

Lecture Slides for

INTRODUCTION TO MACHINE LEARNING

3RD EDITION

ETHEM ALPAYDIN

© The MIT Press, 2014

alpaydin@boun.edu.tr

<http://www.cmpe.boun.edu.tr/~ethem/i2ml3e>

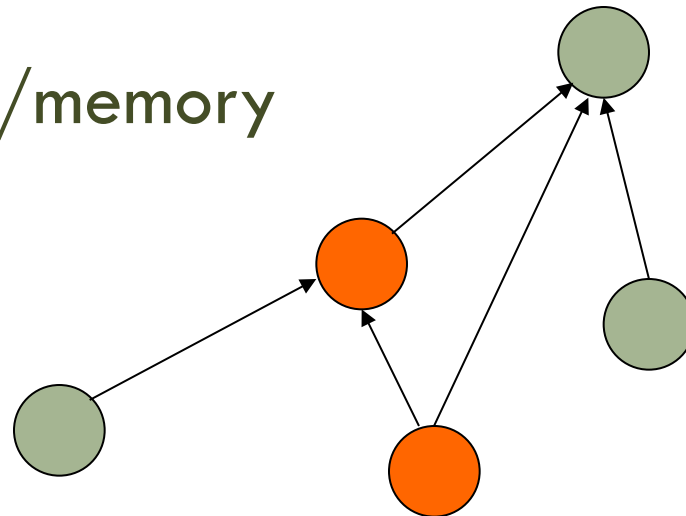
CHAPTER 11:

MULTILAYER PERCEPTRONS

Neural Networks

3

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



Understanding the Brain

4

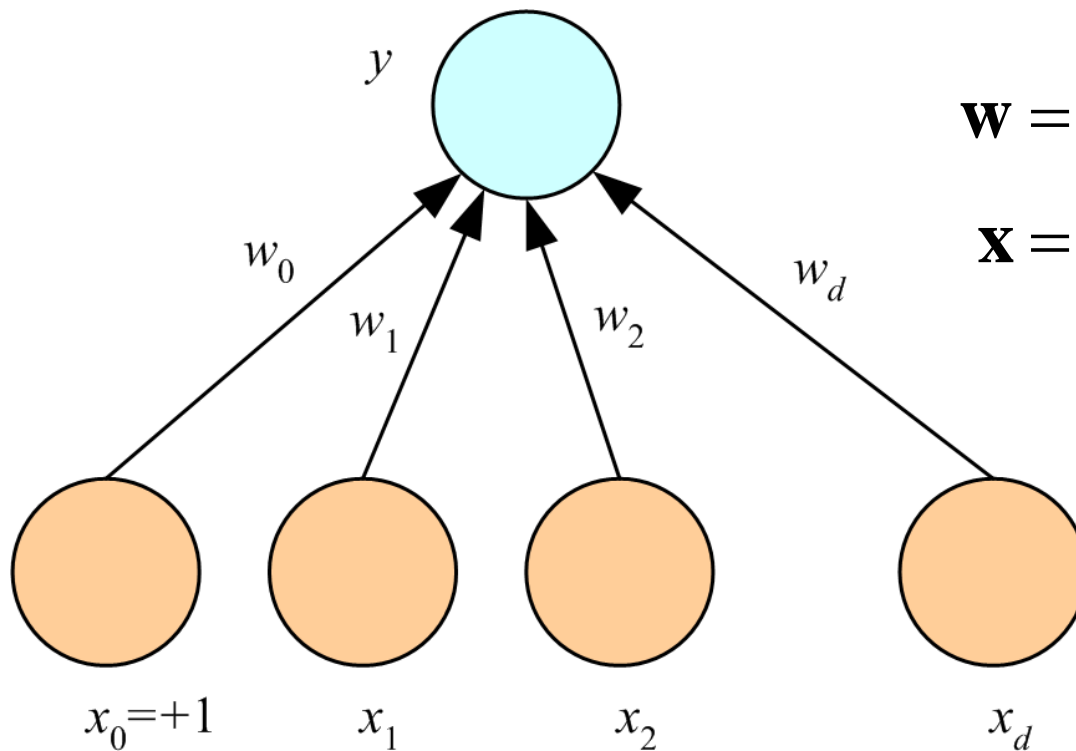
- Levels of analysis (Marr, 1982)
 1. Computational theory
 2. Representation and algorithm
 3. Hardware implementation
- Reverse engineering: From hardware to theory
- Parallel processing: SIMD vs MIMD

