

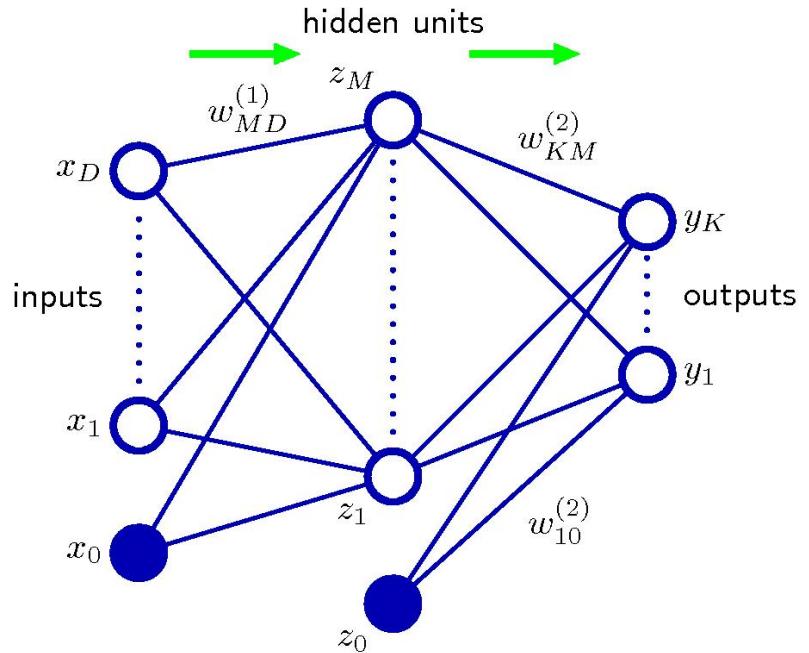
Multilayer Perceptron

Rui Kuang

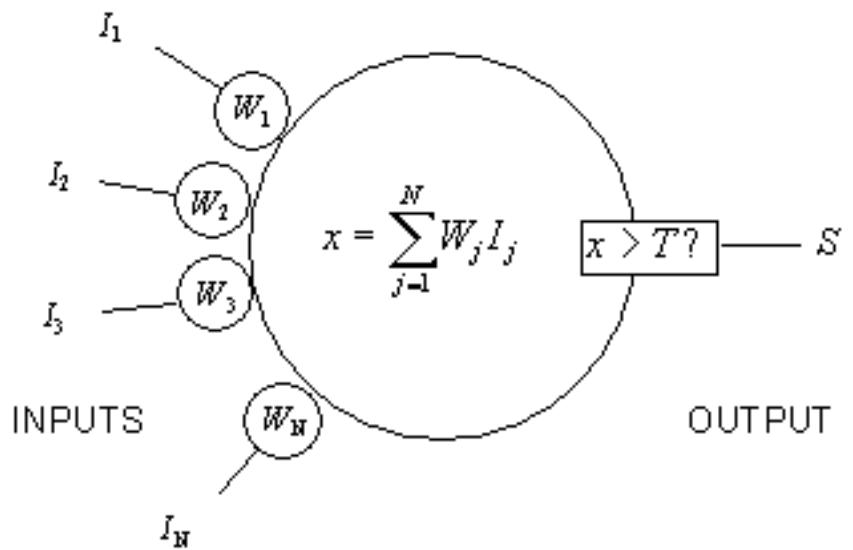
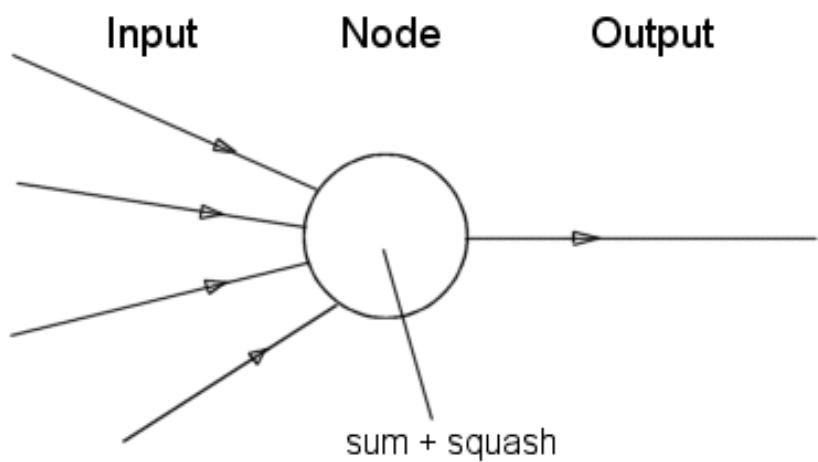
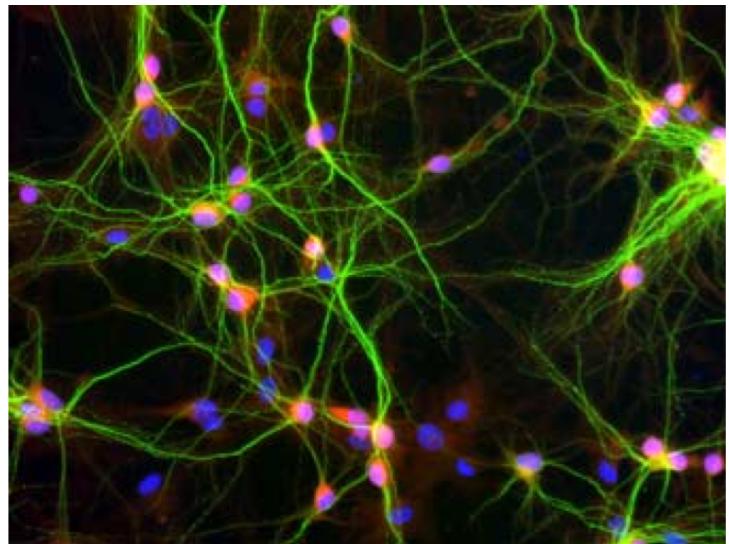
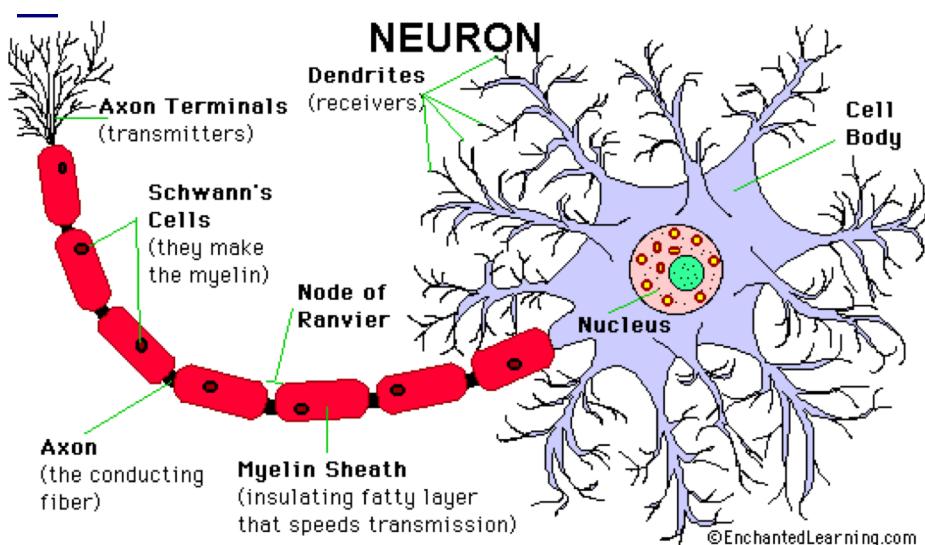
Department of Computer Science and Engineering
University of Minnesota

Artificial Neural Networks

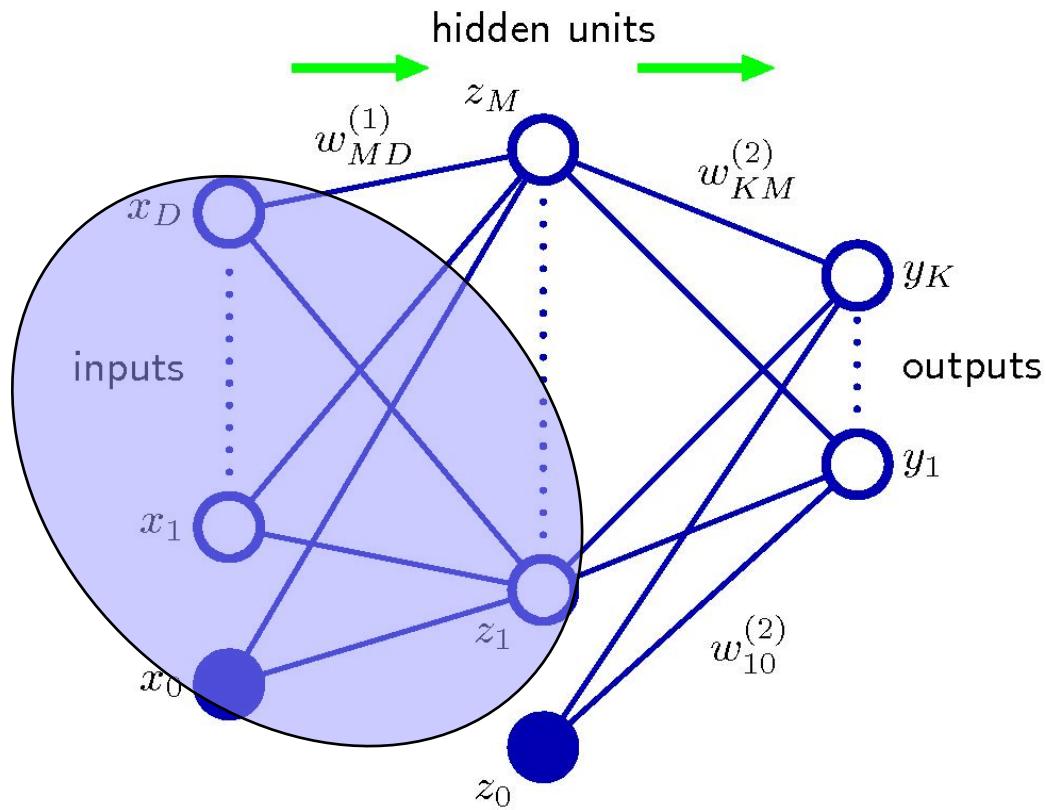
- Network models with processing units (variables) with connections (linear coefficients) between them.
- Motivated by human brain for information processing.
- Structures of hidden units incorporate prior knowledge.
- Basic units are each perceptron.



■ Neurone vs. Node

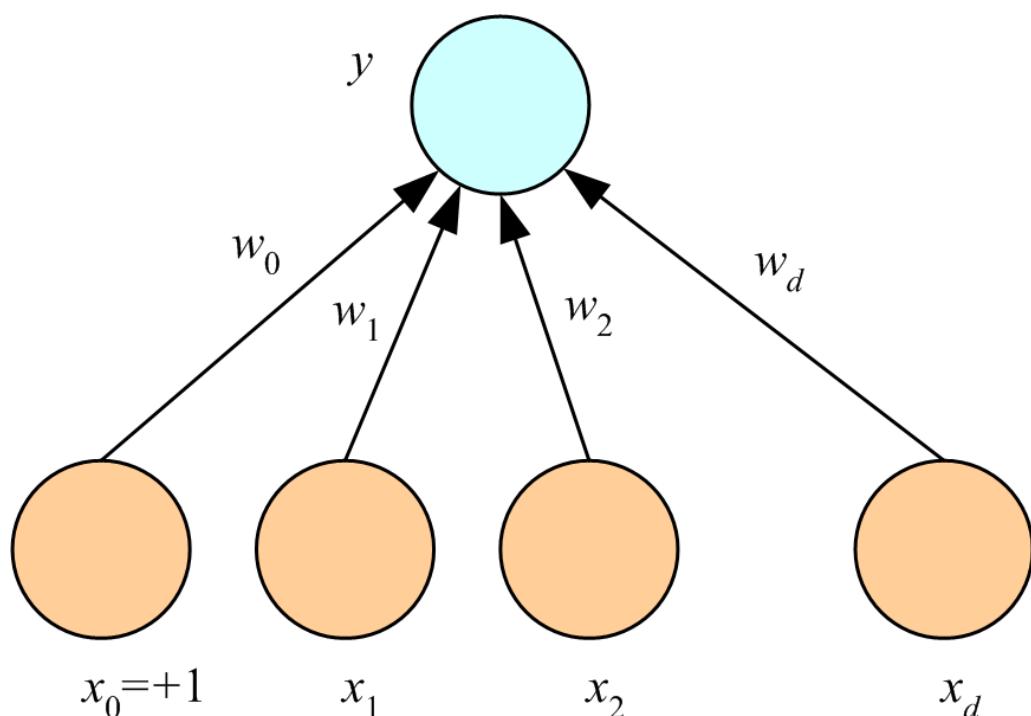


Artificial Neural Networks



- Basic processing units are called perceptrons.
- The perceptrons are connected to form multilayers.
- When evaluating the likelihood, errors are propagated through the network.

