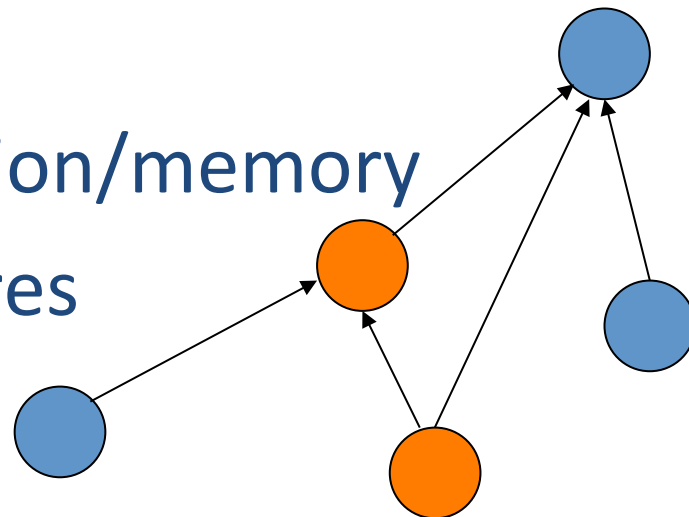


- ▶ Basic Concepts
 - ▶ supervised, unsupervised and semi-supervised learning
 - ▶ density estimation, classification, regression
 - ▶ error measures
 - ▶ optimization
 - ▶ linear and nonlinear models
- ▶ Perceptrons and Multilayer Perceptrons
 - ▶ perceptron
 - ▶ multilayer perceptron
 - ▶ training
 - ▶ overfitting and its prevention
 - ▶ learning from data and hints
 - ▶ neural networks for time series
- ▶ Deep Neural Networks
 - ▶ shallow networks' problems
 - ▶ optimization
 - ▶ regularization
 - ▶ convolution
 - ▶ hardware

Neural Networks

- 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

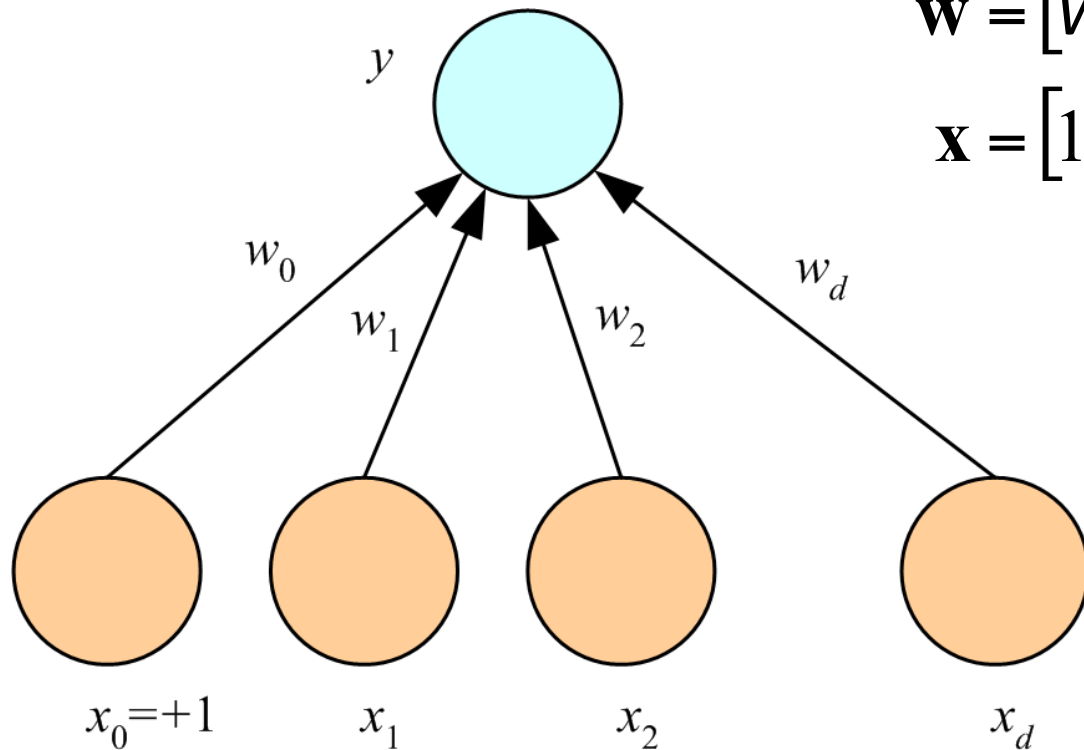


Perceptron

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

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

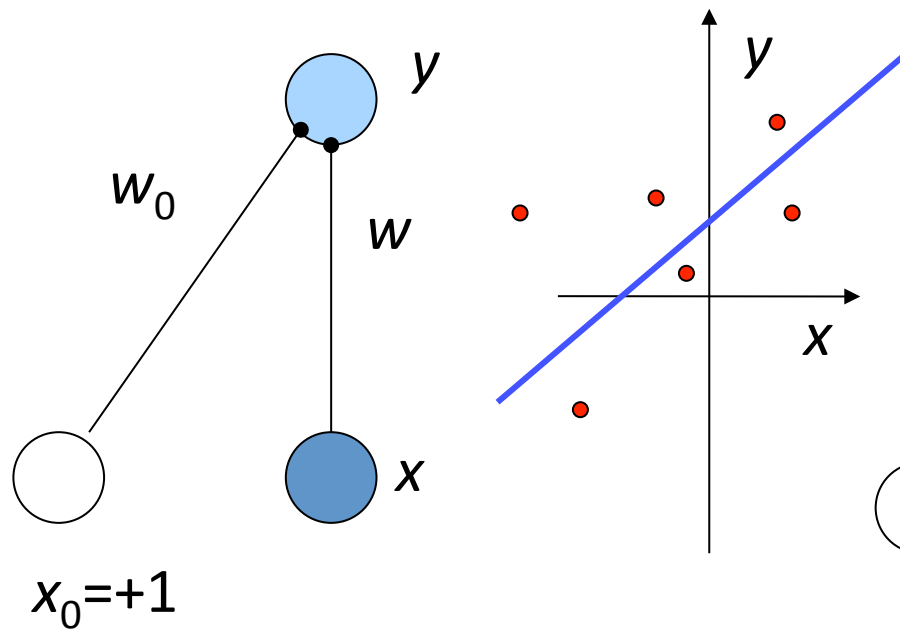
$$\mathbf{x} = [1, x_1, \dots, x_d]^T$$



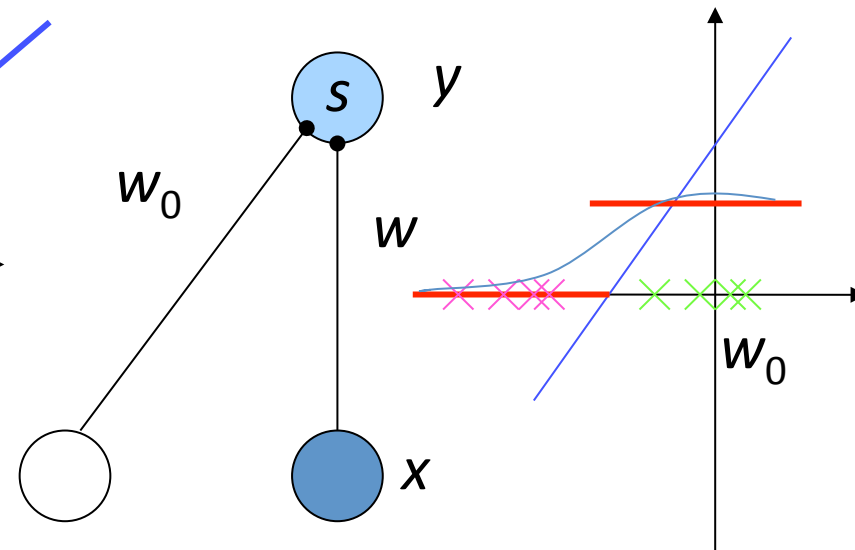
(Rosenblatt, 1962)

What a Perceptron Does

- Regression: $y=wx+w_0$



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



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

Training a Perceptron

- Regression:

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

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

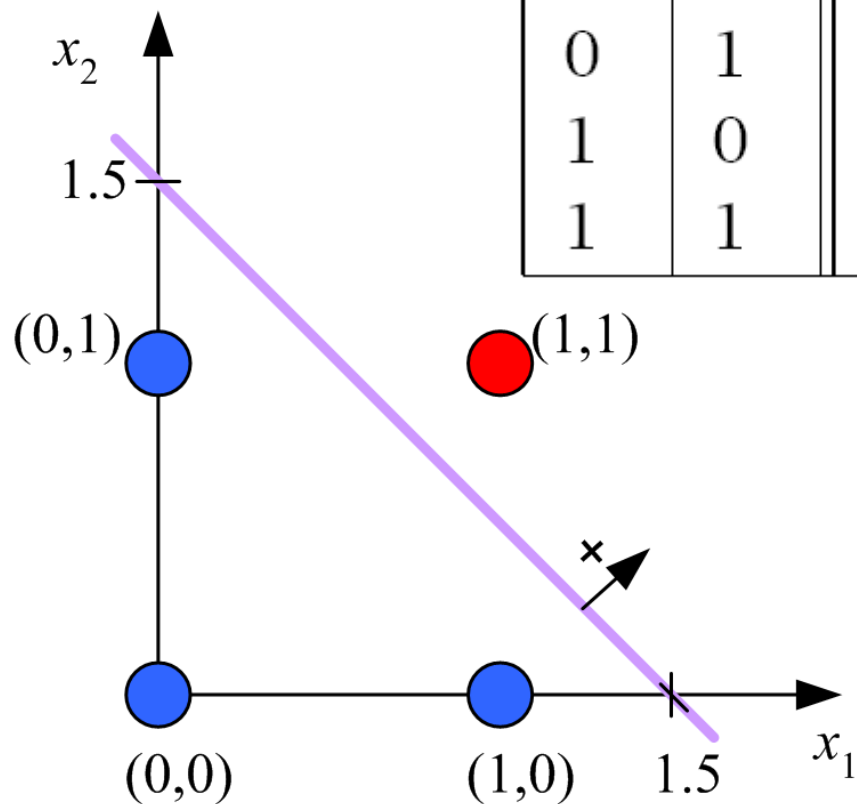
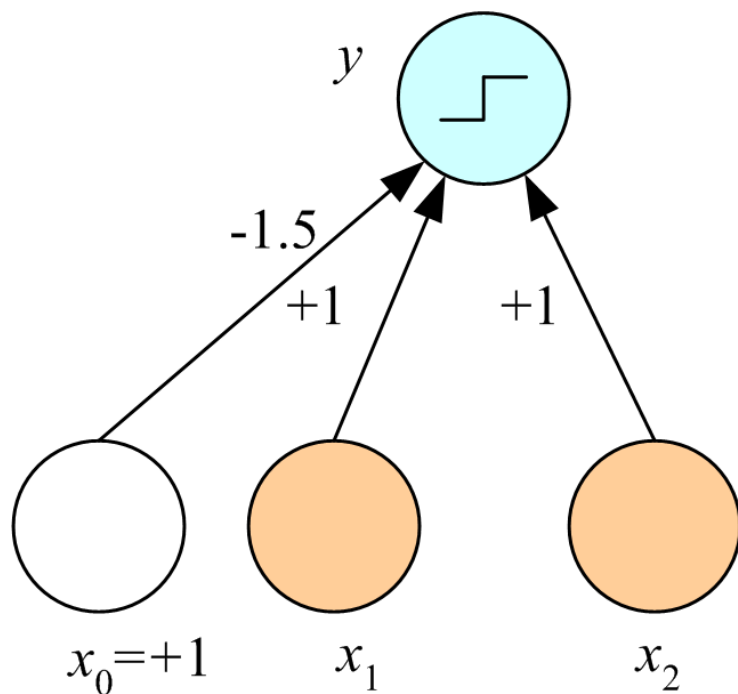
- Classification:

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

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

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

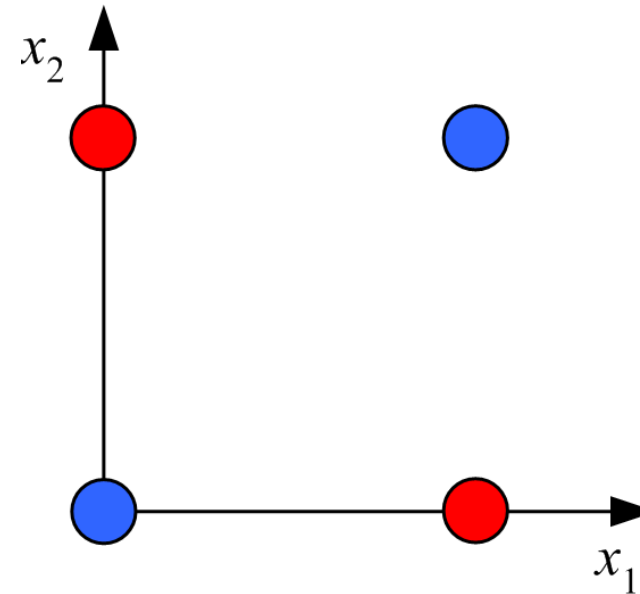
Learning Boolean AND



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:

$$w_0 \leq 0$$

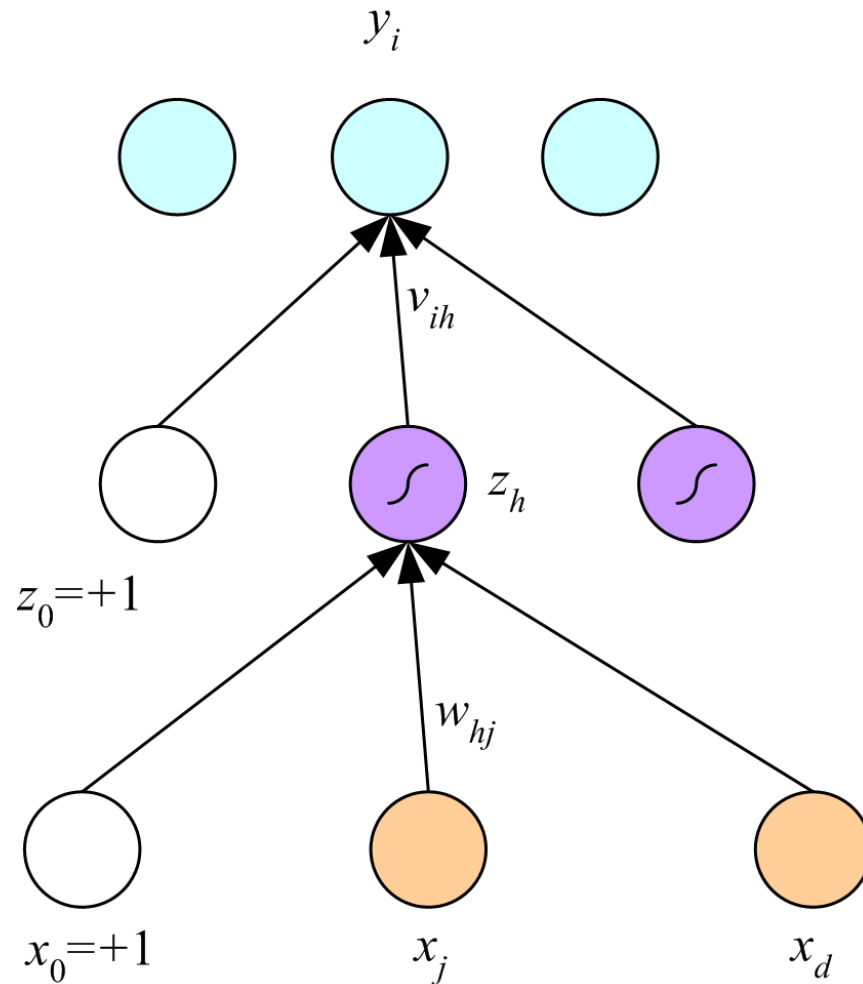
$$w_2 + w_0 > 0$$

$$w_1 + w_0 > 0$$

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

(Minsky and Papert, 1969)