Neural net: SIMD with modifiable local memory

Learning: Update by training/experience

Perceptron

5



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

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

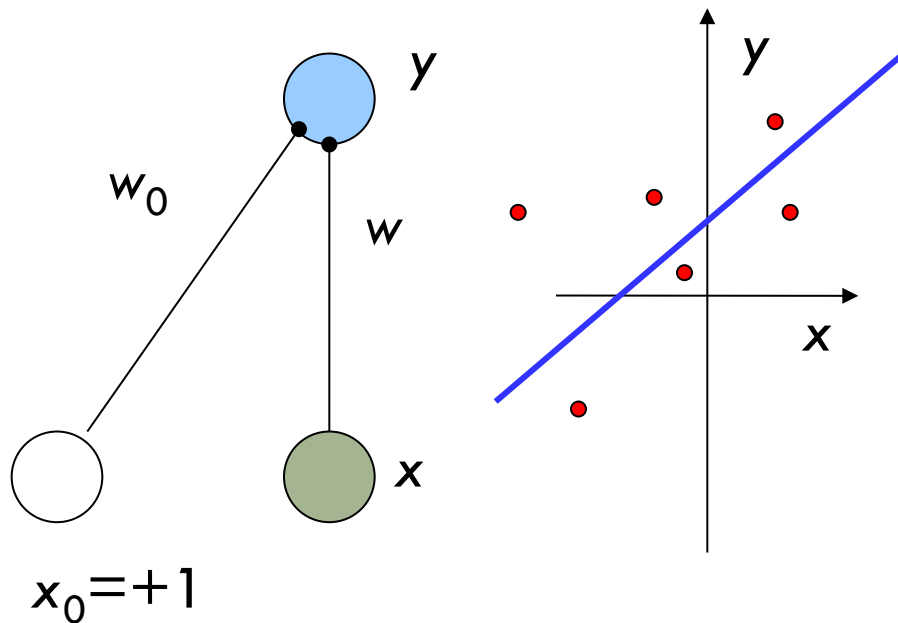
$$\mathbf{x} = [1, x_1, \dots, x_d]^T$$

(Rosenblatt, 1962)

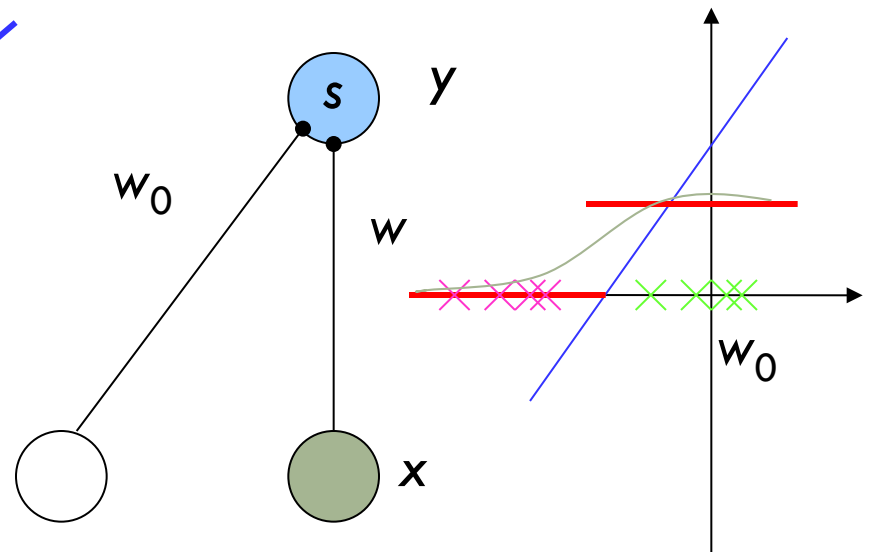
What a Perceptron Does

6

□ Regression: $y = wx + w_0$



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



$$y = \text{sigmoid}(o) = \frac{1}{1 + \exp[-\mathbf{w}^T \mathbf{x}]}$$

K Outputs

7

Regression:

