

CS536: Machine Learning

Artificial Neural Networks

Fall 2005

Ahmed Elgammal

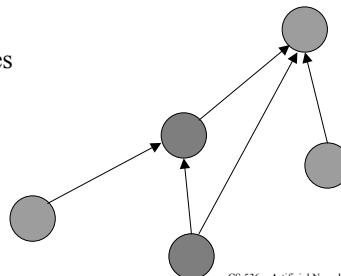
Dept of Computer Science

Rutgers University

Neural Networks

Biological Motivation: Brain

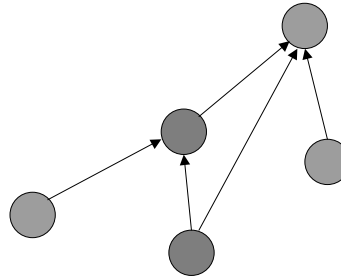
- Networks of processing units (neurons) with connections (synapses) between them
- Large number of neurons: 10^{11}
- Large connectivity: each connected to, on average, 10^4 others
- Switching time 10^{-3} second
- Parallel processing
- Distributed computation/memory
- Processing is done by neurons and the memory is in the synapses
- Robust to noise, failures



Neural Networks

Characteristic of Biological Computation

- Massive Parallelism
- Locality of Computation
- Adaptive (Self Organizing)
- Representation is Distributed



Understanding the Brain

- 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

- ALVINN system

Pomerleau (1993)

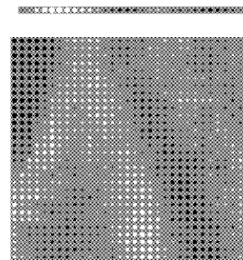
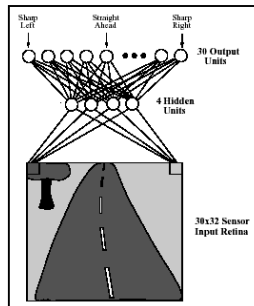
- Many successful examples:

Speech phoneme recognition [Waibel]

Image classification [Kanade, Baluja, Rowley]

Financial prediction

Backgammon [Tesauro]



CS 536 – Artificial Neural Networks - - 5

When to use ANN

- Input is high-dimensional discrete or real-valued (e.g. raw sensor input). Inputs can be highly correlated or independent.
- Output is discrete or real valued
- Output is a vector of values
- Possibly noisy data. Data may contain errors
- Form of target function is unknown
- Long training time are acceptable
- Fast evaluation of target function is required
- Human readability of learned target function is unimportant

⇒ ANN is much like a black-box

CS 536 – Artificial Neural Networks - - 6

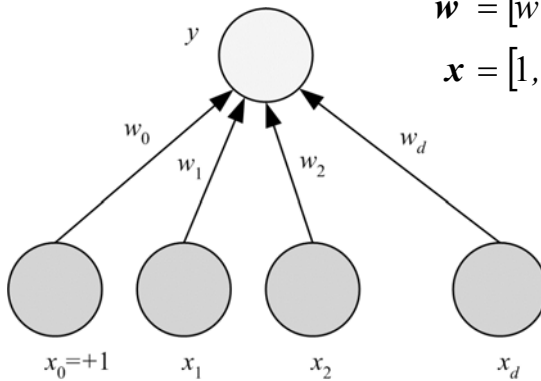
Perceptron

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

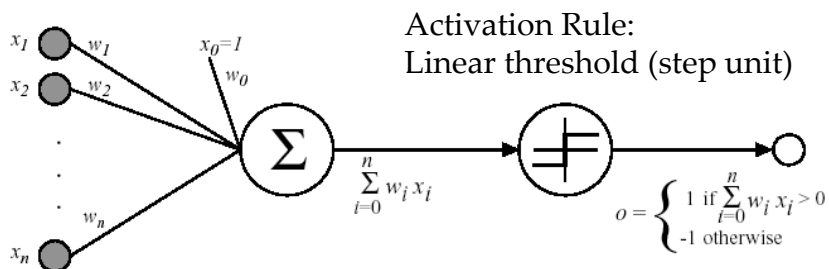
$$\mathbf{w} = [w_0, w_1, \dots, w_d]^T$$

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

(Rosenblatt, 1962)



Perceptron

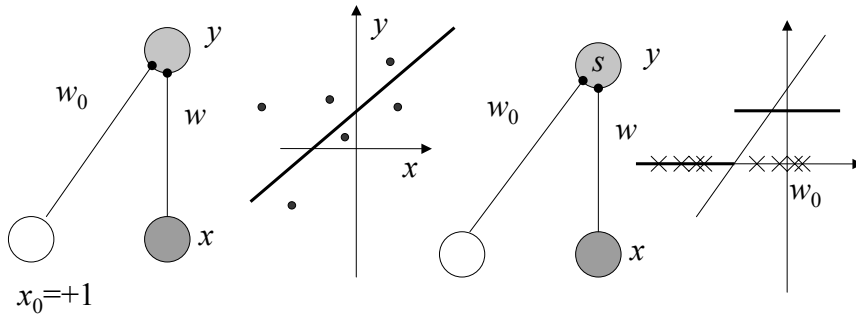


$$o(x_1, \dots, x_n) = \begin{cases} 1 & \text{if } w_0 + w_1 x_1 + \dots + w_n x_n > 0 \\ -1 & \text{otherwise.} \end{cases}$$

