Machine Learning - CMPE462

Week 8

Linear Discrimination Multilayer Perceptrons

Emre Ugur, BM 33 emre.ugur@boun.edu.tr http://www.cmpe.boun.edu.tr/~emre/courses/cmpe462 cmpe462@listeci.cmpe.boun.edu.tr

Exercise (histogram, k-nn, etc.)

Assume you are given the following training sample composed of 10 instances, where the first column corresponds to the feature value (x) and the second column corresponds to the class (c) of the instance

Predict the class of a new instance (x=5) using all the methods below.

x (5 | ?

Projects

- By April 9
- Email to instructor
 - Team members and
 - Description of the project
- Penalty: -5% x late-day
 - If not provided by deadline.

Likelihood- vs. Discriminant-based Classification

• Likelihood-based: Assume a model for $p(\mathbf{x}|C_i)$, use Bayes' rule to calculate $P(C_i|\mathbf{x})$

$$g_i(\mathbf{x}) = \log P(C_i|\mathbf{x})$$

- Discriminant-based: Assume a model for $g_i(\mathbf{x}|\Phi_i)$; no density estimation
- Estimating the boundaries is enough; no need to accurately estimate the densities inside the boundaries

Linear Discriminant

Linear discriminant:

$$g_i(\mathbf{x} | \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0} = \sum_{j=1}^{a} \mathbf{w}_{ij} \mathbf{x}_j + \mathbf{w}_{i0}$$

- Advantages:
 - Simple: O(d) space/computation
 - Knowledge extraction: Weighted sum of attributes; positive/negative weights, magnitudes (credit scoring)

Generalized Linear Model

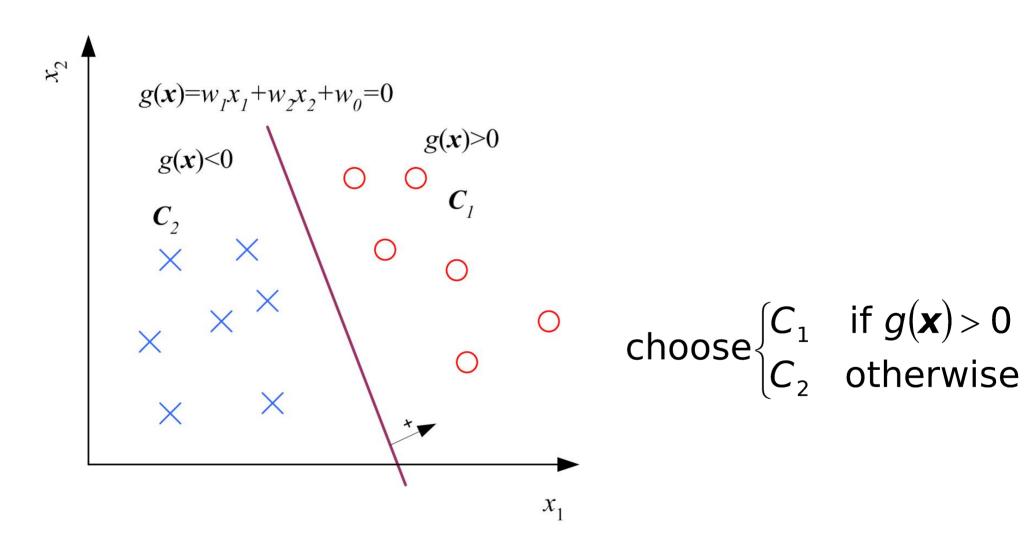
Higher-order (product) terms:

$$Z_1 = X_1$$
, $Z_2 = X_2$, $Z_3 = X_1^2$, $Z_4 = X_2^2$, $Z_5 = X_1X_2$

Map from **x** to **z** using nonlinear basis functions and use a linear discriminant in **z**-space

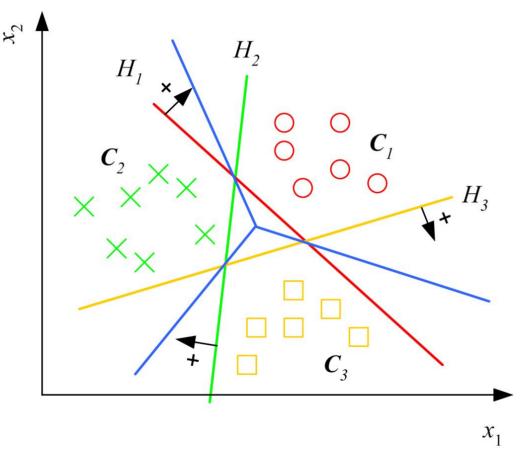
$$g_i(\mathbf{x}) = \sum_{j=1}^k \mathbf{w}_{ij} \phi_j(\mathbf{x})$$

Two Classes



Multiple Classes

$$g_i(\mathbf{x} \mid \mathbf{w}_i, \mathbf{w}_{i0}) = \mathbf{w}_i^T \mathbf{x} + \mathbf{w}_{i0}$$