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

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

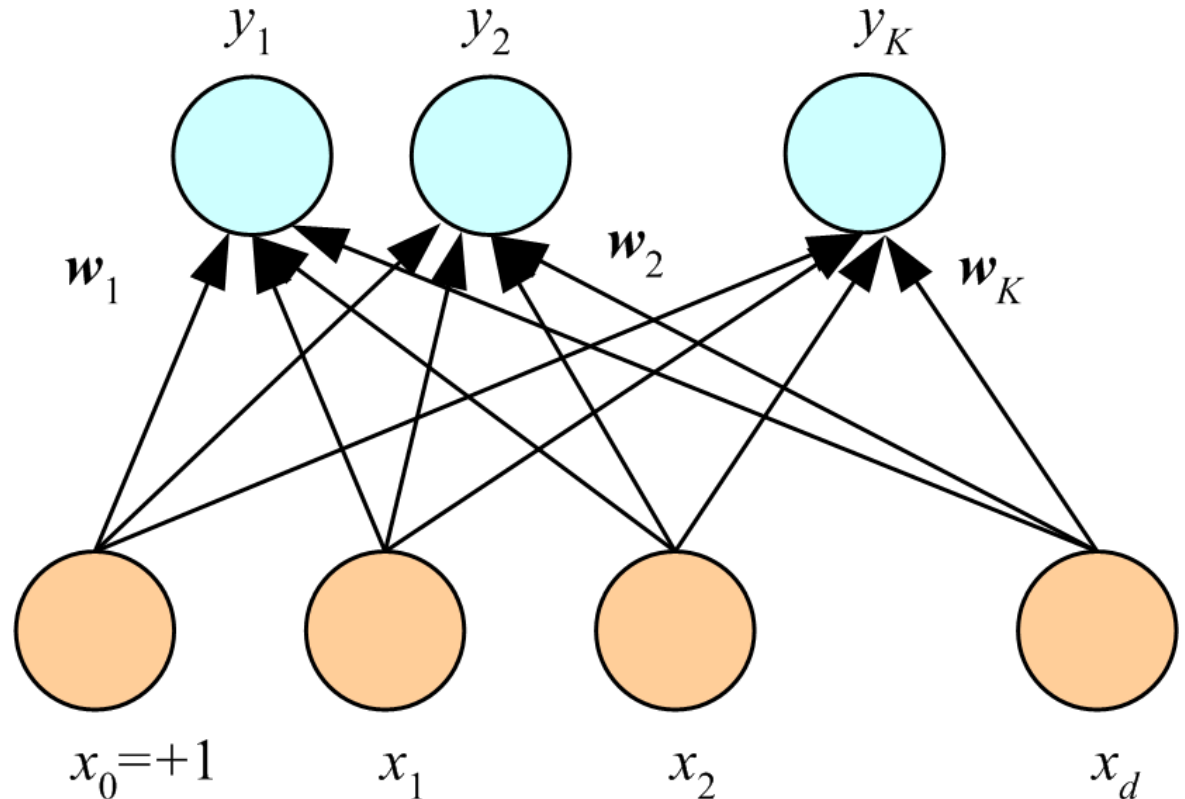
Classification:

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

$$y_i = \frac{\exp o_i}{\sum_k \exp o_k}$$

choose C_i

if $y_i = \max_k y_k$



Training

8

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

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

$$\text{Update} = \text{LearningFactor} \cdot (\text{DesiredOutput} - \text{ActualOutput}) \cdot \text{Input}$$

Training a Perceptron: Regression

- Regression (Linear output):

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

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

- Stochastic optimization, online learning
 - ▣ Update weights based on one sample at a time
 - ▣ Stochastic gradient descent (SGD)