Or, more succinctly: $o(\mathbf{x}) = \text{sgn}(\mathbf{w} \cdot \mathbf{x})$

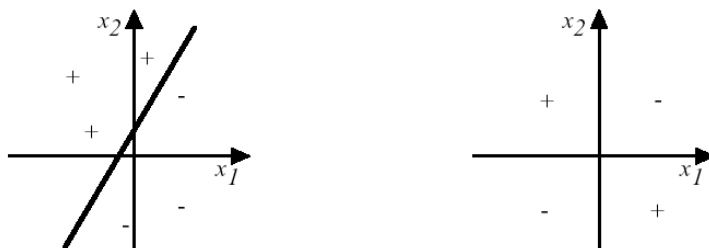
What a Perceptron Does

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



CS 536 – Artificial Neural Networks - 9

Perceptron Decision Surface



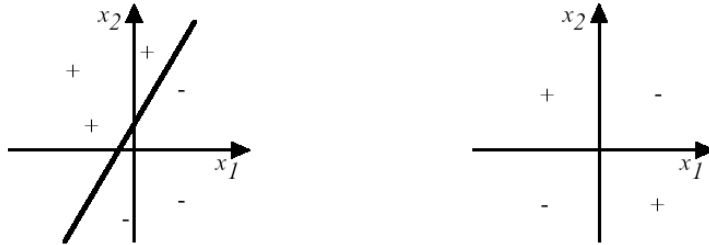
Perceptron as hyperplane decision surface in the n -dimensional input space

The perceptron outputs 1 for instances lying on one side of the hyperplane and outputs -1 for instances on the other side

Data that can be separated by a hyperplane: *linearly separable*

CS 536 – Artificial Neural Networks - 10

Perceptron Decision Surface



A single unit can represent some useful functions

- What weights represent

$g(x_1, x_2) = \text{AND}(x_1, x_2)$? Majority, Or

But some functions not representable

- e.g., not linearly separable
- Therefore, we'll want networks of these...

Perceptron training rule

$$w_i \leftarrow w_i + \Delta w_i$$

where

$$\Delta w_i = \eta (t - o) x_i$$

Where:

- $t = c(\mathbf{x})$ is target value
- o is perceptron output
- η is small constant (e.g., .1) called the *learning rate* (or *step size*)

Perceptron training rule

Can prove it will converge

- If training data is linearly separable
- and η sufficiently small
- Perceptron Convergence Theorem (Rosenblatt): if the data are linearly separable then the perceptron learning algorithm converges in finite time.

Gradient Descent – Delta Rule

Also known as LMS (least mean squares) rule or Widrow-Hoff rule.

To understand, consider simpler *linear unit*, where

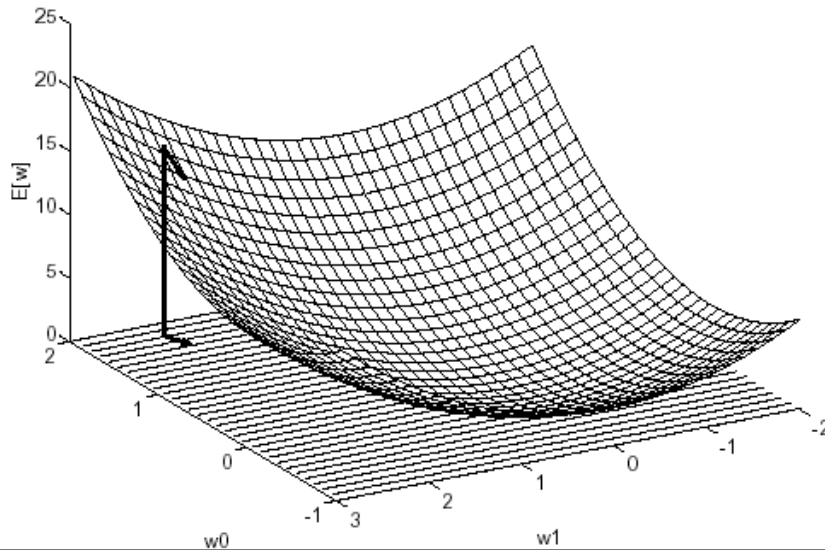
$$o = w_0 + w_1x_1 + \dots + w_nx_n$$

Let's learn w_i 's to minimize squared error

$$E[\mathbf{w}] \equiv 1/2 \sum_{d \in D} (t_d - o_d)^2$$

Where D is set of training examples

Error Surface



Gradient Descent

Gradient

$$\nabla E[\mathbf{w}] = [\partial E / \partial w_0, \partial E / \partial w_1, \dots, \partial E / \partial w_n]$$

When interpreted as a vector in weight space, the gradient specifies the direction that produces the steepest increase in E

Training rule:

$$\Delta \mathbf{w} = -\eta \nabla E[\mathbf{w}]$$

in other words:

$$\Delta w_i = -\eta \partial E / \partial w_i$$

This results in the following update rule:

$$\Delta w_i = \eta \sum_d (t_d - o_d) (x_{i,d})$$

Gradient of Error

$$\partial E / \partial w_i$$

$$= \partial / \partial w_i \frac{1}{2} \sum_d (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d \partial / \partial w_i (t_d - o_d)^2$$

$$= \frac{1}{2} \sum_d 2 (t_d - o_d) \partial / \partial w_i (t_d - o_d)$$

$$= \sum_d (t_d - o_d) \partial / \partial w_i (t_d - \mathbf{w} \cdot \mathbf{x}_d)$$

$$= \sum_d (t_d - o_d) (-x_{i,d})$$

Learning Rule:

$$\Delta w_i = -\eta \partial E / \partial w_i$$

$$\Rightarrow \Delta w_i = \eta \sum_d (t_d - o_d) (x_{i,d})$$

Stochastic Gradient Descent

Batch mode Gradient Descent:

Do until satisfied

1. Compute the gradient $\nabla E_D [\mathbf{w}]$
2. $\mathbf{w} \leftarrow \mathbf{w} - \nabla E_D [\mathbf{w}]$

Incremental mode Gradient Descent:

Do until satisfied

- For each training example d in D
 1. Compute the gradient $\nabla E_d [\mathbf{w}]$
 2. $\mathbf{w} \leftarrow \mathbf{w} - \nabla E_d [\mathbf{w}]$

More Stochastic Grad. Desc.

$$E_D[\mathbf{w}] \equiv 1/2 \sum_{d \in D} (t_d - o_d)^2$$

$$E_d[\mathbf{w}] \equiv 1/2 (t_d - o_d)^2$$

Incremental Gradient Descent can approximate Batch Gradient Descent arbitrarily closely if η set small enough

Incremental Learning Rule: $\Delta w_i = \eta (t_d - o_d) (x_{i,d})$

Delta Rule: $\Delta w_i = \eta (t - o) (x_i)$

$\delta = (t - o)$

Gradient Descent Code

GRADIENT-DESCENT(training examples, η)

Each training example is a pair of the form $\langle \mathbf{x}, t \rangle$, where \mathbf{x} is the vector of input values, and t is the target output value. η is the learning rate (e.g., .05).