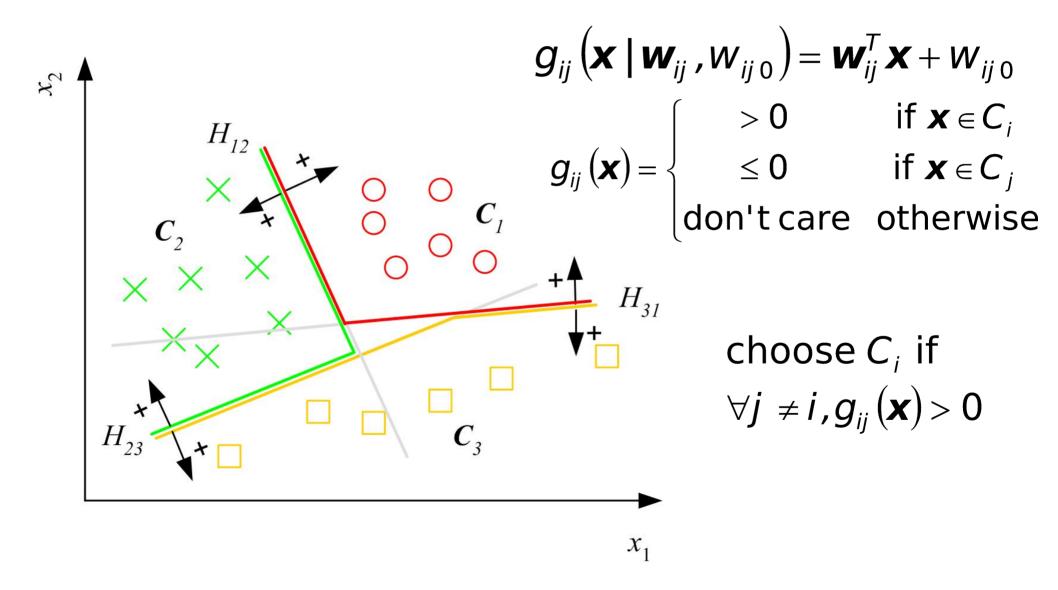


Choose C_i if

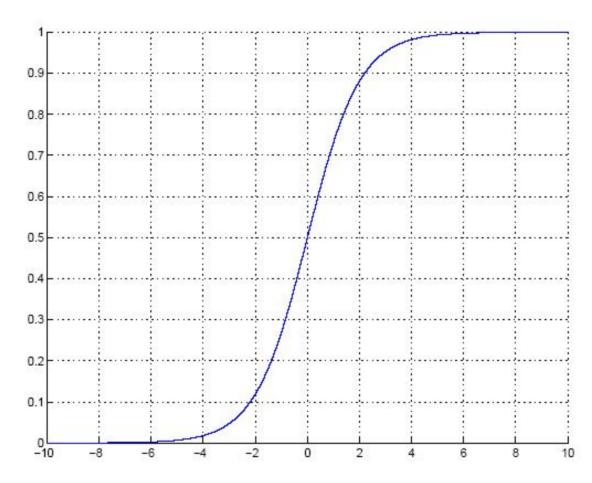
$$g_i(\mathbf{x}) = \max_{j=1}^K \mathbf{x} g_j(\mathbf{x})$$

Classes are linearly separable

Pairwise Separation



Sigmoid (Logistic) Function



- 1. Calculate $g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$ and choose C_1 if $g(\mathbf{x}) > 0$, or
- 2. Calculate $y = \text{sigmoid}(\mathbf{w}^T \mathbf{x} + \mathbf{w}_0)$ and choose C_1 if y > 0.5

Gradient-Descent

• E(w|X) is error with parameters w on sample X w^* =arg min $_w E(w|X)$

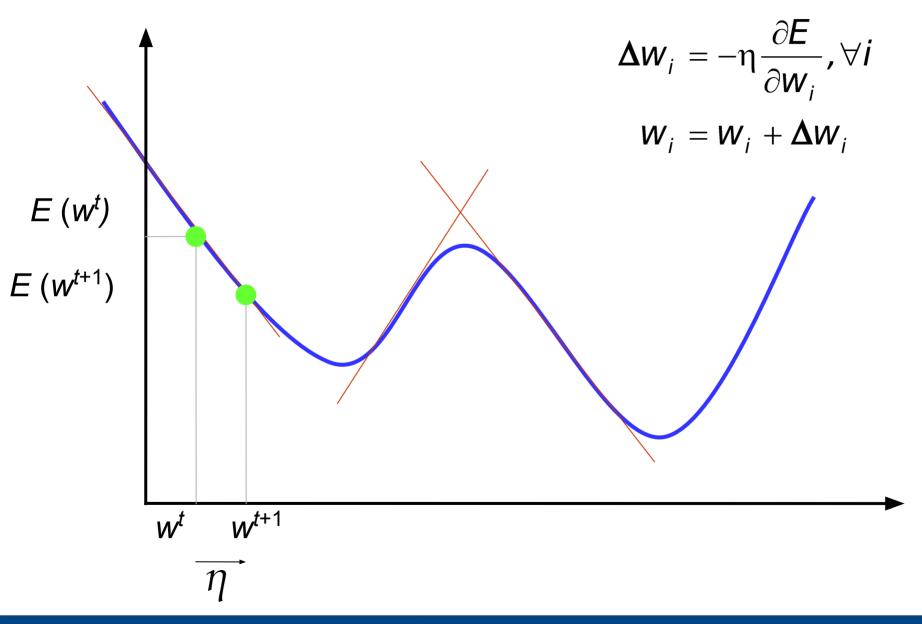
Gradient

$$\nabla_{w} E = \left[\frac{\partial E}{\partial w_{1}}, \frac{\partial E}{\partial w_{2}}, \dots, \frac{\partial E}{\partial w_{d}} \right]'$$

Gradient-descent:

Starts from random **w** and updates **w** iteratively in the negative direction of gradient

Gradient-Descent



Brain

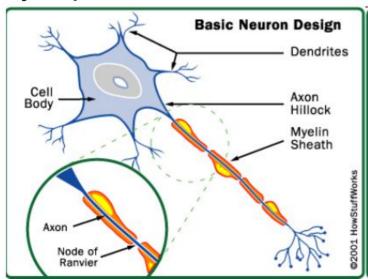
- Large number of neurons: 1010
- Large connectitivity: 105
- Parallel processing
- Distributed computation/memory
- Robust to noise, failures



