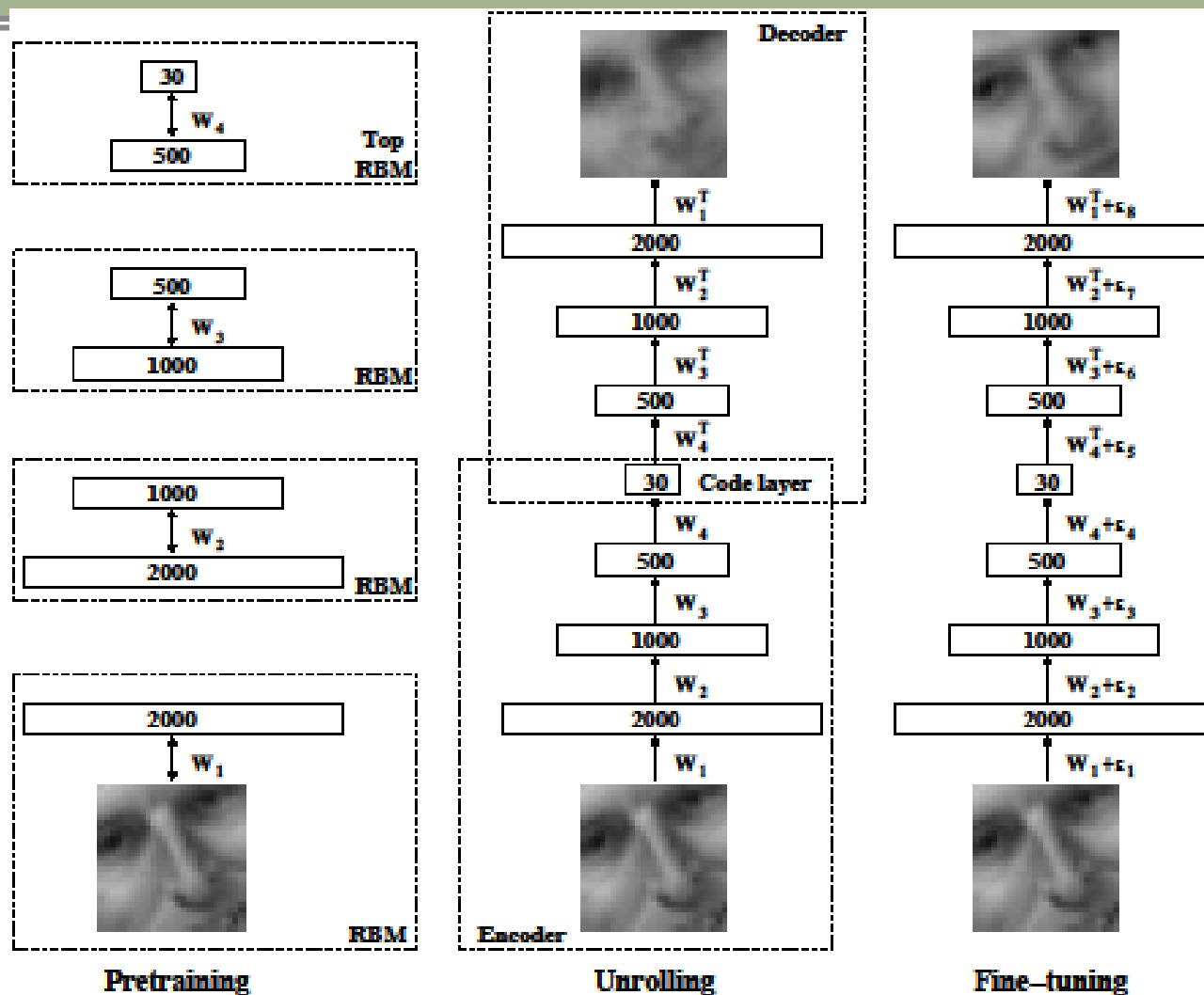