Classification

10

□ Single sigmoid output

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

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

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

□ $K > 2$ softmax outputs

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

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

□ Stochastic gradient descent: one sample at a time

Classification using Square Loss

11

□ Single sigmoid output

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

$$E^t(\mathbf{w} \mid \mathbf{x}^t, \mathbf{r}^t) = (r^t - y^t)^2$$

$$\Delta w_j^t = \eta (r^t - y^t) y^t (1 - y^t) x_j^t$$

□ Training: Square loss vs. Cross entropy

- ▣ Cross entropy leads to a convex optimization problem

- ▣ Square loss leads to a non-convex optimization problem

□ Classification with one layer: Logistic regression

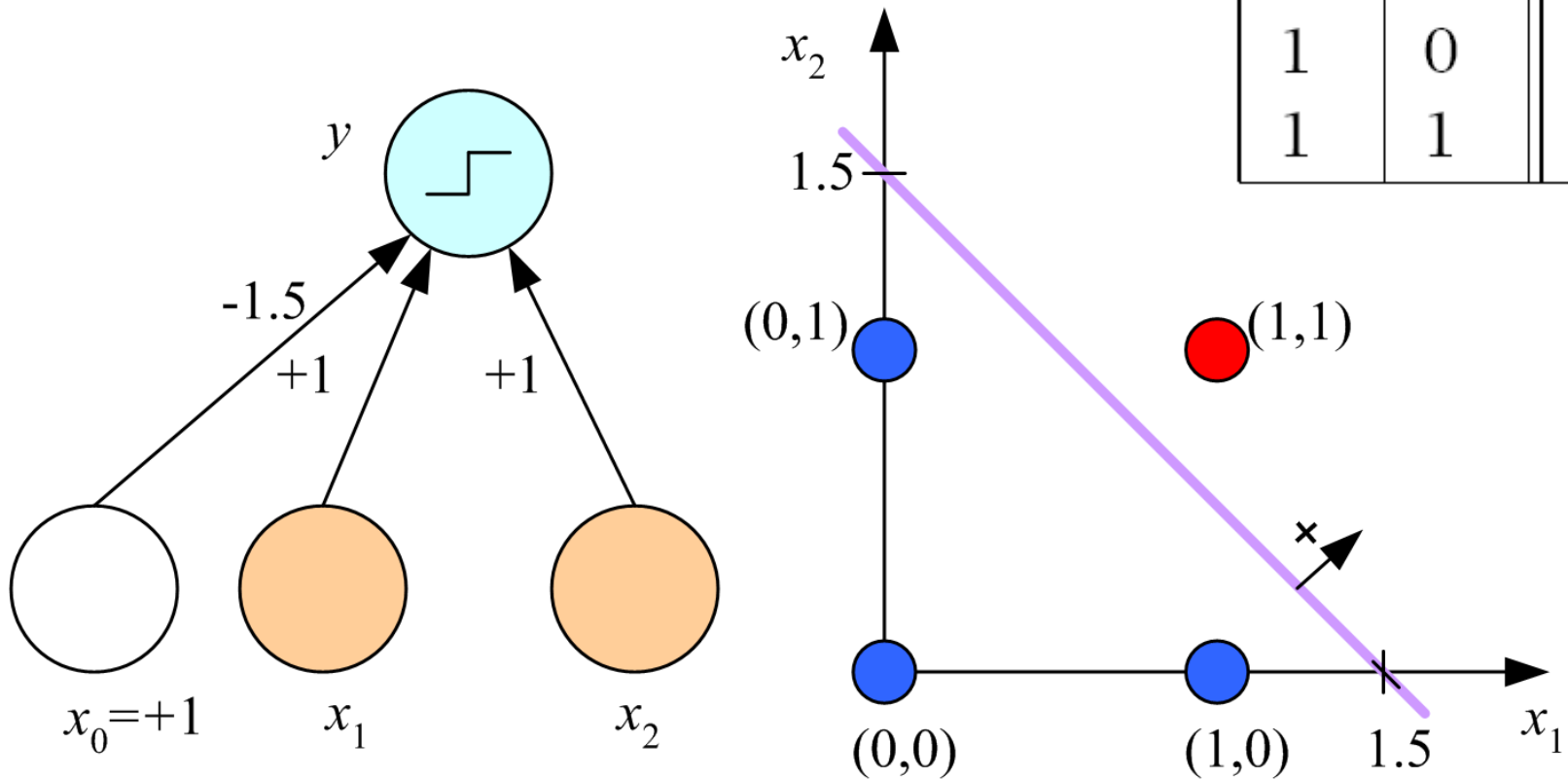
- ▣ Soft-max/sigmoid transfer function

- ▣ Cross entropy as loss function

Learning Boolean AND

13

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

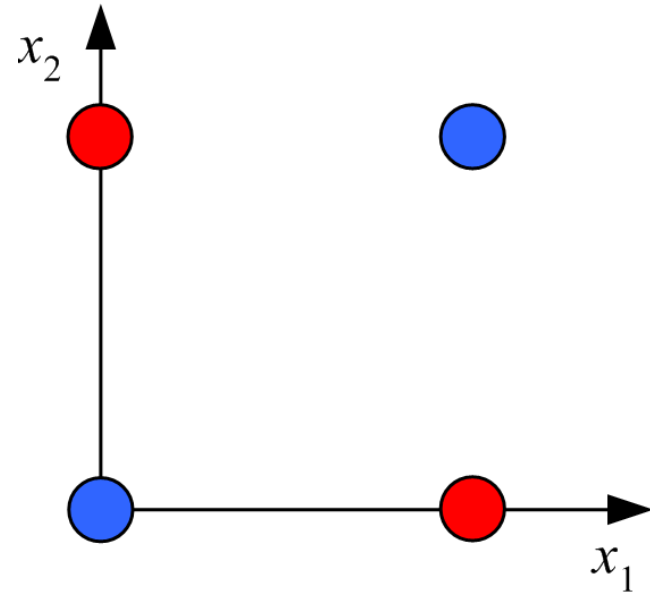
□ No w_0, w_1, w_2 satisfy:

$$w_0 \leq 0$$

$$w_2 + w_0 > 0$$

$$w_1 + w_0 > 0$$

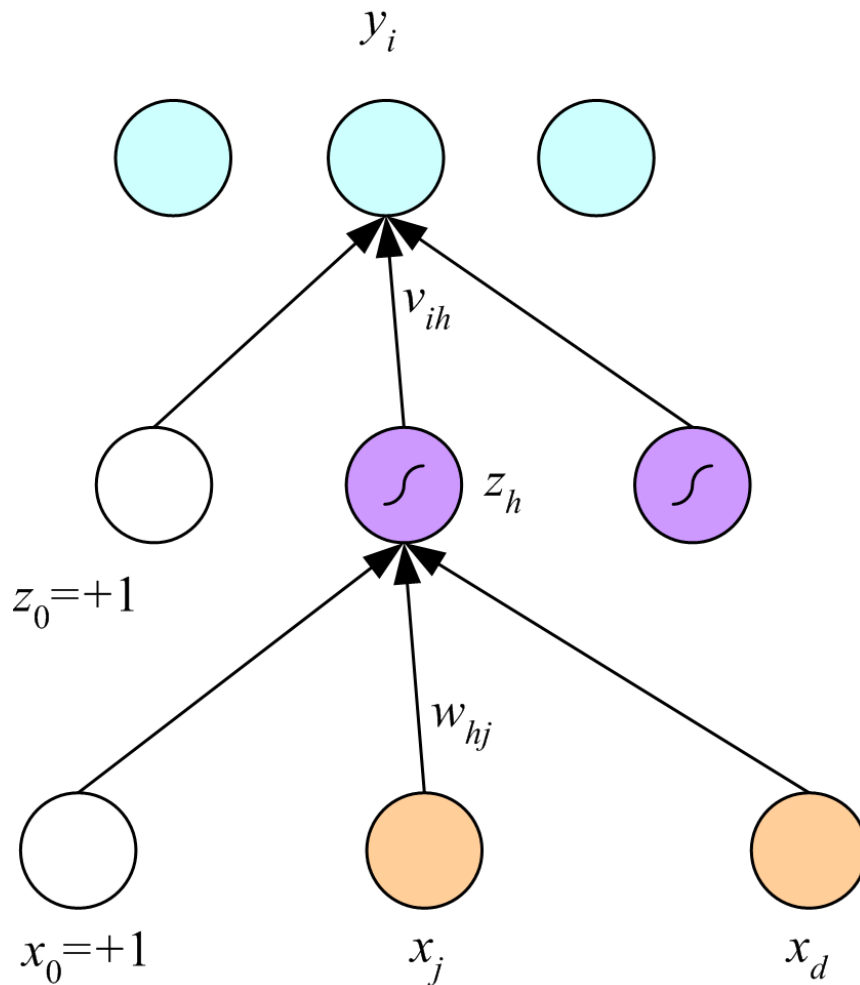
$$w_1 + w_2 + w_0 \leq 0$$



(Minsky and Papert, 1969)

Multilayer Perceptrons

15



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

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

(Rumelhart et al., 1986)

MLP Transfer: Sigmoid vs ReLU

16

- Traditional transfer function: sigmoid (or tanh)

