# Neural Networks

1

# Outline

- The Brain
- Perceptrons
- Gradient descent
- Multi-layer networks
- Backpropagation

# Artificial Neural Networks

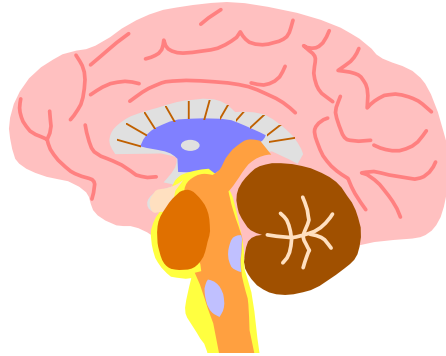- ## Other terms/names
    - connectionist
    - parallel distributed processing
    - neural computation
    - adaptive networks..

- ## History
    - 1943-McCulloch & Pitts are generally recognised as the designers of the first neural network
    - 1949-First learning rule
    - 1969-Minsky & Papert - perceptron limitation - Death of ANN
    - 1980's - Re-emergence of ANN - multi-layer networks

3

# The biological inspiration



- The brain has been extensively studied by scientists.
- Vast complexity prevents all but rudimentary understanding.
- Even the behaviour of an individual neuron is extremely complex
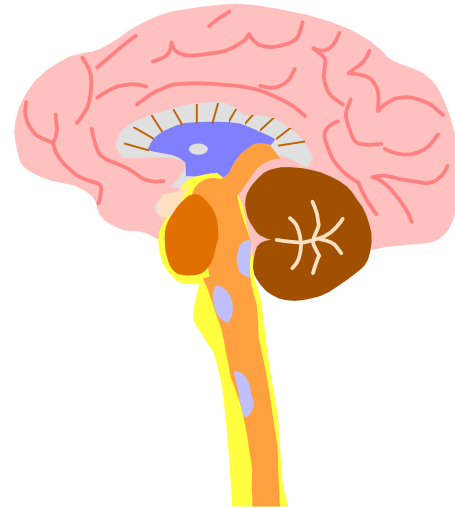
# Features of the Brain

- Ten billion ($10^{10}$) neurons
- Neuron switching time $>10^{-3}$secs
- Face Recognition ~0.1secs
- On average, each neuron has several thousand connections
- Hundreds of operations per second
- High degree of parallel computation
- Distributed representations
- Die off frequently (never replaced)
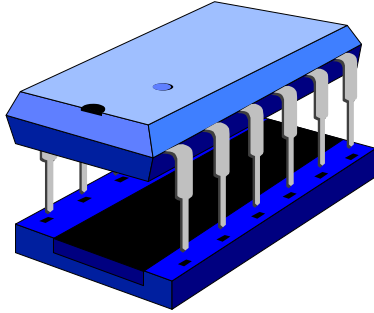- Compensated for problems by massive parallelism

# Brain and Machine

- The Brain
  - Pattern Recognition
  - Association
  - Complexity
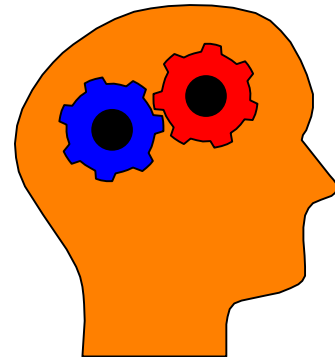  - Noise Tolerance

- The Machine
  - Calculation
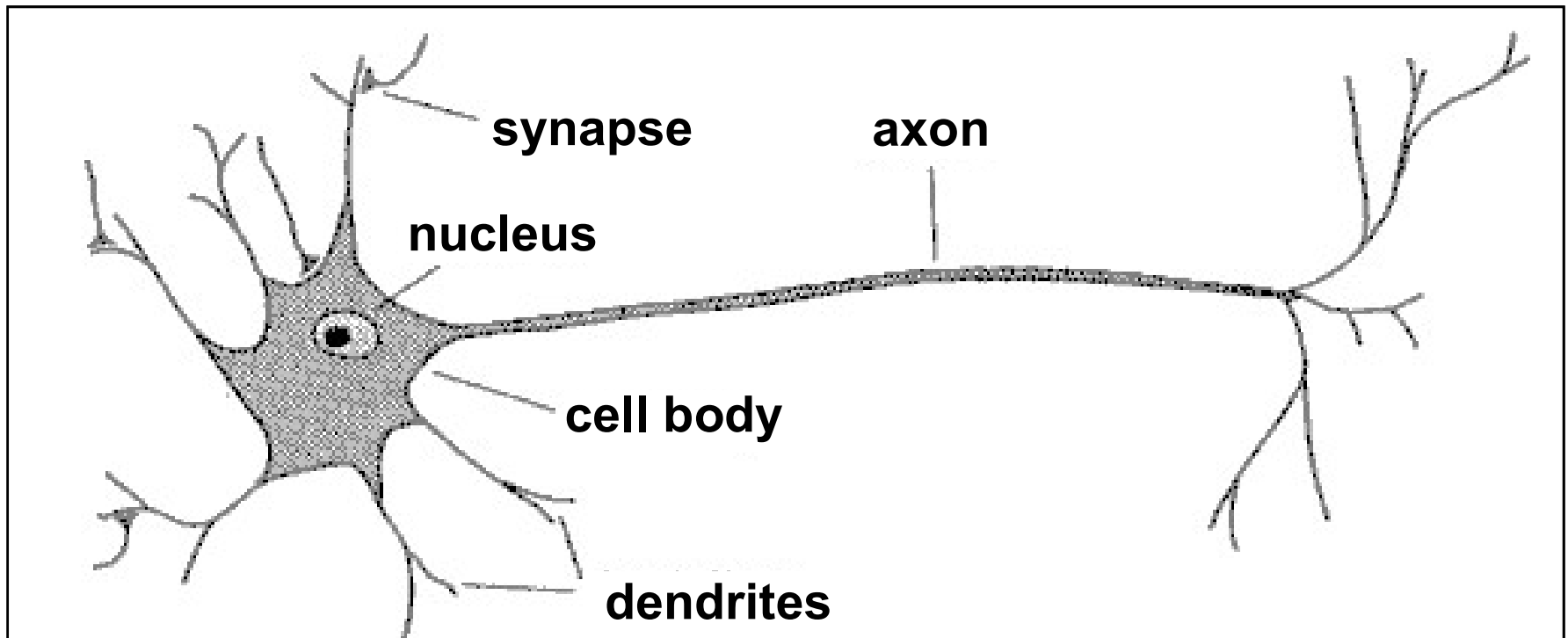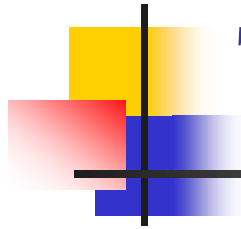  - Precision
  - Logic

# The contrast in architecture

- The Von Neumann architecture uses a single processing unit;
  - Tens of millions of operations per second
  - Absolute arithmetic precision

- The brain uses many slow unreliable processors acting in parallel

# The Structure of Neurons



synapse

axon

nucleus

cell body

dendrites

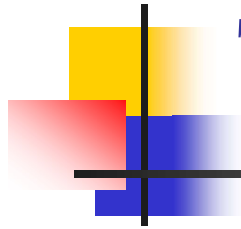# The Structure of Neurons

- A neuron only fires if its input signal exceeds a certain amount (the <span style="color:red">threshold</span>) in a short time period.

- Synapses vary in strength
  - Good connections allowing a large signal
  - Slight connections allow only a weak signal.
  - Synapses can be either <span style="color:blue">excitatory</span> or <span style="color:red">inhibitory</span>.