Perceptron



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

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

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

(Rosenblatt, 1962)

Perceptron vs Gradient Descent

- The same update rule in batch model or online model.

Gradient Decent :

$$\mathbf{w} = \mathbf{0};$$

Repeat until converge

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

$$\mathbf{w} = \mathbf{w} + \eta \mathbf{X} (\mathbf{r} - \mathbf{y})^T;$$

Perceptron :

$$\mathbf{w} = \mathbf{0};$$

Repeat until converge

for each example $(\mathbf{x}^t, \mathbf{r}^t)$:

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

$$\mathbf{w} = \mathbf{w} + \eta (r^t - y^t) \mathbf{x}^t;$$

Perceptron Training

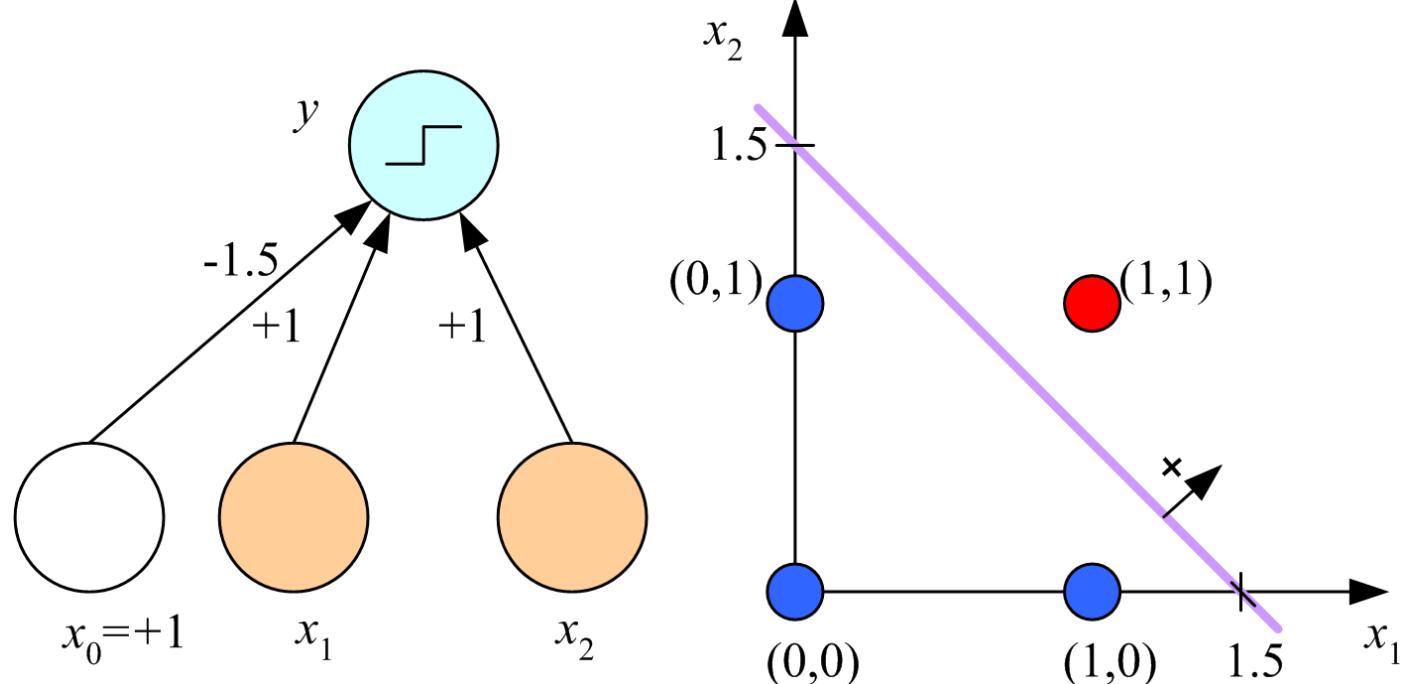
- Online training: instances seen one by one. (vs batch learning):
 - No need to store the whole sample
 - Problem may change in time
 - Wear and degradation in system components
- Stochastic gradient-descent: Update after a single pattern (or a subset of samples)

$$\Delta \mathbf{w}^t = \eta (r^t - y^t) \mathbf{x}^t$$

Update=LearningFactor·(DesiredOutput–ActualOutput)·Input

Learning Boolean AND

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

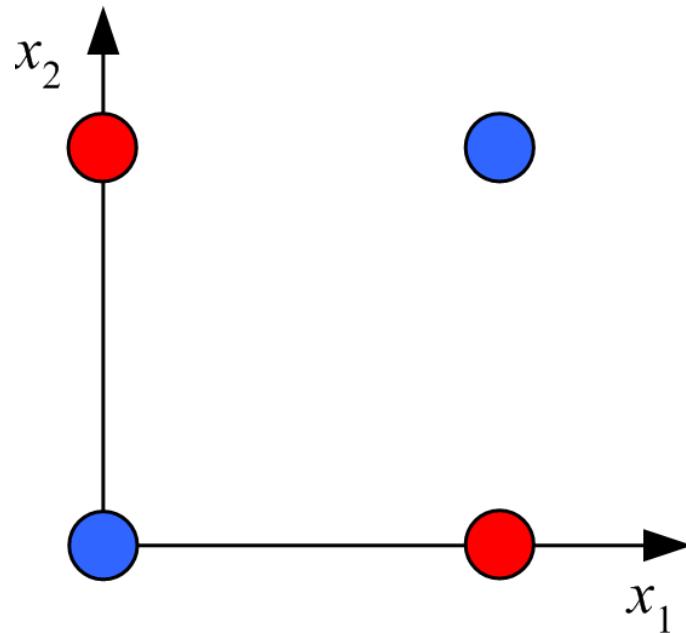


$$y = \text{step}(w^T x) = \text{step}(w_0 + w_1 x_1 + w_2 x_2)$$

$$\text{step}(a) = \begin{cases} 1 & a > 0 \\ 0 & a \leq 0 \end{cases}$$

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:

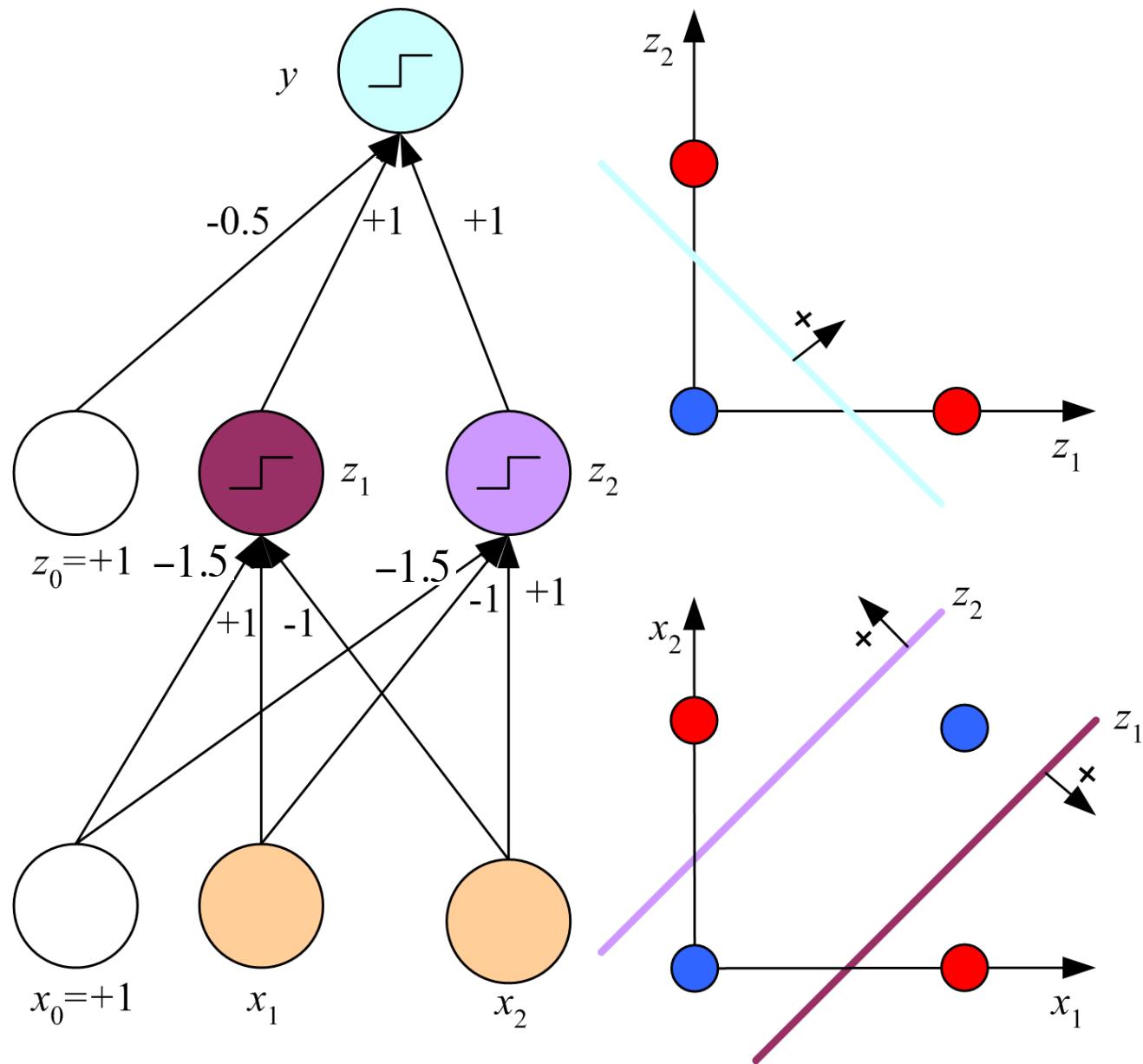
$$w_0 \leq 0$$

$$w_2 + w_0 > 0$$

$$w_1 + w_0 > 0$$

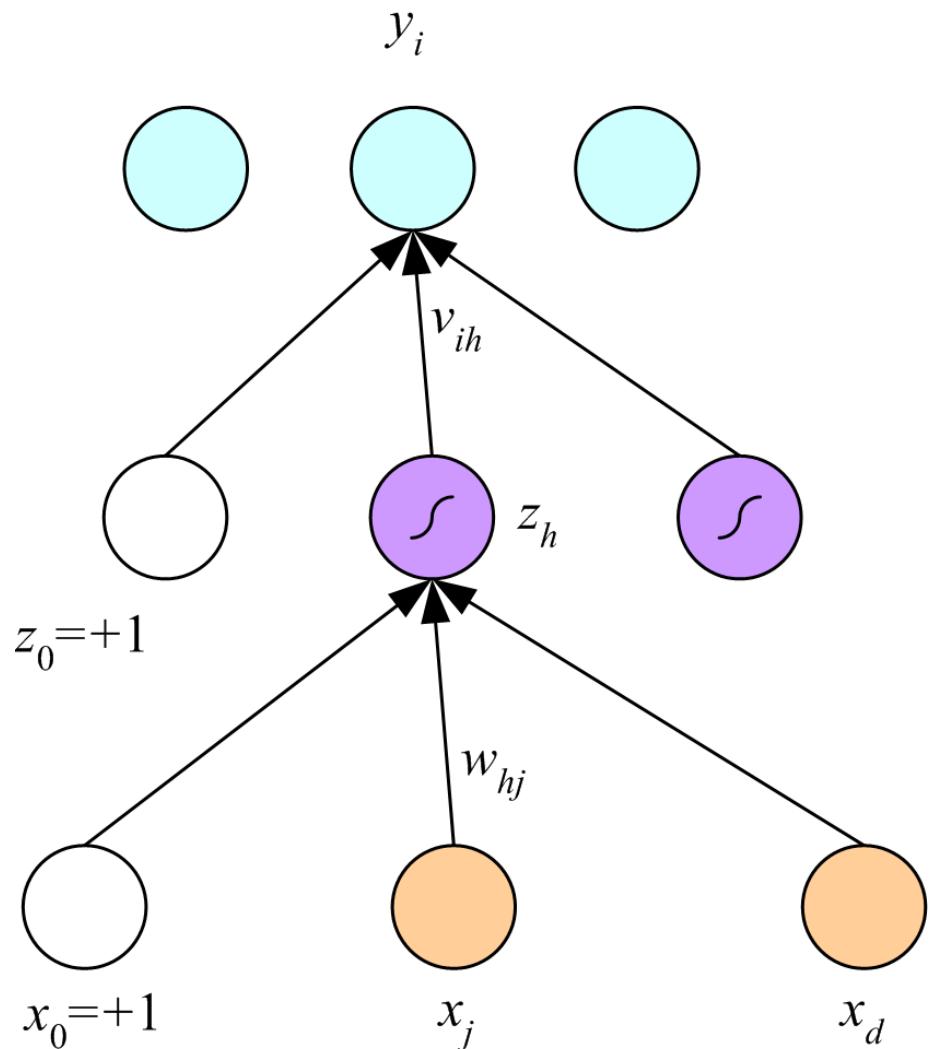
$$w_1 + w_2 + w_0 \leq 0$$

No w is a solution!