Multilayer Perceptrons

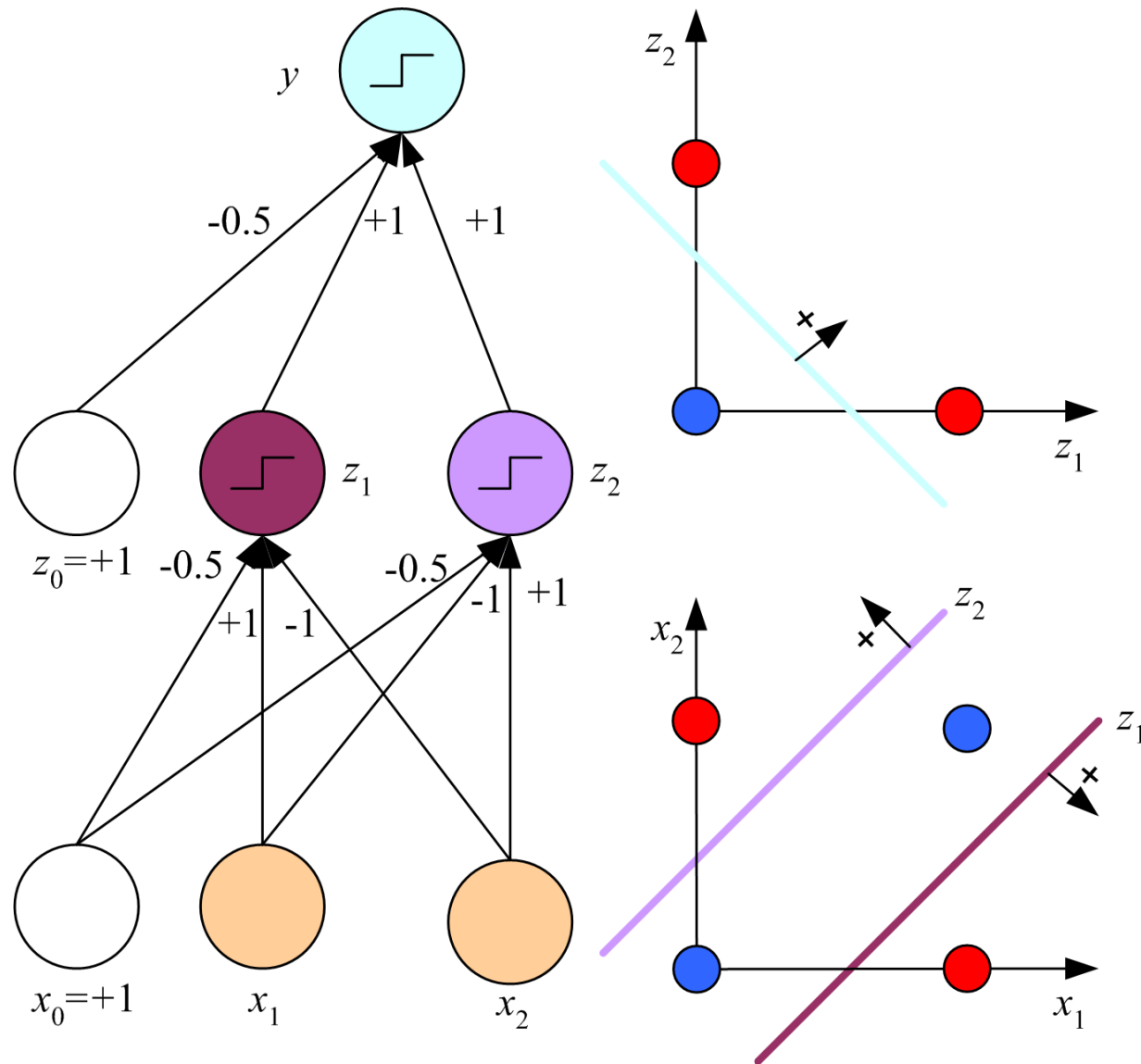


$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

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

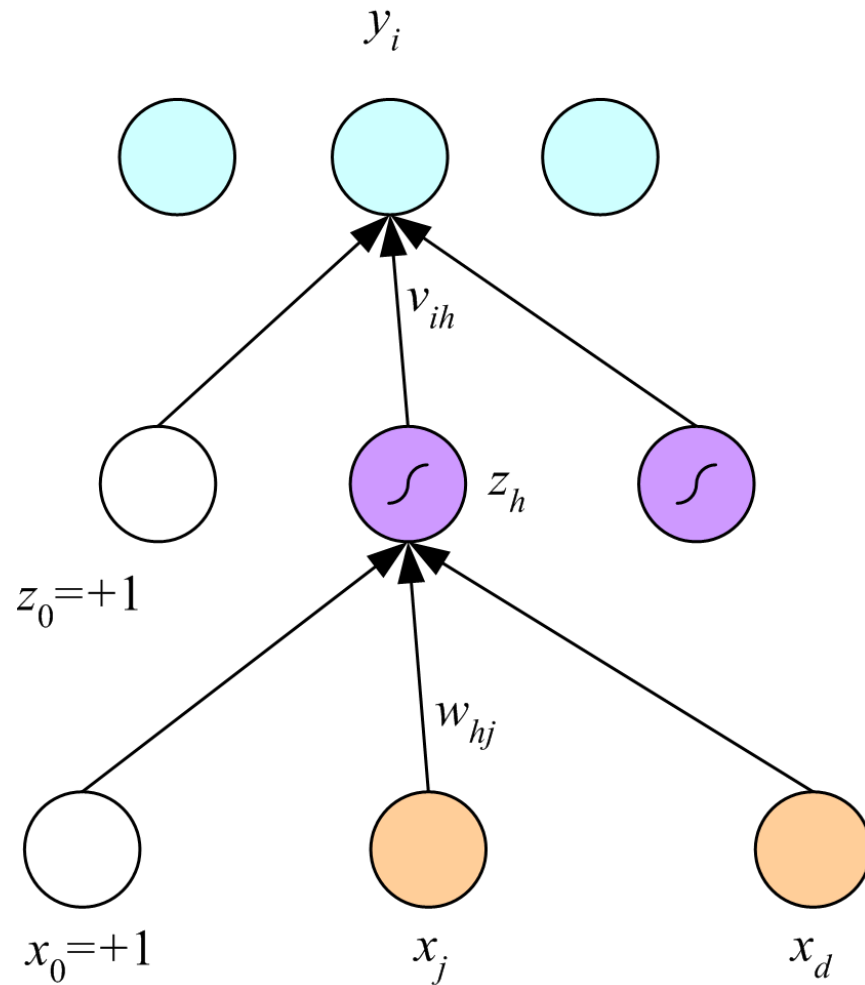
$$= \frac{1}{1 + \exp \left[- \left(\sum_{j=1}^d w_{hj} x_j + w_{h0} \right) \right]}$$

(Rumelhart et al., 1986)



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

Backpropagation



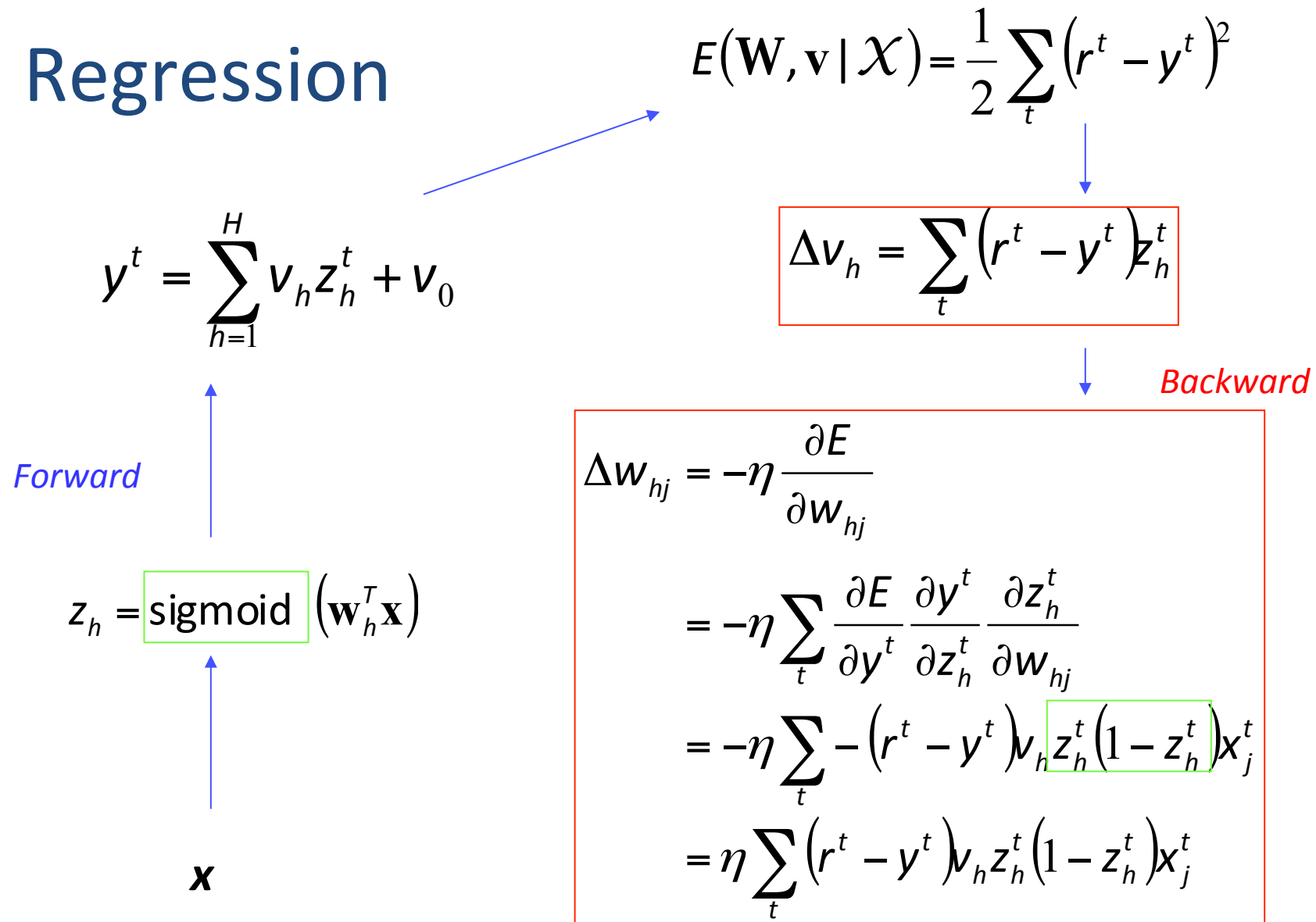
$$y_i = \mathbf{v}_i^T \mathbf{z} = \sum_{h=1}^H v_{ih} z_h + v_{i0}$$