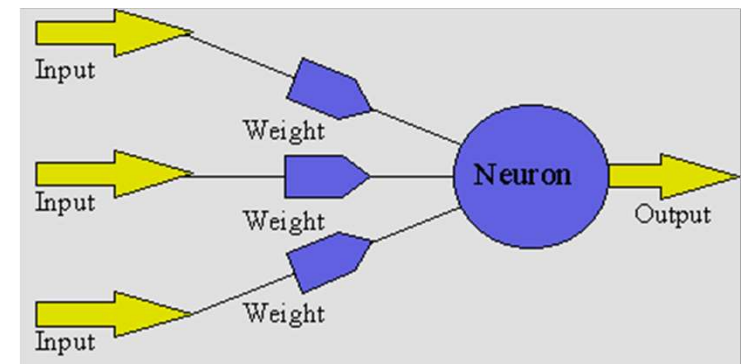# The Structure of Neurons

A neuron has a cell body, a branching **i**nput structure (the dendr**I**te) and a branching **o**utput structure (the ax**O**n)

- Axons connect to dendrites via synapses.
- Electro-chemical signals are propagated from the dendritic input, through the cell body, and down the axon to other neurons

# Properties of Artificial Neural Nets (ANNs)

# Properties of Artificial Neural Nets (ANNs)

- Many simple neuron-like threshold switching units

- Many weighted interconnections among units

- Highly parallel, distributed processing

- Learning by tuning the connection weights

# Appropriate Problem Domains for Neural Network Learning

- Input is high-dimensional discrete or real-valued (e.g. raw sensor input)
- Output is discrete or real valued
- Output is a vector of values
- Form of target function is unknown
- Humans do not need to interpret the results (black box model)

# Perceptron

- Linear treshold unit (LTU)

$x_{0=1}$

$x_1$   $w_1$

$w_0$

$x_2$   $w_2$

$\Sigma$

$w_n$

$x_n$

$\sum_{i=0}^{n} w_i \ x_i$

$o$

$$o(x_i)= \begin{cases} 1 \text{ if } \sum_{i=0}^{n} w_i \ x_i > 0 \\ -1 \text{ otherwise} \end{cases}$$

# Activation functions

- Transforms neuron's input into output.
- Features of activation functions:
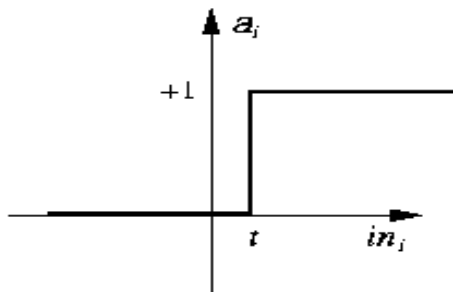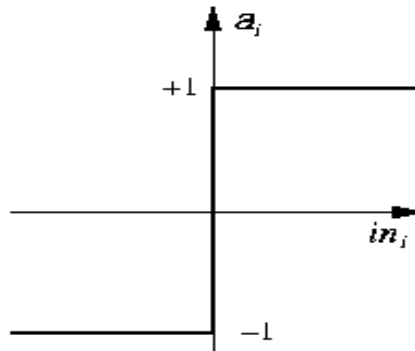  - A squashing effect is required
    - Prevents accelerating growth of activation levels through the network.



(a) Step function          (b) Sign function          (c) Sigmoid function

# Standard activation functions

- The hard-limiting threshold function
  - Corresponds to the biological paradigm
    - either fires or not

- Sigmoid functions ('S'-shaped curves)

  - The logistic function $\longrightarrow$  $\phi(x) = \dfrac{1}{1 + e^{-ax}}$

  - The hyperbolic tangent (symmetrical)
  - Both functions have a simple differential
  - Only the shape is important

# ARTIFICIAL NEURON

**Topics:** connection weights, bias, activation function

range determined
by $g(\cdot)$



bias $b$ only
changes the
position of
the riff

(from Pascal Vincent's slides)

17

# Perceptron Learning Rule

$w_i = w_i + \Delta w_i$

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

$t = c(x)$ is the target value

$o$ is the perceptron output

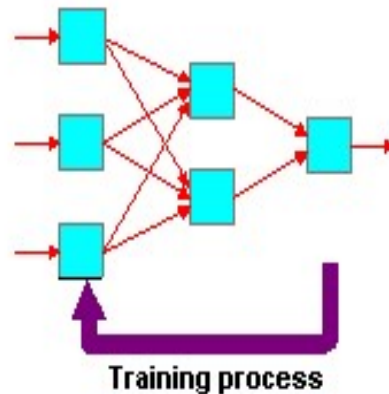$\eta$ Is a small constant (e.g. 0.1) called *learning rate*

- If the output is correct ($t = o$) the weights $w_i$ are not changed
- If the output is incorrect ($t \neq o$) the weights $w_i$ are changed such that the output of the perceptron for the new weights is *closer* to t.
- The algorithm converges to the correct classification
    - if the training data is linearly separable
    - and $\eta$ is sufficiently small

# Supervised Learning

- Training and test data sets
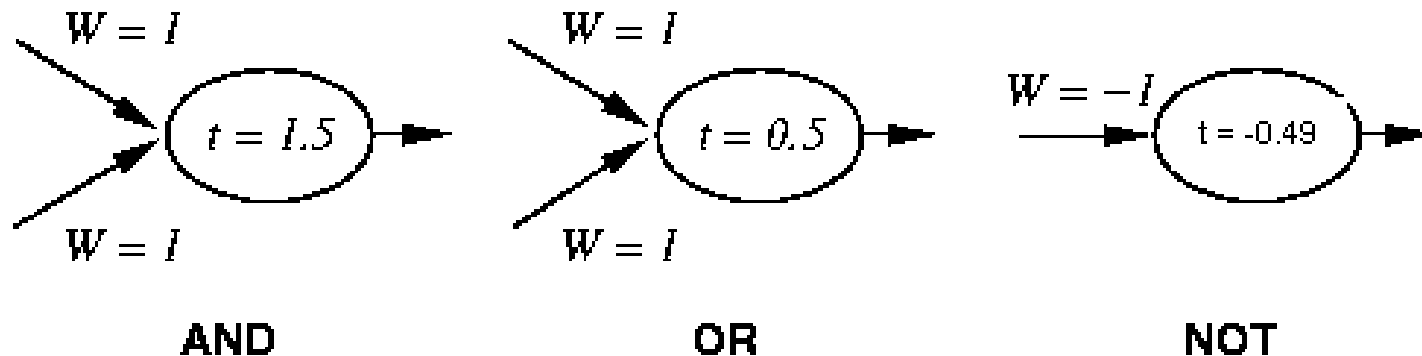- Training set; input & target



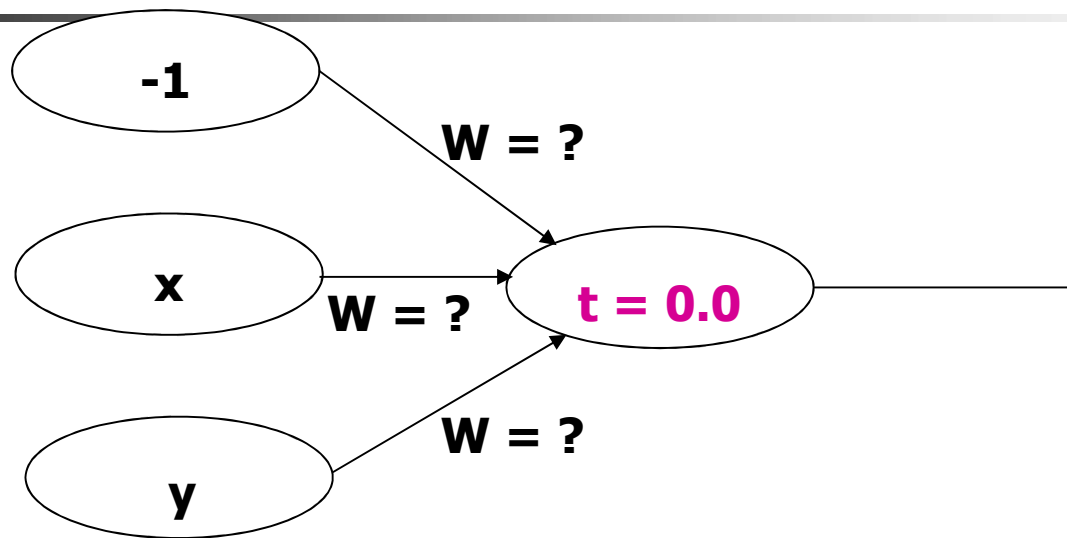| Sepal length | Sepal width | Petal length | Petal width | Class |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | 0 |
| 4.9 | 3.0 | 1.4 | 0.2 | 2 |
| 4.7 | 3.2 | 1.3 | 0.2 | 0 |
| 4.6 | 3.1 | 1.5 | 0.2 | 1 |

# Perceptron Training



$$\text{Output} = \begin{cases} 1 \text{ if } \sum_{i=0} w_i\, x_i > t \\ 0 \text{ otherwise} \end{cases}$$

- Linear threshold is used.
- W - weight value
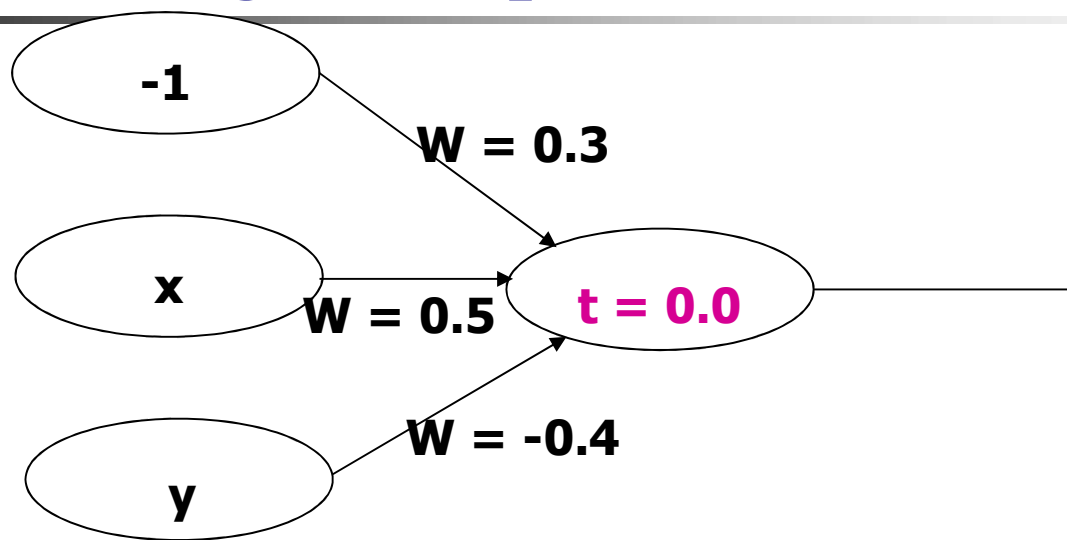- t - threshold value

# Training Perceptrons



**For AND**

| A B | Output |
|-----|--------|
| 0 0 | 0 |
| 0 1 | 0 |
| 1 0 | 0 |
| 1 1 | 1 |

•What are the weight values?