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

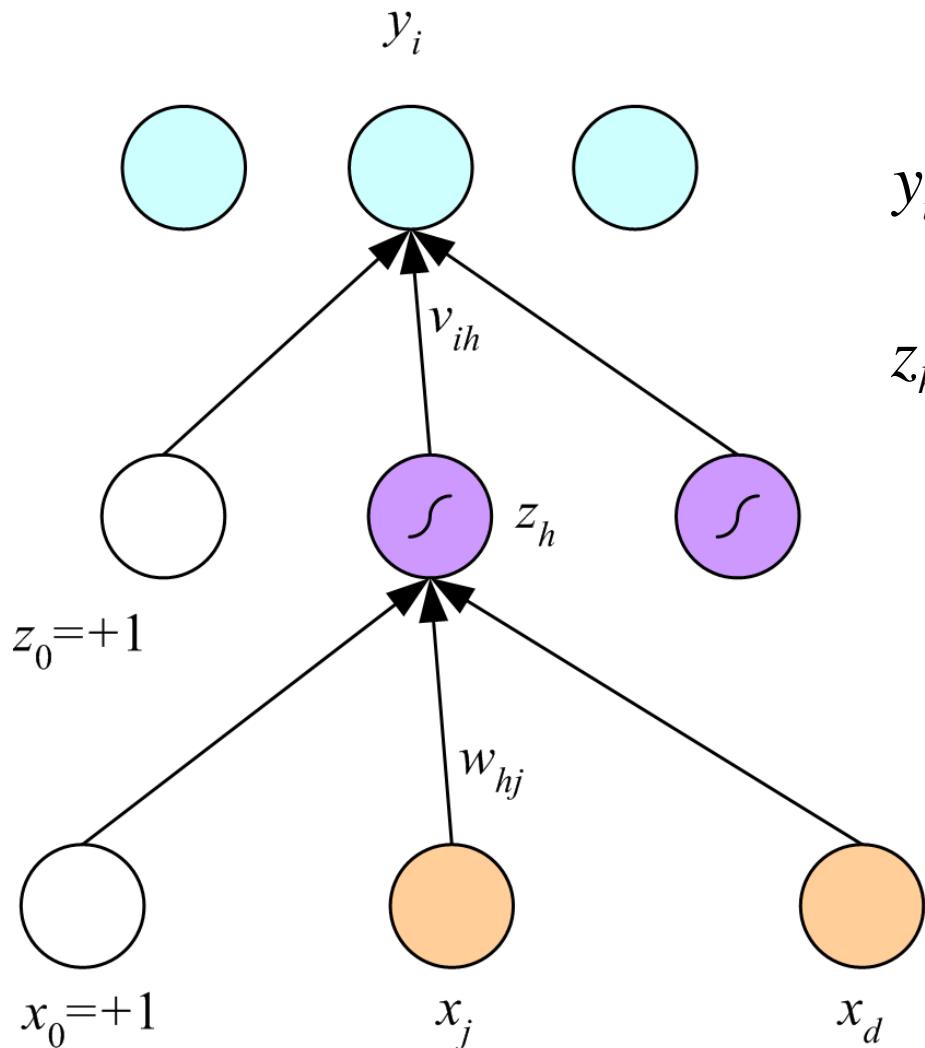
Multilayer Perceptrons



$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$
$$z_h = \text{sigmoid} (\mathbf{w}_h^T \mathbf{x})$$
$$= \frac{1}{1 + \exp \left[- \left(\sum_{j=1}^d w_{hj} x_j + w_{h0} \right) \right]}$$

(Rumelhart et al., 1986)

Backpropagation



$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$
$$z_h = \text{sigmoid} \left(\mathbf{w}_h^T \mathbf{x} \right)$$