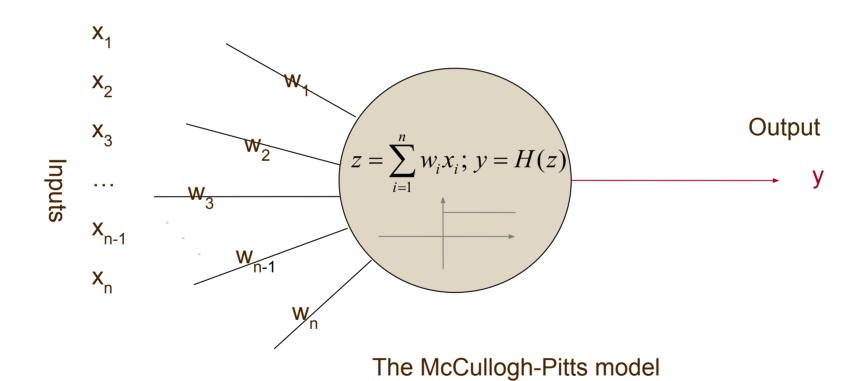


Biological inspiration

- The spikes travelling along the axon of the pre-synaptic neuron trigger the release of neurotransmitter substances at the synapse.
- The neurotransmitters cause excitation or inhibition in the dendrite of the post-synaptic neuron.
- The integration of the excitatory and inhibitory signals may produce spikes in the post-synaptic neuron.
- The contribution of the signals depends on the strength of the synaptic connection.



Artificial neurons



History

- 1943: McCulloch—Pitts "neuron"
 - Started the field
- 1962: Rosenblatt's perceptron
 - Learned its own weight values; convergence proof
- 1969: Minsky & Papert book on perceptrons
 - Proved limitations of single-layer perceptron networks
- 1982: Hopfield and convergence in symmetric networks
 - Introduced energy-function concept
- 1986: Backpropagation of errors
 - Method for training multilayer networks
- Present: Probabilistic interpretations, Bayesian and spiking networks,
- Deep-Networks

Linear Classifiers (Perceptrons)

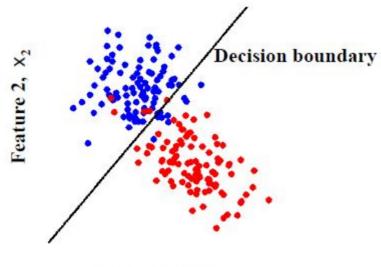
- Linear Classifiers
 - a linear classifier is a mapping which partitions feature space using
 - a linear function (a straight line, or a hyperplane)
 - separates the two classes using a straight line in feature space

Linearly separable data

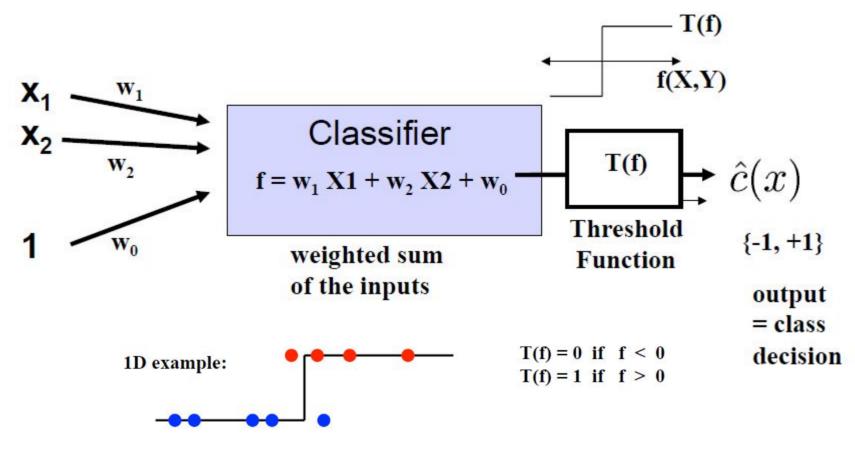
Feature 1, X1

Leature 2, x, Decision boundary

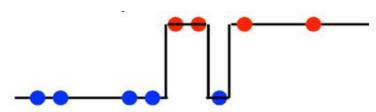
Linearly non-separable data



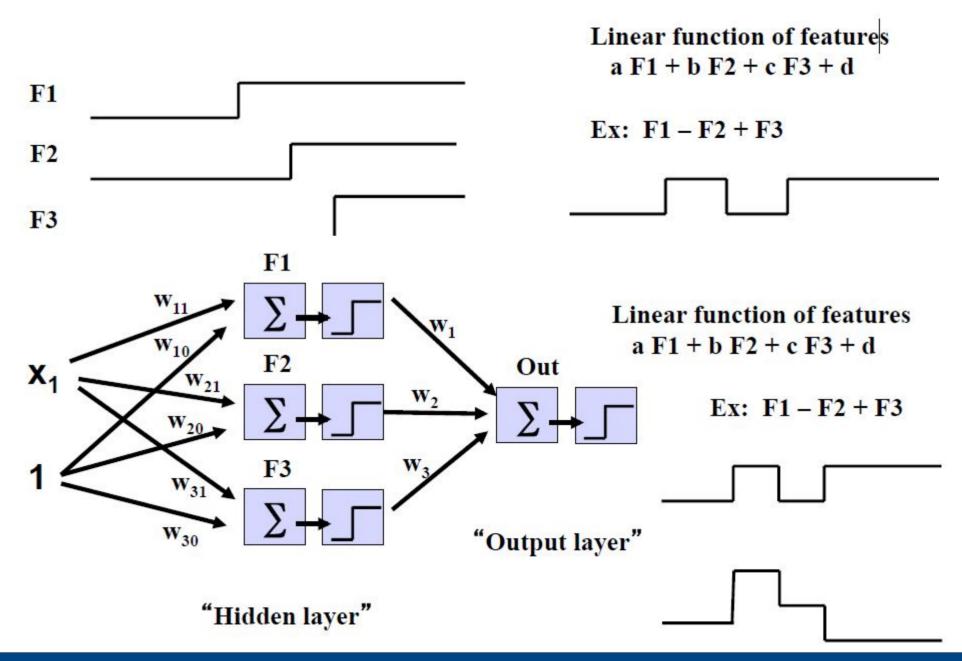
Perceptron Classifier (2 features)



