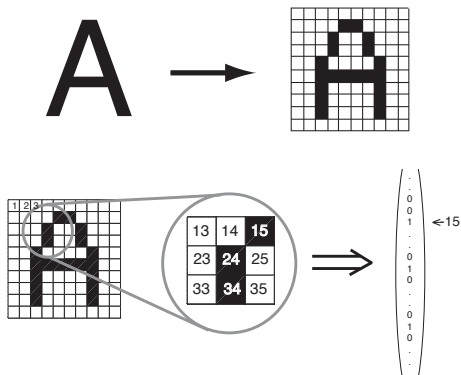# Neuroinformatics

April 5, 2012

Lecture 8: Feed Forward Networks

# Digital representation of a letter



**Optical character recognition**: Predict meaning from features.
E.g., given features **x**, what is the character **y**

$$f : \mathbf{x} \in \mathbf{S}_1^n \to \mathbf{y} \in \mathbf{S}_2^m$$

# Examples given by lookup table

Boolean AND function

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 0 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

Look-up table for a non-boolean example function

| $x_1$ | $x_2$ | $y$ |
|-------|-------|-----|
| 1 | 2 | -1 |
| 2 | 1 | 1 |
| 3 | -2 | 5 |
| -1 | -1 | 7 |
| ... | ... | ... |

# The population node as perceptron

**Update rule:** $\mathbf{r}^{\text{out}} = g(\mathbf{w}\mathbf{r}^{\text{in}})$ (component-wise: $r_i^{\text{out}} = g(\sum_j w_{ij} r_j^{\text{in}})$)
For example: $r_i^{\text{in}} = x_i$, $\tilde{y} = r^{\text{out}}$, linear grain function $g(x) = x$:
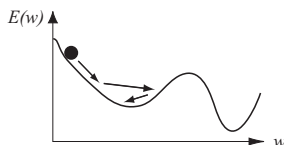
$$\tilde{y} = w_1 x_1 + w_2 x_2$$

# How to find the right weight values?

**Objective (error) function**, for example: mean square error (MSE)

$$E = \frac{1}{2} \sum_i (r_i^{\text{out}} - y_i)^2$$

**Gradient descent** method:
$$w_{ij} \leftarrow w_{ij} - \epsilon \frac{\partial E}{\partial w_{ij}}$$
$$= w_{ij} - \epsilon (y_i - r_i^{\text{out}}) r_j^{\text{in}} \qquad \text{for MSE, linear gain}$$



Initialize weights arbitrarily
Repeat until error is sufficiently small
    Apply a sample pattern to the input nodes: $r_i^0 = r_i^{\text{in}} = \xi_i^{\text{in}}$
    Calculate rate of the output nodes: $r_i^{\text{out}} = g(\sum_j w_{ij} r_j^{\text{in}})$
    Compute the delta term for the output layer: $\delta_i = g'(h_i^{\text{out}})(\xi_i^{\text{out}} - r_i^{\text{out}})$
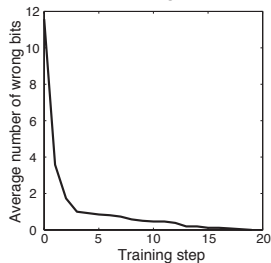    Update the weight matrix by adding the term: $\Delta w_{ij} = \epsilon \delta_i r_j^{\text{in}}$
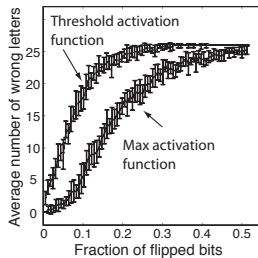
# Example: OCR



A. Training pattern

```
>> displayLetter(1)
      +++
      +++
     +++++
     ++ ++
    ++   ++
   +++   +++
   +++++++++
  +++++++++++
  +++       +++
  +++       +++
  +++       +++
  +++       +++
```
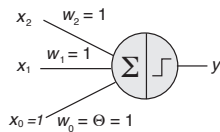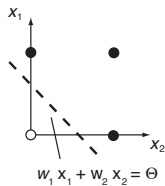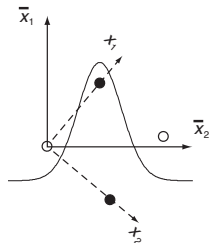
B. Learning curve