- Initialize each w_i to some small random value
- Until the termination condition is met, Do
 - Initialize each Δw_i to zero.
 - For each $\langle \mathbf{x}, t \rangle$ in training examples, Do
 - Input the instance \mathbf{x} to the unit and compute the output o
 - For each linear unit weight w_i , Do
$$\Delta w_i \leftarrow \Delta w_i + \eta (t - o) x_i$$
 - For each linear unit weight w_i , Do
$$w_i \leftarrow w_i + \Delta w_i$$

Summary

Perceptron training rule will succeed if

- Training examples are linearly separable
- Sufficiently small learning rate η

Linear unit training uses gradient descent (delta rule)

- Guaranteed to converge to hypothesis with minimum squared error
- Given sufficiently small learning rate η
- Even when training data contains noise
- Even when training data not H separable

K Outputs

Regression:

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

$$\mathbf{y} = \mathbf{W}\mathbf{x} \quad \text{Linear Map from } \mathbb{R}^d \Rightarrow \mathbb{R}^k$$

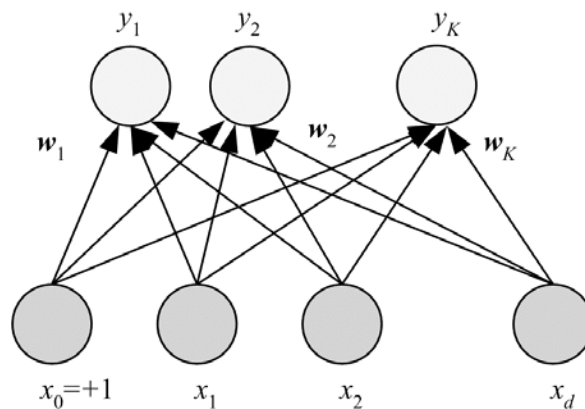
Classification:

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

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

choose C_i

if $y_i = \max_k y_k$



Training

- 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):

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

Update = Learning Factor \cdot (Desired Output - Actual Output) \cdot Input

Training a Perceptron: Regression

- Regression (Linear output):

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

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

Classification

- Single sigmoid output

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

- $K > 2$ softmax outputs

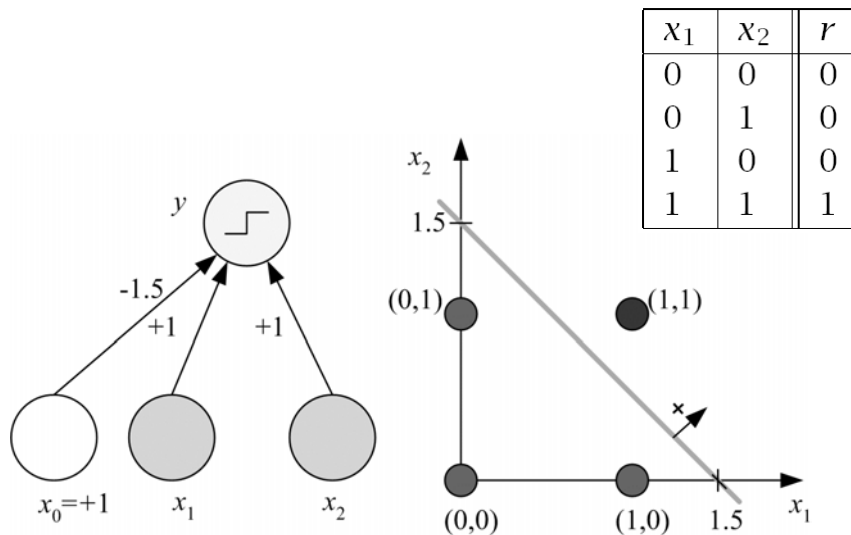
$$E^t(\mathbf{w} | \mathbf{x}^t, \mathbf{r}^t) = -r^t \log y^t - (1 - r^t) \log (1 - y^t)$$

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

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

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

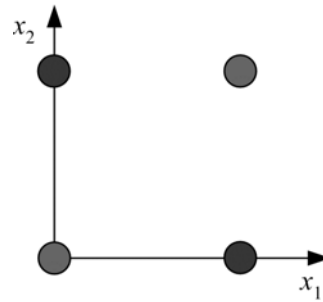
Learning Boolean AND



XOR

• No v

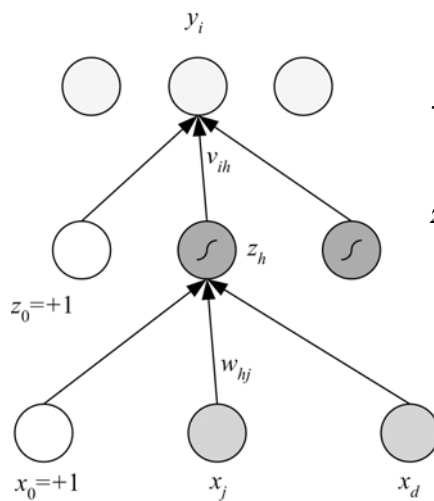
x_1	x_2	r
0	0	0
0	1	1
1	0	1
1	1	0



$$\begin{aligned}
 w_0 &\leq 0 \\
 w_2 + w_0 &> 0 \\
 w_1 + w_0 &> 0 \\
 w_1 + w_2 + w_0 &\leq 0
 \end{aligned}$$

(Minsky and Papert, 1969)

