# Neural Networks

# Perceptron

- ❖ The perceptron algorithm was invented in 1957 at the Cornell Aeronautical Laboratory by Frank Rosenblatt.
- ❖ Perceptron is an algorithm for supervised classification.
- ❖ It is a type of linear classifier.
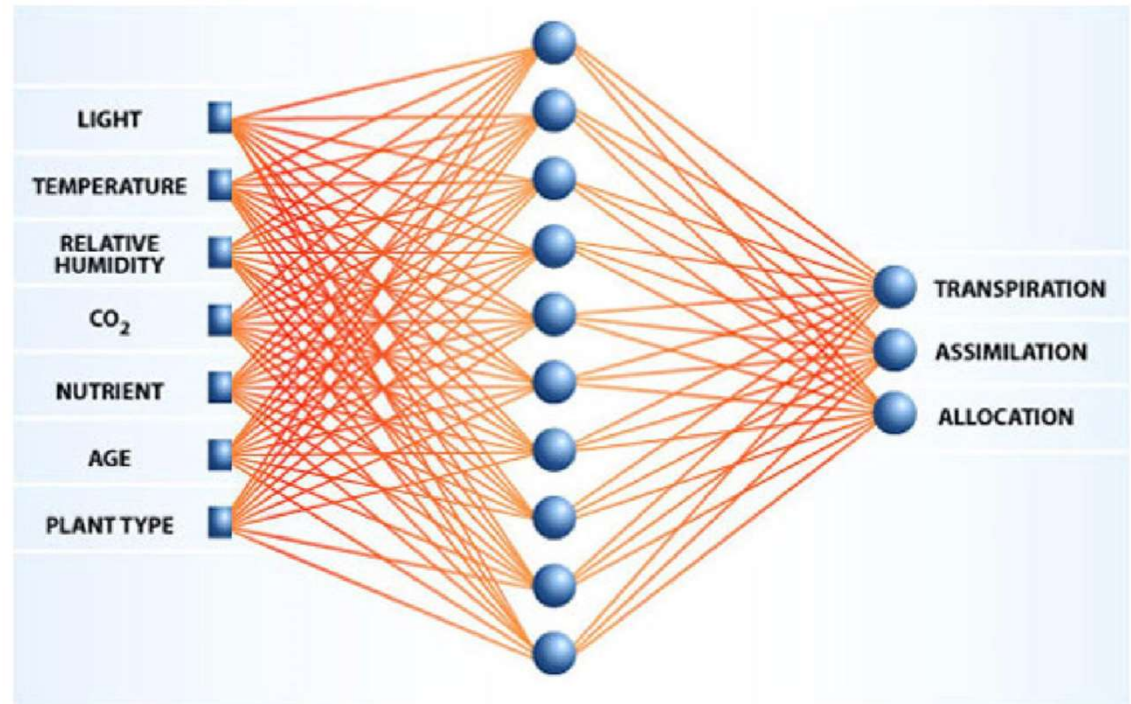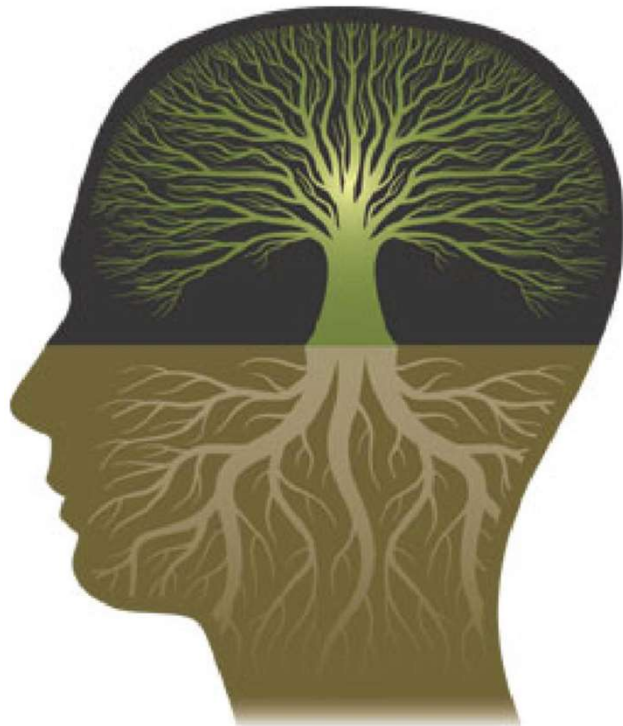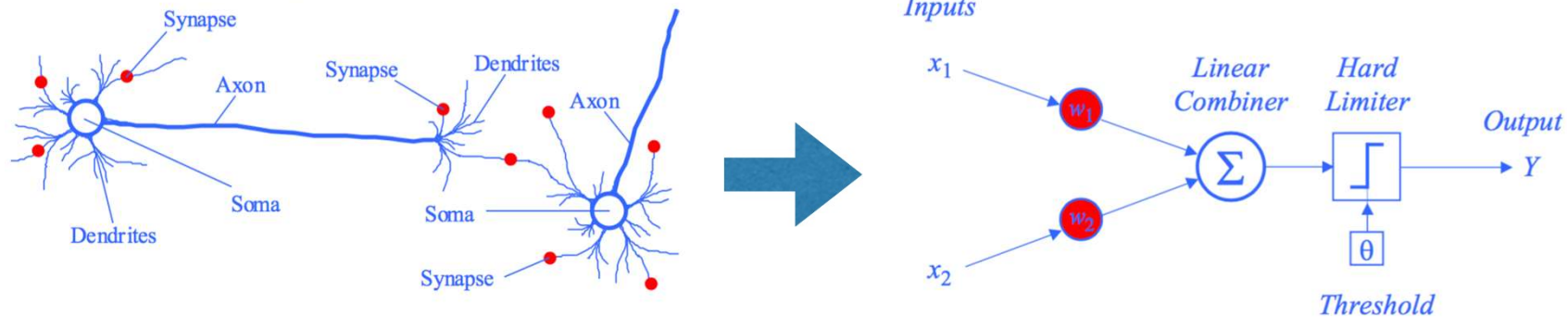- ❖ It lays the foundation of artificial neural networks (ANN).
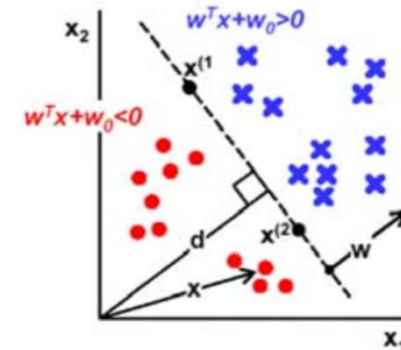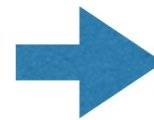
# Inspired from Neural Networks