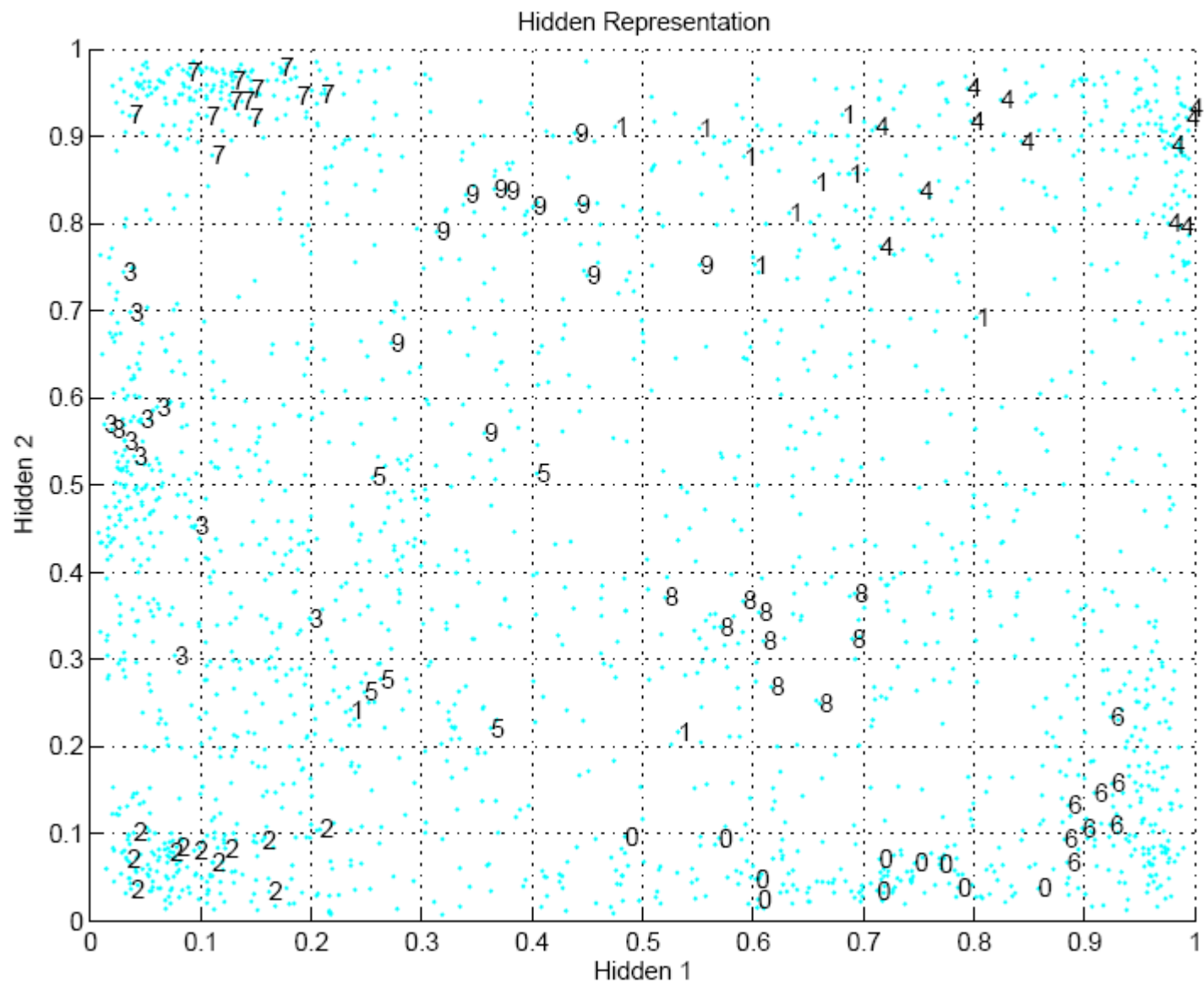