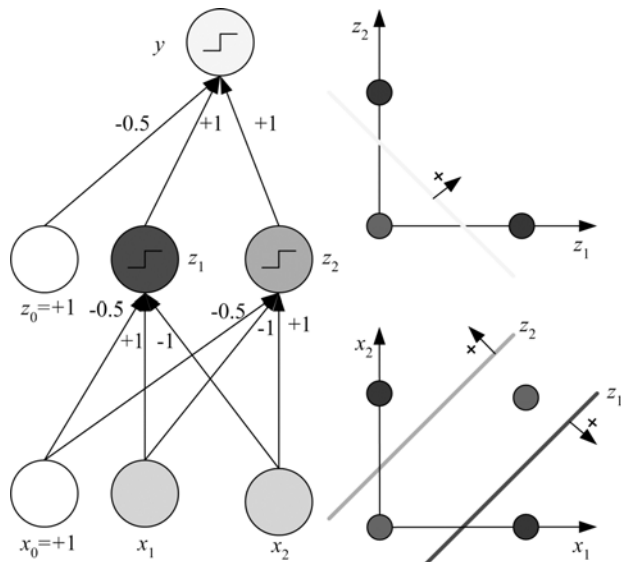
Multilayer Perceptrons



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

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

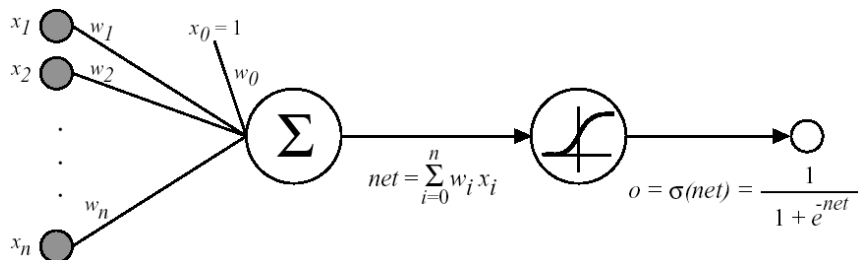
(Rumelhart et al., 1986)



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

CS 536 – Artificial Neural Networks - - 29

Sigmoid Activation



$\sigma(x)$ is the sigmoid (s-like) function

$$\frac{1}{1 + e^{-x}}$$

Nice property:

$$d \sigma(x)/dx = \sigma(x) (1 - \sigma(x))$$

Other variations:

$\sigma(x) = 1/(1 + e^{-k \cdot x})$ where k is a positive constant that determines the steepness of the threshold.

CS 536 – Artificial Neural Networks - - 30

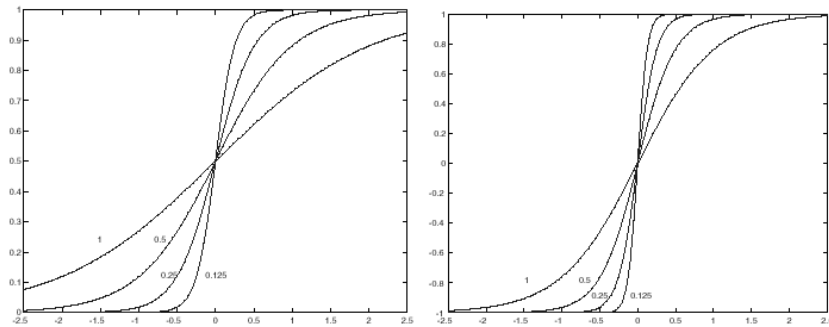


Figure 22.14. On the **left**, a series of squashing functions obtained using $\phi(x; \nu) = \frac{e^x}{1+e^x}$, for different values of ν indicated on the figure. On the **right**, a series of squashing functions obtained using $\phi(x; \nu, A) = A \tanh(x/\nu)$ for different values of ν indicated on the figure. Generally, for x close to the center of the range, the squashing function is linear; for x small or large, it is strongly non-linear.

CS 536 – Artificial Neural Networks - - 31

Error Gradient for Sigmoid

$$\begin{aligned}
 \partial E / \partial w_i &= \partial / \partial w_i \frac{1}{2} \sum_d (t_d - o_d)^2 \\
 &= \frac{1}{2} \sum_d \partial / \partial w_i (t_d - o_d)^2 \\
 &= \frac{1}{2} \sum_d 2 (t_d - o_d) \partial / \partial w_i (t_d - o_d) \\
 &= \sum_d (t_d - o_d) (-\partial o_d / \partial w_i) \\
 &= - \sum_d (t_d - o_d) (\partial o_d / \partial net_d \partial net_d / \partial w_i)
 \end{aligned}$$

CS 536 – Artificial Neural Networks - - 32

Even more...

But we know:

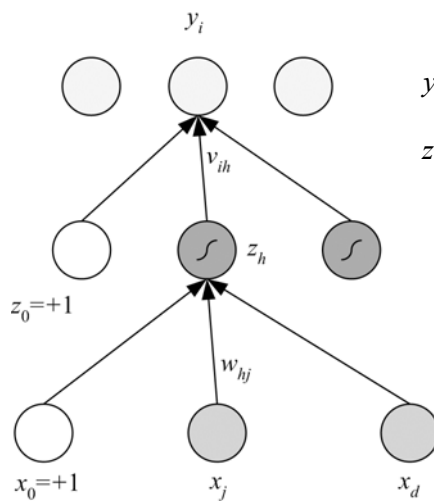
$$\begin{aligned} \partial o_d / \partial net_d &= \partial \sigma(net_d) / \partial net_d = o_d (1 - o_d) \end{aligned}$$

$$\partial net_d / \partial w_i = \partial (\mathbf{w} \cdot \mathbf{x}_d) / \partial w_i = x_{i,d}$$

So:

$$\partial E / \partial w_i = - \sum_d (t_d - o_d) o_d (1 - o_d) x_{i,d}$$

Backpropagation



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

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

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

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

Backpropagation Algorithm

Initialize all weights to small random numbers.