C. Generalization ability

Threshold activation function

Max activation function

# Example: Boolean function

## A. Boolean OR function

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

$w_1 x_1 + w_2 x_2 = \Theta$

$x_2 \quad w_2 = 1$

$x_1 \quad w_1 = 1$

$\Sigma \quad \quad y$

$x_0 = 1 \quad w_0 = \Theta = 1$

## B. Boolean XOR function

| $x_1$ | $x_2$ | $y$ |
|---|---|---|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

?

$\overline{x}_1$

$+_1$

$\overline{x}_2$

$+_2$

## perceptronTrain.m

```
1   %% Letter recognition with threshold perceptron
2    clear; clf;
3    nIn=12*13; nOut=26;
4    wOut=rand(nOut,nIn)-0.5;
5
6   % training vectors
7    load pattern1;
8    rIn=reshape(pattern1', nIn, 26);
9    rDes=diag(ones(1,26));
10
11  % Updating and training network
12   for training_step=1:20;
13       % test all pattern
14        rOut=(wOut*rIn)>0.5;
15        distH=sum(sum((rDes-rOut).^2))/26;
16        error(training_step)=distH;
17       % training with delta rule
18        wOut=wOut+0.1*(rDes-rOut)*rIn';
19   end
20
21   plot(0:19,error)
22   xlabel('Training step')
23   ylabel('Average Hamming distance')
```

# Percepton as Linear Classifier: ML approach

- Assume a binary classification problem, i.e. $S = \{s_1, s_2\}$.
- One discriminant function $g(\vec{x})$ enough: classify
  $$y = \begin{cases} s_1, & \text{if } g(\vec{x}) > 0; \\ s_2, & \text{otherwise.} \end{cases}$$
- we will estimate $\vec{b}, c$ directly from the given sample
  $D = \{(\vec{x}_1, y_1), (\vec{x}_2, y_2) \dots (\vec{x}_m, y_m)\}$.
- We want $\left( \vec{b}^t \vec{x}_i + c \right) > 0$ if $y_i = s_1$ and $\left( \vec{b}^t \vec{x}_i + c \right) < 0$ otherwise.
- Same as requesting $\left( \vec{b}^t \vec{z}_i + c \right) > 0$ for all $z_i$, where $z_i = x_i$ if
  $y_i = s_1$ and $z_i = -x_i$ otherwise.
- Let formally $z_i^{n+1} = 1 \; \forall i$ and $\vec{w} = [\vec{b}, c]$ (add $c$ as the last
  component of $\vec{w}$).
- Thus we can write simply $g(\vec{z}) = \vec{w}^t \vec{z}$ and request $\vec{w}^t \vec{z}_i > 0$ for all
  $z_i$.
- Let
  $$E(\vec{w}) = \sum_{\vec{z}_i \in M} -\vec{w}^t \vec{z}_i$$

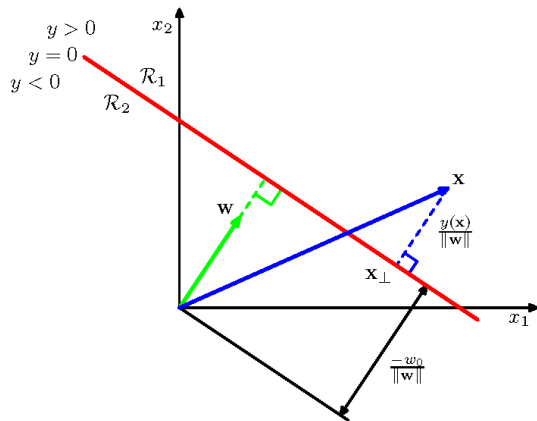  where $M$ is the set $\vec{z}_i$ that are misclassified.