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

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

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

Regression



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 $h = 1, \dots, H$

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

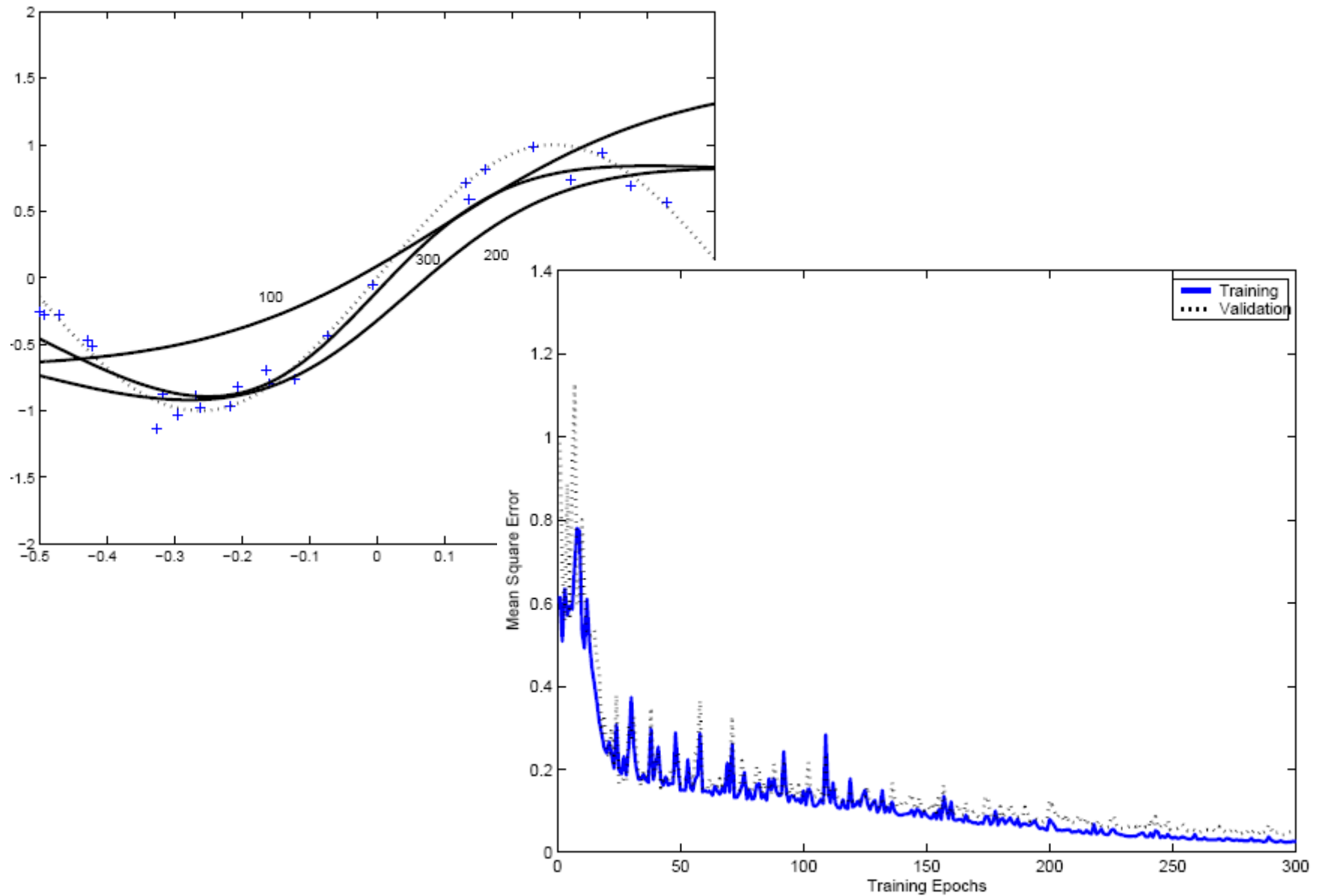
For $i = 1, \dots, K$

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

For $h = 1, \dots, H$

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

Until convergence



Improving Convergence

- Momentum

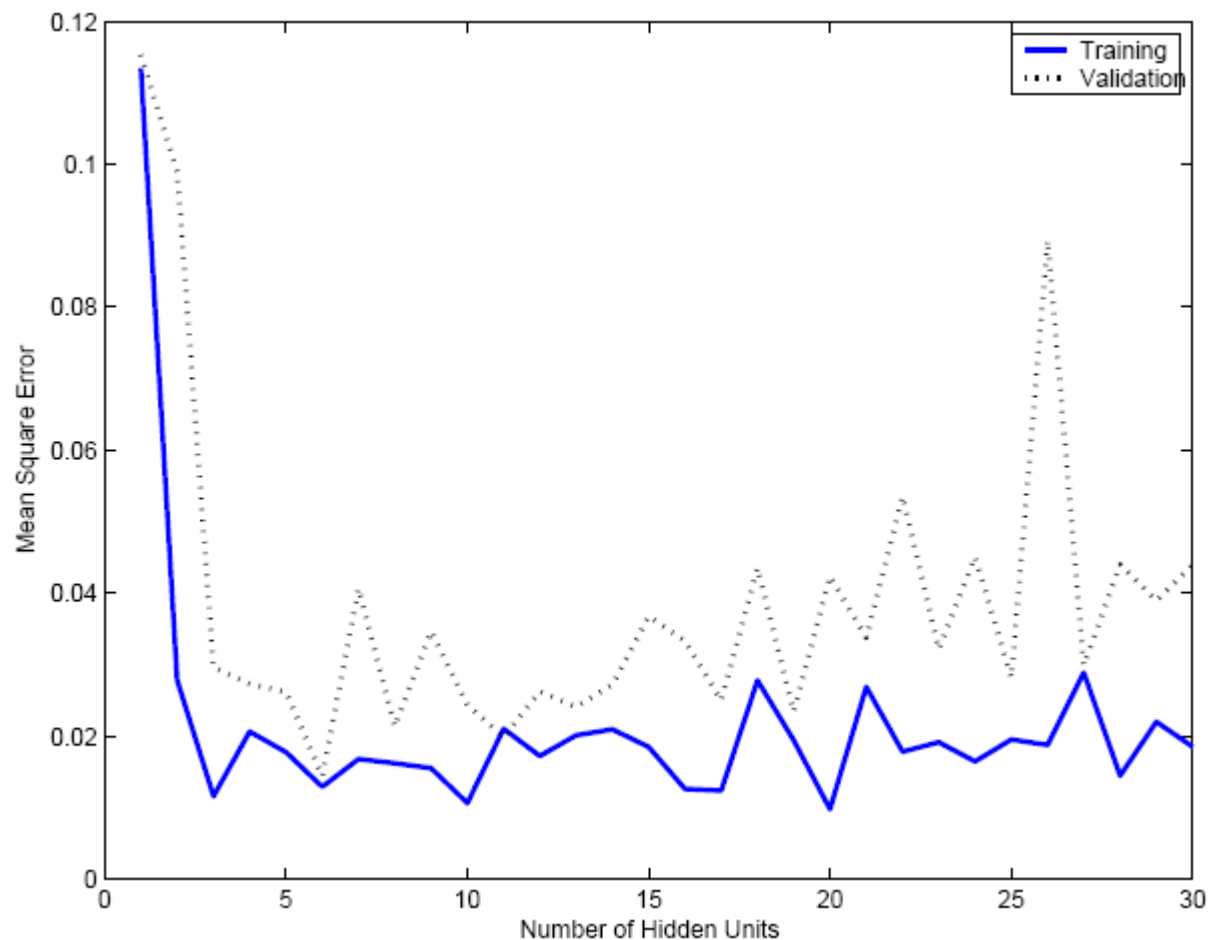
$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \alpha \Delta w_i^{t-1}$$