Until satisfied, Do

- For each training example, Do
 - Input the training example to the network and compute the network outputs
 - For each output unit k

$$\delta_k = o_k(1 - o_k)(t_k - o_k)$$
 - For each hidden unit h

$$\delta_h = o_h(1 - o_h) \sum_{k \text{ in outputs}} w_{h,k} \delta_k$$

- Update each network weight $w_{i,j}$

$$w_{i,j} \leftarrow w_{i,j} + \Delta w_{i,j} \text{ where } \Delta w_{i,j} = \eta \delta_j a_i$$

Regression

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

Forward

$$z_h = \text{sigmoid}(w_h^T x)$$

Backward

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

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

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

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

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

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

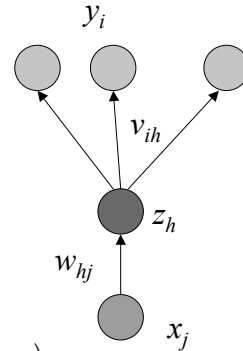
Regression with Multiple Outputs

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

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

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

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



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

Repeat

For all $(\mathbf{x}^t, \mathbf{r}^t) \in \mathcal{X}$ 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 $h = 1, \dots, H$

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

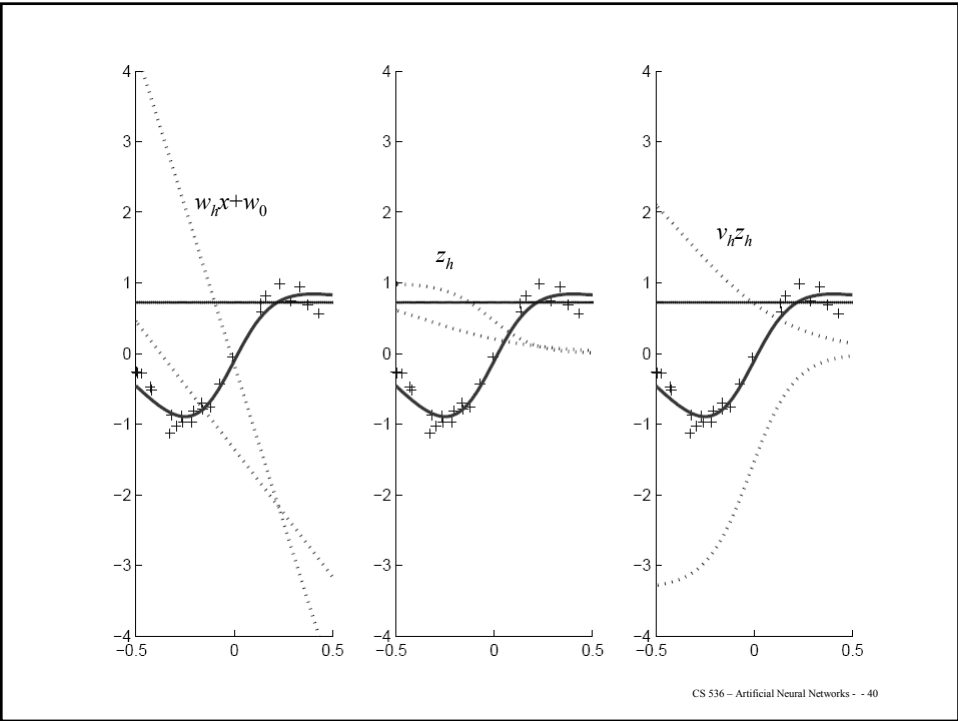
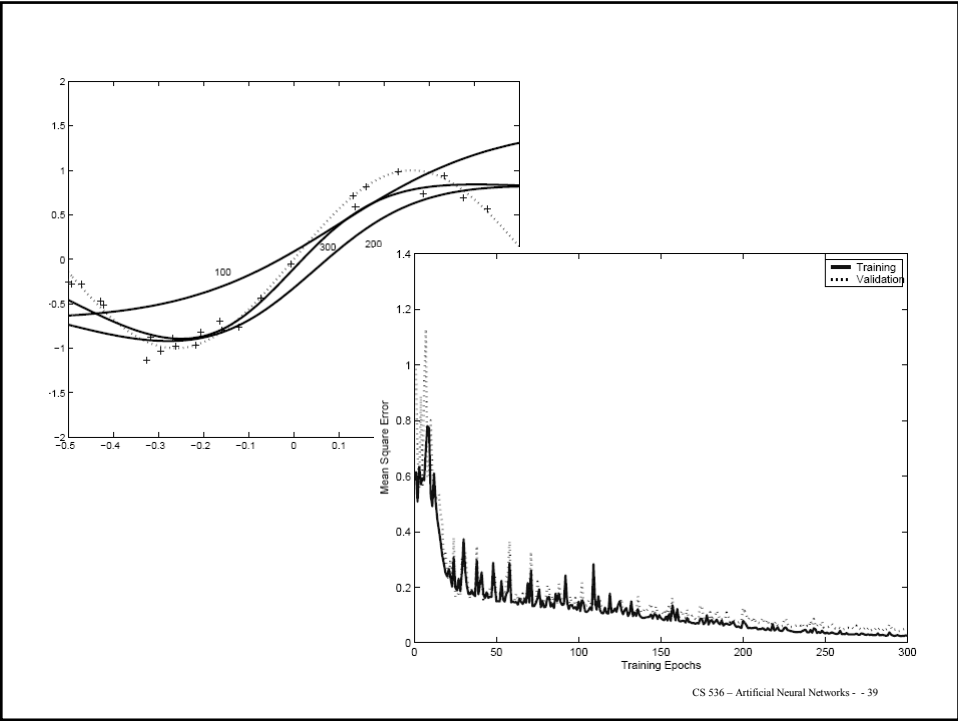
For $i = 1, \dots, K$