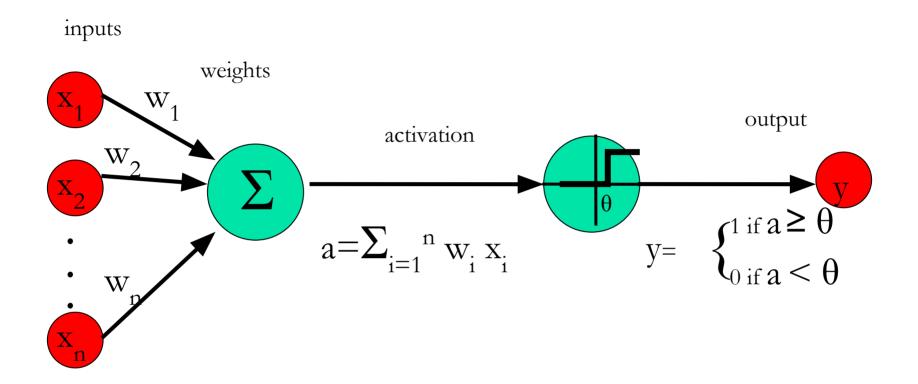
- If features change, it will find non-linear boundaries.
- What features can produce the following decision rule?



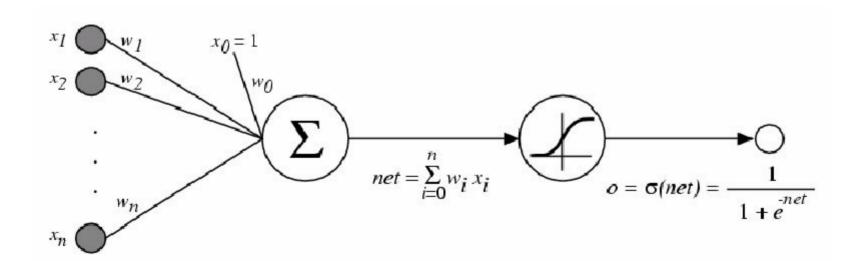
Features and perceptrons



Threshold Unit



Sigmoid Unit



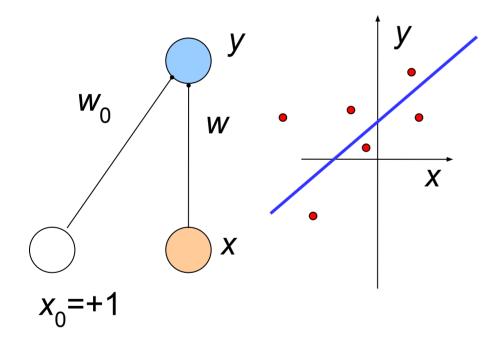
 $\sigma(x)$ is the sigmoid function

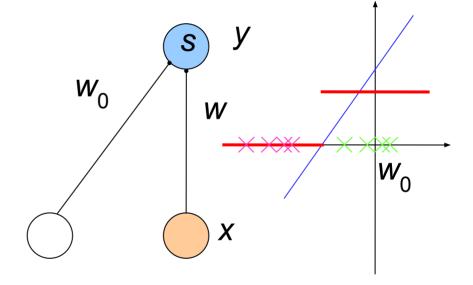
Nice property:
$$\frac{d\sigma(x)}{dx} = \sigma(x)(1-\sigma(x))^{\frac{1.0}{20.5}}$$

What a perceptron does

Regression: y=wx+w₀

• Classification: $y=1(wx+w_0>0)$





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

- If we need posterior probability

K outputs

Regression

$$\mathbf{y}_i = \sum_{j=1}^d \mathbf{w}_{ij} \mathbf{x}_j + \mathbf{w}_{i0} = \mathbf{w}_i^T \mathbf{x}$$

$$y = Wx$$

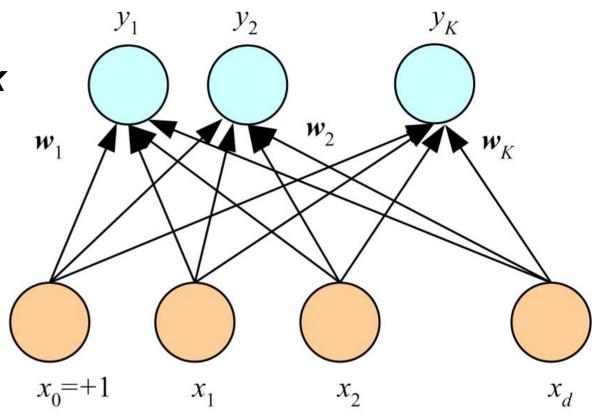
Classification

$$O_i = \mathbf{W}_i^T \mathbf{X}$$

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

 $chooseC_i$

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=LearningFactor(DesiredOutput-ActualOutput).Input

Training a Perceptron

Regression (linear output)