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

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

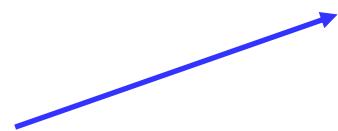
Regression

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

Forward

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

\mathbf{x}



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

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

Backward

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

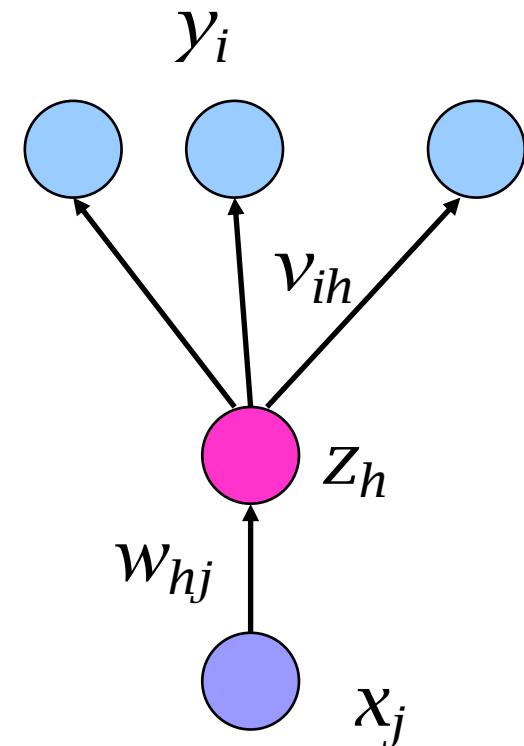
Regression with Multiple Outputs

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

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

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

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



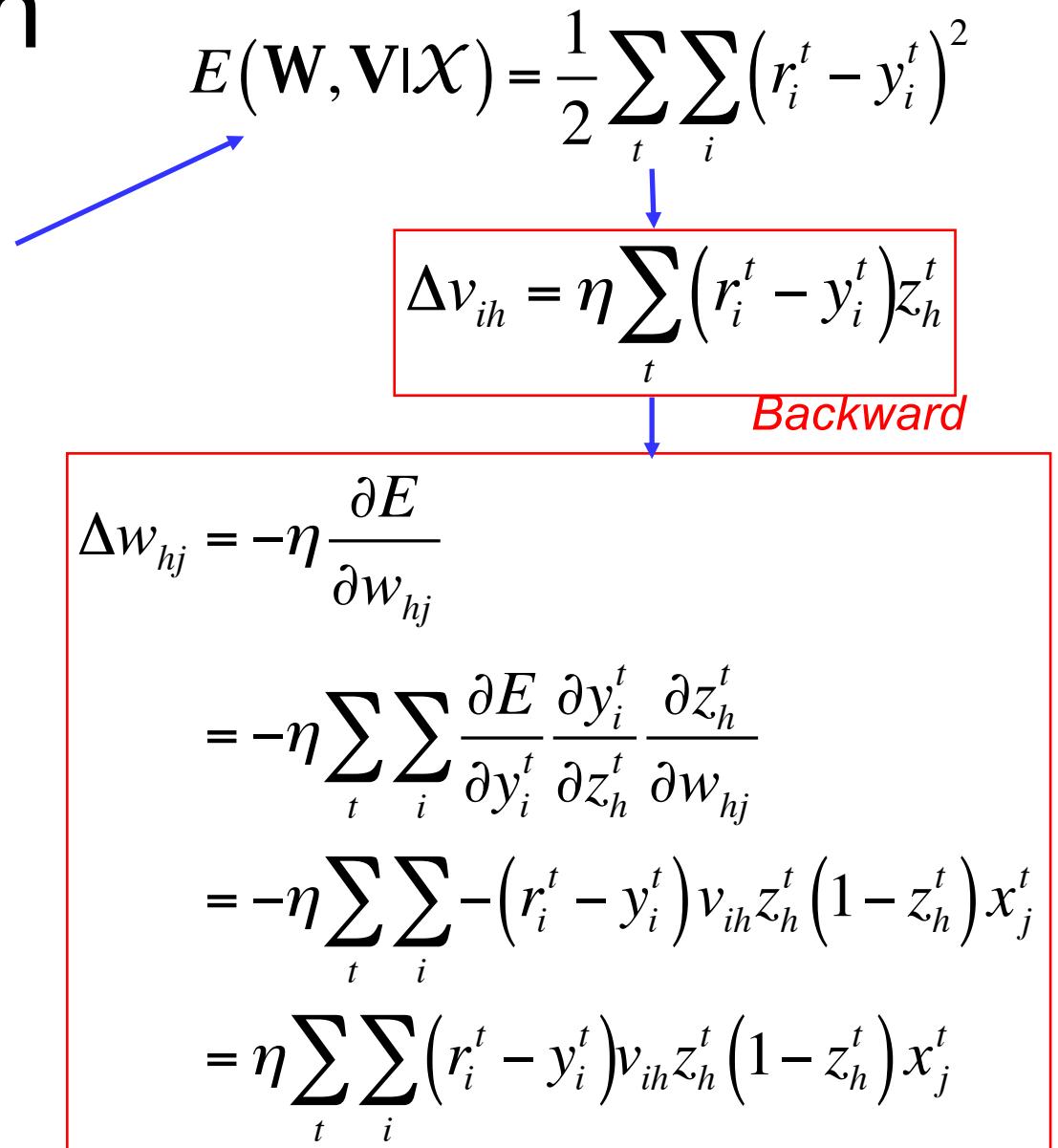
Multi-regression

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

Forward

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

\mathbf{x}



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 $h = 1, \dots, H$

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

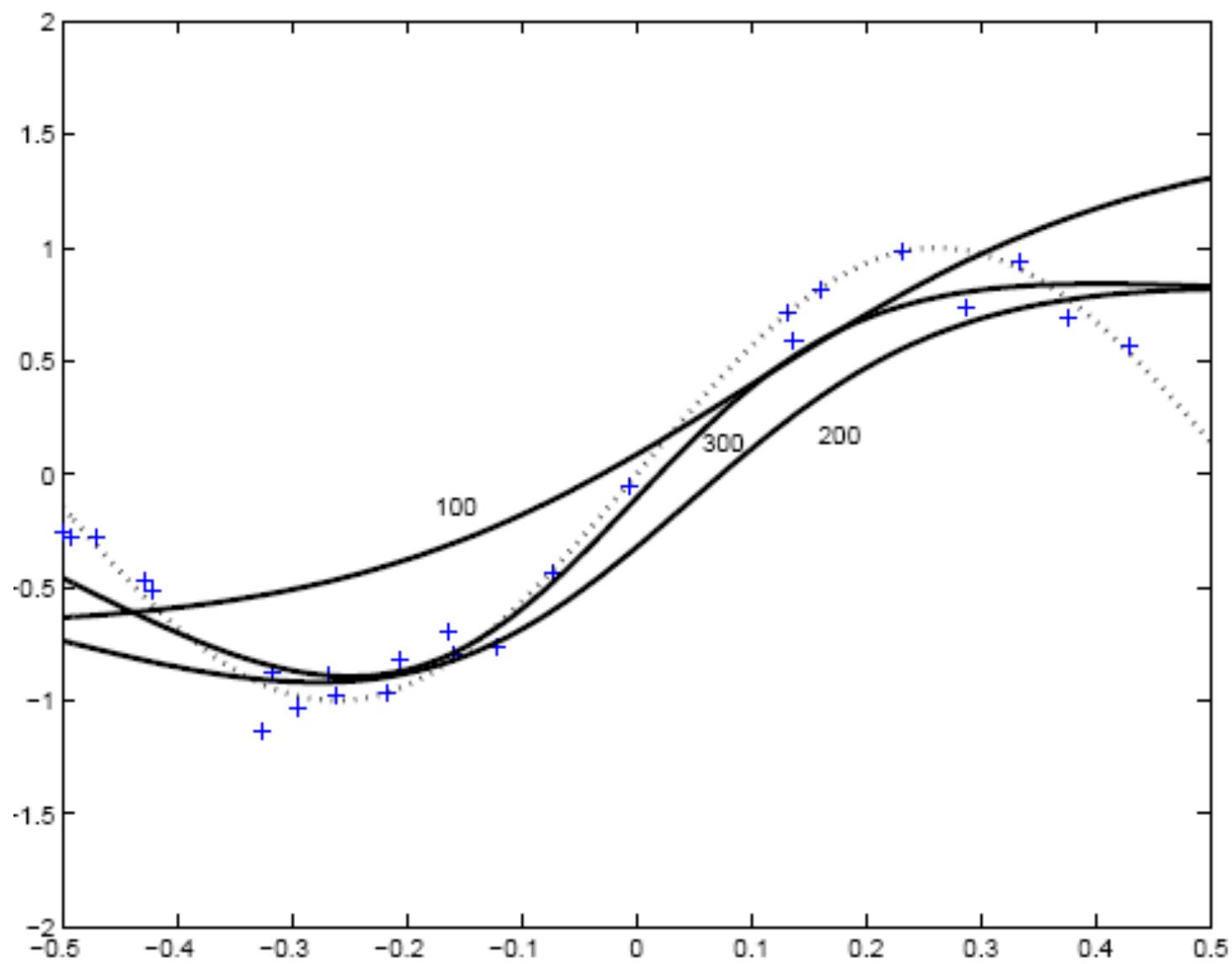
For $i = 1, \dots, K$

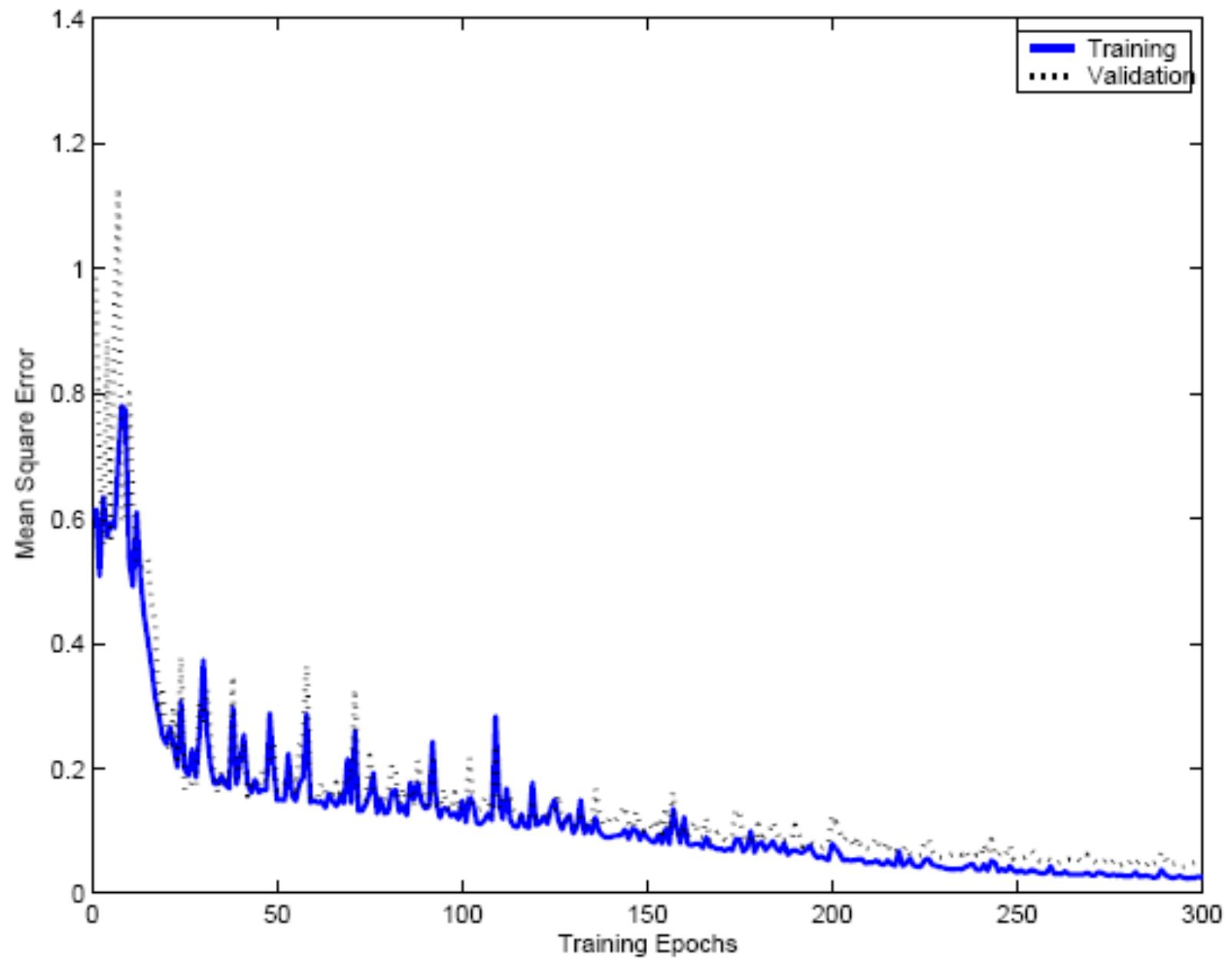
$$\mathbf{v}_i \leftarrow \mathbf{v}_i + \Delta \mathbf{v}_i$$

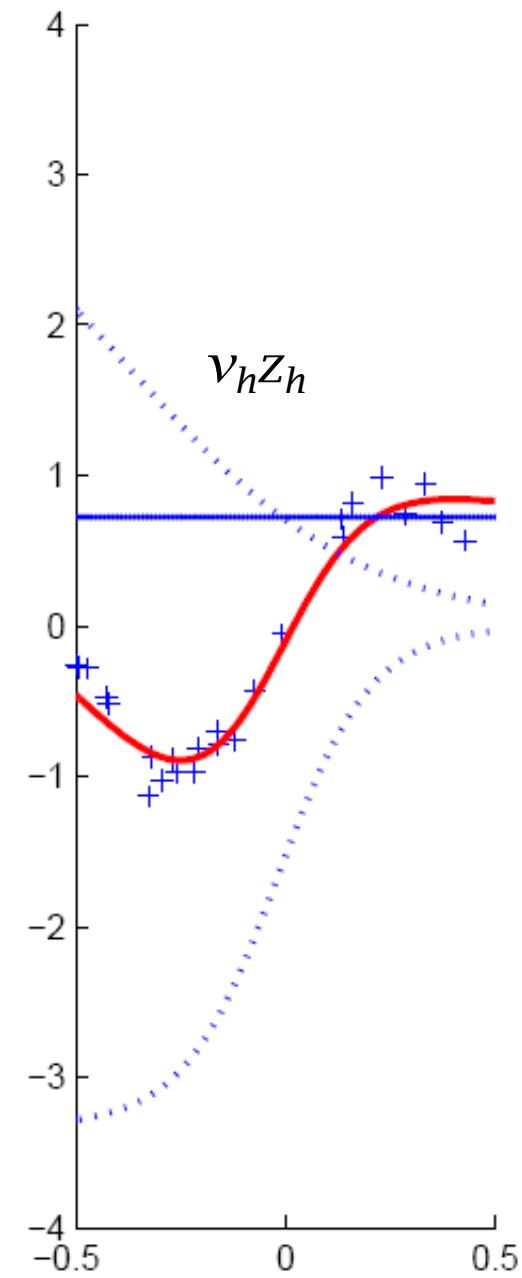
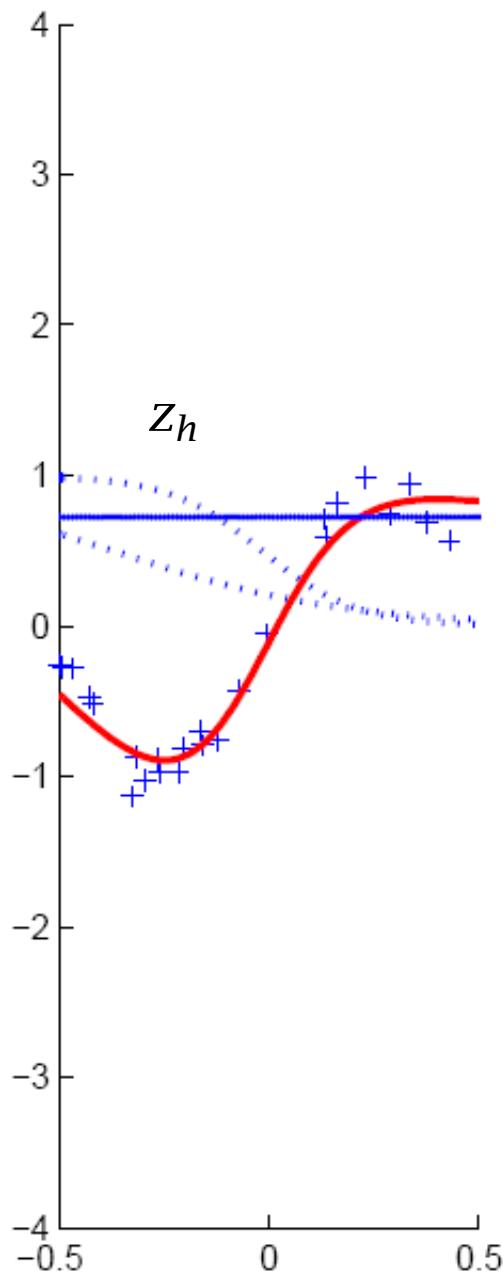
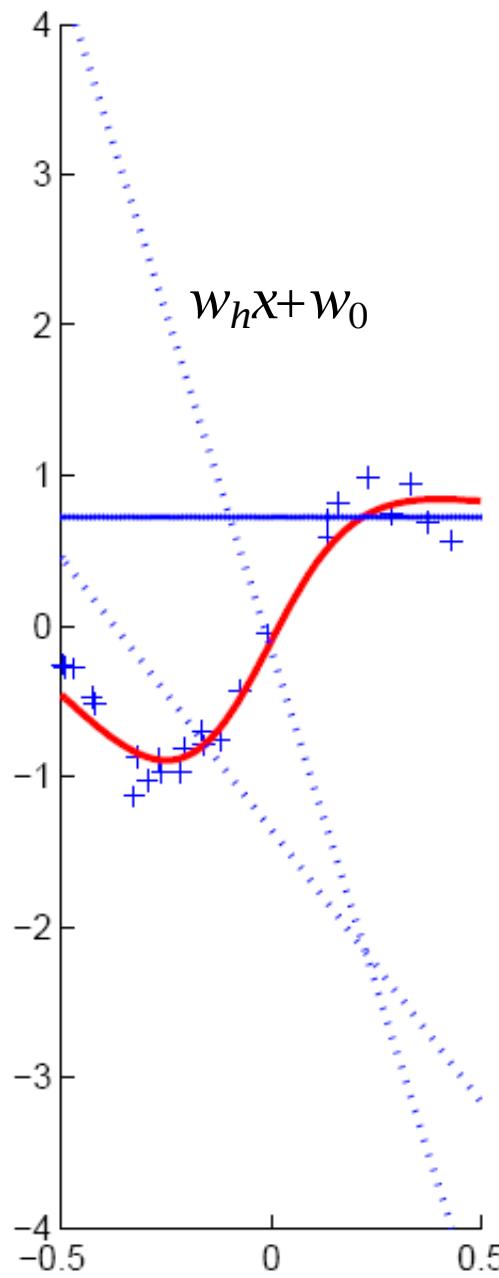
For $h = 1, \dots, H$