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

For $h = 1, \dots, H$

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

Until convergence



Two-Class Discrimination

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

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

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

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

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

K>2 Classes

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

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

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

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

Multiple Hidden Layers

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

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

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

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

Improving Convergence

- Momentum

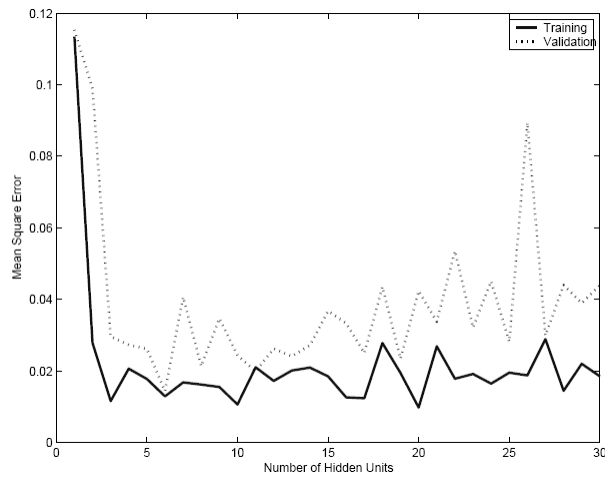
$$\Delta w_i^t = -\eta \frac{\partial E^t}{\partial w_i} + \alpha \Delta w_i^{t-1}$$

- Adaptive learning rate

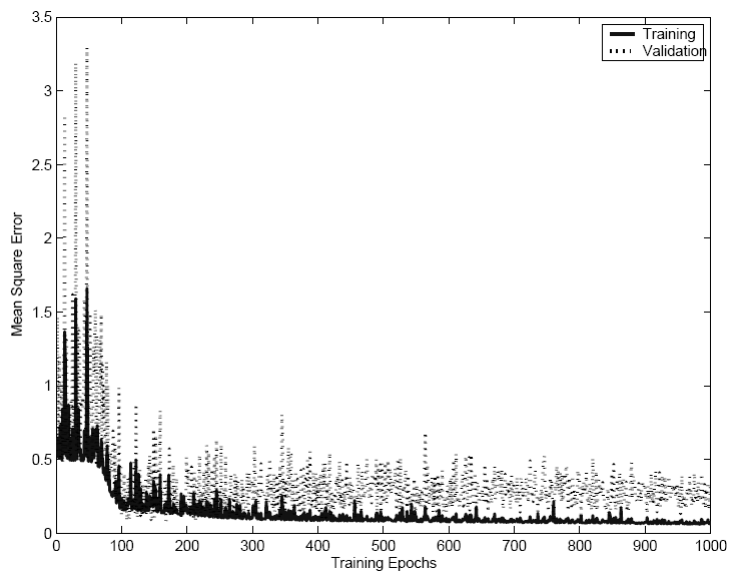
$$\Delta \eta = \begin{cases} +a & \text{if } E^{t+\tau} < E^t \\ -b\eta & \text{otherwise} \end{cases}$$

Overfitting/Overtraining

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

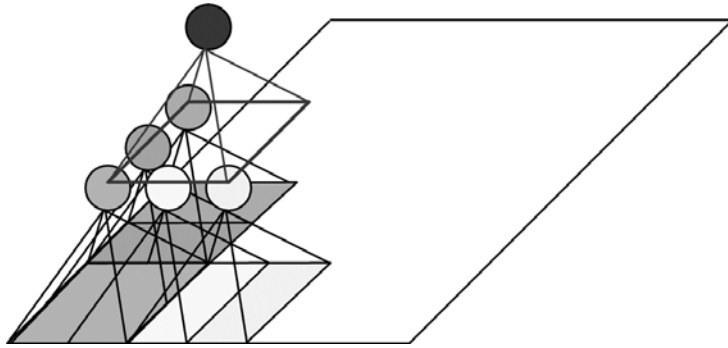


CS 536 – Artificial Neural Networks - - 45



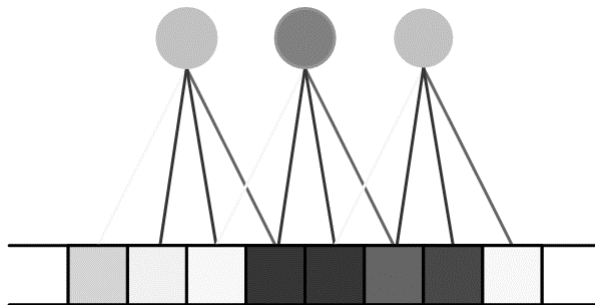
CS 536 – Artificial Neural Networks - - 46

Structured MLP



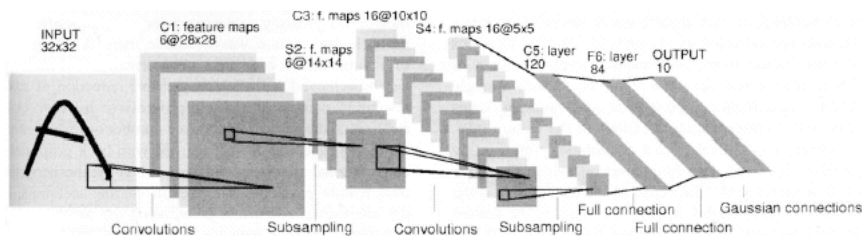
(Le Cun et al, 1989)

Weight Sharing



Convolutional neural networks

- Also known as gradient-based learning
- Template matching using NN classifiers seems to work
- Natural features are filter outputs
 - probably, spots and bars, as in texture
 - but why not learn the filter kernels, too?
- a perceptron approximates convolution.
- Network architecture: Two types of layers
 - Convolution layers: convolving the image with filter kernels to obtain filter maps
 - Subsampling layers: reduce the resolution of the filter maps
 - The number of filter maps increases as the resolution decreases



A convolutional neural network, LeNet; the layers filter, subsample, filter, subsample, and finally classify based on outputs of this process.

Figure from “Gradient-Based Learning Applied to Document Recognition”, Y. Lecun et al Proc. IEEE, 1998 copyright 1998, IEEE



Fig. 4. Size-normalized examples from the MNIST database.