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

$$\Delta w_{j}^{t} = \eta(r^{t} - y^{t})x_{j}^{t}$$

Classification with single sigmoid output

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

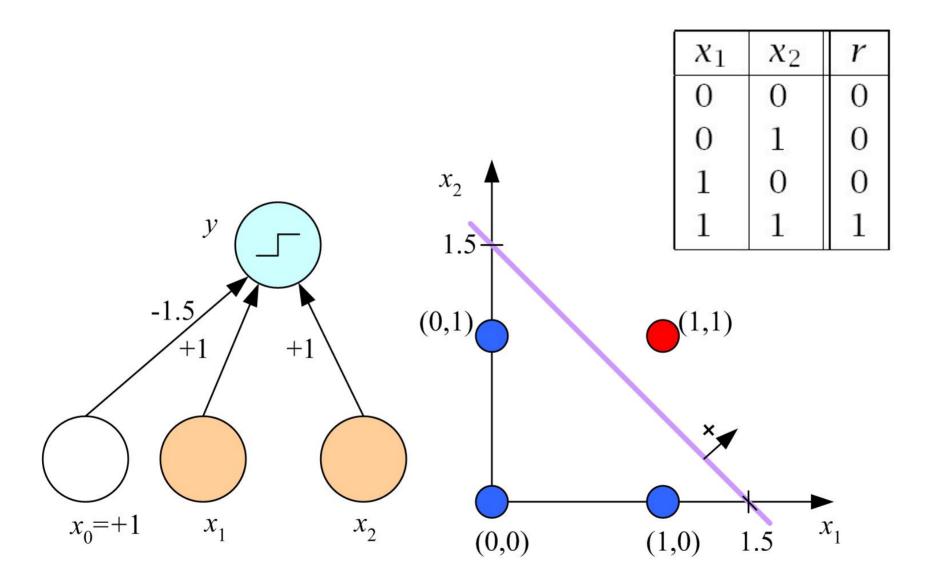
$$E^{t}(\mathbf{w} \mid \mathbf{x}^{t}, \mathbf{r}^{t}) = -r^{t} \log y^{t} - (1 - r^{t}) \log (1 - y^{t})$$

$$\Delta w_{i}^{t} = \eta (r^{t} - y^{t}) x_{i}^{t}$$

Convergence

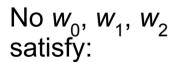
- The algorithm converges to the correct classification
 - if the training data is linearly separable, and
 - learning rate is sufficiently small

Learning Boolean AND

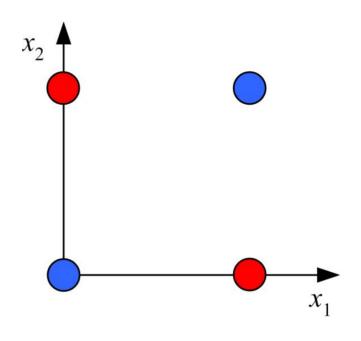


Learning Boolean XOR

x_1	<i>x</i> ₂	r
0	0	0
0	1	1
1	0	1
1	1	0



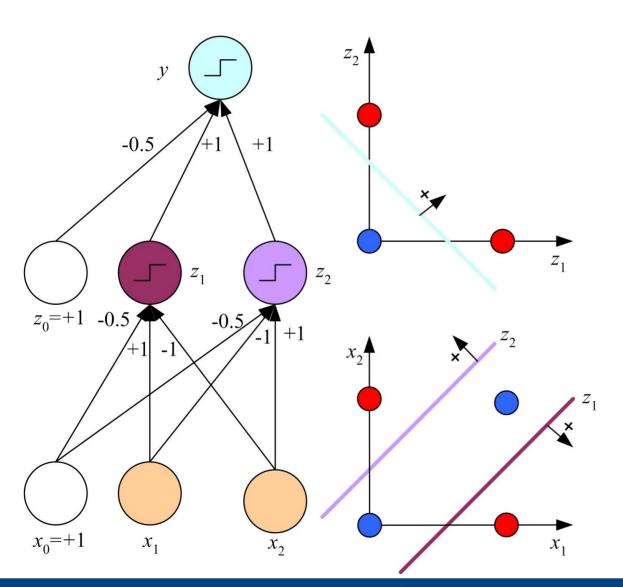
$$w_0 \le 0$$
 $w_2 + w_0 > 0$
 $w_1 + w_0 > 0$
 $w_1 + w_2 + w_0 \le 0$



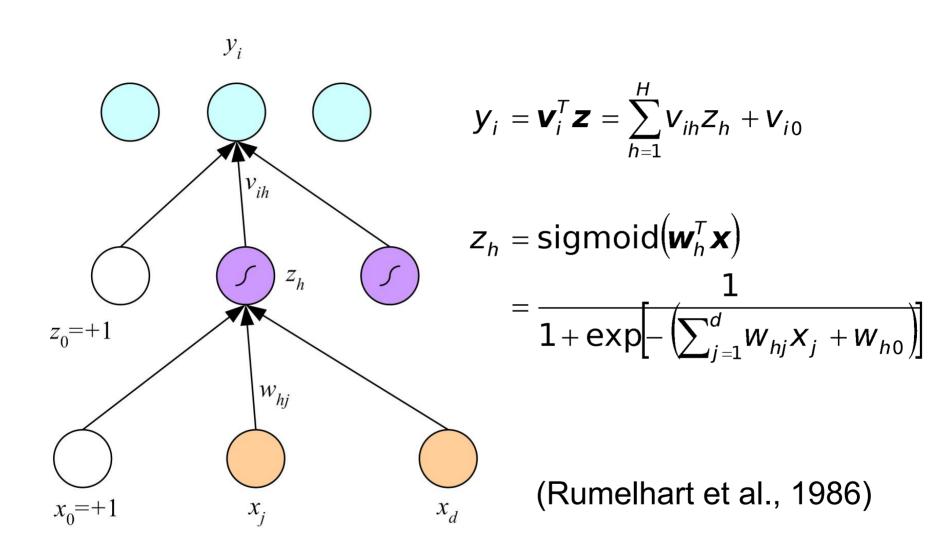
(Minsky and Papert, 1969)