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

Until convergence

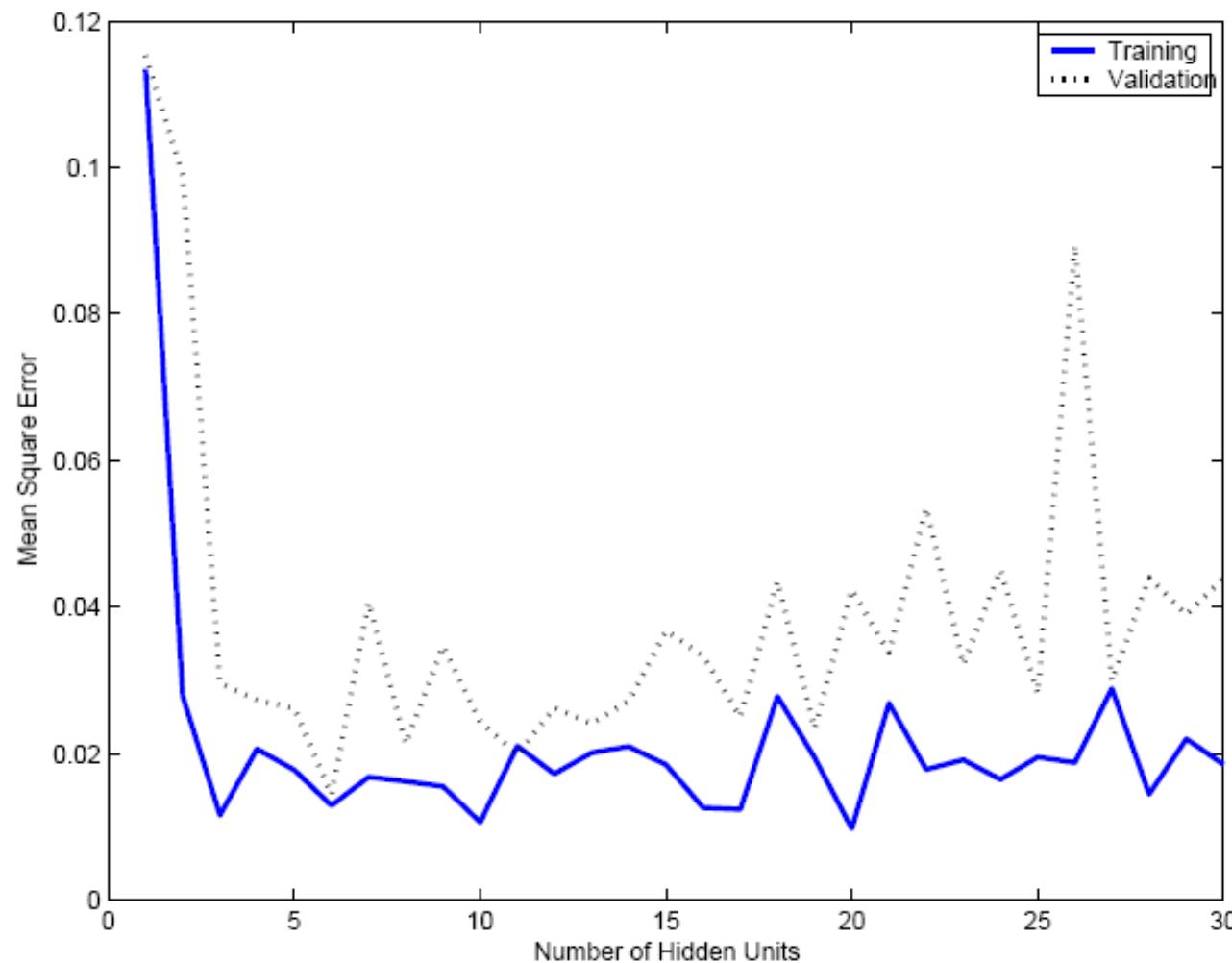


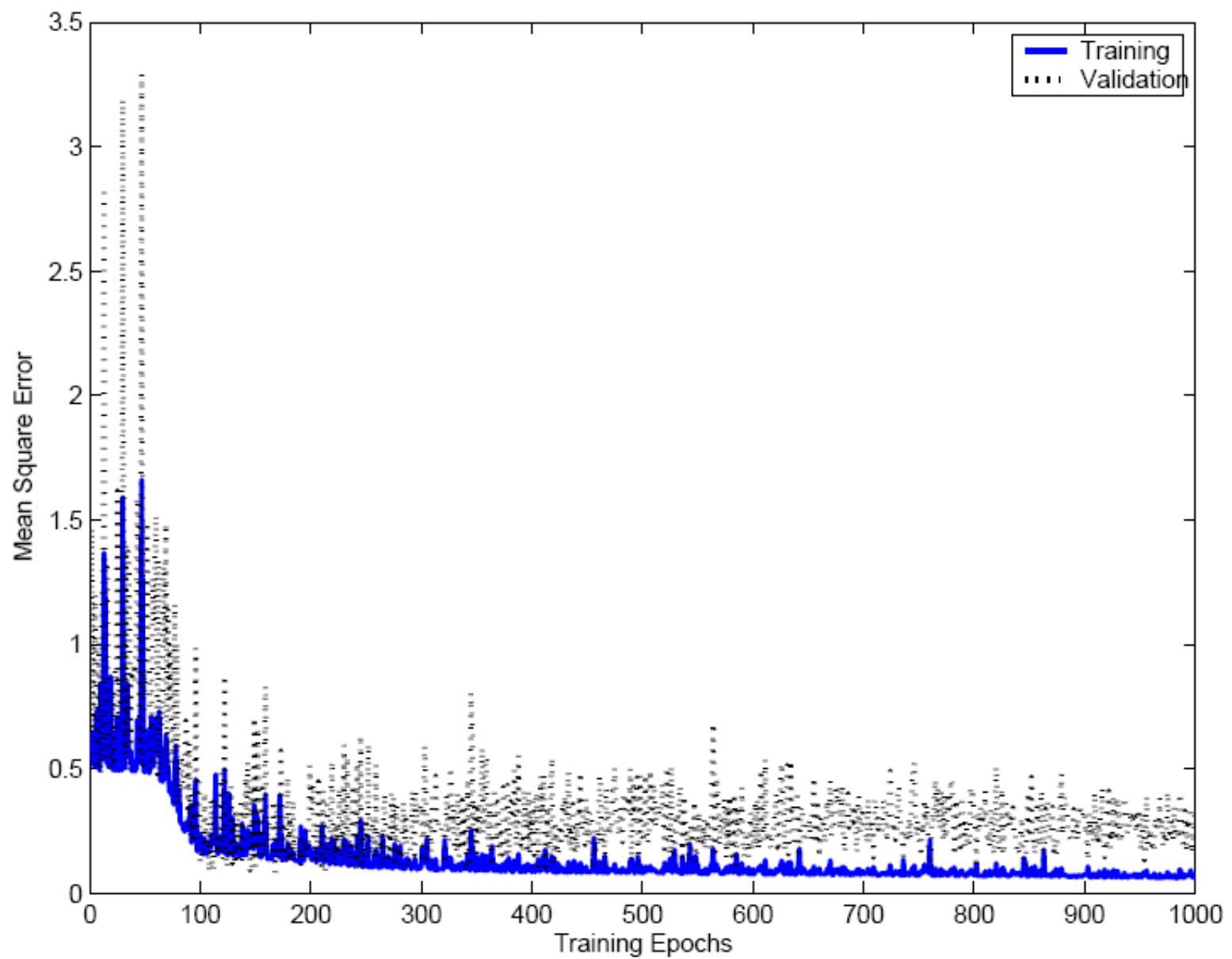




Overfitting/Overtraining

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





Two-Class Classification

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

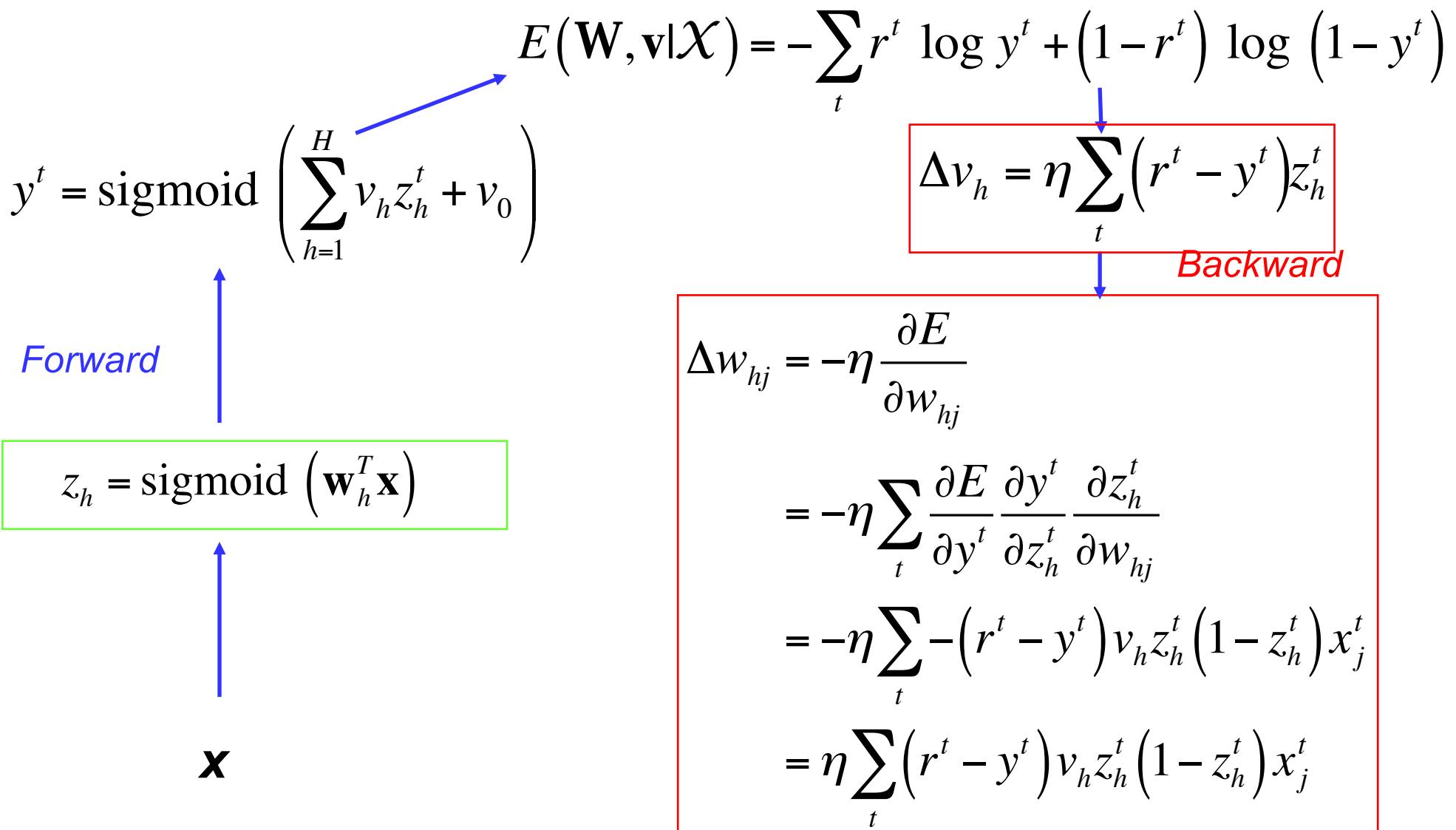
$$y^t = \text{sigmoid} \left(\sum_{h=1}^H v_h z_h^t + v_0 \right)$$

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

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

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

Classification



$K > 2$ Classes

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

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

$$\Delta v_{ih} = \eta \sum_t (r_i^t - y_i^t) z_h^t \quad (\text{See softmax logistic regression})$$

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

Multi-classification

$$o_i^t = \sum_{h=1}^H v_{ih} z_h^t + v_{i0} \quad y_i^t = \frac{\exp o_i^t}{\sum_k \exp o_k^t}$$

Forward

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

\mathbf{x}

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

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

Backward

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

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

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

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

Multiple Hidden Layers

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

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

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

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

Improving Convergence

■ Momentum

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

■ Adaptive learning rate

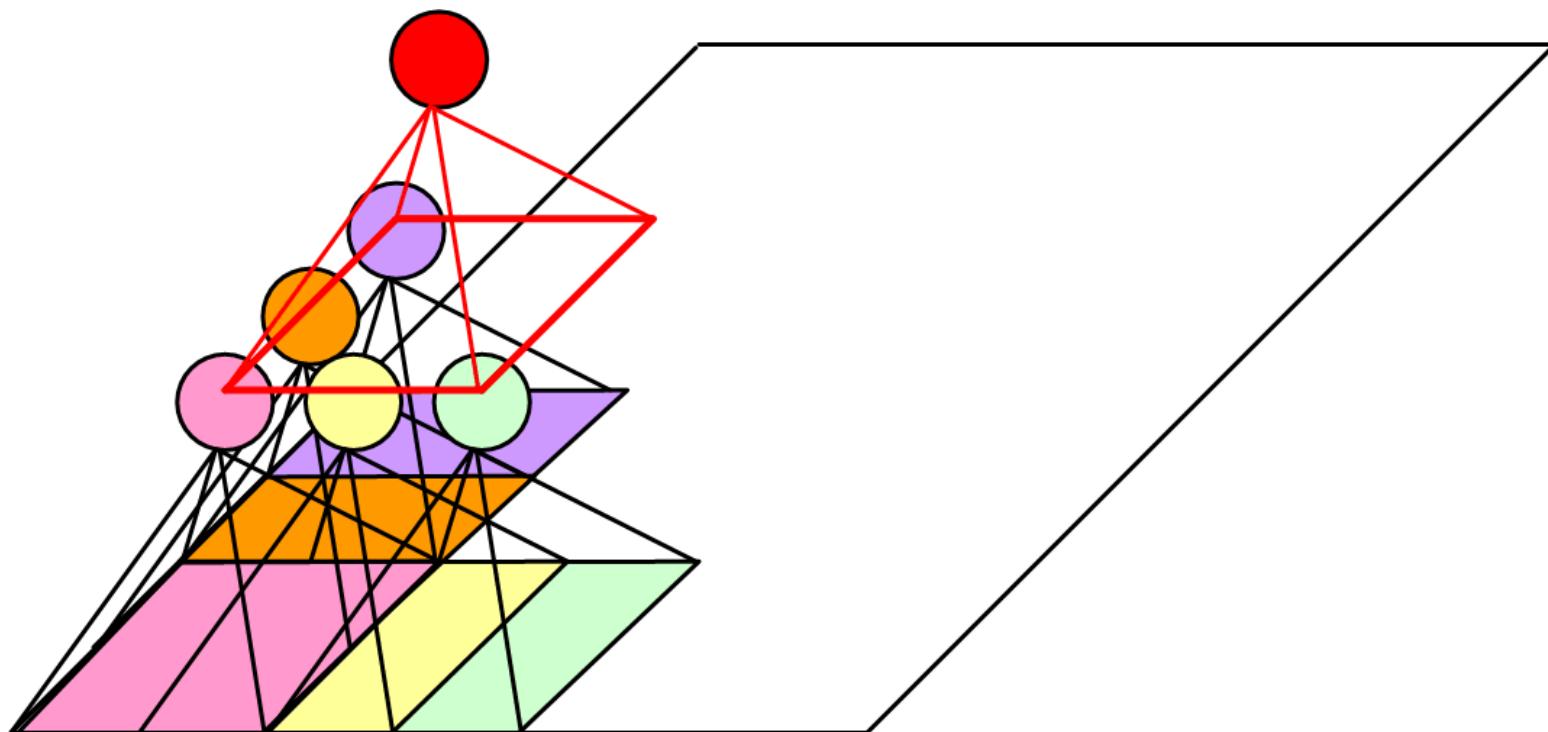
$$\Delta\eta = \begin{cases} +a & \text{if } E^{t+\tau} < E^t \\ -b\eta & \text{otherwise} \end{cases}$$

Go faster!