Figure from Hinton and Salakhutdinov, Science, 2006



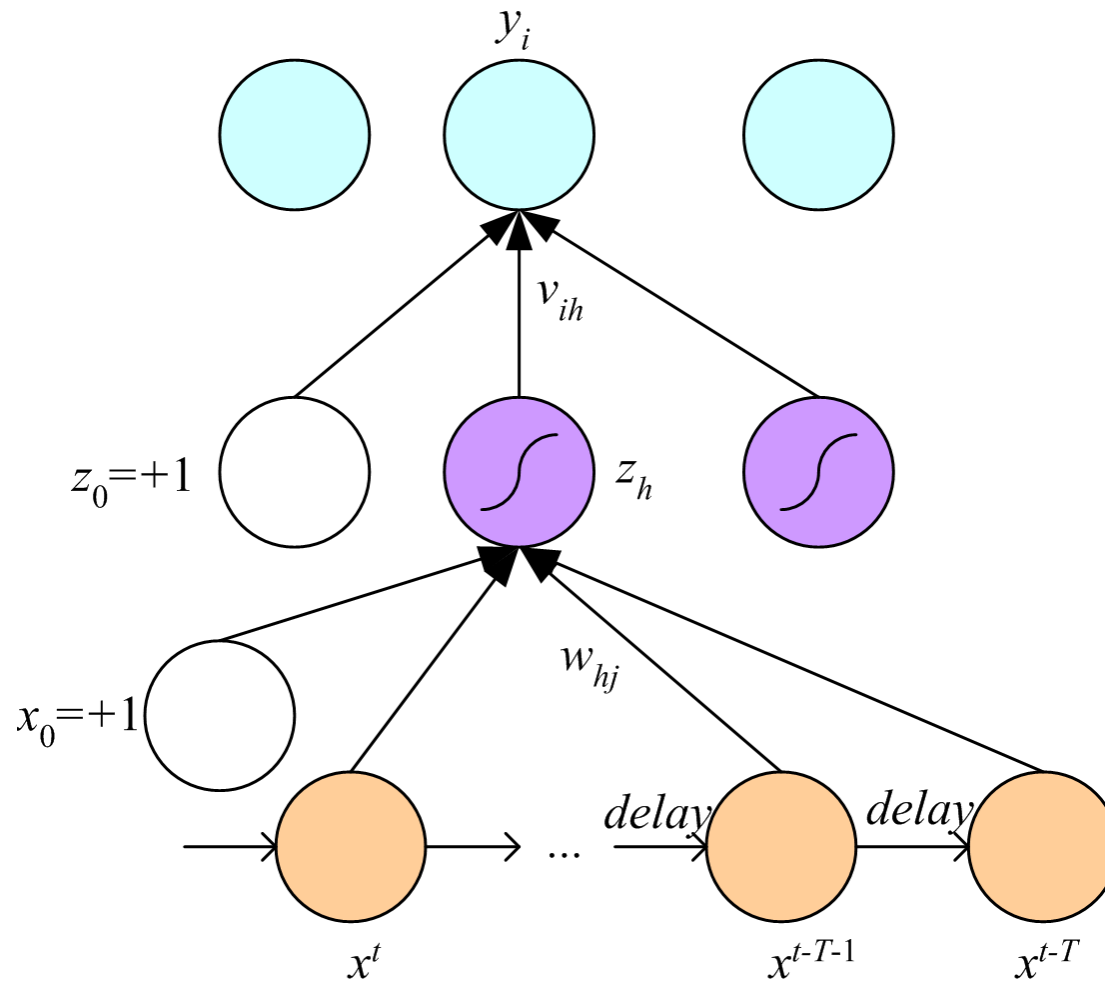
Learning Time

37

- Applications:
 - ▣ Sequence recognition: Speech recognition
 - ▣ Sequence reproduction: Time-series prediction
 - ▣ Sequence association
- Network architectures
 - ▣ Time-delay networks (Waibel et al., 1989)
 - ▣ Recurrent networks (Rumelhart et al., 1986)