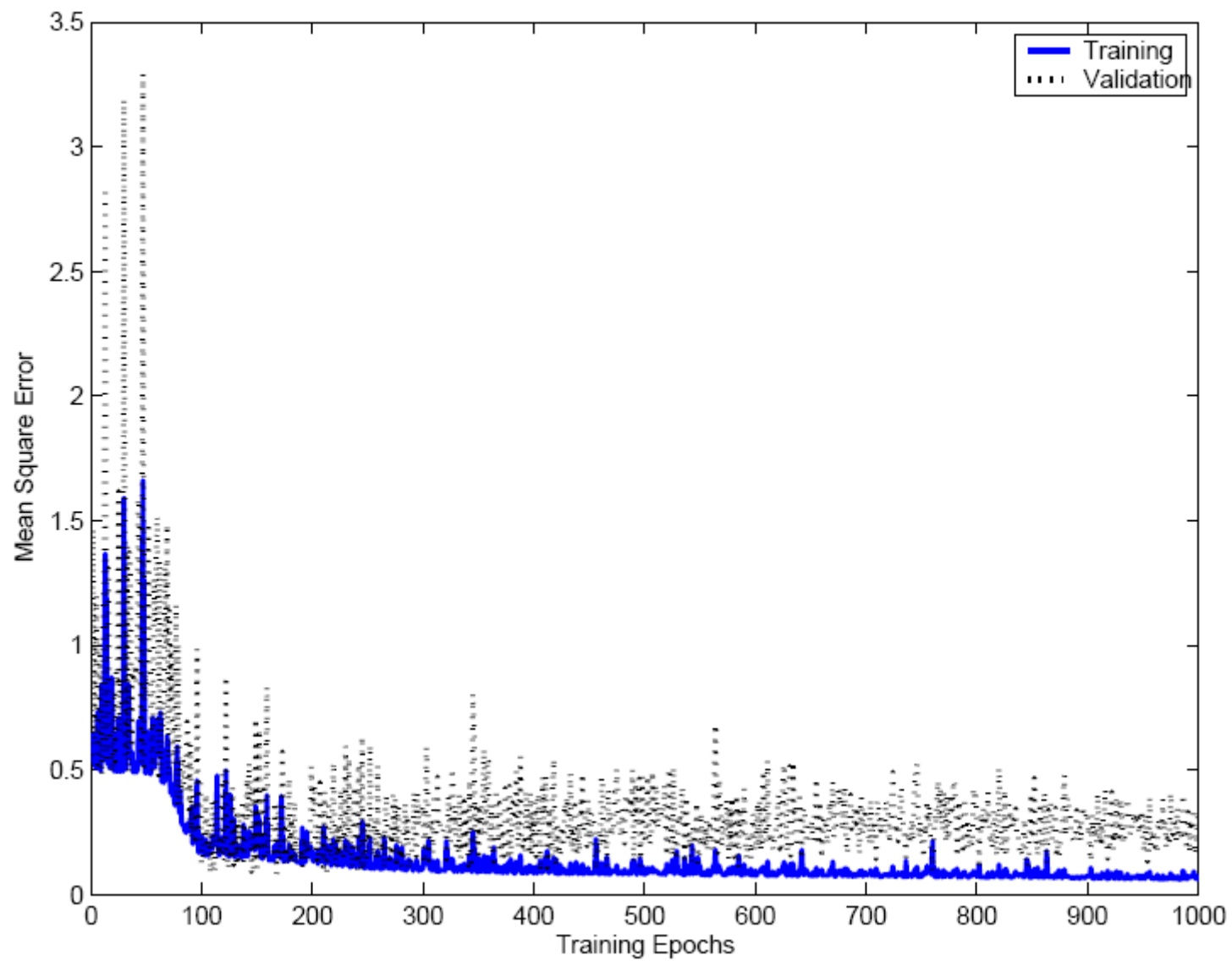
- Adaptive learning rate

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

Overfitting/Overtraining

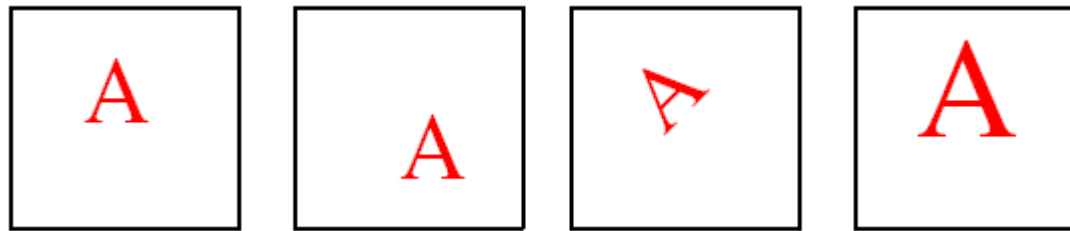
Number of weights: $H(d+1)+(H+1)K$





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”: $E_h = [g(\mathbf{x} | \theta) - g(\mathbf{x}' | \theta)]^2$

Approximation hint:

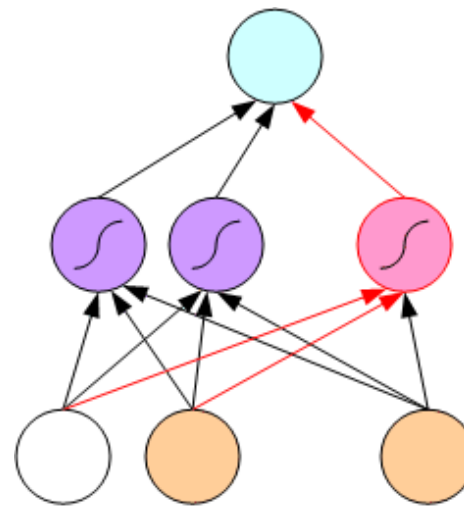
$$E_h = \begin{cases} 0 & \text{if } g(\mathbf{x} | \theta) \in [a_x, b_x] \\ (g(\mathbf{x} | \theta) - a_x)^2 & \text{if } g(\mathbf{x} | \theta) < a_x \\ (g(\mathbf{x} | \theta) - b_x)^2 & \text{if } g(\mathbf{x} | \theta) > b_x \end{cases}$$

Tuning the Network Size

- Destructive
- Weight decay:
- Constructive
- Growing networks

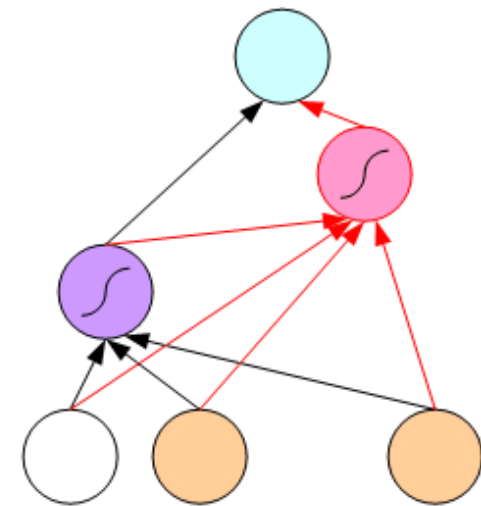
$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i$$