Multilayer perceptrons

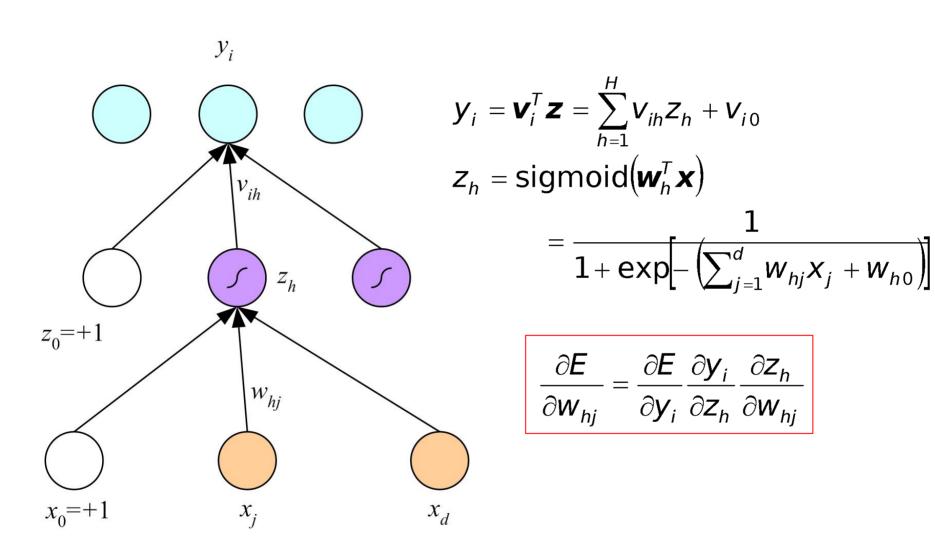
 $x_1 \text{ XOR } x_2 = (x_1 \text{ AND } \sim x_2) \text{ OR } (\sim x_1 \text{ AND } x_2)$



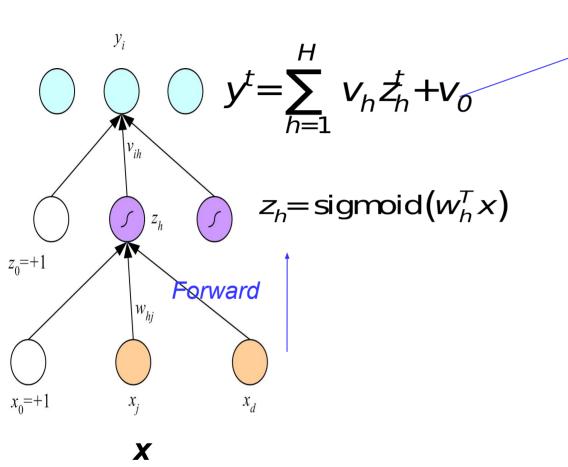
Multilayer Perceptrons



Training multilayer perceptron



Training multilayer perceptron



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

$$\Delta v_{h} = \sum_{t} (r^{t} - y^{t}) z_{h}^{t}$$

$$Backward$$

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

$$-\eta \sum_{t} \frac{\partial E}{\partial y^{t}} \frac{\partial y^{t}}{\partial z_{h}^{t}} \frac{\partial z_{h}^{t}}{\partial w_{hj}}$$

$$-\eta \sum_{t} -(r^{t} - y^{t}) v_{h} z_{h}^{t} (1 - z_{h}^{t}) x_{j}^{t}$$

$$\eta \sum_{t} (r^{t} - y^{t}) v_{h} z_{h}^{t} (1 - z_{h}^{t}) x_{j}^{t}$$

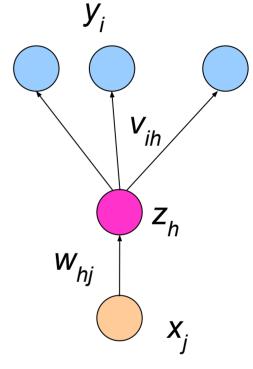
Regression with Multiple Outputs

$$E(\mathbf{W}, \mathbf{V} \mid \mathbf{X}) = \frac{1}{2} \sum_{t} \sum_{i} (r_{i}^{t} - y_{i}^{t})^{2}$$

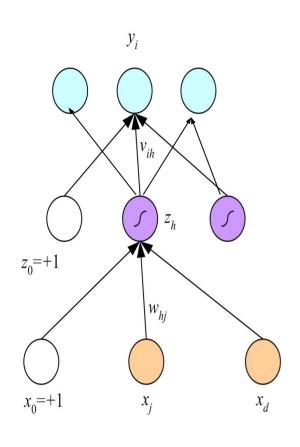
$$y_{i}^{t} = \sum_{h=1}^{H} V_{ih} Z_{h}^{t} + V_{i0}$$

$$\Delta V_{ih} = \eta \sum_{t} (r_{i}^{t} - y_{i}^{t}) Z_{h}^{t}$$

$$\Delta W_{hj} = \eta \sum_{t} \left[\sum_{i} (r_{i}^{t} - y_{i}^{t}) V_{ih} \right] Z_{h}^{t} (1 - Z_{h}^{t}) X_{j}^{t}$$