•Initialize with random weight values

# Training Perceptrons



**For AND**

| A B | Output |
|-----|--------|
| 0 0 | 0 |
| 0 1 | 0 |
| 1 0 | 0 |
| 1 1 | 1 |

| $I_1$ | $I_2$ | $I_3$ | Summation | Output |
|-------|-------|-------|-----------|--------|
| -1 | 0 | 0 | (-1*0.3) + (0*0.5) + (0*-0.4) = -0.3 | 0 |
| -1 | 0 | 1 | (-1*0.3) + (0*0.5) + (1*-0.4) = -0.7 | 0 |
| -1 | 1 | 0 | (-1*0.3) + (1*0.5) + (0*-0.4) = 0.2 | 1 |
| -1 | 1 | 1 | (-1*0.3) + (1*0.5) + (1*-0.4) = -0.2 | 0 |

# Simple network

$$\text{output} = \begin{cases} 1 \text{ if } \sum_{i=0} w_i\, x_i > t \\ 0 \text{ otherwise} \end{cases}$$

**For AND**

A B Output

0 0    0

0 1    0

1 0    0

1 1    1

-1

W = 1.5

X

W = 1

t = 0.0

Y

W = 1

# Learning algorithm

Epoch : Presentation of the entire training set to the neural network.
In the case of the AND function an epoch consists of four sets of inputs being presented to the network (i.e. [0,0], [0,1], [1,0], [1,1])

Error: The error value is the amount by which the value output by the network differs from the target value. For example, if we required the network to output 0 and it output a 1, then   Error = -1
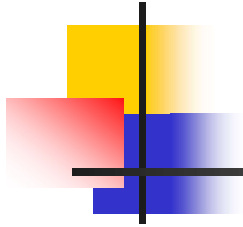
# Learning algorithm

**Target Value, T** : When we are training a network we not only present it with the input but also with a value that we require the network to produce. For example, if we present the network with [1,1] for the AND function the training value will be 1

**Output , O** : The output value from the neuron

$\underline{\mathbf{I_{i-}}}$ : Inputs being presented to the neuron
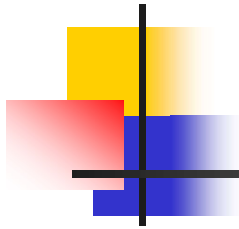
$\underline{\mathbf{W_i}}$ : Weight from input neuron ($I_j$) to the output neuron

**LR** : The learning rate. This dictates how quickly the network converges. It is set by a matter of experimentation. It is typically 0.1
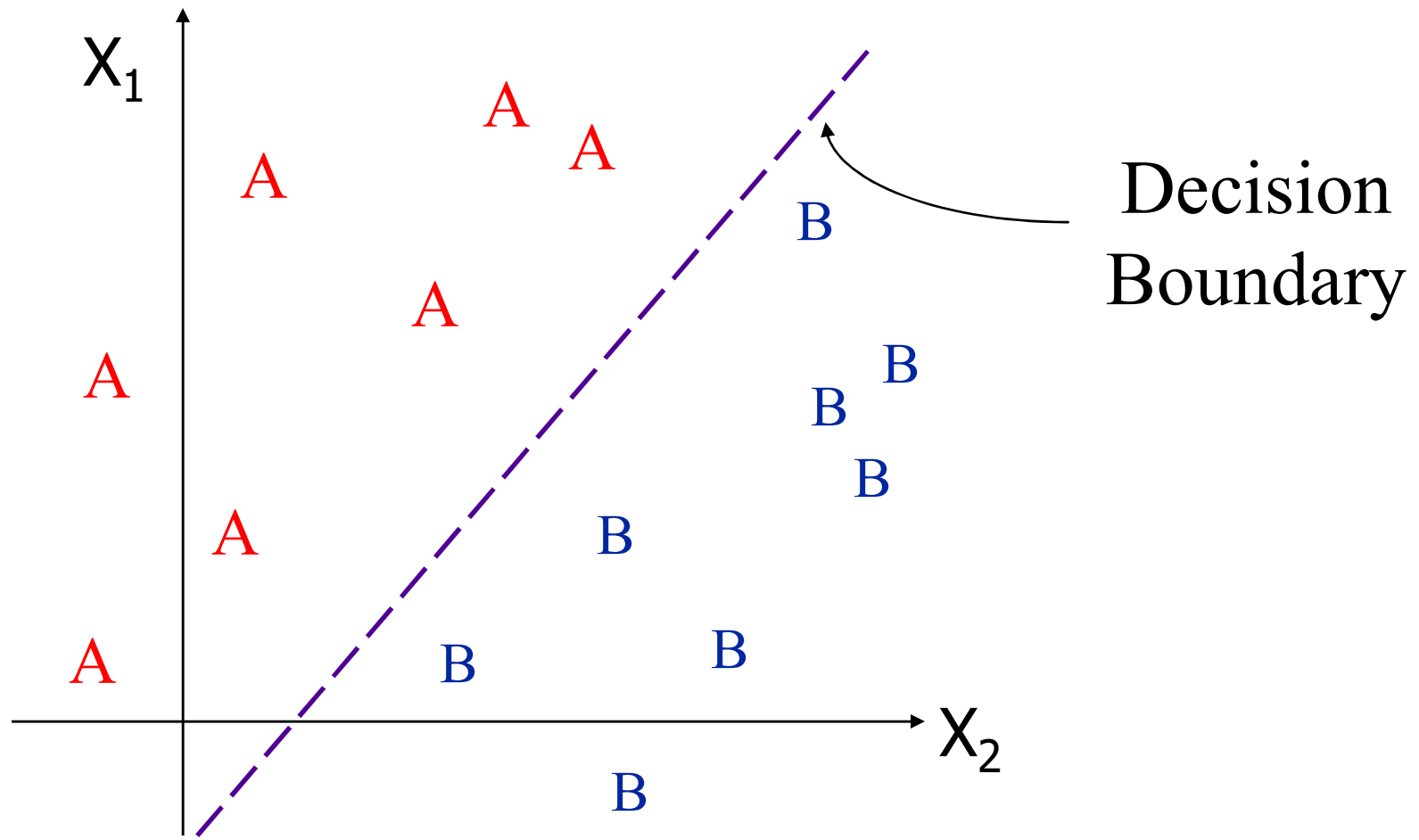
# Decision boundaries

- In simple cases, divide feature space by drawing a hyperplane across it.

- Known as a decision boundary.

- Discriminant function: returns different values on opposite sides. (straight line)

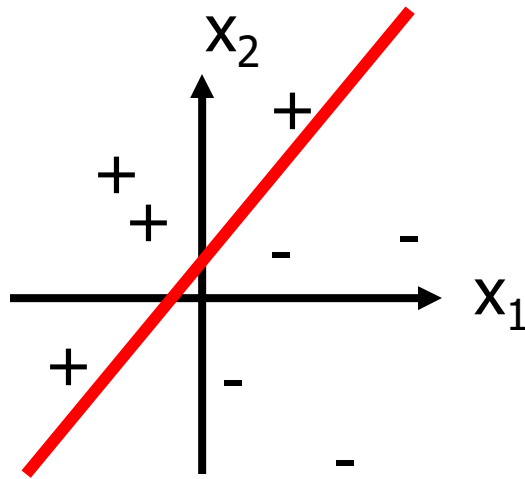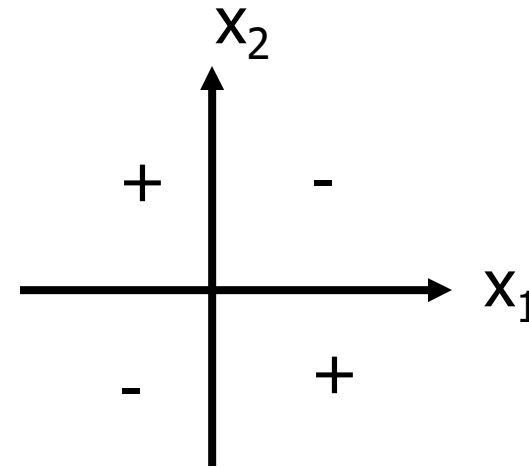- Problems which can be thus classified are linearly separable.