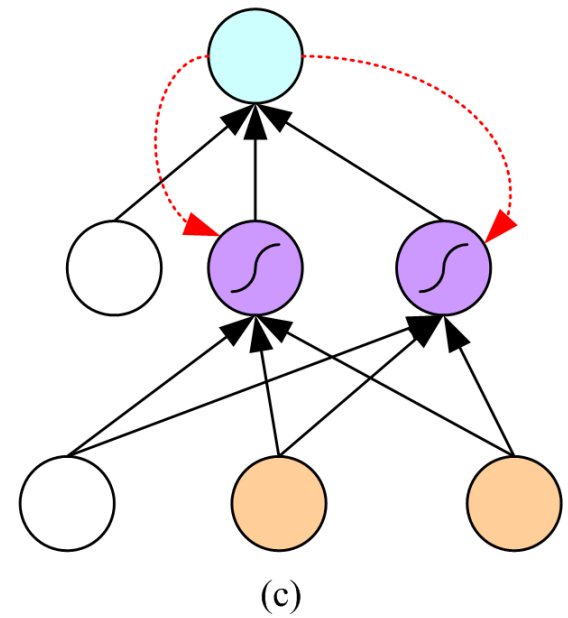
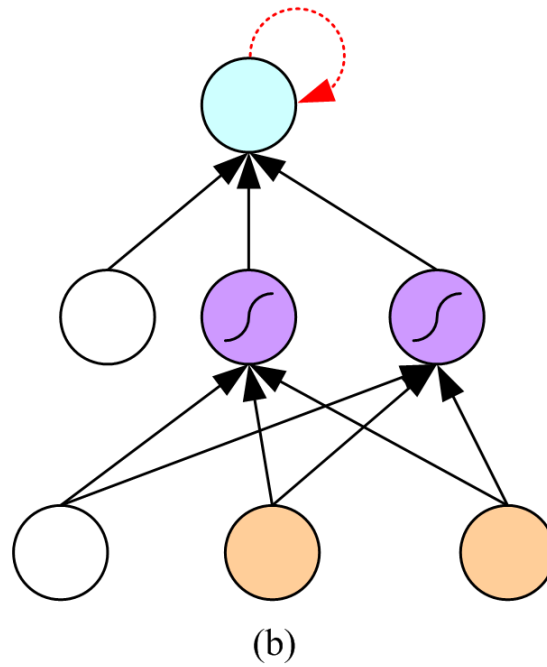
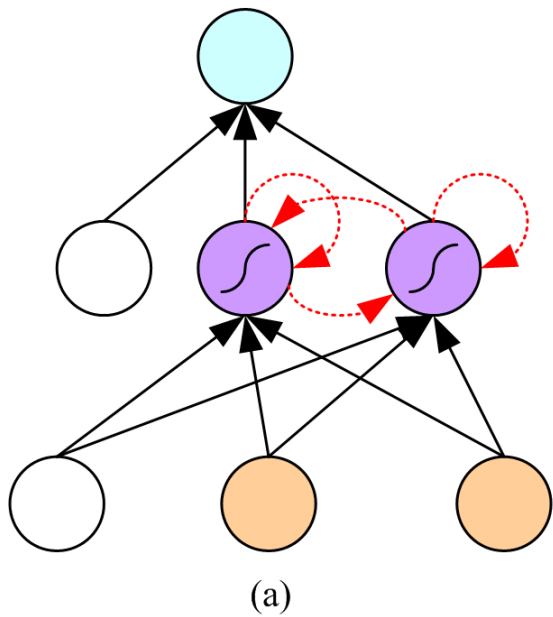
Time-Delay Neural Networks