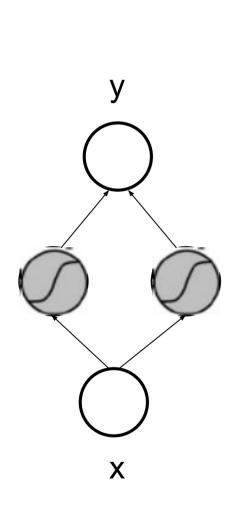


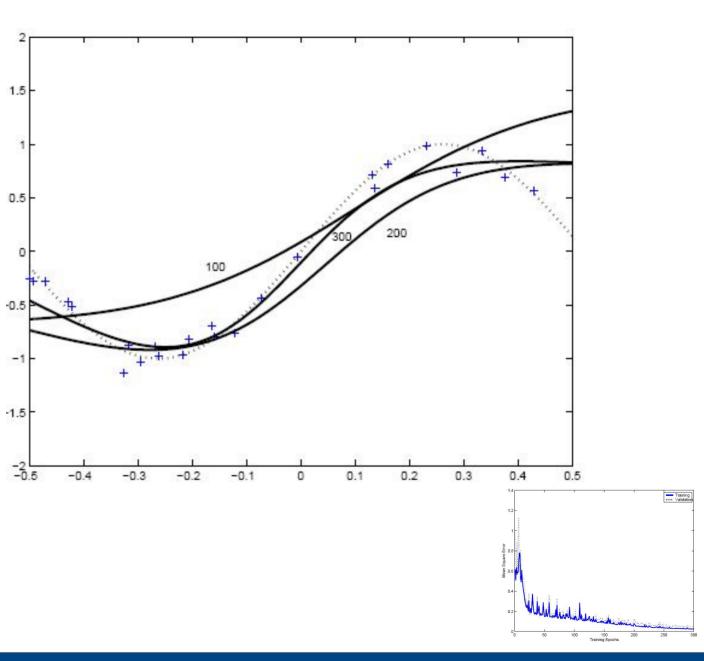
Algorithm for backpropagation



Initialize all v_{ih} and w_{hj} to rand(-0.01, 0.01)Repeat For all $(\boldsymbol{x}^t, r^t) \in \mathcal{X}$ in random order For $h = 1, \ldots, H$ $z_h \leftarrow \operatorname{sigmoid}(\boldsymbol{w}_h^T \boldsymbol{x}^t)$ For $i = 1, \ldots, K$ $y_i = \boldsymbol{v}_i^T \boldsymbol{z}$ For $i = 1, \ldots, K$ $\Delta \boldsymbol{v}_i = \eta(r_i^t - y_i^t)\boldsymbol{z}$ For $h = 1, \ldots, H$ $\Delta \boldsymbol{w}_h = \eta \left(\sum_i (r_i^t - y_i^t) v_{ih} \right) z_h (1 - z_h) \boldsymbol{x}^t$ For $i = 1, \ldots, K$ $\boldsymbol{v}_i \leftarrow \boldsymbol{v}_i + \Delta \boldsymbol{v}_i$ For $h = 1, \ldots, H$ $\boldsymbol{w}_h \leftarrow \boldsymbol{w}_h + \Delta \boldsymbol{w}_h$ Until convergence

Sample training data, $y^t = \sin(6x) + N(0,0.1)$





Sample training data, $y^t = \sin(6x) + N(0,0.1)$

$$y^{t} = \sum_{h=1}^{H} V_{h} Z_{h}^{t} + V_{0}$$

$$y$$

$$Z_{h} = sigmoid(\mathbf{w}_{h}^{T} \mathbf{x})$$

$$\mathbf{w}_{h}^{T} \mathbf{x}$$

$$\mathbf{x}$$

0.5

-0.5

Two-Class Discrimination

One sigmoid output

$$y^{t} = \operatorname{sigmoid}\left(\sum_{h=1}^{H} v_{h} z_{h}^{t} + v_{0}\right)$$

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

$$\Delta v_{h} = \eta \sum_{t} (r^{t} - y^{t}) z_{h}^{t}$$

$$\Delta w_{hj} = \eta \sum_{t} (r^{t} - y^{t}) v_{h} z_{h}^{t} (1 - z_{h}^{t}) x_{j}^{t}$$

Parametric Discrimination Revisited

In chapter 5, we saw that if the class densities, $p(x|C_i)$, are Gaussian and share a common covariance matrix, the discriminant function is linear

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where the parameters can be analytically calculated

Let us again see the special case where there are two classes: We define $y \equiv P(C_1|x)$ and $p(C_2|x) = 1 - y$. Then in classification, we

choose
$$C_1$$
 if $\begin{cases} y > 0.5 \\ \frac{y}{1-y} > 1 \\ \log \frac{y}{1-y} > 0 \end{cases}$ and C_2 otherwise

 $\log y/(1-y)$ is known as the *logit* transformation or *log odds* of y....

$$logit(P(C_1|x)) = log \frac{P(C_1|x)}{1 - P(C_1|x)} = w^T x + w_0$$

The inverse of logit

$$\log \frac{P(C_1|\mathbf{x})}{1 - P(C_1|\mathbf{x})} = \mathbf{w}^T \mathbf{x} + w_0$$

is the logistic function, also called the sigmoid function (see figure 10.5):

$$P(C_1|\mathbf{x}) = \text{sigmoid}(\mathbf{w}^T\mathbf{x} + \mathbf{w}_0) = \frac{1}{1 + \exp[-(\mathbf{w}^T\mathbf{x} + \mathbf{w}_0)]}$$