# Percepton ML view

- $E(\vec{b}, c)$ is always non-negative.
- If $E(\vec{w}) = 0$ then all examples in $D$ are correctly classified and $D$ is linearly separable. We want to find the minimum of $E(\vec{w})$.
- $E(\vec{w})$ is piece-wise linear. A gradient algorithm can be used to search a minimum.
- Gradient algorithm: go towards a minimum by making discrete steps in $\Re^{n+1}$ in the direction opposite to the gradient of $E(\vec{w})$.

$$\nabla(E(\vec{w})) = \left( \frac{\partial E(\vec{w})}{\partial w_1}, \frac{\partial E(\vec{w})}{\partial w_2}, \ldots \frac{\partial E(\vec{w})}{\partial w_{n+1}} \right) = \sum_{z_i \in M} -\vec{z}$$
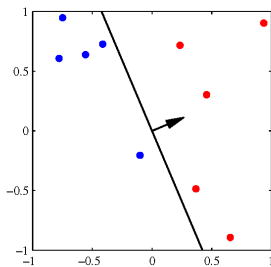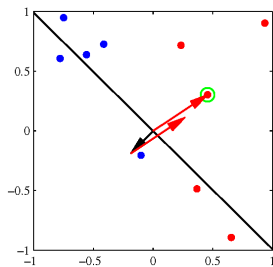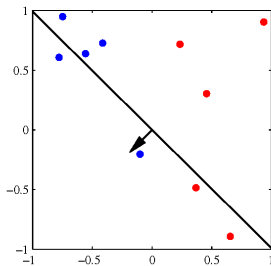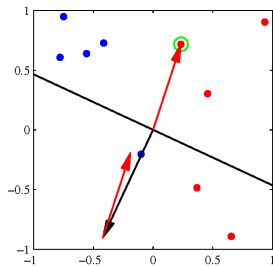
- The perceptron gradient algorithm:
  1. $k = 0$. Choose a random $\vec{w}$.
  2. $k \leftarrow k + 1$
  3. $\vec{w} \leftarrow \vec{w} + \eta(k) \sum_{z_i \in M_k} \vec{z}$
  4. if $|\nu(k) \sum_{z_i \in M_k} \vec{z}| > \theta$ go to 2
  5. return $\vec{w}$
- $\eta$ - the learning rate, $\theta$ - an error threshold.

# Percepton graphical representation

- $y(\vec{x}) = \vec{w}^t \vec{x} + w_0$, $y(\vec{x_a}) = y(\vec{x_b})$
- $\vec{x_a}$ a $\vec{x_b}$ is on decision surface, hence $\vec{w}^t(\vec{x_a} - \vec{x_b}) = 0$)
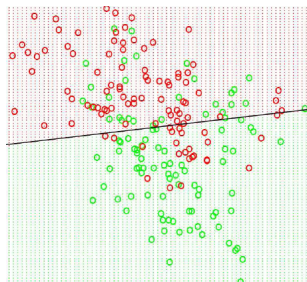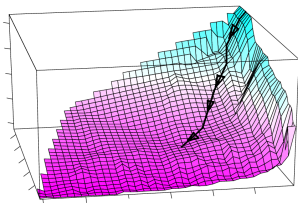- w is orthonomal to decision surface, $\omega_0(b)$ is translation [Bishop]

# Percepton learning

# Percepton - linear separability

▶ If the two classes are linearly separable, the perceptron algorithm will terminate in a finite number of steps with zero training error.

▶ A problem that is linearly non-separable in $\Re^n$ may be separable after being transformed to $\Re^{n'}$ $n' > n$. For example, new coordinates may contain all quadratic terms:

$$[x(1), \ldots x(n), x^2(1), x(1)x(2), x(1)x(3), \ldots x^2(n)]$$

▶ A linear separation method such as the perceptron may be applied in the extended space, generating nonlinear separation in the original space.
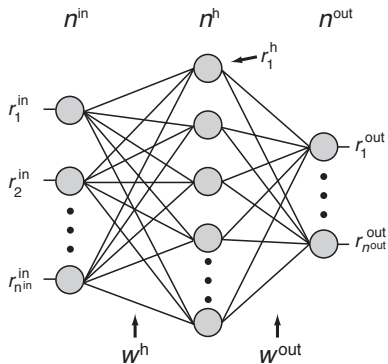
# Percepton - history