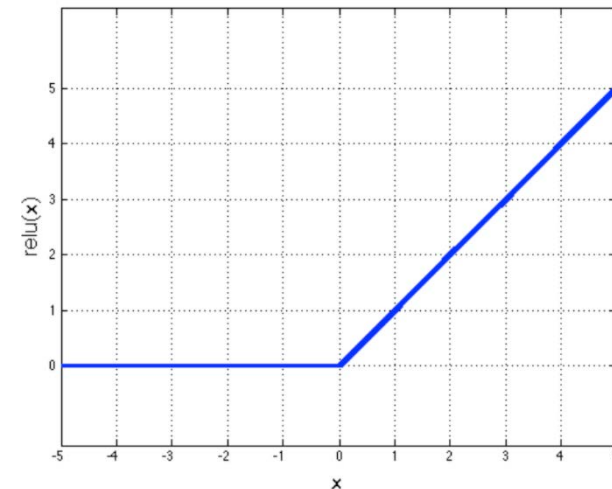
$$z_h = \text{sigmoid}(\mathbf{w}_h^T \mathbf{x}) = \frac{1}{1 + \exp\left[-\left(\sum_{j=1}^d w_{hj} x_j + w_{h0}\right)\right]}$$

- Rectified Linear Unit (ReLU)

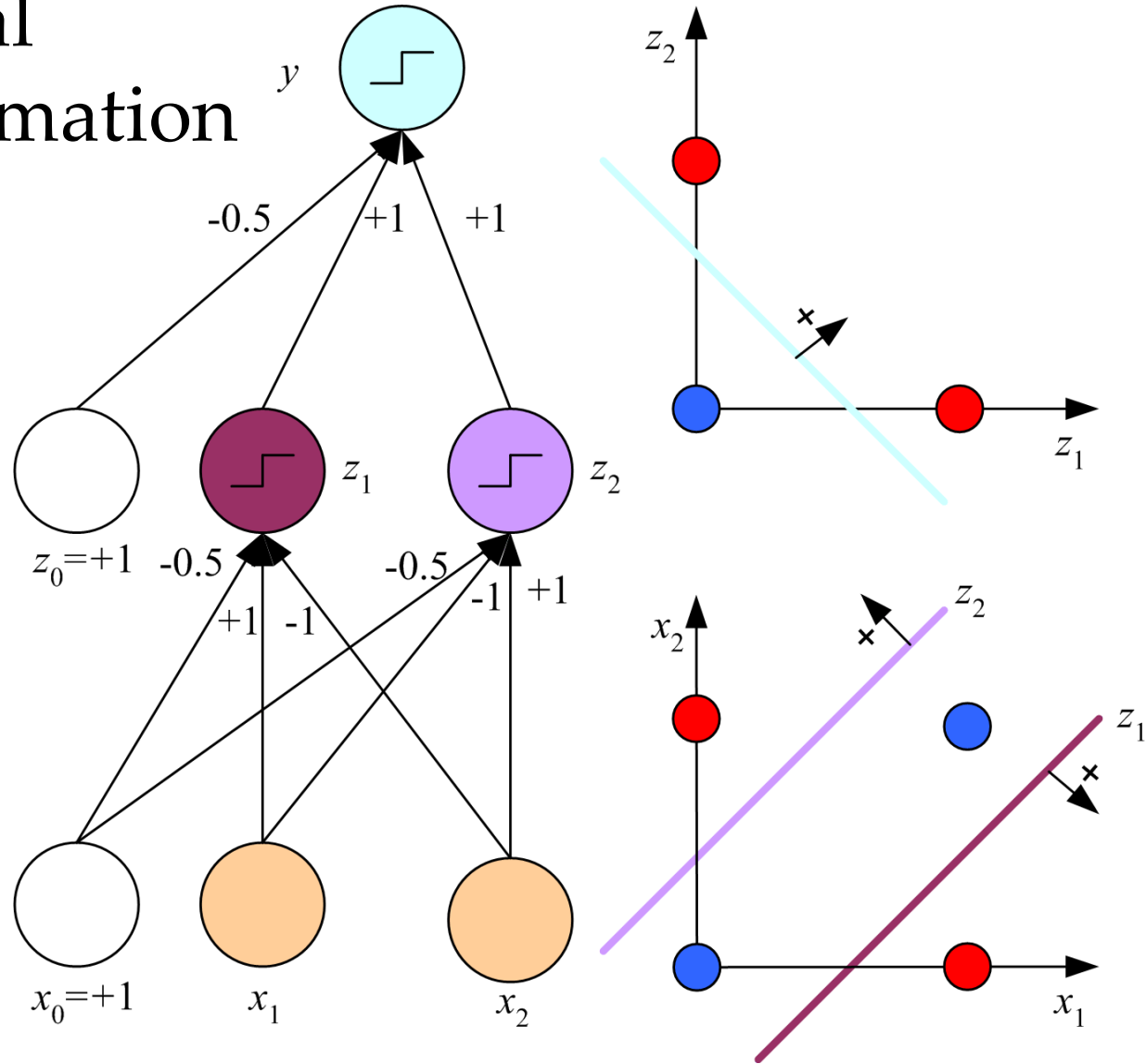
$$z_h = \text{ReLU}(\mathbf{w}_h^T \mathbf{x} + w_{h,0}) \quad \text{ReLU}(a) = \max(0, a)$$