# Linear Separability

# Decision Surface of a Perceptron

$x_2$

+

+

+

-

-

$x_1$

+

-

-

**Linearly separable**

$x_2$

+

-

$x_1$

-

+
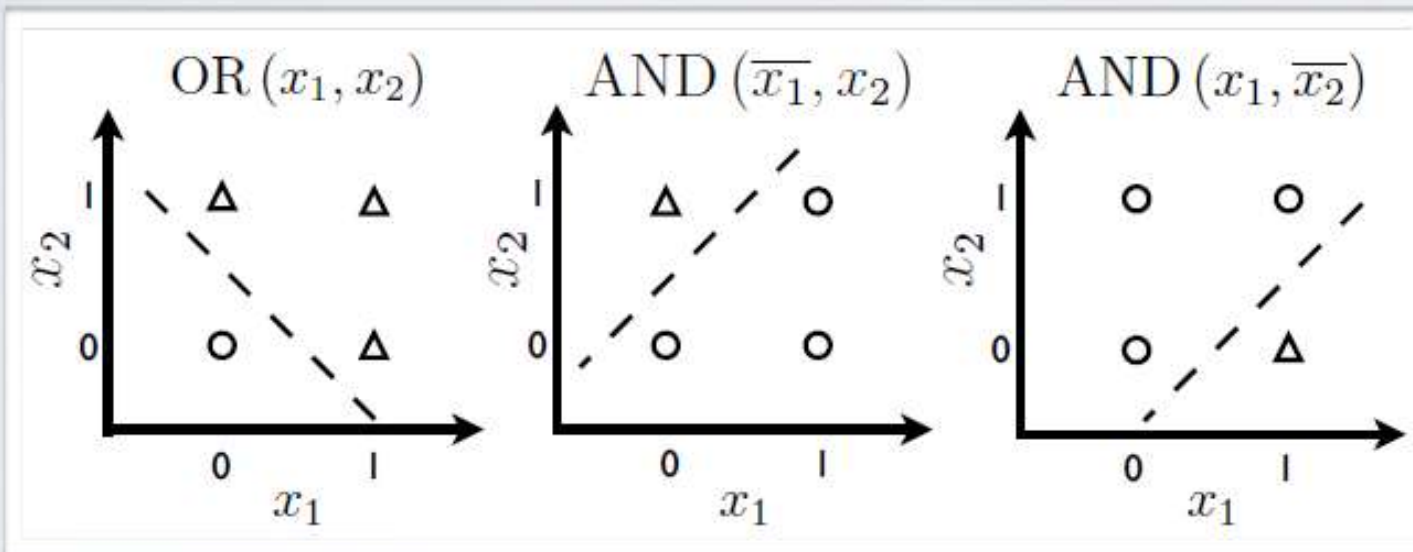
Non-Linearly separable

- Perceptron is able to represent some useful functions
- AND($x_1$,$x_2$) choose weights $w_0$=-1.5, $w_1$=1, $w_2$=1
- But functions that are not linearly separable (e.g. XOR) are not representable

# ARTIFICIAL NEURON

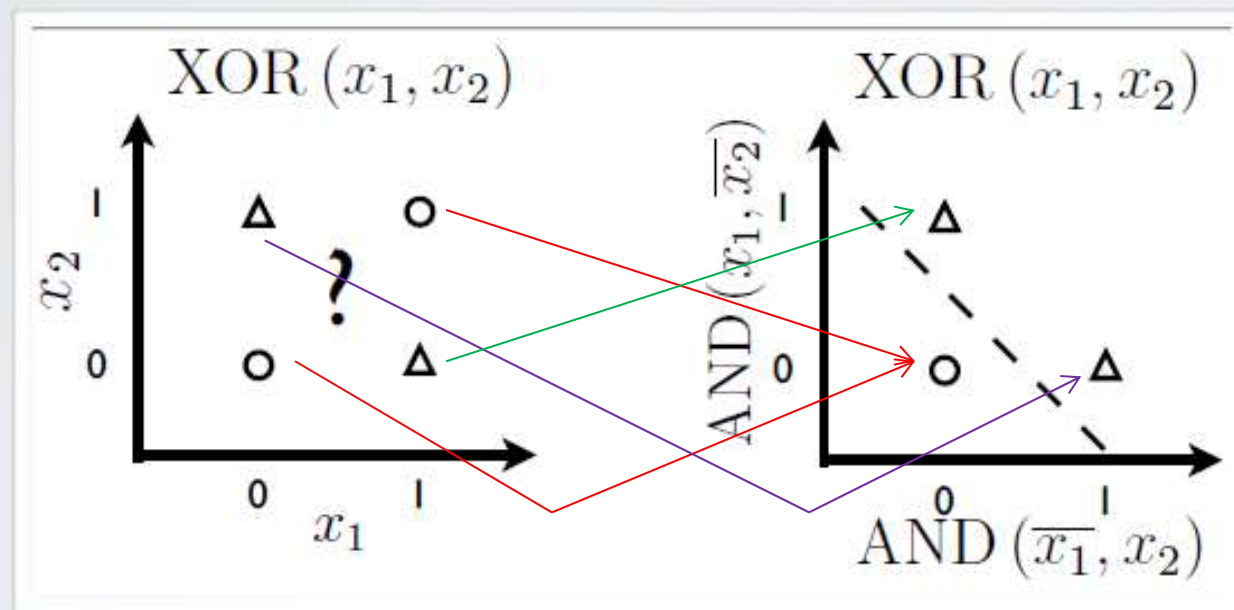**Topics:** capacity of single neuron

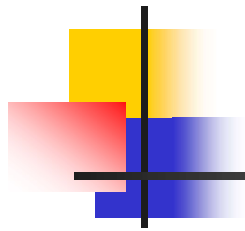• Can solve linearly separable problems

# ARTIFICIAL NEURON

**Topics:** capacity of single neuron

• Can't solve non linearly separable problems...



• ... unless the input is transformed in a better representation

# Hyperplane partitions

- An extra layer models a convex hull
  - "An area with no dents in it"
  - Perceptron models, but can't learn
  - Sigmoid function learning of convex hulls
  - Two layers add convex hulls together
  - Sufficient to classify anything "sane".
- In theory, further layers add nothing
- In practice, extra layers may be better

# Different Non-Linearly Separable Problems

| Structure | Types of Decision Regions | Exclusive-OR Problem | Classes with Meshed regions | Most General Region Shapes |
|---|---|---|---|---|
| *Single-Layer* | *Half Plane Bounded By Hyperplane* | A B B A | B A | |
| *Two-Layer* | *Convex Open Or Closed Regions* | A B B A | B A | |
| *Three-Layer* | Arbitrary (Complexity Limited by No. of Nodes) | A B B A | B A | |

# to continue...

# Multi Layer Perceptron (MLP)

**Topics:** single hidden layer neural network

- Hidden layer pre-activation:

$$\mathbf{a(x)} = \mathbf{b}^{(1)} + \mathbf{W}^{(1)}\mathbf{x}$$

$$\left( a(\mathbf{x})_i = b_i^{(1)} + \sum_j W_{i,j}^{(1)} x_j \right)$$

- Hidden layer activation:

$$\mathbf{h(x)} = \mathbf{g(a(x))}$$

- Output layer activation:

$$f(\mathbf{x}) = o\left( b^{(2)} + \mathbf{w}^{(2)\top} \mathbf{h}^{(1)}\mathbf{x} \right)$$

output activation function



34

# Multi Layer Perceptron (MLP)

**Topics:** softmax activation function

- For multi-class classification:
  - ▸ we need multiple outputs (1 output per class)
  - ▸ we would like to estimate the conditional probability $p(y = c|\mathbf{x})$

- We use the softmax activation function at the output:

$$\mathbf{o}(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \left[ \frac{\exp(a_1)}{\sum_c \exp(a_c)} \cdots \frac{\exp(a_C)}{\sum_c \exp(a_c)} \right]^{\top}$$

  - ▸ strictly positive
  - ▸ sums to one

- Predicted class is the one with highest estimated probability

# Multi Layer Perceptron (MLP)

**Topics:** multilayer neural network

- Could have $L$ hidden layers:
  - layer pre-activation for $k>0$   $(\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x})$

  $$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

  - hidden layer activation ($k$ from 1 to $L$):

  $$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

  - output layer activation ($k=L+1$):

  $$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$

# Types of Layers

- The input layer.
  - Introduces input values into the network.
  - No activation function or other processing.
- The hidden layer(s).
  - Perform classification of features
  - Two hidden layers are sufficient to solve any problem
  - Features imply more layers may be better
- The output layer.
  - Functionally just like the hidden layers
  - Outputs are passed on to the world outside the neural network.