Oscillating!
Slow down!

Convolutional Neural Network

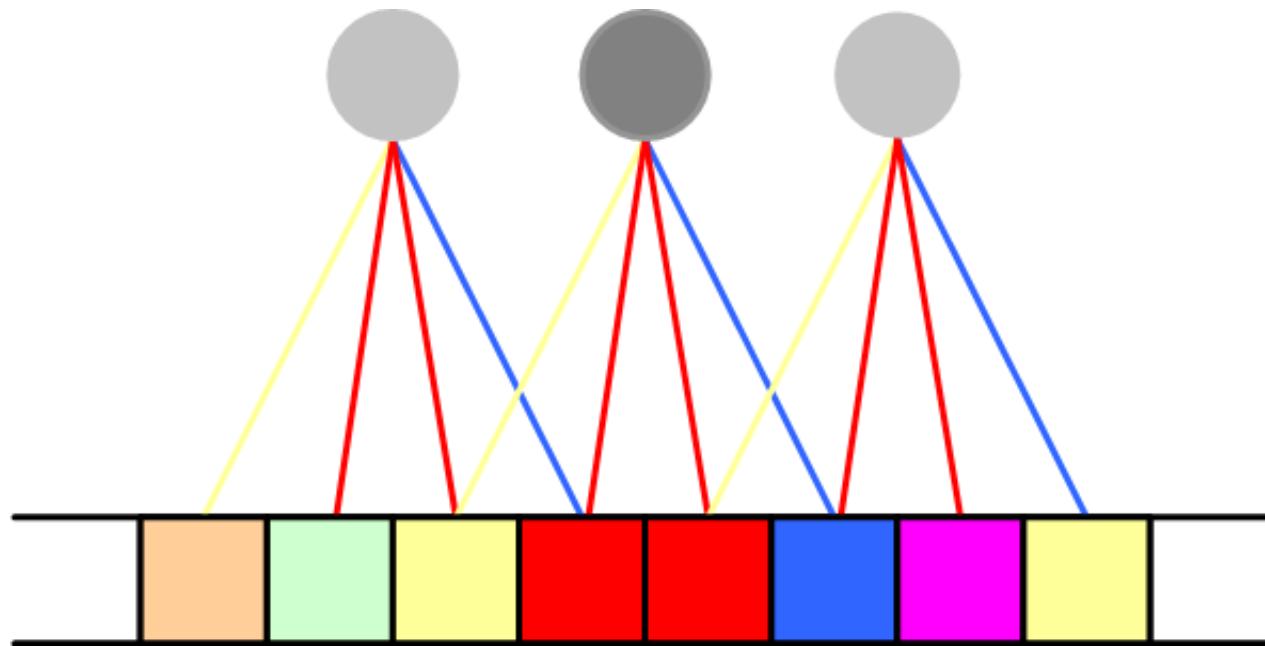
- Hidden units are connected to subsets of inputs



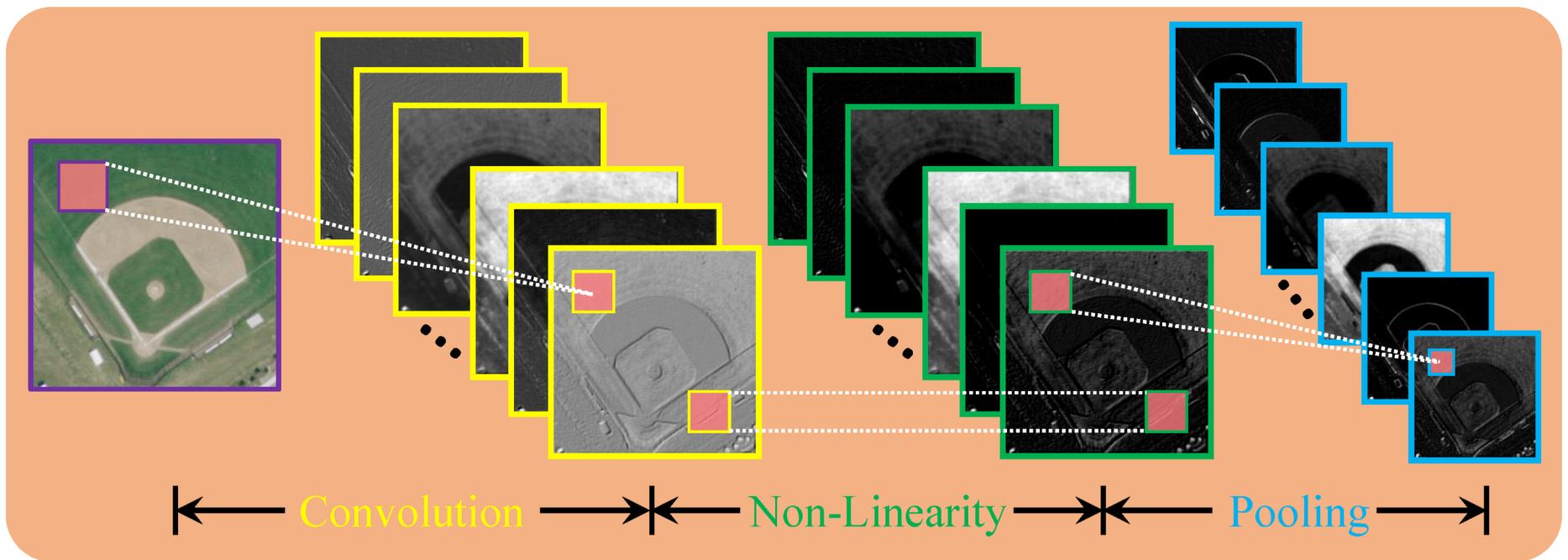
(Le Cun et al, 1989)

Weight Sharing

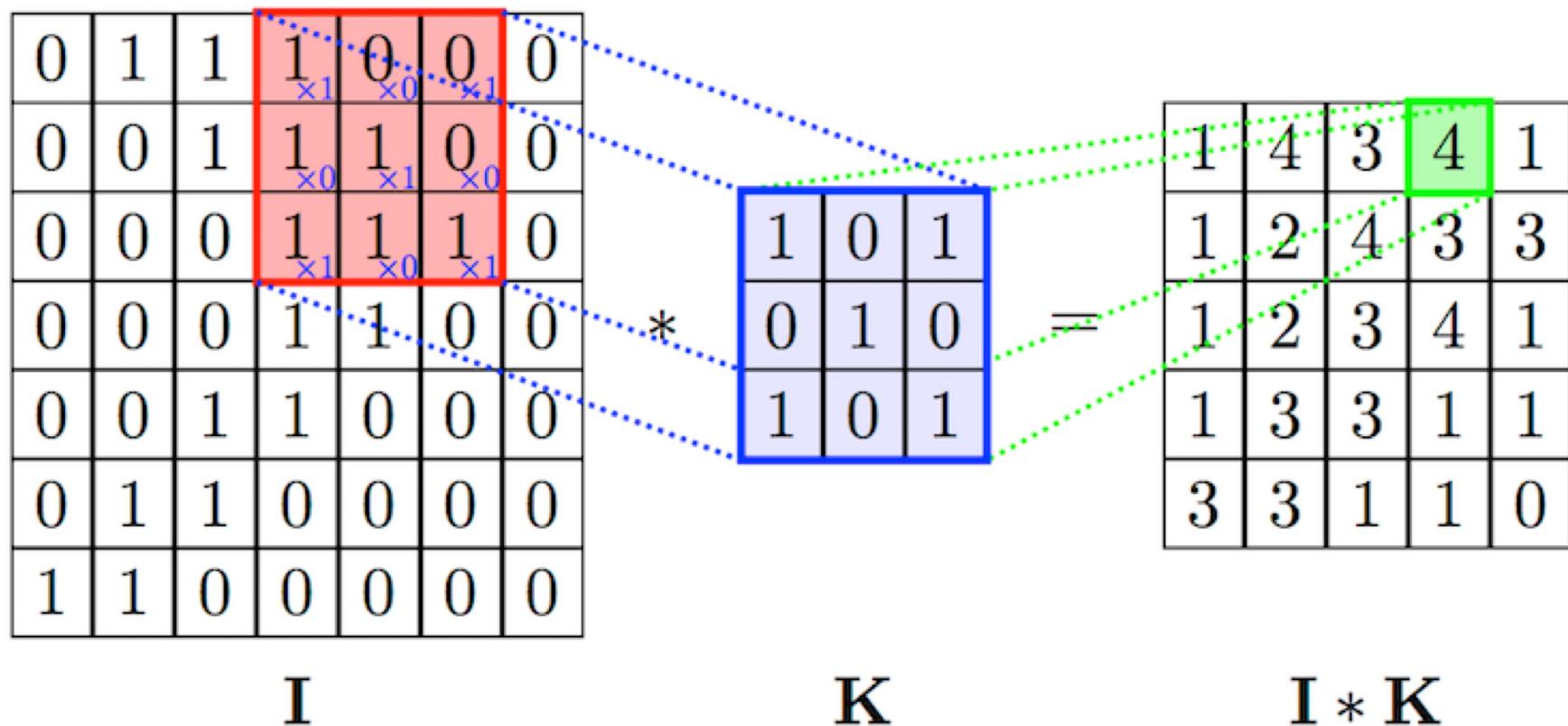
- Some of the weights are shared among several units.