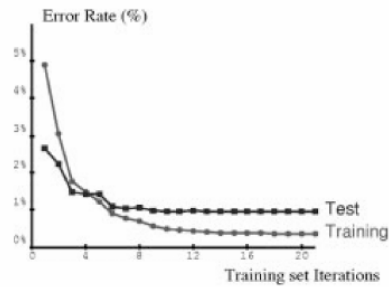


Fig. 5. Training and test error of LeNet-5 as a function of the number of passes through the 60,000 pattern training set (without distortions). The average training error is measured on-the-fly as training proceeds. This explains why the training error appears to be larger than the test error initially. Convergence is attained after 10-12 passes through the training set.

LeNet is used to classify handwritten digits. Notice that the test error rate is not the same as the training error rate, because the test set consists of items not in the training set. Not all classification schemes necessarily have small test error when they have small training error. Error rate 0.95% on MNIST database

Figure from "Gradient-Based Learning Applied to Document Recognition", Y. Lecun et al
Proc. IEEE, 1998 copyright 1998, IEEE

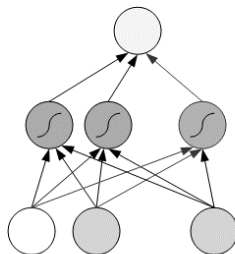
CS 536 – Artificial Neural Networks - - 51

Tuning the Network Size

- Destructive
- Weight decay:
- Constructive
- Growing networks

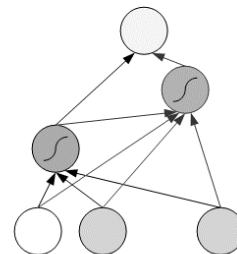
$$\Delta w_i = -\eta \frac{\partial E}{\partial w_i} - \lambda w_i$$

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



Dynamic Node Creation

(Ash, 1989)



Cascade Correlation

(Fahlman and Lebiere, 1989)

CS 536 – Artificial Neural Networks - - 53

Bayesian Learning

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

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

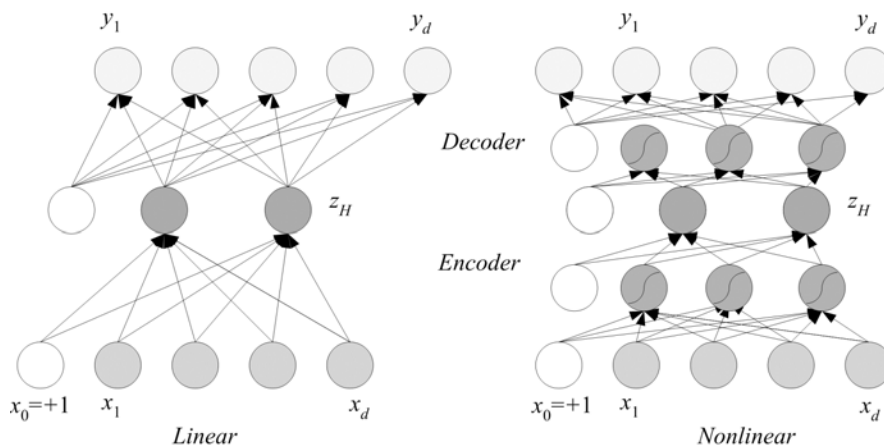
- Consider weights w_i as random vars, prior $p(w_i)$

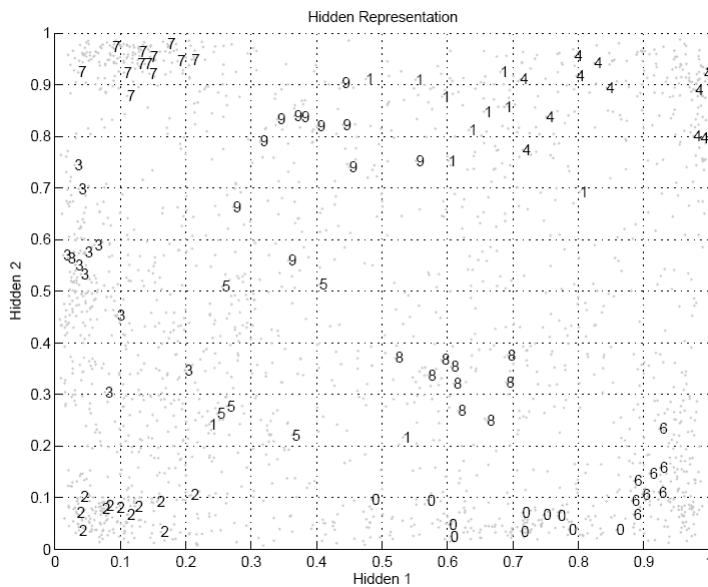
$$p(\mathbf{w}) = \prod_i p(w_i) \text{ where } p(w_i) = c \cdot \exp \left[-\frac{w_i^2}{2(1/2\lambda)} \right]$$

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

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

Dimensionality Reduction





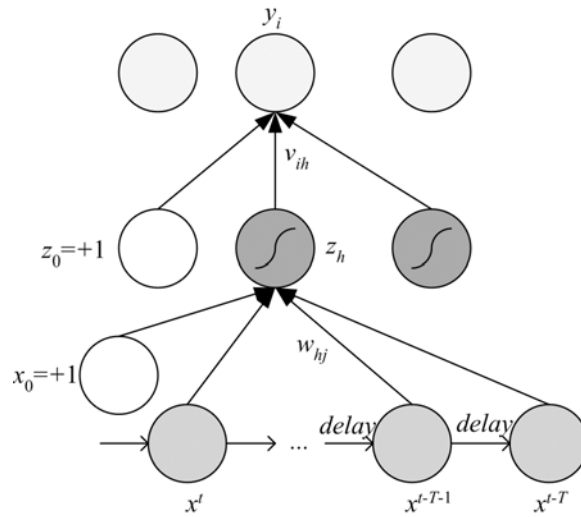
CS 536 – Artificial Neural Networks - - 56

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)