# Capacity of MLP

**Topics:** capacity of single neuron

• Can't solve non linearly separable problems...

$$\text{XOR}\ (x_1, x_2)$$

$$\text{XOR}\ (x_1, x_2)$$

• ... unless the input is transformed in a better representation

38

# Capacity of MLP

**Topics:** multilayer neural network

- Could have $L$ hidden layers:

  ‣ layer pre-activation for $k>0$  ($\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$)
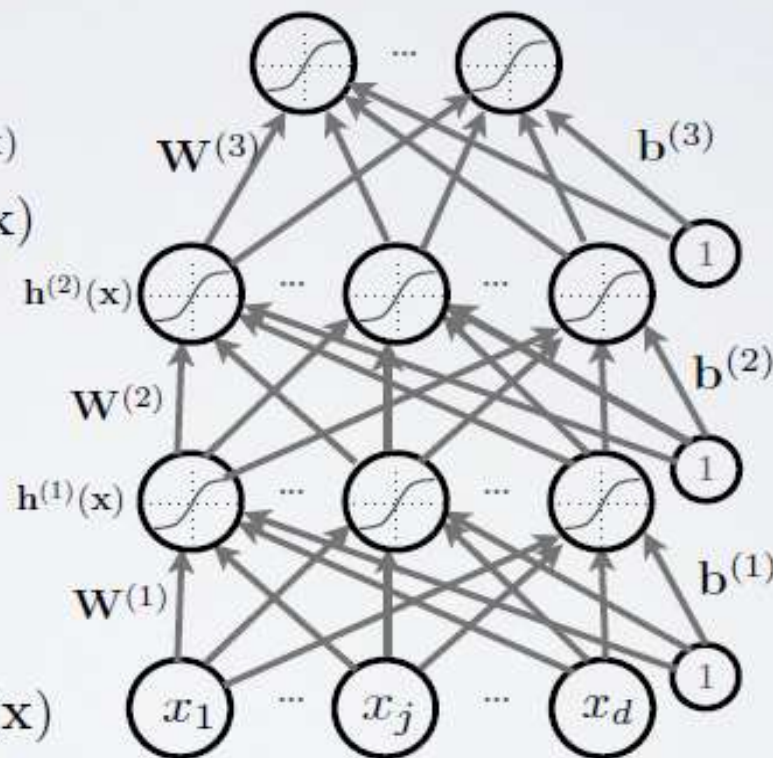
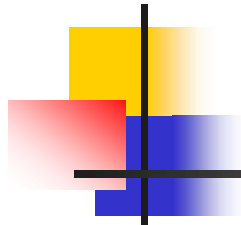  $$\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$$

  ‣ hidden layer activation ($k$ from 1 to $L$):

  $$\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$$

  ‣ output layer activation ($k=L+1$):

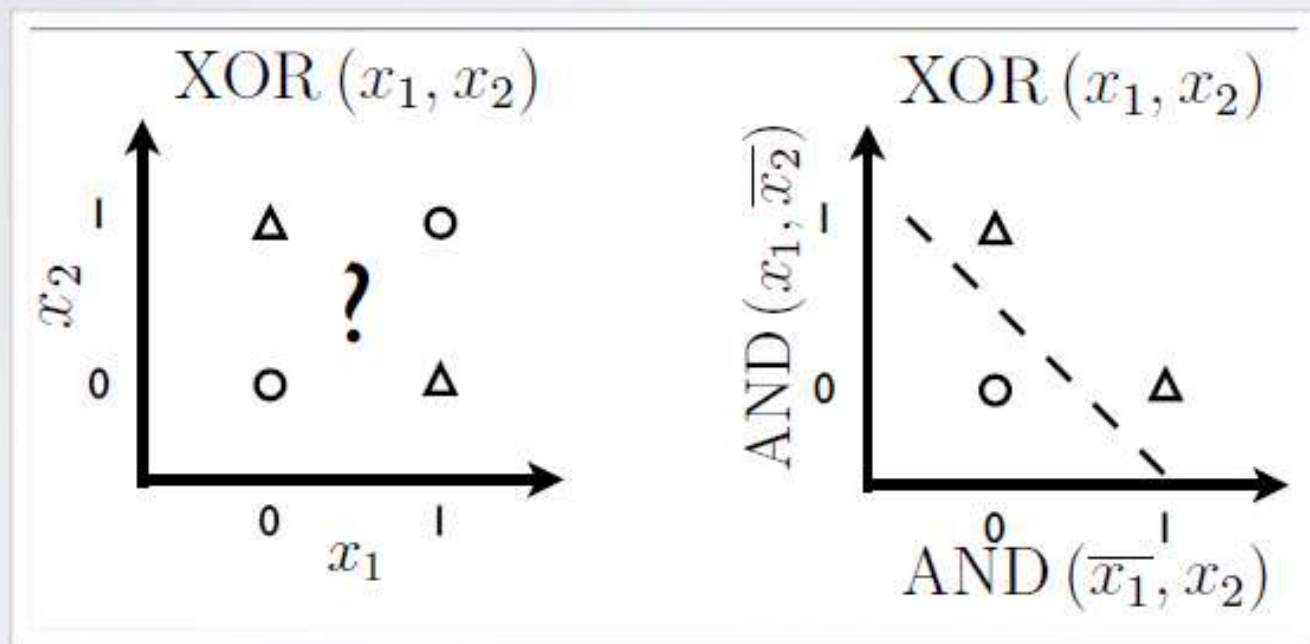  $$\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$$

# Capacity of MLP



**Topics:** single hidden layer neural network

(from Pascal Vincent's slides)

# Capacity of MLP



**Topics:** single hidden layer neural network

(from Pascal Vincent's slides)

# Capacity of MLP



**Topics:** single hidden layer neural network

(from Pascal Vincent's slides)

# Capacity of MLP

**Topics:** universal approximation

- Universal approximation theorem (Hornik, 1991):
  - ‣ "a single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units"

- The result applies for sigmoid, tanh and many other hidden layer activation functions

- This is a good result, but it doesn't mean there is a learning algorithm that can find the necessary parameter values!

# to continue...

# Training Algorithms

- Adjust neural network weights to map inputs to outputs.

- Use a set of sample patterns where the desired output (given the inputs presented) is known.

- The purpose is to learn to generalize
  - Recognize features which are common to good and bad exemplars

# Back-Propagation

- A training procedure which allows multi-layer feedforward Neural Networks to be trained;
- Can theoretically perform "any" input-output mapping;
- Can learn to solve linearly inseparable problems.

# Activation functions and training

- For feed-forward networks:
  - A continuous function can be differentiated allowing gradient-descent.
  - Back-propagation is an example of a gradient-descent technique.
  - Reason for prevalence of sigmoid

# Gradient Descent Learning Rule

- Consider linear unit without threshold and continuous output o (not just −1,1)

  - $o = w_0 + w_1 x_1 + \ldots + w_n x_n$

- Train the $w_i$'s such that they minimize the squared error

  - $E[w_1, \ldots, w_n] = \frac{1}{2} \sum_{d \in D} (t_d - o_d)^2$

  where D is the set of training examples

# Gradient Descent

$D=\{<(1,1),1>,<(-1,-1),1>,$
$<(1,-1),-1>,<(-1,1),-1>\}$

Gradient:
$\nabla E[w]=[\partial E/\partial w_0,\dots \partial E/\partial w_n]$

$\Delta w=-\eta \ \nabla E[w]$

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

$=\partial/\partial w_i \ 1/2\Sigma_d(t_d-o_d)^2$

$= \partial/\partial w_i \ 1/2\Sigma_d(t_d-\Sigma_i w_i x_i)^2$

$= \Sigma_d(t_d- o_d)(-x_i)$

# Gradient Descent

Gradient-Descent(*training_examples*, $\eta$)

Each training example is a pair of the form $<(x_1,\ldots x_n),t>$ where $(x_1,\ldots,x_n)$ is the vector of input values, and t is the target output value, $\eta$ is the learning rate (e.g. 0.1)

- Initialize each $w_i$ to some small random value
- Until the termination condition is met, Do
    - Initialize each $\Delta w_i$ to zero
    - For each $<(x_1,\ldots x_n),t>$ in *training_examples* Do
        - Input the instance $(x_1,\ldots,x_n)$ to the linear unit and compute the output o
        - For each linear unit weight $w_i$ Do
            - $\Delta w_i = \Delta w_i + \eta\ (t\text{-}o)\ x_i$
    - For each linear unit weight wi Do
        - $w_i = w_i + \Delta w_i$

# Incremental Stochastic Gradient Descent

- Batch mode : gradient descent

  $w = w - \eta \, \nabla E_D[w]$ over the entire data D

  $E_D[w] = 1/2 \Sigma_d (t_d - o_d)^2$

- Incremental mode: gradient descent

  $w = w - \eta \, \nabla E_d[w]$ over individual training examples d

  $E_d[w] = 1/2 \, (t_d - o_d)^2$

Incremental Gradient Descent can approximate Batch Gradient Descent arbitrarily closely if $\eta$ is small enough