- Use ReLU (and variants)

- Simplifies gradient descent
- Makes learning faster
- Avoids (sigmoid) saturation



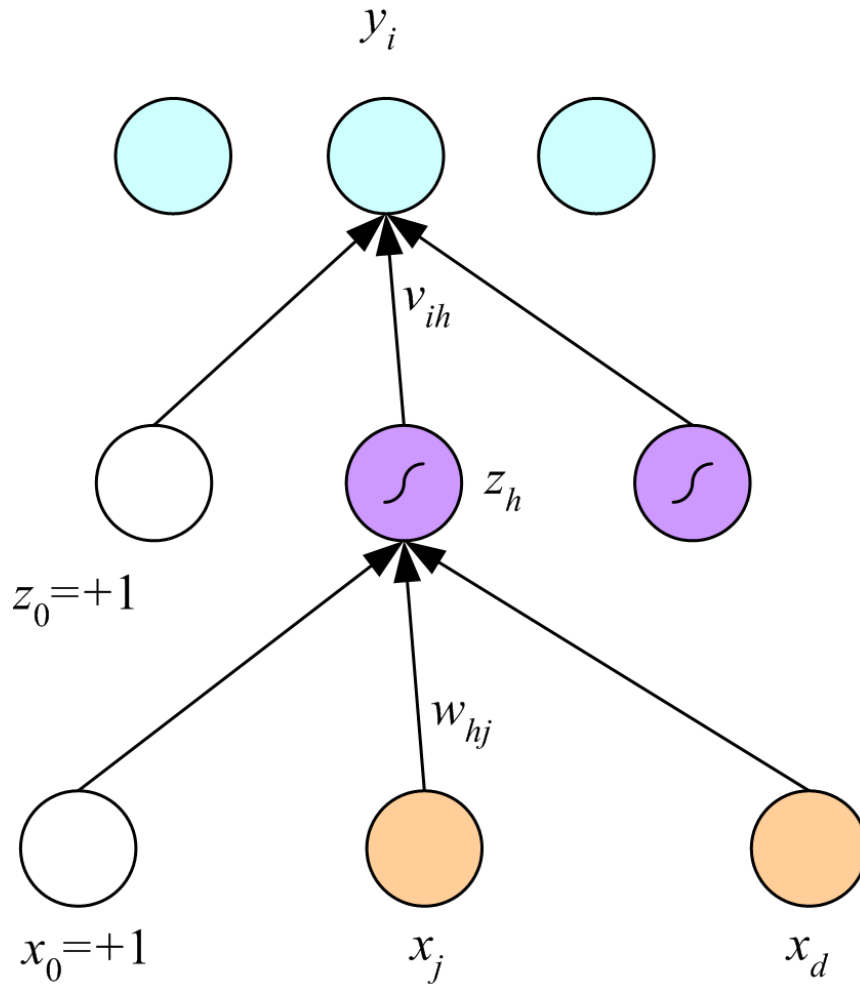
Universal Approximation



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

Regression, Backpropagation

18



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

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

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

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

E is the loss at the output,
e.g., $(r^t - y^t)^2$

Regression with Multiple Outputs

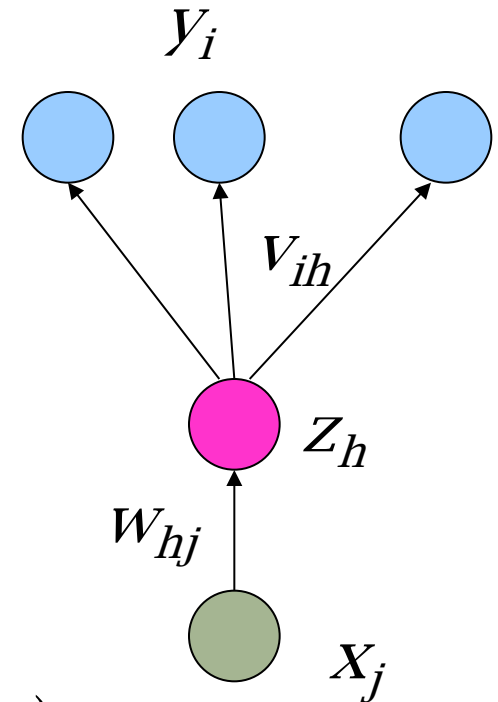
19

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

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

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

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



Backpropagation

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

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

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

Forward

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

\mathbf{x}

Backward

$$\begin{aligned} \Delta \mathbf{w}_{hj} &= -\eta \frac{\partial E}{\partial \mathbf{w}_{hj}} \\ &= -\eta \sum_t \frac{\partial E}{\partial y^t} \frac{\partial y^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial \mathbf{w}_{hj}} \\ &= -\eta \sum_t -(r^t - y^t) v_h \boxed{z_h^t (1 - z_h^t)} x_j^t \\ &= \eta \sum_t (r^t - y^t) v_h z_h^t (1 - z_h^t) x_j^t \end{aligned}$$