Convolution and Pooling



Convolution



More on Filters

1/9	1/9	1/9
1/9	1/9	1/9
1/9	1/9	1/9

0	-1	0
-1	5	-1
0	-1	0



Blurry

Sharpen

<https://www.saama.com/blog/different-kinds-convolutional-filters/>

More on Filters

-1	-2	-1
0	0	0
1	2	1

Horizontal

-1	0	1
-2	0	2
-1	0	1

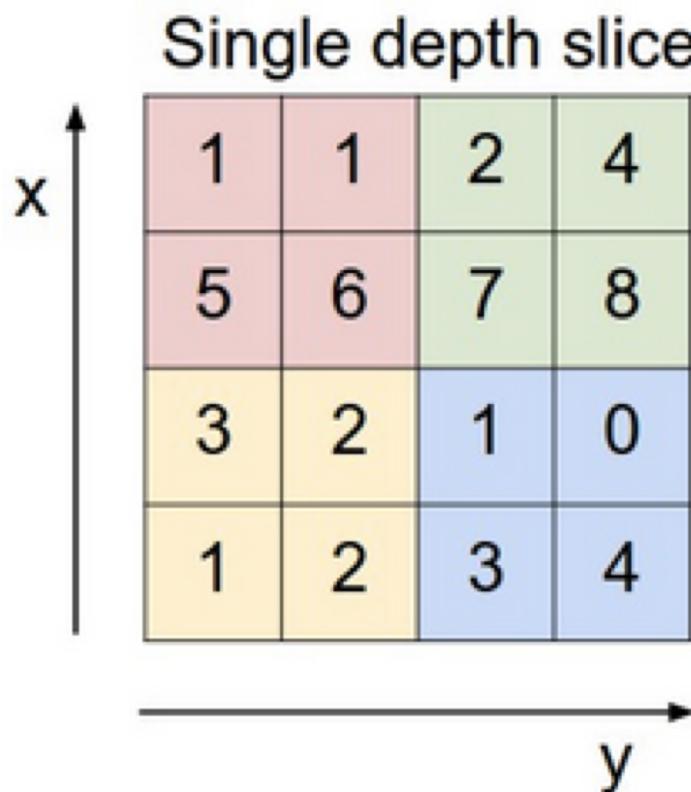
Vertical



Edge detection

<https://www.saama.com/blog/different-kinds-convolutional-filters/>

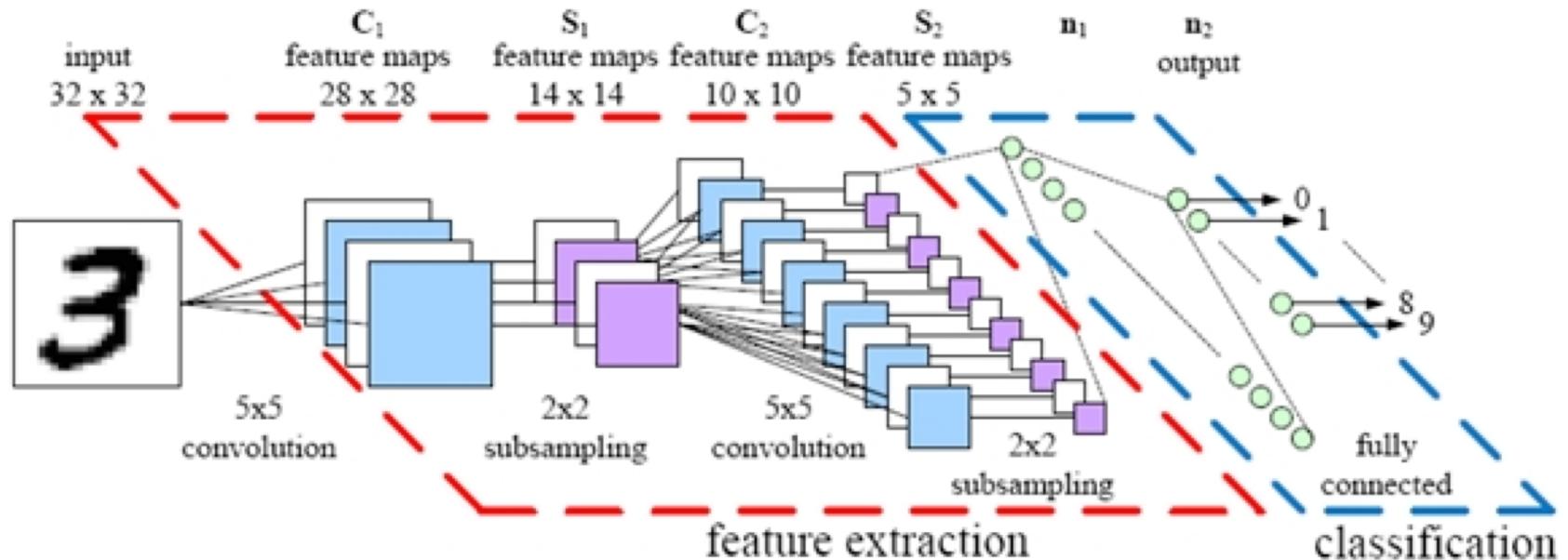
Pooling (subsampling)



max pool with 2x2 filters
and stride 2

6	8
3	4

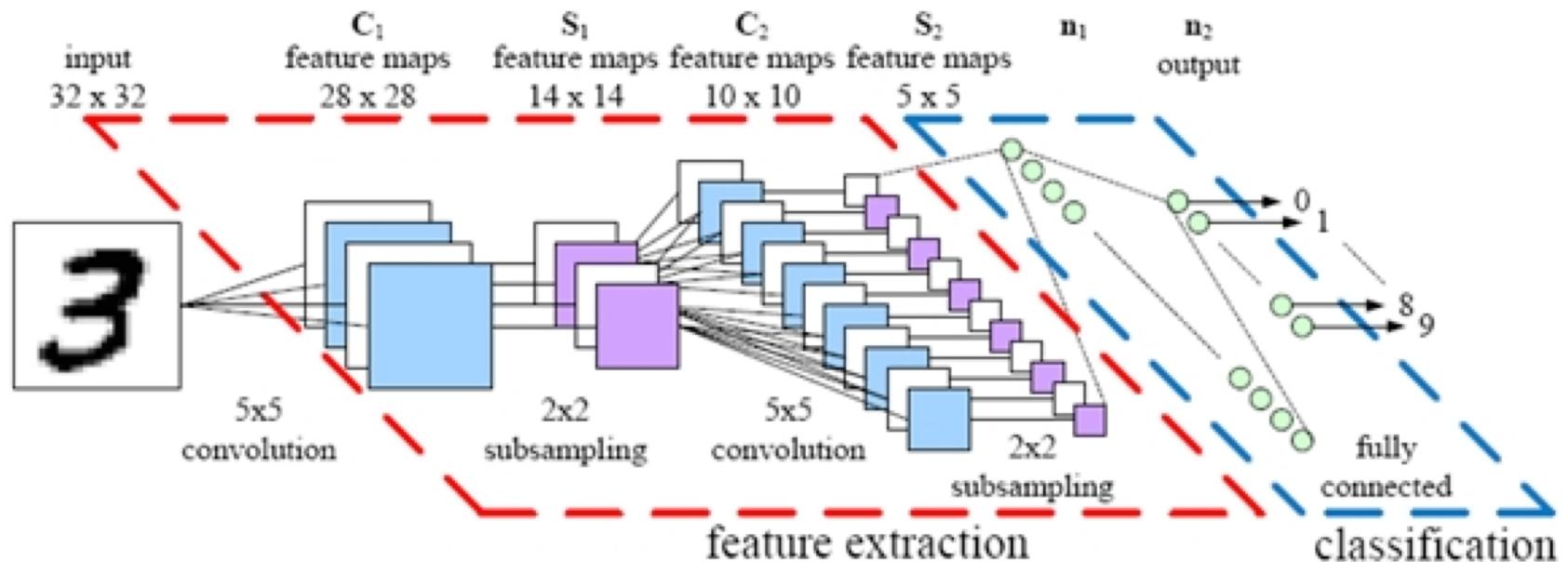
CNN Example



■ Convolutional layer

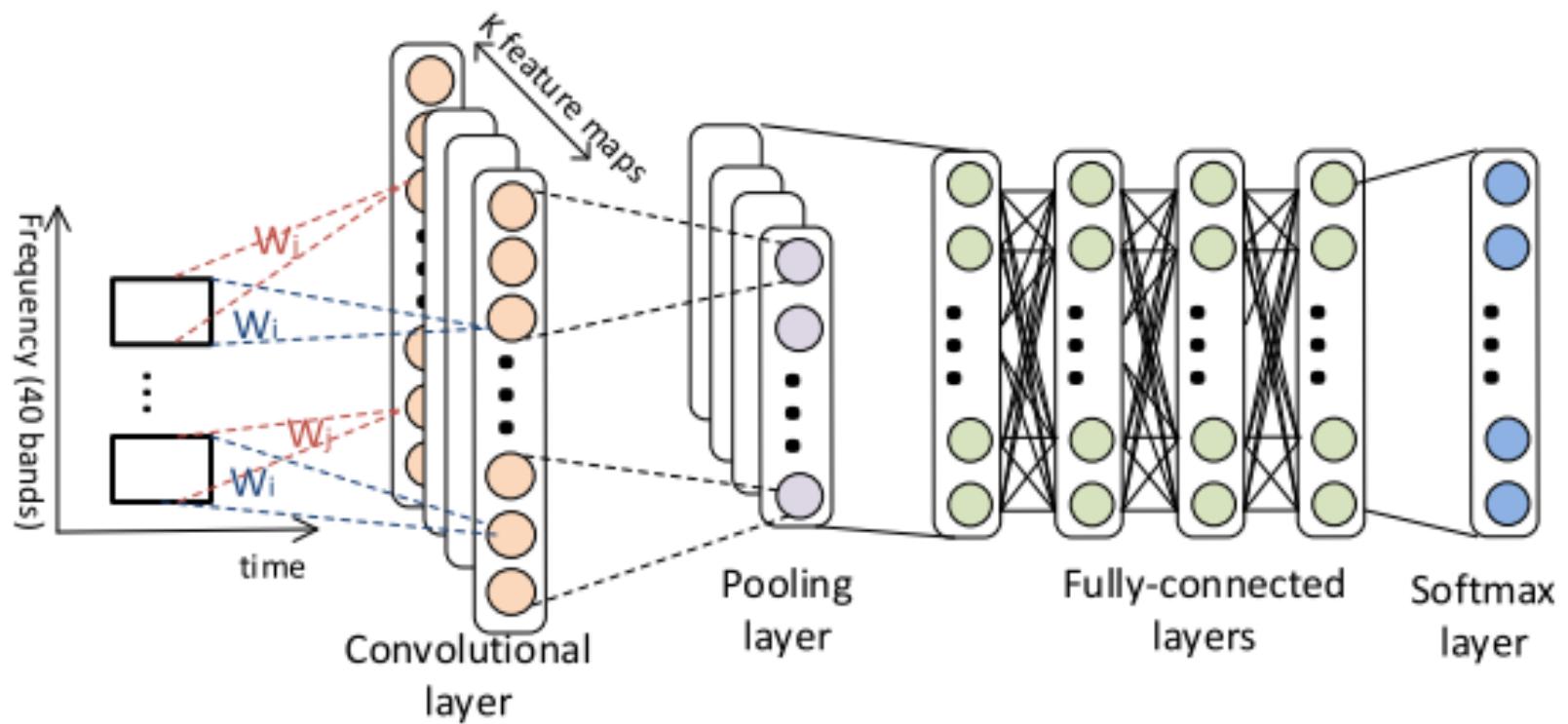
- Each unit takes subregion as input.
- Shared weights across all the units in a feature map.
- There are multiple feature maps in the convolutional layer.

CNN Example



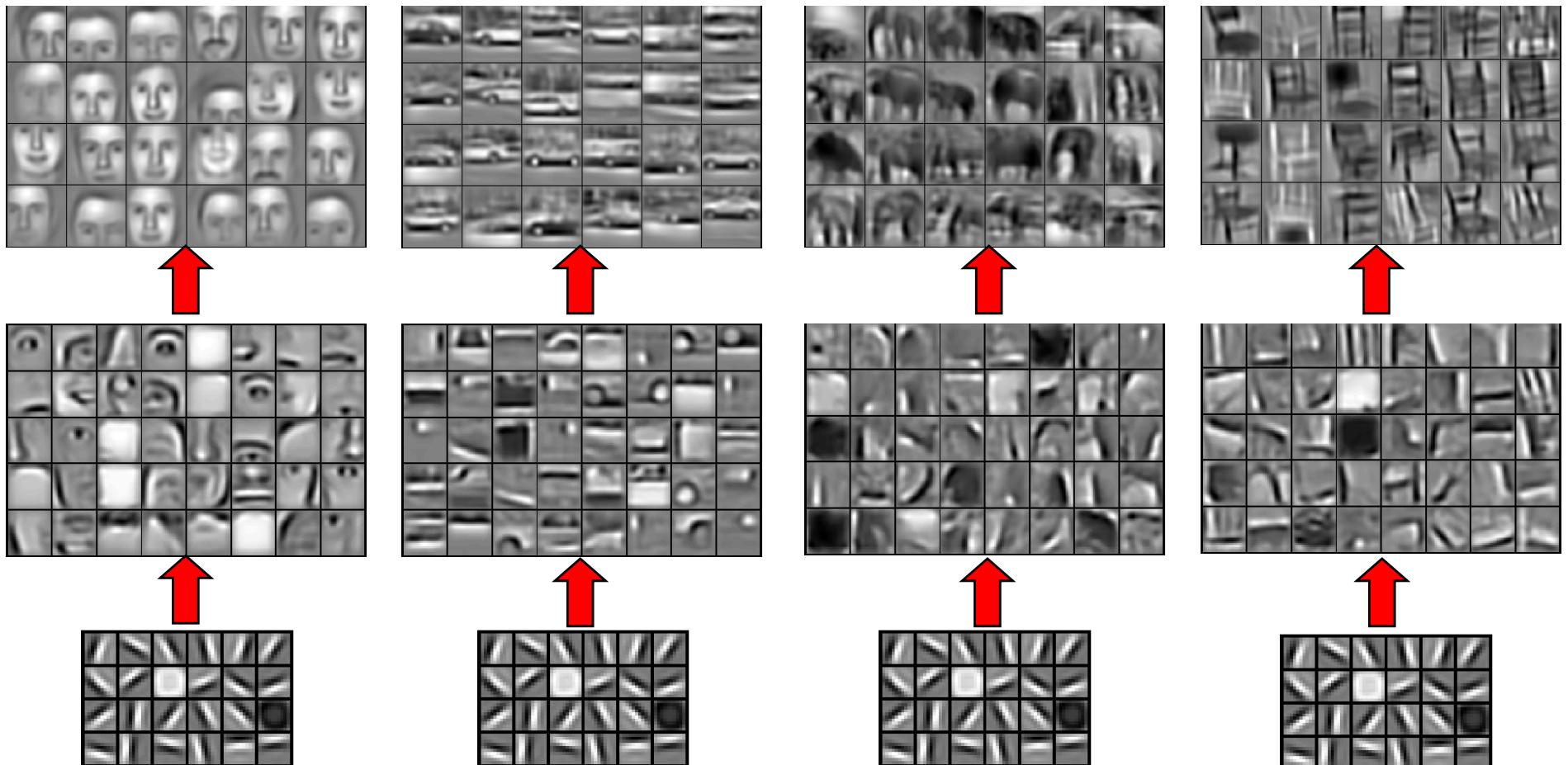
- Nonoverlapping subsampling layer.
 - Units in one plane connect to 2x2 region (receptive field) in a feature map.
 - Taking the mean or max of the inputs in the region.
- Final layers are fully connected NN.