# Comparison Perceptron and Gradient Descent Rule

Perceptron learning rule guaranteed to succeed if

- Training examples are linearly separable
- Sufficiently small learning rate $\eta$

Linear unit training rules uses gradient descent

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

# Multi-Layer Networks

output layer

hidden layer

input layer

53

# Sigmoid Unit



$x_{0=1}$

$x_1$ $w_1$

$w_0$ $net = \sum_{i=0}^{n} w_i x_i$ $o = \sigma(net) = 1/(1+e^{-net})$

$x_2$ $w_2$

$\Sigma$
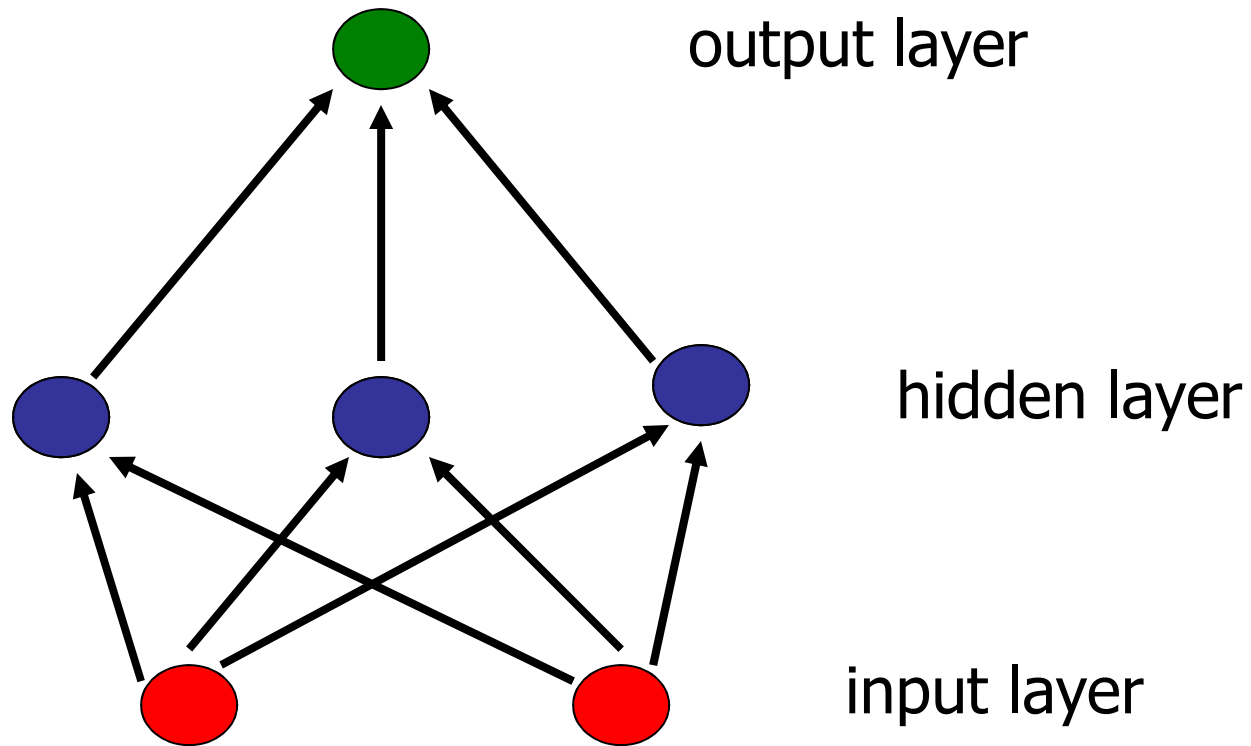
$w_n$

$x_n$

$\sigma(x)$ is the sigmoid function: $1/(1+e^{-x})$

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

Derive gradient decent rules to train:
• one sigmoid function

$\partial E/\partial w_i = -\sum_d (t_d - o_d) o_d (1-o_d) x_i$

• Multilayer networks of sigmoid units
backpropagation:

54

# Backpropagation Algorithm

- Initialize each $w_i$ to some small random value
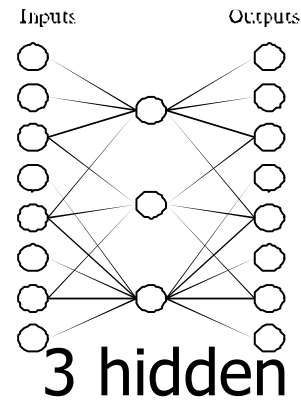- Until the termination condition is met, Do
  - For each training example $<(x_1,...x_n),t>$ Do
    - Input the instance $(x_1,...,x_n)$ to the network and compute the network outputs $o_k$
    - For each output unit k
      - $\delta_k = o_k(1-o_k)(t_k-o_k)$
    - For each hidden unit h

      - $\delta_h = o_h(1-o_h) \sum_k w_{h,k} \delta_k$
    - For each network weight $w_{,j}$ Do
    - $w_{i,j} = w_{i,j} + \Delta w_{i,j}$ where
      $\Delta w_{i,j} = \eta \ \delta_j \ x_{i,j}$

55

# Backpropagation

- Gradient descent over entire *network* weight vector
- Easily generalized to arbitrary directed graphs
- Will find a local, not necessarily global error minimum
  -in practice often works well (can be invoked multiple times with different initial weights)
- Often include weight *momentum* term

$$\Delta w_{i,j}(t)= \eta \ \delta_j \ x_{i,j} + \alpha \ \Delta w_{i,j}\ (t\text{-}1)$$

- Minimizes error training examples
  - Will it generalize well to unseen instances (over-fitting)?
- Training can be slow typical 1000-10000 iterations
  (use Levenberg-Marquardt instead of gradient descent)
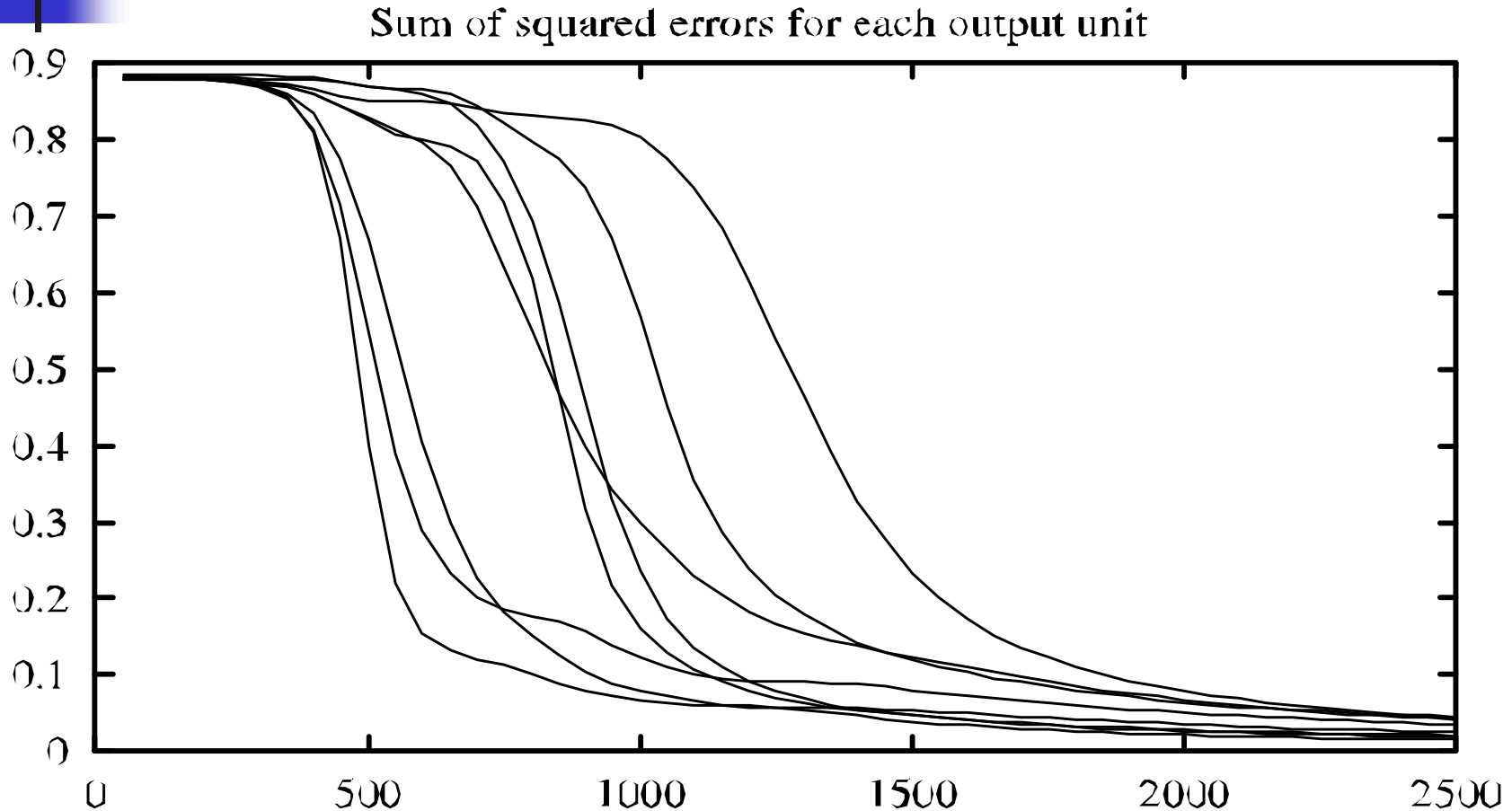- Using network after training is fast