Logistic Discrimination / Logistic Regression

In *logistic discrimination*, we do not model the class-conditional densities, $p(x|C_i)$, but rather their ratio. Let us again start with two classes and assume that the log likelihood ratio is linear:

$$\log \frac{p(\boldsymbol{x}|C_1)}{p(\boldsymbol{x}|C_2)} = \boldsymbol{w}^T \boldsymbol{x} + w_0^o$$

Rearranging terms, we get the sigmoid function

$$y = \hat{P}(C_1|\mathbf{x}) = \frac{1}{1 + \exp[-(\mathbf{w}^T\mathbf{x} + \mathbf{w}_0)]}$$

We assume r^t , given x^t , is Bernoulli with probability $y^t \equiv P(C_1|x^t)$ as calculated in equation 10.21:

$$r^t | \mathbf{x}^t \sim \text{Bernoulli}(y^t)$$

 $l(\mathbf{w}, w_0 | \mathcal{X}) = \prod_t (y^t)^{(r^t)} (1 - y^t)^{(1 - r^t)}$

$$E(\boldsymbol{w}, w_o | \mathcal{X}) = -\sum_t r^t \log y^t + (1 - r^t) \log (1 - y^t)$$

$$\Delta w_{j} = -\eta \frac{\partial E}{\partial w_{j}} = \eta \sum_{t} \left(\frac{r^{t}}{y^{t}} - \frac{1 - r^{t}}{1 - y^{t}} \right) y^{t} (1 - y^{t}) x_{j}^{t}$$

$$= \eta \sum_{t} (r^{t} - y^{t}) x_{j}^{t}, j = 1, \dots, d$$

$$\Delta w_{0} = -\eta \frac{\partial E}{\partial w_{0}} = \eta \sum_{t} (r^{t} - y^{t})$$

Improving Convergence: Momentum

- Successive Δw values may be so different that large oscillations may occur and slow convergence
- Take a running average by incorporating the previous update

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

t is time index (epoch number or iteration number

 get an effect of averaging and smooth the trajectory during convergence.

Improving Convergence: Adaptive Learning Rate

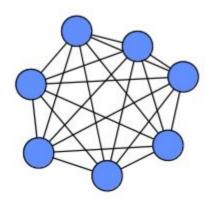
- adaptive for faster convergence, where it is kept large when learning takes place and is decreased when learning slows down
- Increase learning rate by a constant amount if the error on the training set decreases and decrease it geometrically if it increases

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

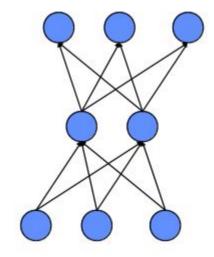
Practice

Scale inputs and outputs

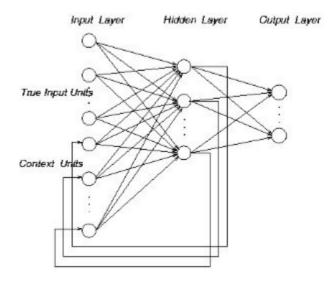
Topologies of Neural Networks



completely connected



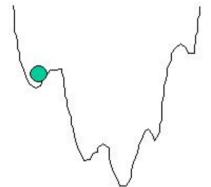
feedforward (directed, a-cyclic)



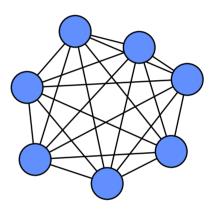
recurrent (feedback connections)

Hopfield Networks

- Act as "autoassociative" memories to store patterns
- Network converges to local minima which store different patterns.

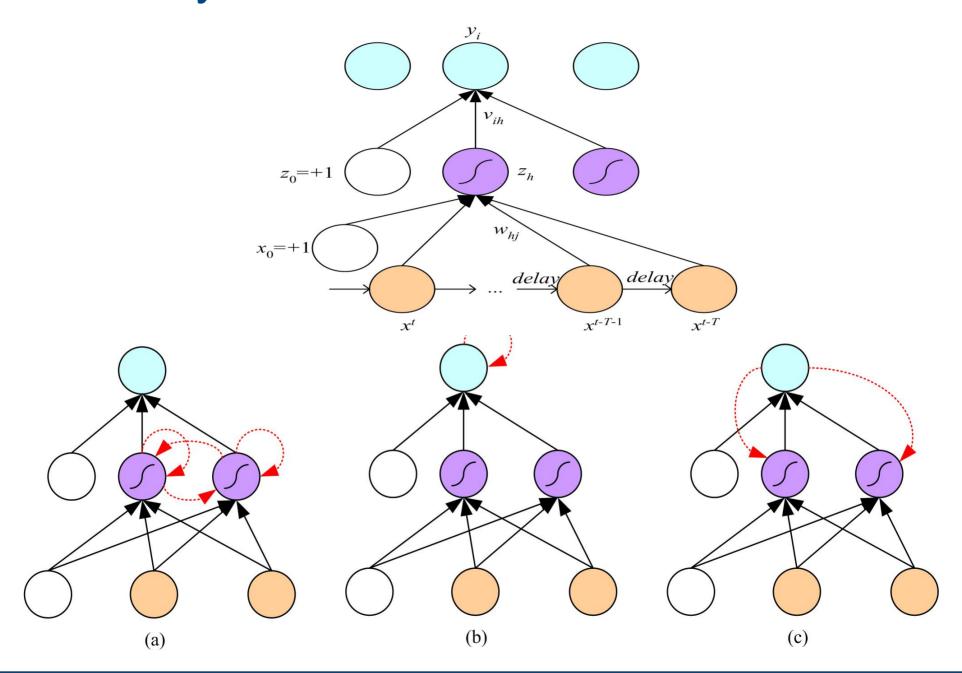






completely connected

Time-delay Neural Networks



Recurrent Networks