38



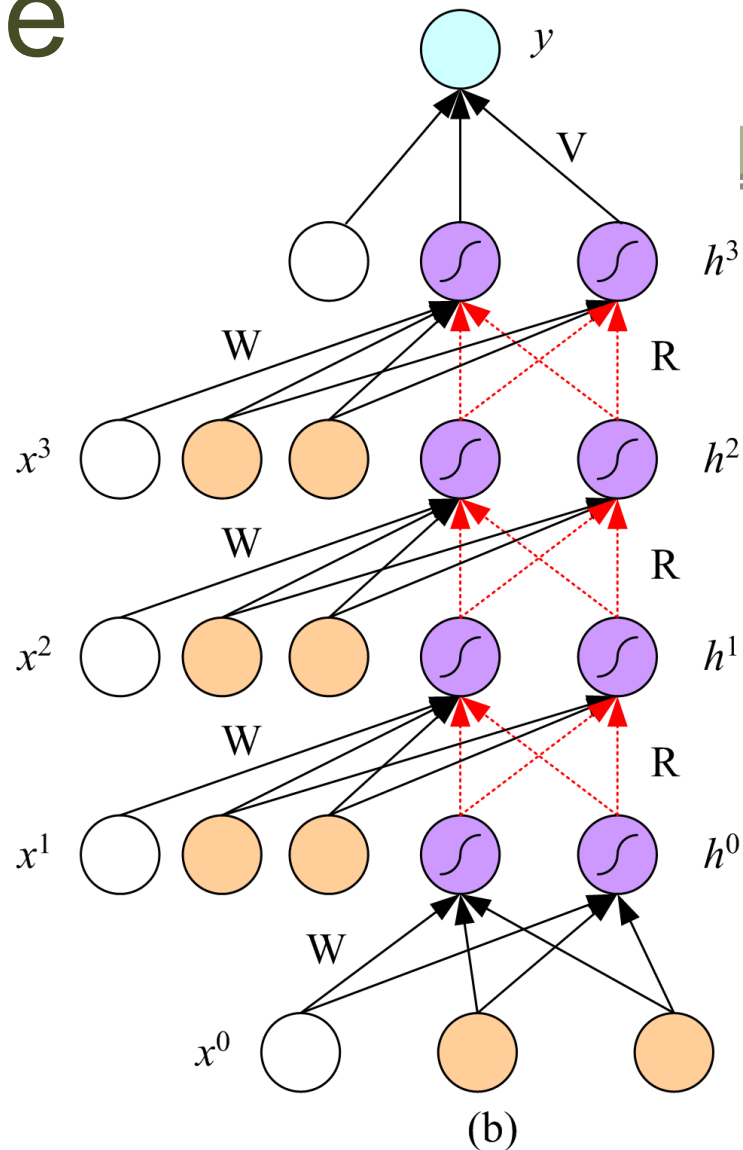
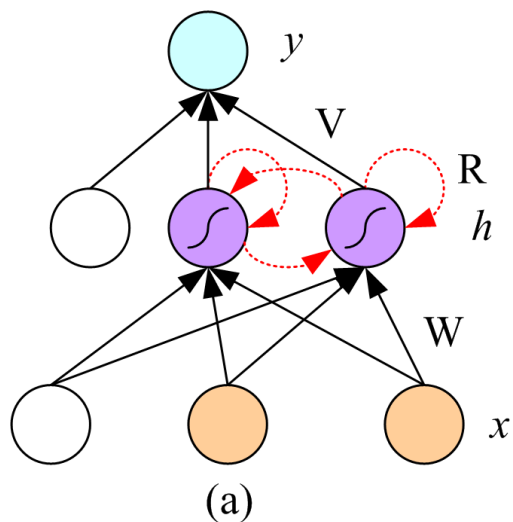
Recurrent Networks

39



Unfolding in Time

40



Deep Networks

41

- Layers of feature extraction units
- Can have local receptive fields as in convolution networks, or can be fully connected
- Can be trained layer by layer using an autoencoder in an unsupervised manner
- No need to craft the right features or the right basis functions or the right dimensionality reduction method; learns **multiple layers of abstraction** all by itself given a lot of data and a lot of computation
- Applications in vision, language processing, ...