# 8-3-8 Binary Encoder-Decoder



8 inputs          3 hidden          8 outputs

A target function:

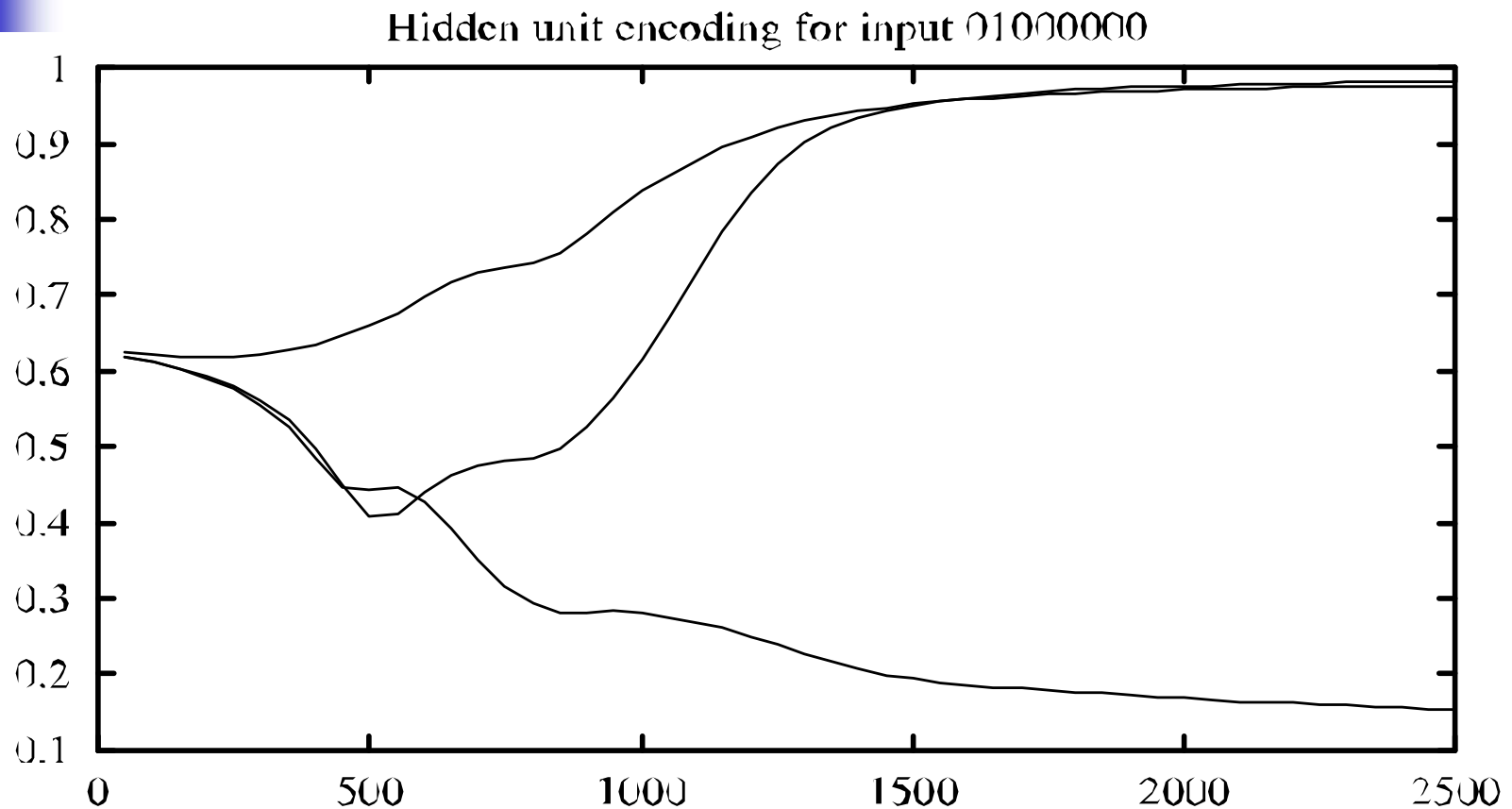| Input | Output |
|---|---|
| 10000000 → | 10000000 |
| 01000000 → | 01000000 |
| 00100000 → | 00100000 |
| 00010000 → | 00010000 |
| 00001000 → | 00001000 |
| 00000100 → | 00000100 |
| 00000010 → | 00000010 |
| 00000001 → | 00000001 |

Hidden values
.89 .04 .08
.01 .11 .88
.01 .97 .27
.99 .97 .71
.03 .05 .02
.22 .99 .99
.80 .01 .98
.60 .94 .01

Can this be learned??

# Sum of Squared Errors for the Output Units



Sum of squared errors for each output unit

# Hidden Unit Encoding for Input 0100000



Hidden unit encoding for input 01000000

# Convergence of Backprop

Gradient descent to some local minimum
- Perhaps not global minimum
- Add momentum
- Stochastic gradient descent
- Train multiple nets with different initial weights

Nature of convergence
- Initialize weights near zero
- Therefore, initial networks near-linear
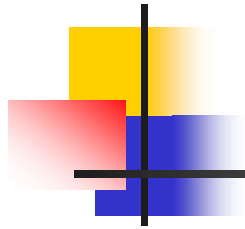- Increasingly non-linear functions possible as training progresses

# Expressive Capabilities of ANN

Boolean functions

- Every boolean function can be represented by network with single hidden layer
- But might require exponential (in number of inputs) hidden units

Continuous functions

- Every bounded continuous function can be approximated with arbitrarily small error, by network with one hidden layer [Cybenko 1989, Hornik 1989]
- Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988]
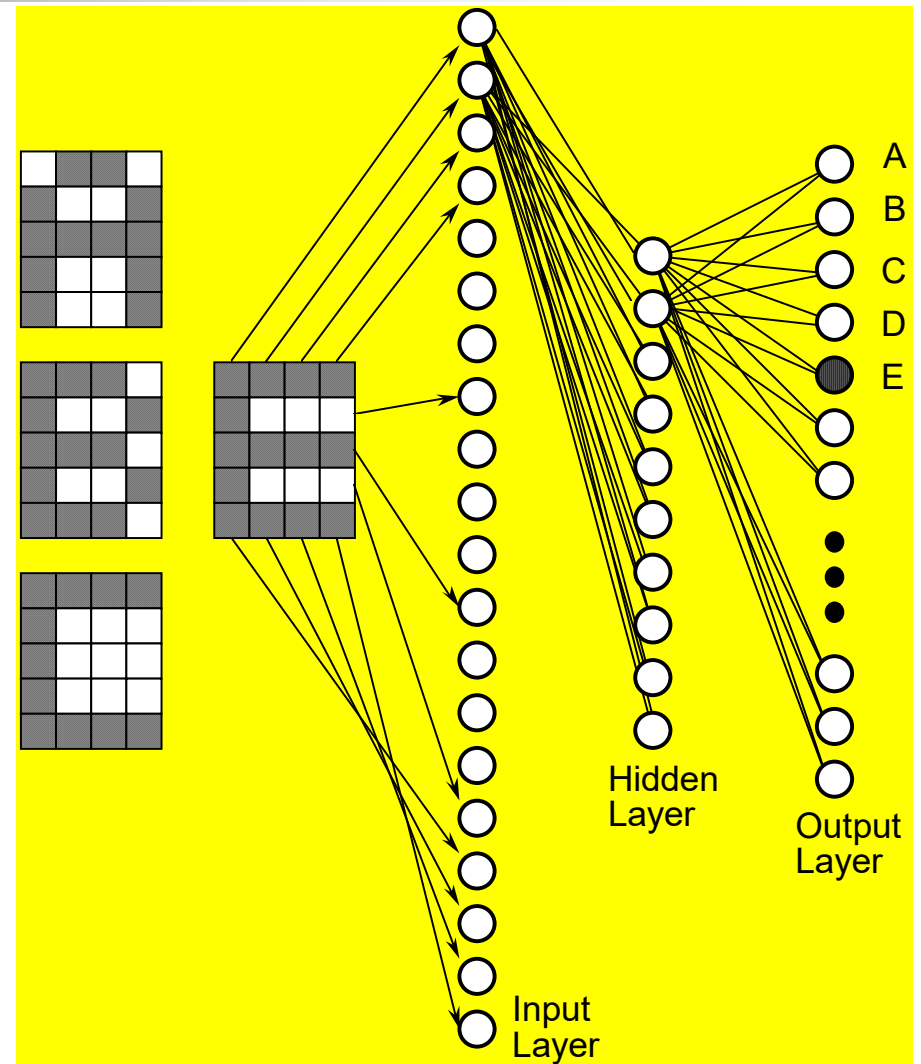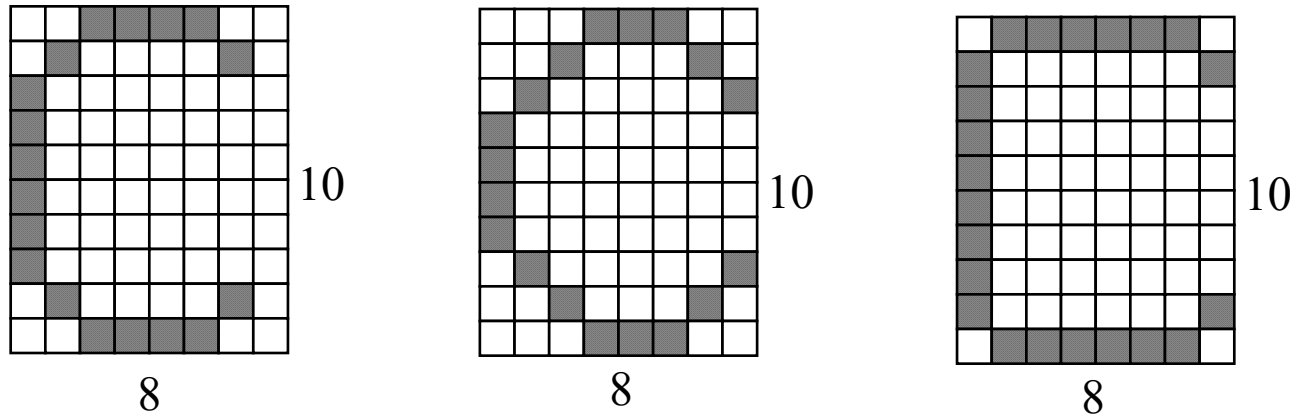
# Applications