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



Dynamic Node Creation

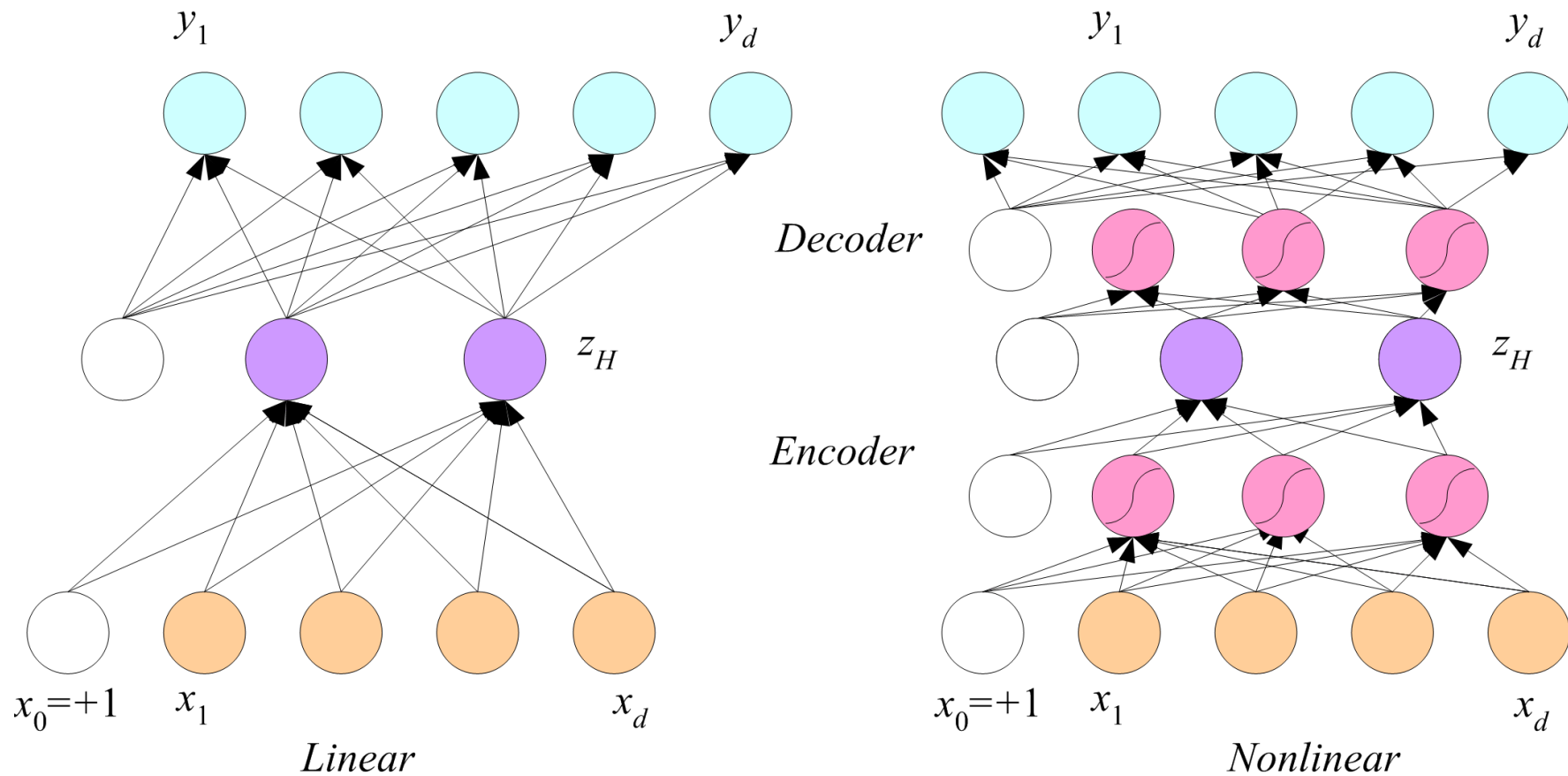
(Ash, 1989)



Cascade Correlation

(Fahlman and Lebiere, 1989)

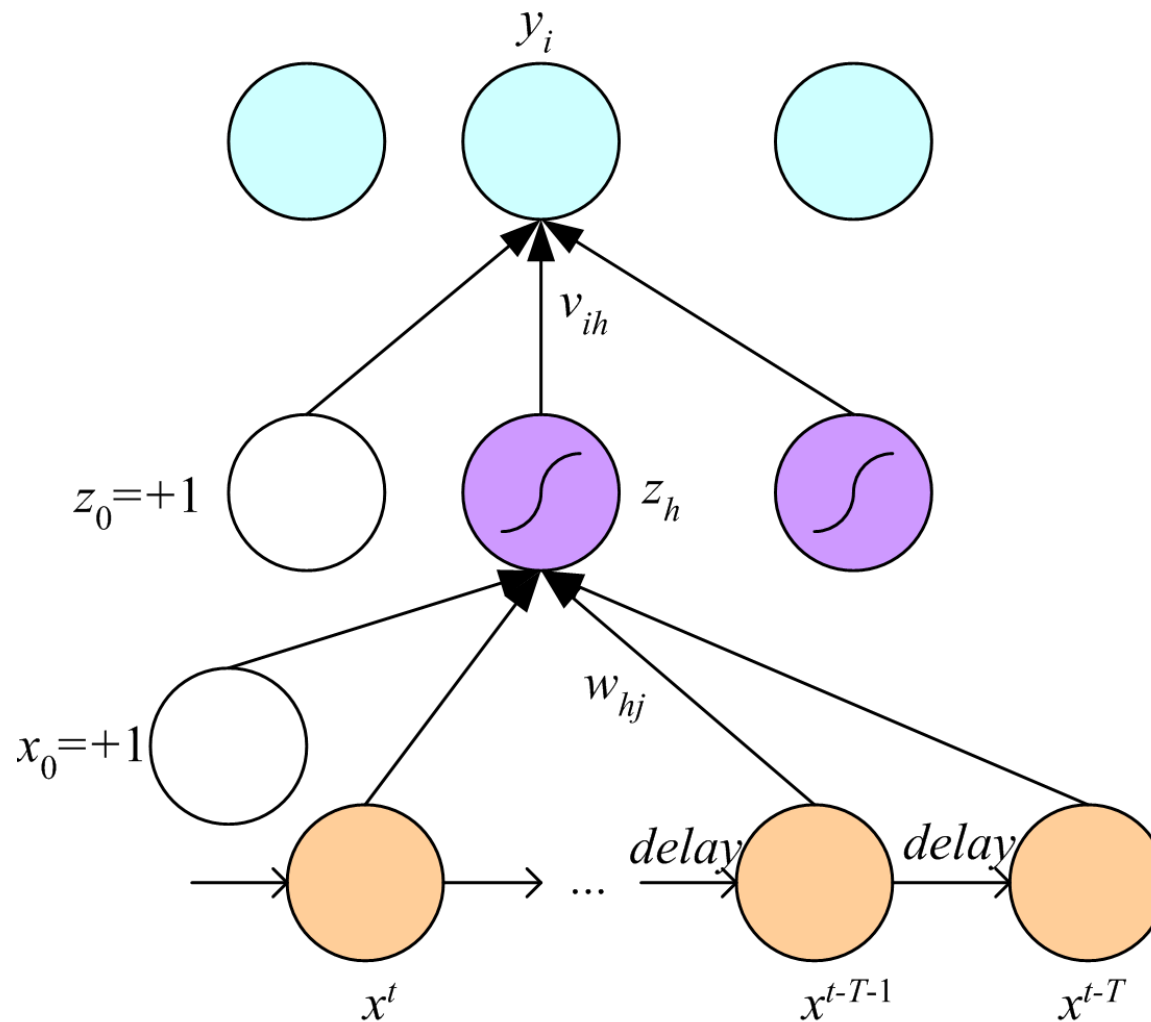
Dimensionality Reduction



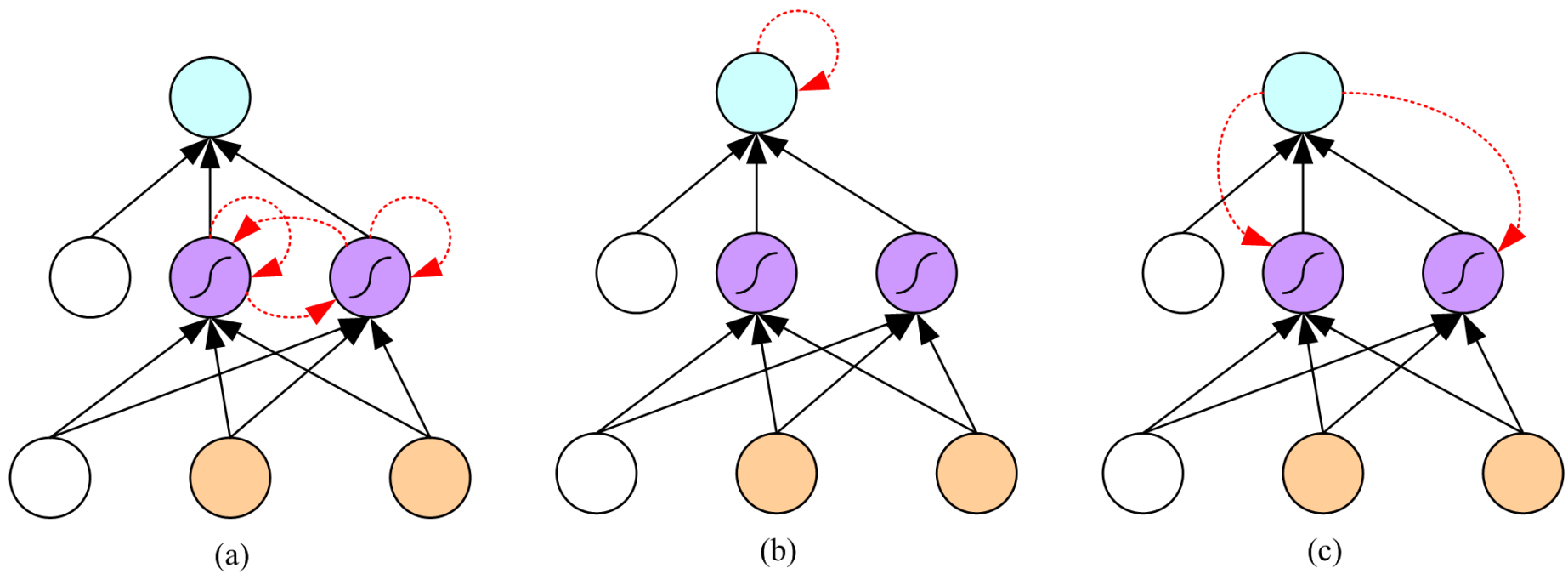
Learning Time

- 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)