- ► Frank Rosenblatt - HW realization of percepton in 1958



- ► Learning of simple symbols and alphabet - inspiration by brain nets
- ► Character was illuminated by powerful lights, image focused onto 20 x 20 array of cadmium sulphide photocells giving 400 pixel image
- ► Patch board - different configuration of input features
- ► Rack of adaptive weights, each weight rotary variable resistor driven by electric motor - weights were adjusted automatically by the elarning algorithm
- ► MARK 1 computer (Harvard - IBM): 765000 parts, 16 m long, 2.4 m height, 2 m wide, 3 operation per second, multiplication took 6 sec

# The multilayer Perceptron (MLP)



Update rule: $\mathbf{r}^{\text{out}} = g^{\text{out}}(\mathbf{w}^{\text{out}} g^{\text{h}}(\mathbf{w}^{\text{h}} \mathbf{r}^{\text{in}}))$

Learning rule (error backpropagation): $w_{ij} \leftarrow w_{ij} - \epsilon \frac{\partial E}{\partial w_{ij}}$

# The error-backpropagation algorithm

Initialize weights arbitrarily

Repeat until error is sufficiently small

> Apply a sample pattern to the input nodes: $r_i^0 := r_i^{\text{in}} = \xi_i^{\text{in}}$
>
> Propagate input through the network by calculating the rates of nodes in successive layers $l$: $r_i^l = g(h_i^l) = g(\sum_j w_{ij}^l r_j^{l-1})$
>
> Compute the delta term for the output layer: $\delta_i^{\text{out}} = g'(h_i^{\text{out}})(\xi_i^{\text{out}} - r_i^{\text{out}})$
>
> Back-propagate delta terms through the network: $\delta_i^{l-1} = g'(h_i^{l-1}) \sum_j w_{ji}^l \delta_j^l$
>
> Update weight matrix by adding the term: $\Delta w_{ij}^l = \epsilon \delta_i^l r_j^{l-1}$

# MLP as universal approximator

- Hidden layer enables realization of complicated non-linear fces
- Each neuron can have its own activation fce
- We suppose that we have only ONE type of activation fce
- **QUESTION: Can 3-forward layer approximate any non-linear function?**
- **ANSWER: YES- thanks to A.Kolmogorov Any continuous fce can be implemented by 3-layes net under assumption of sufficient number of $n_H$ hidden neurons, suitable non-linearities and weights $w$.**
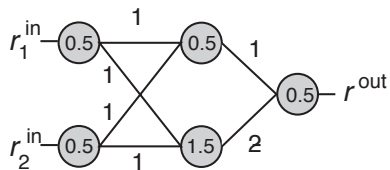
# Andrej Kolmogorov

- ▶ He constructed p̈erpetuum mobileïn high school, his teacher could not discover the trick
- ▶ First he studied history in Moscow university
- ▶ He published the first scientific work on realities in Novgorod area during 15. a 16. centurary
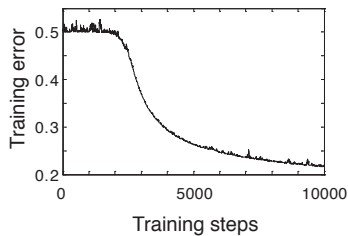- ▶ The biggest contribution in probability field