- The properties of neural networks define where they are useful.
  - Can learn complex mappings from inputs to outputs, based solely on samples
  - Difficult to analyse: firm predictions about neural network behaviour difficult;
    - Unsuitable for safety-critical applications.
  - Require limited understanding from trainer, who can be guided by heuristics.

# Neural network for OCR

- feedforward network
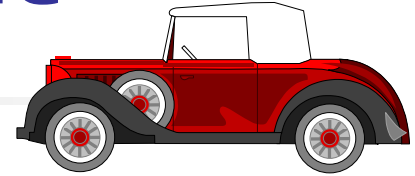- trained using Back- propagation
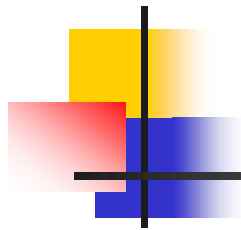
# OCR for 8x10 characters



- NN are able to generalise
- learning involves generating a partitioning of the input space
- for single layer network input space must be linearly separable
- what is the dimension of this input space?
- how many points in the input space?
- this network is binary(uses binary values)

# Engine management

- The behaviour of a car engine is influenced by a large number of parameters
  - temperature at various points
  - fuel/air mixture
  - lubricant viscosity.
- Major companies have used neural networks to dynamically tune an engine depending on current settings.
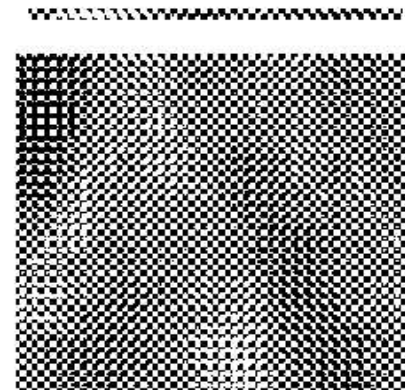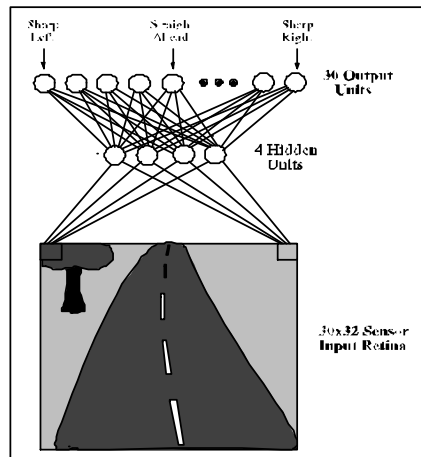
# ALVINN
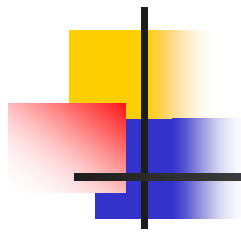
Drives 70 mph on a public highway



30 outputs for steering
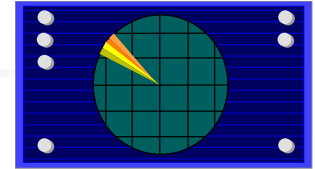
4 hidden units

30x32 pixels as inputs



30x32 weights into one out of four hidden unit



66

# Signature recognition

- Each person's signature is different.
- There are structural similarities which are difficult to quantify.
- One company has manufactured a machine which recognizes signatures to within a high level of accuracy.
  - Considers speed in addition to gross shape.
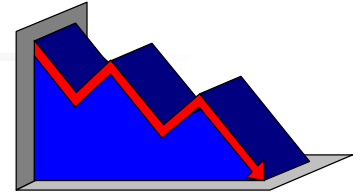  - Makes forgery even more difficult.

# Sonar target recognition

- Distinguish mines from rocks on sea-bed
- The neural network is provided with a large number of parameters which are extracted from the sonar signal.
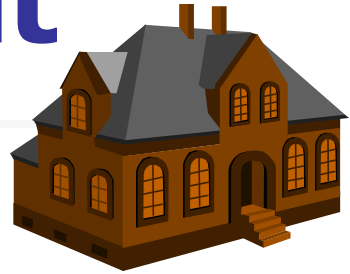- The training set consists of sets of signals from rocks and mines.

# Stock market prediction

- "Technical trading" refers to trading based solely on known statistical parameters; e.g. previous price

- Neural networks have been used to attempt to predict changes in prices.

- Difficult to assess success since companies using these techniques are reluctant to disclose information.
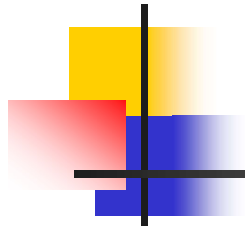
# Mortgage assessment

- Assess risk of lending to an individual.

- Difficult to decide on marginal cases.

- Neural networks have been trained to make decisions, based upon the opinions of expert underwriters.

- Neural network produced a 12% reduction in delinquencies compared with human experts.

# Neural Network Problems

- Many Parameters to be set
- Overfitting
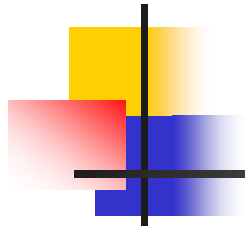- long training times
- ...

# Parameter setting

- Number of layers
- Number of neurons
  - too many neurons, require more training time
- Learning rate
  - from experience, value should be small ~0.1
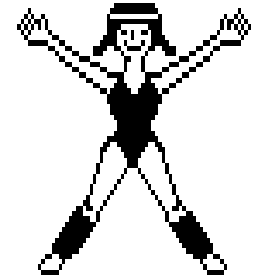- Momentum term
- ..

# Over-fitting

- With sufficient nodes can classify any training set exactly

- May have poor generalisation ability.

- Cross-validation with some patterns
  - Typically 30% of training patterns
  - Validation set error is checked each epoch
  - Stop training if validation error goes up

# **Training time**

- How many epochs of training?
  - Stop if the error fails to improve (has reached a minimum)
  - Stop if the rate of improvement drops below a certain level
  - Stop if the error reaches an acceptable level
  - Stop when a certain number of epochs have passed

# Literature & Resources

- Textbook:
  - ”Neural Networks for Pattern Recognition”, Bishop, C.M., 1996
- Software:
  - Neural Networks for Face Recognition

  http://www.cs.cmu.edu/afs/cs.cmu.edu/user/mitchell/ftp/faces.html
  - SNNS Stuttgart Neural Networks Simulator

  http://www-ra.informatik.uni-tuebingen.de/SNNS
  - Neural Networks at your fingertips

  http://www.geocities.com/CapeCanaveral/1624/

  http://www.stats.gla.ac.uk/~ernest/files/NeuralAppl.html