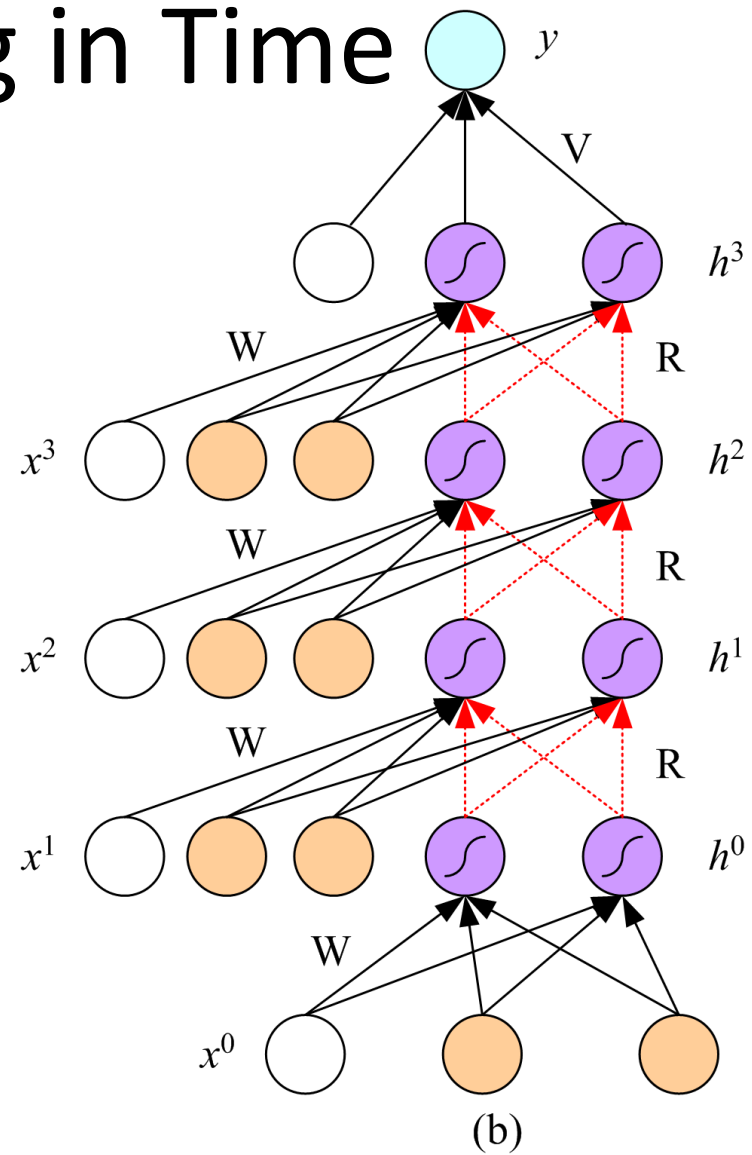
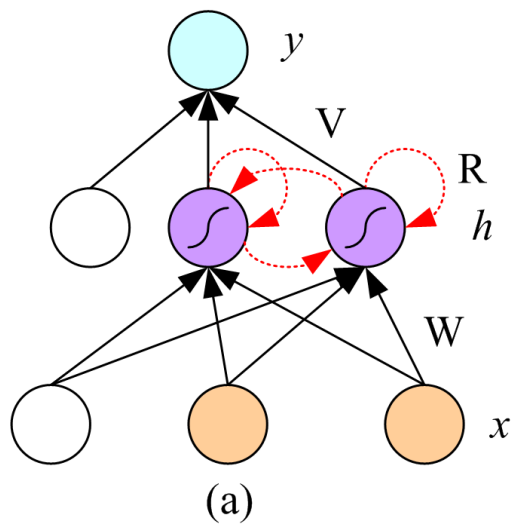
Time-Delay Neural Networks



Recurrent Networks



Unfolding in Time

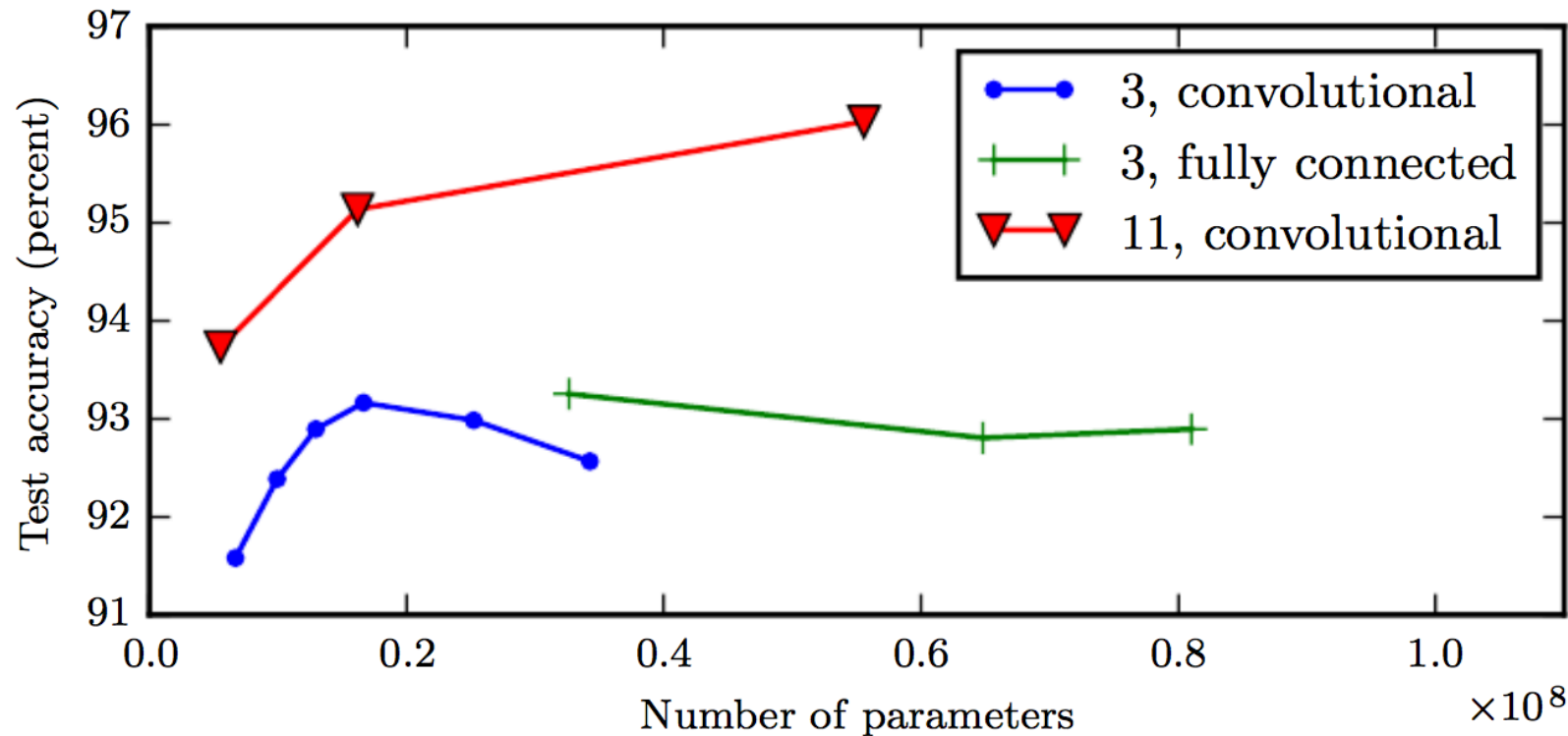


- ▶ Basic Concepts
 - ▶ supervised, unsupervised and semi-supervised learning
 - ▶ density estimation, classification, regression
 - ▶ error measures
 - ▶ optimization
 - ▶ linear and nonlinear models
- ▶ Perceptrons and Multilayer Perceptrons
 - ▶ perceptron
 - ▶ multilayer perceptron
 - ▶ training
 - ▶ overfitting and its prevention
 - ▶ learning from data and hints
 - ▶ neural networks for time series
- ▶ Deep Neural Networks
 - ▶ shallow networks' problems
 - ▶ optimization
 - ▶ regularization
 - ▶ convolution
 - ▶ hardware

Universal Approximation Theorem

- MLP with one hidden layer is a **universal approximator** (Hornik et al., 1989)
- But using multiple layers may lead to simpler networks
- There is an MLP that fits the data, but it might not be possible to find that MLP using the training algorithms we have.

How to Train a MLP with Lots of Data



Having a shallow network with a lot of hidden units may not learn as easy as a deep network.

Deep network may partition tasks into subtasks and learn easier.