## mlp.m

```
1    %% MLP with backpropagation learning on XOR problem
2    clear; clf;
3    N_i=2; N_h=2; N_o=1;
4    w_h=rand(N_h,N_i)-0.5; w_o=rand(N_o,N_h)-0.5;
5
6    % training vectors (XOR)
7    r_i=[0 1 0 1 ; 0 0 1 1];
8    r_d=[0 1 1 0];
9
10   % Updating and training network with sigmoid activation function
11   for sweep=1:10000;
12     % training randomly on one pattern
13       i=ceil(4*rand);
14       r_h=1./(1+exp(-w_h*r_i(:,i)));
15       r_o=1./(1+exp(-w_o*r_h));
16       d_o=(r_o.*(1-r_o)).*(r_d(:,i)-r_o);
17       d_h=(r_h.*(1-r_h)).*(w_o'*d_o);
18       w_o=w_o+0.7*(r_h*d_o')';
19       w_h=w_h+0.7*(r_i(:,i)*d_h')';
20     % test all pattern
21       r_o_test=1./(1+exp(-w_o*(1./(1+exp(-w_h*r_i)))));
22       d(sweep)=0.5*sum((r_o_test-r_d).^2);
23   end
24   plot(d)
```
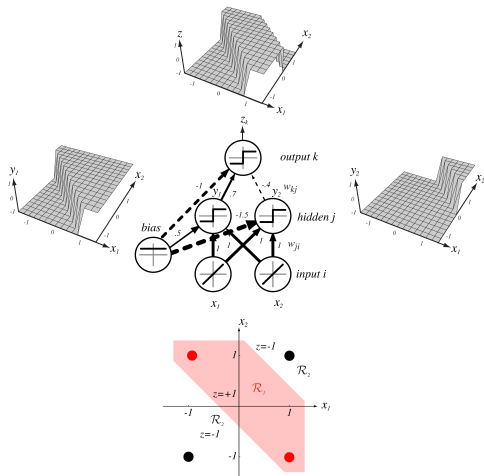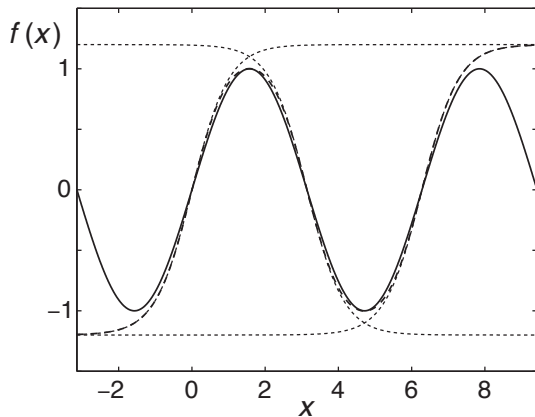
# MLP for XOR function



Learning curve for XOR problem

# Example 3-layer neural net - XOR problem

- $0 \bigoplus 0 = 0,\ 1 \bigoplus 1 = 0,\ 1 \bigoplus 0 = 1,\ 0 \bigoplus 1 = 1$
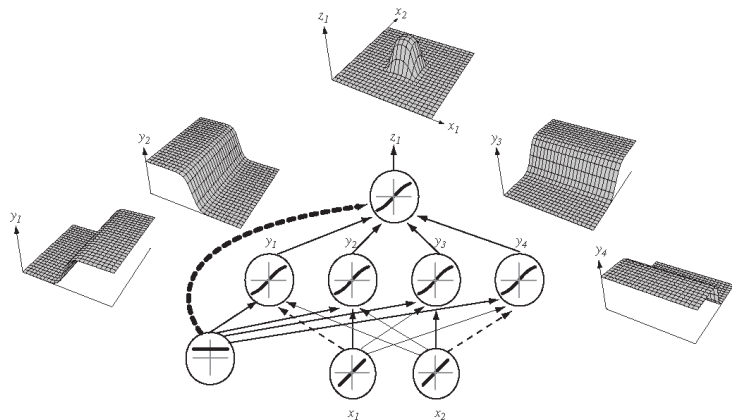- $-1 \bigoplus -1 = -1,\ 1 \bigoplus 1 = -1,\ 1 \bigoplus -1 = 1,\ -1 \bigoplus 1 = 1$
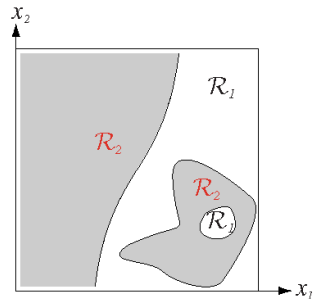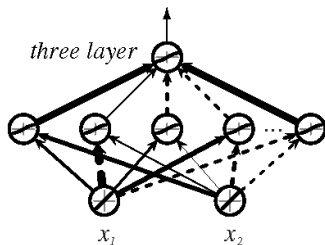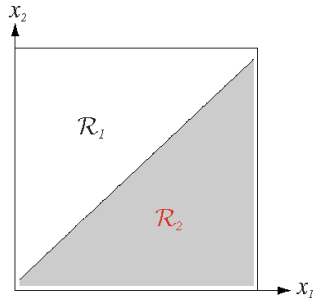
# MLP approximating sine function