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

Backward

$$\Delta v_{i,h} = \eta (r_i^t - y_i^t) z_h^t$$

Backward

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

$$\Delta v_{i,h} = \eta \Delta_i^t z_h^t = \eta \times \text{error} \times \text{input}$$

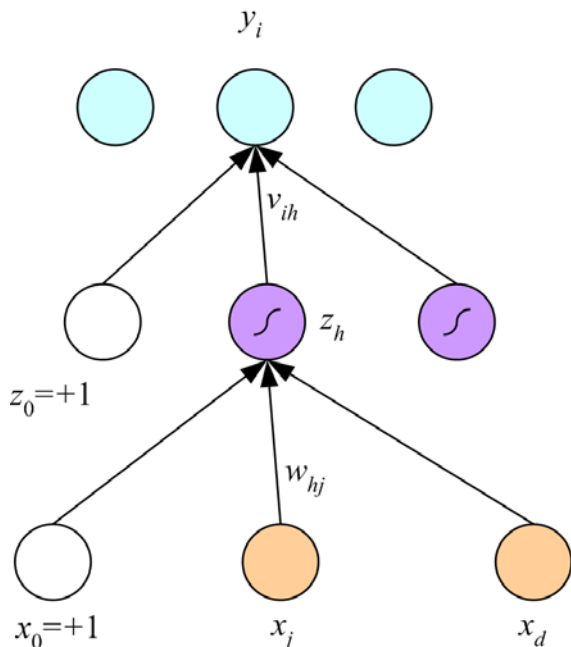
$$\Delta_i^t = (r_i^t - y_i^t) = \text{error}$$

$$\Delta w_{h,j} = \eta \Delta_h^t x_j^t = \eta \times \text{error} \times \text{input}$$

$$\Delta_h^t = \sum_i \Delta_i^t v_{ih} g'(a_h) = \text{backpropagated error}$$

g : transfer (e.g., sigmoid)

$$a_h = \sum_j w_{hj} x_j + x_0$$



Initialize all v_{ih} and w_{hj} to $\text{rand}(-0.01, 0.01)$

Repeat

For all $(\mathbf{x}^t, r^t) \in \mathcal{X}$ in random order

For $h = 1, \dots, H$

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

For $i = 1, \dots, K$

$$y_i = \mathbf{v}_i^T \mathbf{z}$$

For $i = 1, \dots, K$

$$\Delta \mathbf{v}_i = \eta(r_i^t - y_i^t) \mathbf{z}$$

For $h = 1, \dots, H$

$$\Delta \mathbf{w}_h = \eta\left(\sum_i (r_i^t - y_i^t) v_{ih}\right) z_h (1 - z_h) \mathbf{x}^t$$

For $i = 1, \dots, K$

$$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$$

For $h = 1, \dots, H$

$$\mathbf{w}_h \leftarrow \mathbf{w}_h + \Delta \mathbf{w}_h$$

Until convergence

Deep Networks

42

- 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, ...