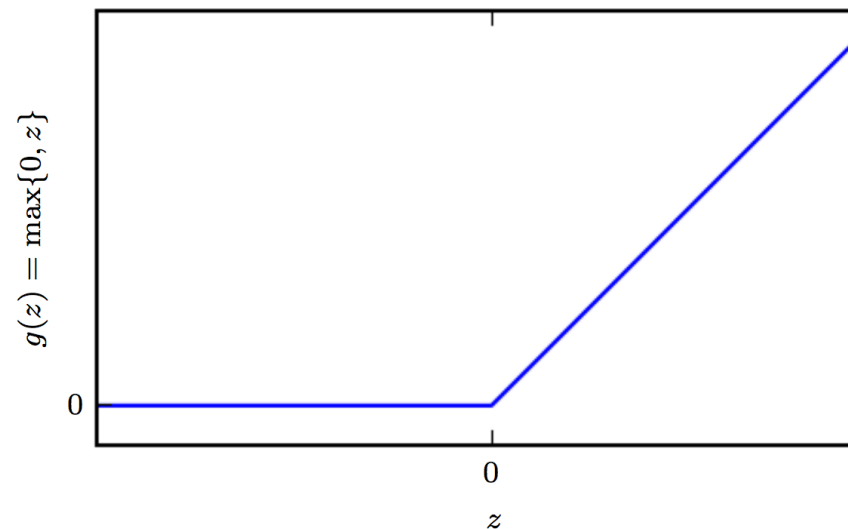
(Goodfellow, 2014)

Vanishing Gradients

- Learning with a lot of layers is difficult.
- Since gradients are computed based on gradients of layers closer to the outputs, using chain rule, the gradients of weights closer to the inputs get very small.

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

- Instead of sigmoid, ReLU (Rectified Linear Activation Function) is used. ReLU outputs are linear for positive inputs and 0 for negative inputs



$$\text{ReLU}(z) = \max(0, z)$$

Deep Learning Libraries



Caffe

theano



- When computing gradients, use graph representation of the MLP.
- Instead of single weights, use vectors/tensors of weights.
- Different libraries (Torch, Caffe vs Theano, TensorFlow) approaches to the gradient computation.

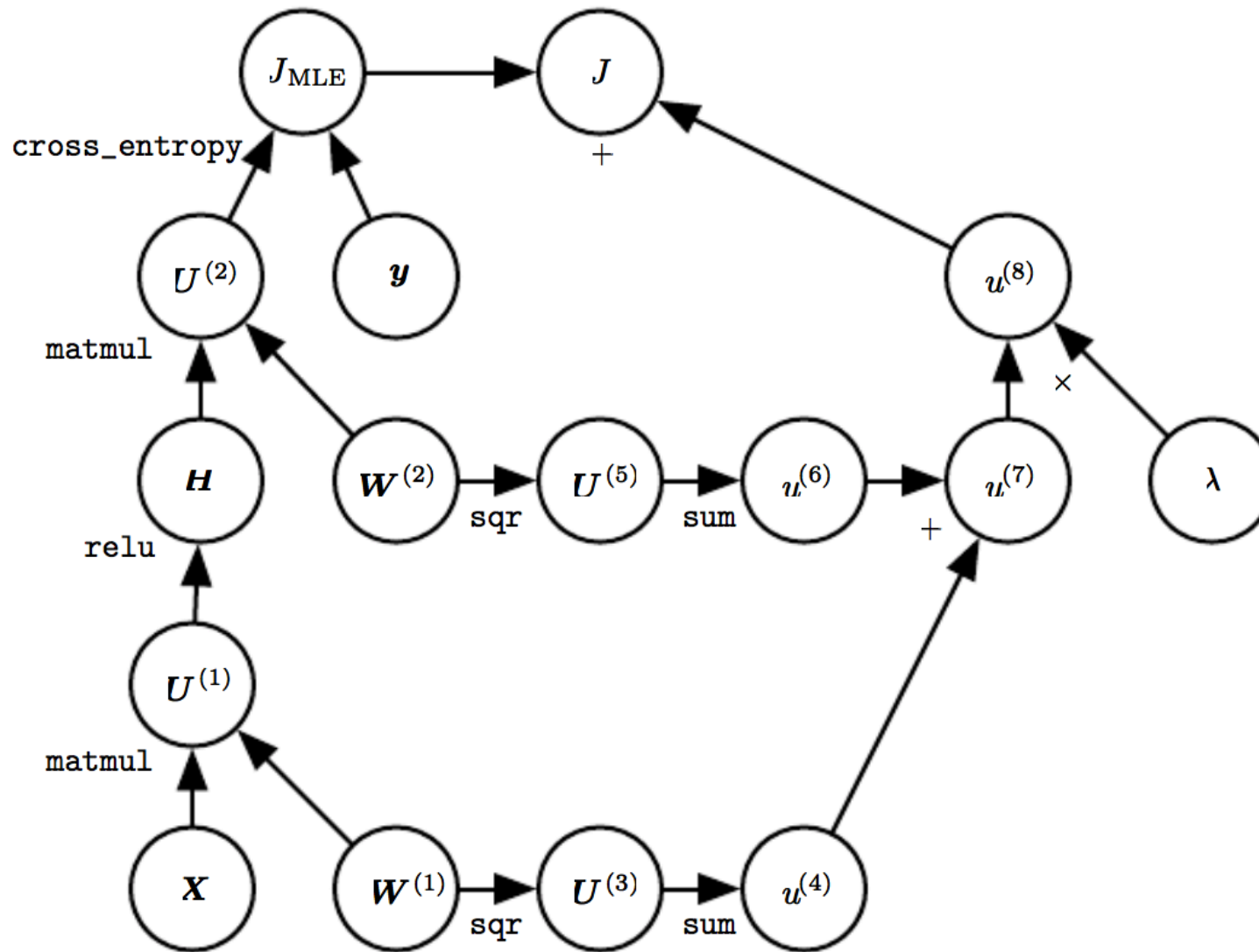


Figure 6.11: The computational graph used to compute the cost used to train our example of a single-layer MLP using the cross-entropy loss and weight decay.

Optimization

- Theano and Tensorflow approximate the Hessian Matrix (higher order derivatives) using Krylov Methods. Matrix inversions (therefore computations of eigenvectors/values are also approximated).
- Nesterov Momentum: similar to momentum, but evaluate gradient after applying the current velocity.
- Minibatches (in order to optimize for the hardware also)
- Conjugate gradient
- Polyak averaging
- Greedy supervised pretraining

Regularization for Deep Learning

- L1 and L2 (ridge regression/Tikhonov regularization) regularization
- Feature selection (results in weight elimination)
- Early stopping
- Dataset augmentation: e.g. add rotated, blurred, scaled etc. images

Regularization for Deep Learning

- Ability to penalize weights at different layers differently (Srebro 2005, constrain norm of each column of the weight matrix of a neural network, rather than constrain the whole weight matrix)
- Noise robustness: dropout algorithm (approximates bagging, zero random some weights at each minibatch)
- Parameter tying and parameter sharing: reduce the number of parameters used by clustering/forcing similar valued parameters together and representing them using smaller number of bits. CNNs already do parameter sharing
- Adversarial training: training on perturbed images

Convolution

- Sparse interactions, parameter sharing and equivariant representations
- Parallel convolutions (feature extraction units), followed by ReLU, followed by pooling to to obtain summary statistics of nearby cells.

Hardware: GPUs and TPUs to make matrix/weight operations faster

LAMBDA BLADE

If you're looking for an enterprise solution for cutting-edge Deep Learning performance and parallelization, the [Lambda Blade](#) is a necessity. Experience how profitable of an investment Deep Learning hardware can be for your team.

Designed for deployment in data centers, our 8 GPU configurations significantly boost Deep Learning performance and are ideal for remote access from multiple concurrent users.

\$16,500

Includes FREE shipping

Ships in 3 days



LAMBDA SINGLE

[SINGLE](#) is an introduction to high performance Deep Learning hardware. Pre-installed is every framework you'll need — including Tensorflow, Torch, and Caffe.

Have any questions? Call us: [1 \(650\) 479-5530](tel:16504795530).

GPU

1x 1080 Ti

NVIDIA GeForce GTX 1080 Ti or Titan Xp. Pascal architecture

STORAGE

250GB SSD

Configurable to 500 GB SATA SSD

PROCESSOR

Intel i5

Intel™ Core i5-6500, 6M Cache, up to 3.60 GHz

RAM

16GB DDR4

Configurable to 32 GB DDR4 RAM.

\$2,599

Includes FREE shipping

Ships in 3 days



source: <https://lambdal.com>

1 at 1

source: <https://lambdal.com>

Teşekkürler. Questions?

cataltepe@itu.edu.tr
zehra@tazi.io



[1] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.

<http://www.deeplearningbook.org>

Alpaydin, Ethem. Introduction to machine learning. MIT press, 2014.