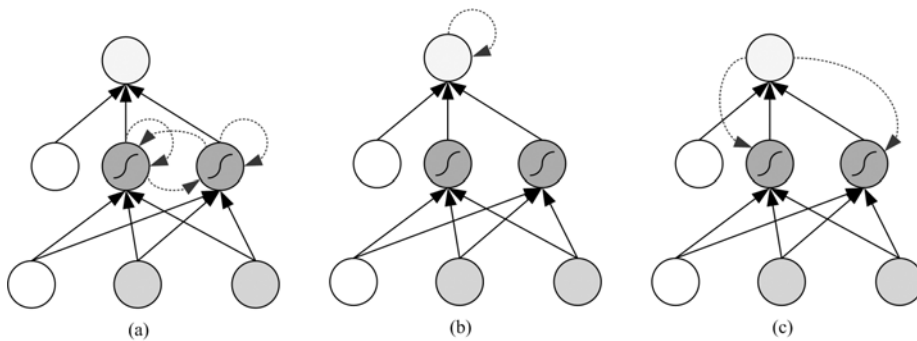
CS 536 – Artificial Neural Networks - - 57

Time-Delay Neural Networks



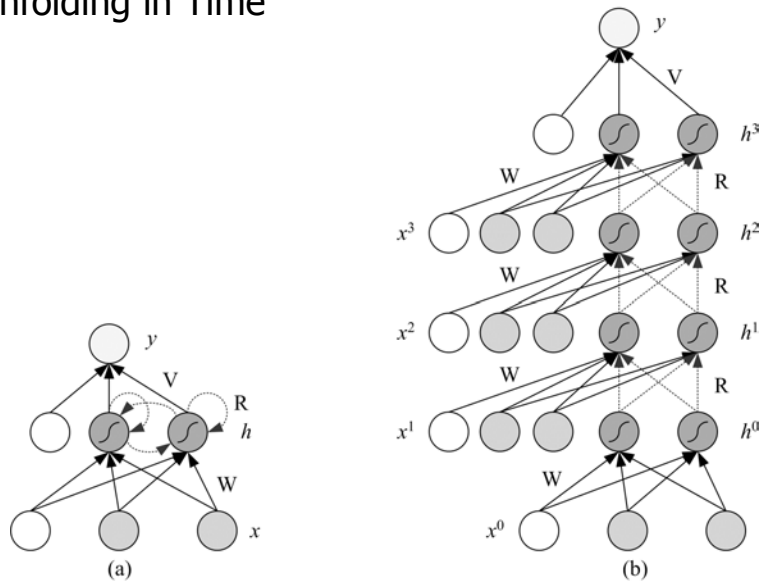
CS 536 – Artificial Neural Networks - - 58

Recurrent Networks

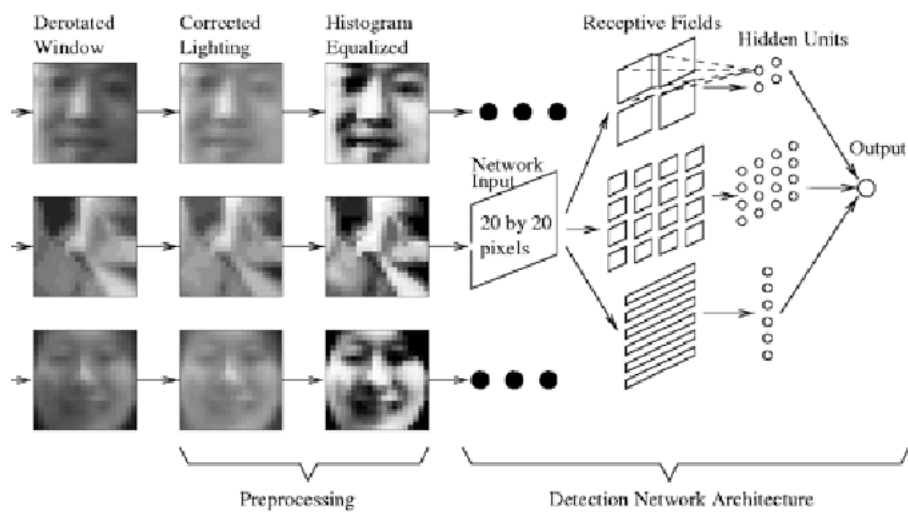


CS 536 – Artificial Neural Networks - - 59

Unfolding in Time



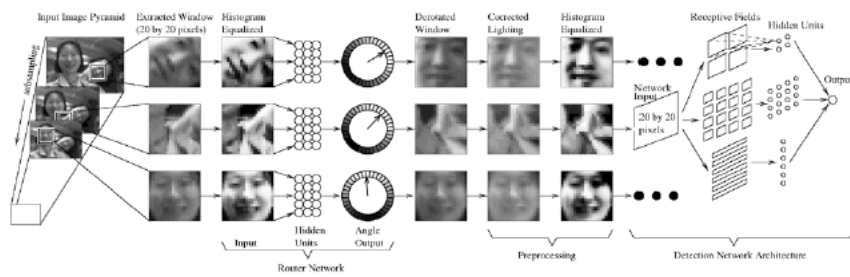
CS 536 – Artificial Neural Networks - - 60



The vertical face-finding part of Rowley, Baluja and Kanade's system

Figure from "Rotation invariant neural-network based face detection," H.A. Rowley, S. Baluja and T. Kanade, Proc. Computer Vision and Pattern Recognition, 1998, copyright 1998, IEEE

CS 536 – Artificial Neural Networks - - 61



Architecture of the complete system: they use another neural net to estimate orientation of the face, then rectify it. They search over scales to find bigger/smaller faces.

Figure from "Rotation invariant neural-network based face detection," H.A. Rowley, S. Baluja and T. Kanade, Proc. Computer Vision and Pattern Recognition, 1998, copyright 1998, IEEE



Figure from "Rotation invariant neural-network based face detection," H.A. Rowley, S. Baluja and T. Kanade, Proc. Computer Vision and Pattern Recognition, 1998, copyright 1998, IEEE

Sources

- Slides by Ethem Elpaydin, “introduction to machine learning” © The MIT Press, 2004
- Slides by Tom M. Mitchell
- Ethem Elpaydin, “introduction to machine learning” Chapter 10
- Tom M. Mitchell “Machine Learning”