CNN Example



AN ANALYSIS OF CONVOLUTIONAL NEURAL NETWORKS FOR SPEECH RECOGNITION,
Huang et al.

Learning of Object Parts



Tuning the Network Size

■ Destructive

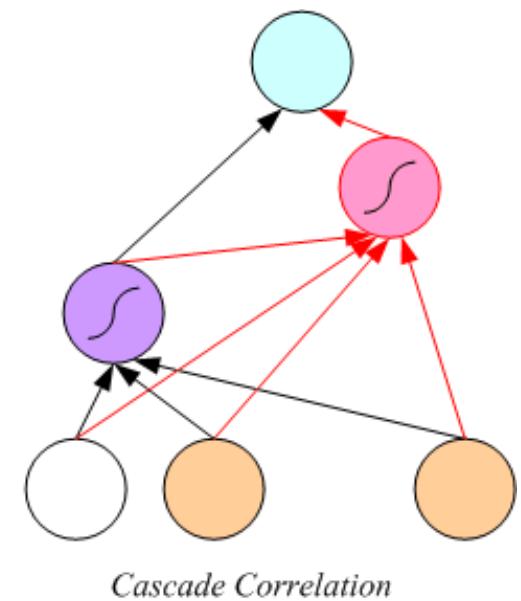
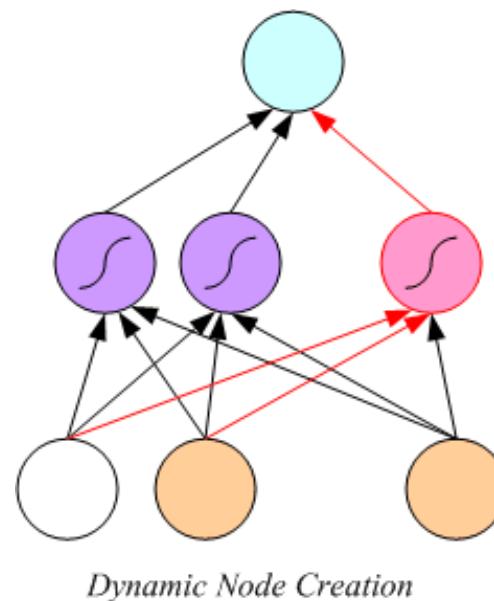
- Weight decay:
penalizing non-zero
parameters
- The same as adding
additional constraints

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

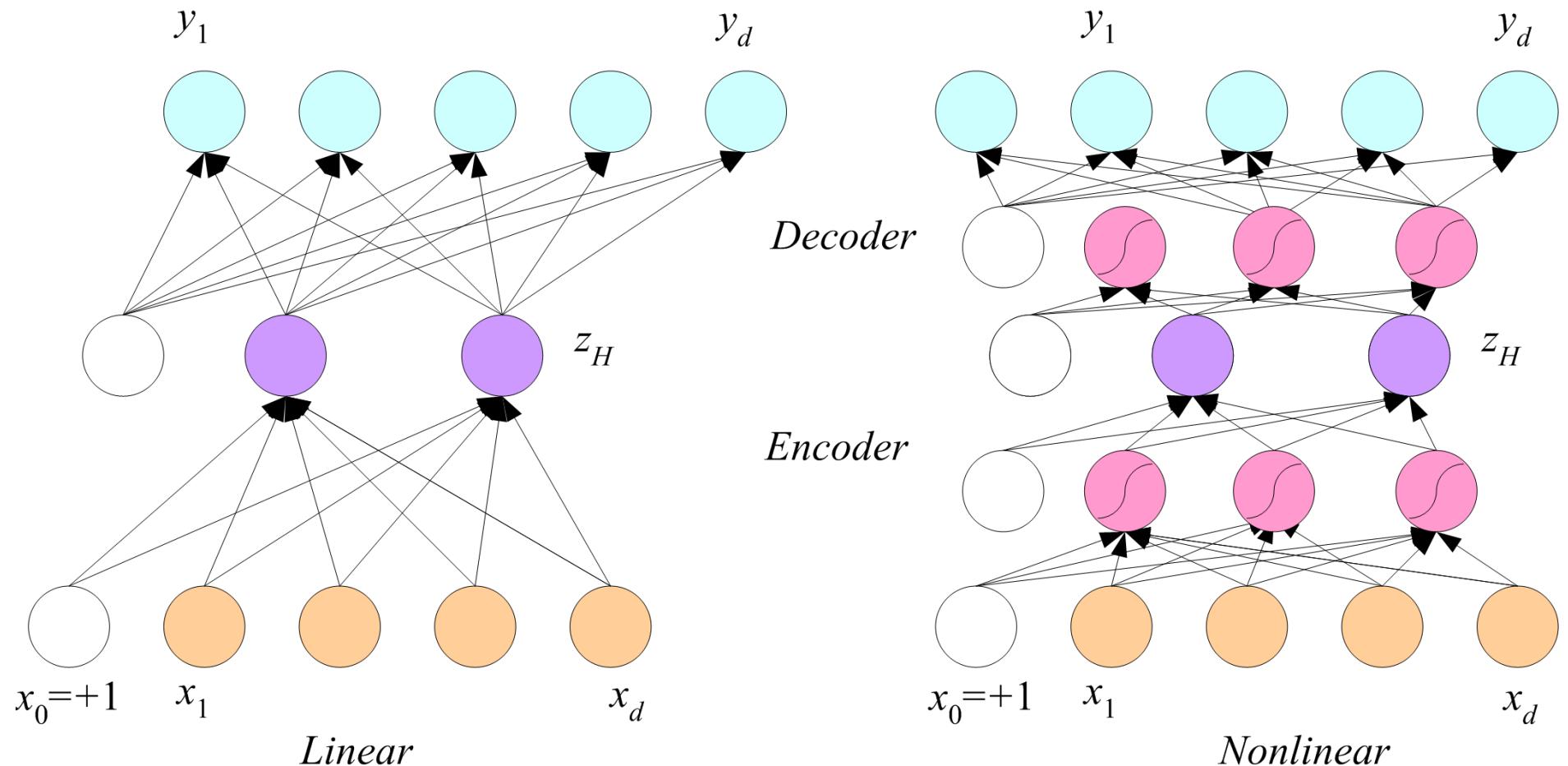
$$E' = E + \frac{\lambda}{2} \sum_i w_i^2$$

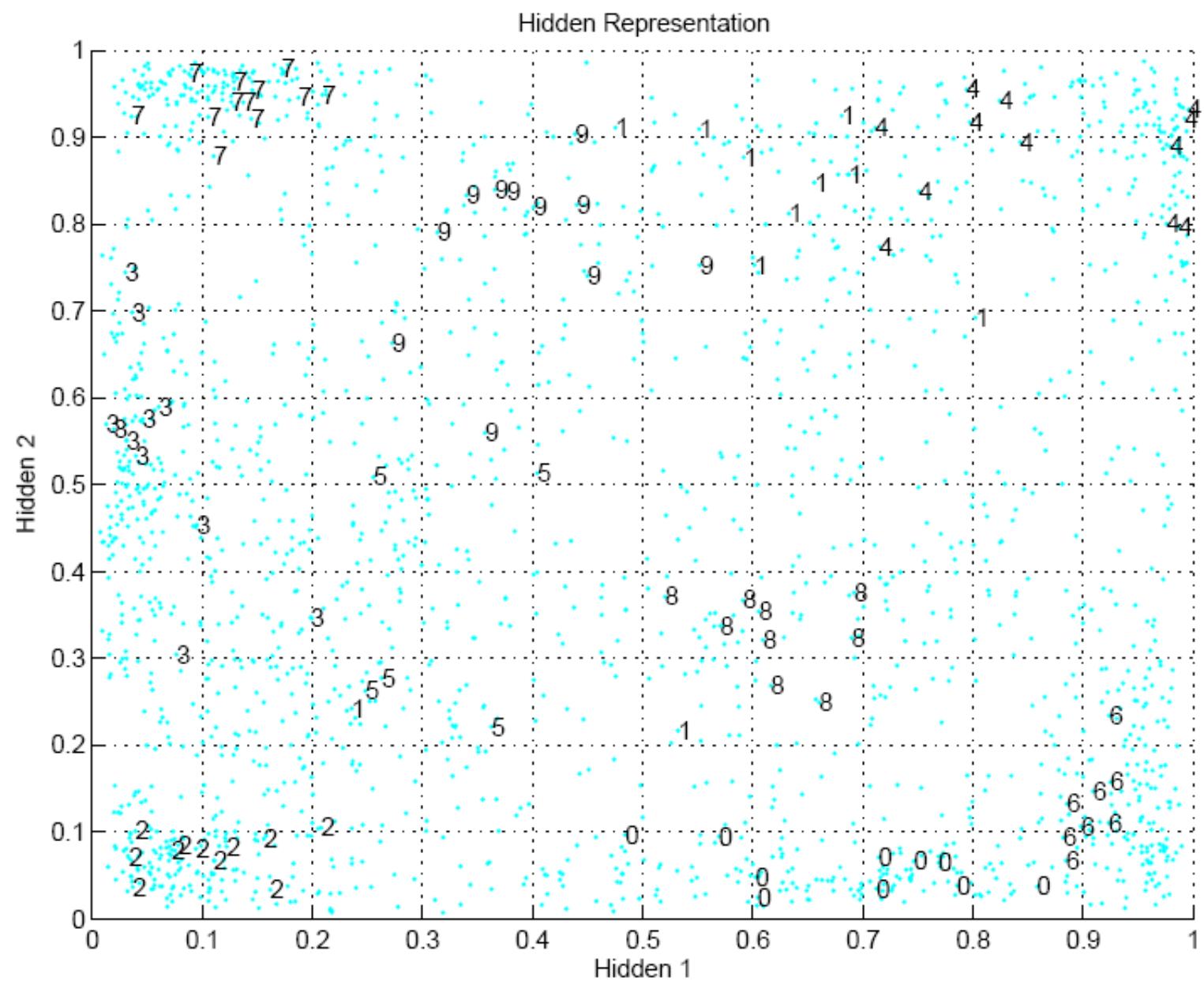
■ Constructive

- Growing networks: if
training error is high, add
more hidden units.



Autoencoder: Dimension Reduction

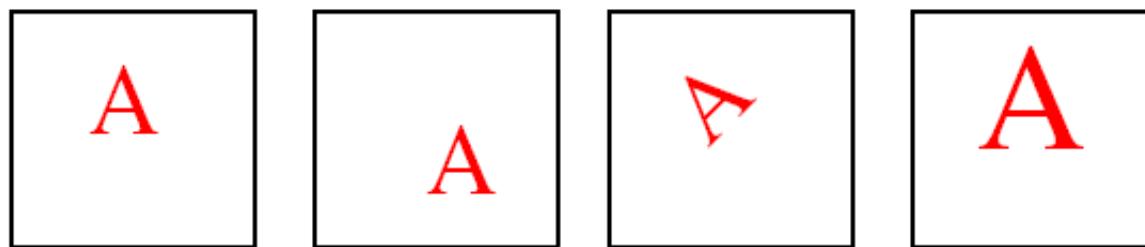




Hints

(Abu-Mostafa, 1995)

- Invariance to translation, rotation, size



- Virtual examples
- Augmented error: $E' = E + \lambda_h E_h$
If \mathbf{x}' and \mathbf{x} are the “same”: $E_h = [g(\mathbf{x}|\theta) - g(\mathbf{x}'|\theta)]^2$