# Non-linear fce approximation

Fourier transform ANALOGY

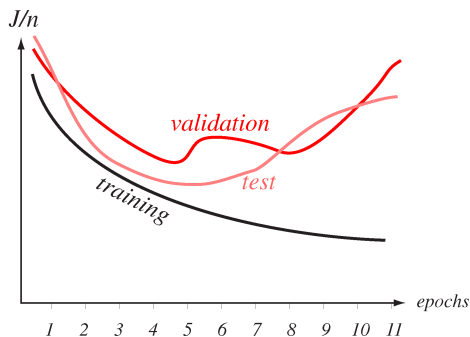*two layer*

$\mathcal{R}_1$

$\mathcal{R}_2$

*three layer*

$\mathcal{R}_1$

$\mathcal{R}_2$

$\mathcal{R}_2$

$\mathcal{R}_1$

# Validation

- error of training set in monotonic-decreasing fce because of gradient algorithm optimization
- we divide data to training and validation set We use validation as stopping criteria (e.g. the first minimum)
- DEMO - Neural Network Toolbox Matlab
  http://www.mathworks.com/products/neuralnet/
- Data are from UCI Machine Learning Repository
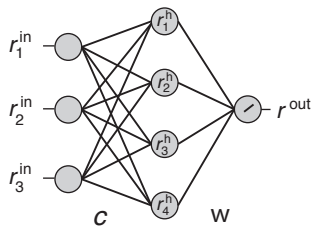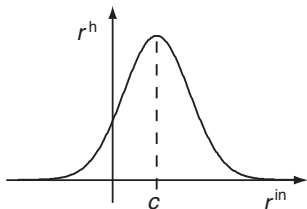  http://mlearn.ics.uci.edu/MLRepository.html

# MLP biological plausibility

1. universal approximator $\rightarrow$ small number of hidden neurons $\rightarrow$ smooth solution & big number of hidden layers in biological systems
2. problematic training with error-back propagation, some exchange between postsynaptic and presynaptic neurons is possible, however
3. inclusion of derivative terms??
4. non-locality of the algorithm, neuron must gather the back-propagated errors from all other nodes to which it projects

# Kernel machine

- ▶ better recognition after transformation of feature space $x_1 x_2, x_i^2$, $x \to \Theta(x)$, $w \to \Theta(w)$
- ▶ the net input of node $h = \sum_i (w_i r_i) = wr$, node in the network, $h = \Theta(w)\Theta(r) = K(w, r)$
- ▶ K is kernel function, special case is Gaussian kernel function $K(w, x) = \frac{(w-x)^2}{2\delta^2}$, FITS tuning curve
- ▶ Radial basis networks

# Advance learning

- shallow part of errro function, very slow convergence, using momentum term

$$\Delta w_{ij}(t+1) = \eta \frac{\partial E}{w_{ij}} + \alpha \Delta w_{ij}(t)$$

- Acceleration of learning process, other fce than MSE: entropic error function

$$E = \frac{1}{2}\sum_{\mu,i}[(1+y_i^{\mu})\log\frac{1+y_i^{\mu}}{1+r_i^{out}} + (1-y_i^{\mu})\log\frac{1-y_i^{\mu}}{1-r_i^{out}}]$$

- measure information content of the output, even less computation of delta term: $g(x) = tanh(x)\delta_i = y)i - r^{out}$

- more sophisticated training using higher-order gradients: in MATLAB Levenberg-Marquardt. The relation of such sophistivate technique to biological learning is,so far, unclear!

- random search $\rightarrow$ stochastic processes, stochastic annealing, genetic algorithms