# Perceptron and Neural Nets

❖ From biological neuron to artificial neuron (perceptron)

Synapse

Axon

Synapse   Dendrites

Axon

Soma

Soma

Dendrites

Synapse

Inputs

$x_1$

$w_1$

*Linear Combiner*

*Hard Limiter*

$\Sigma$

$\theta$

*Output*

$Y$

$x_2$

$w_2$

*Threshold*

❖ Generative Model

$$X = \sum_{i=1}^{n} x_i w_i \qquad y = \begin{cases} +1, if\, X > w_0 \\ -1, if\, X < w_0 \end{cases}$$

$x_2$

$w^T x + w_0 > 0$

$w^T x + w_0 < 0$

$x^{(1}$

$x^{(2}$

$d$

$w$

$x$

$x_1$

❖ Artificial neuron networks

  ❖ supervised learning
  ❖ gradient descent

*Input Signals*

*Output Signals*

*Input Layer*

*Middle Layer*

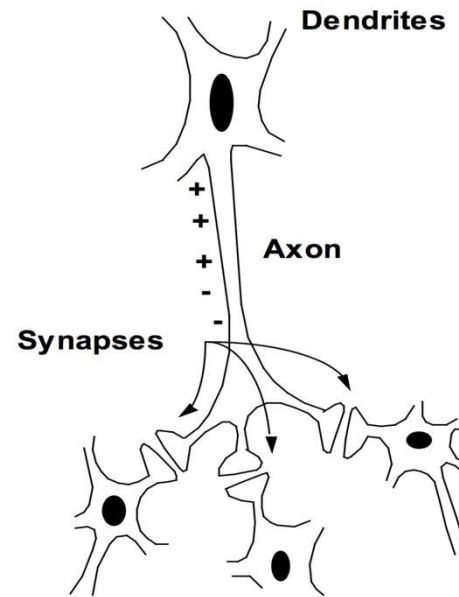*Output Layer*

# Connectionist Models
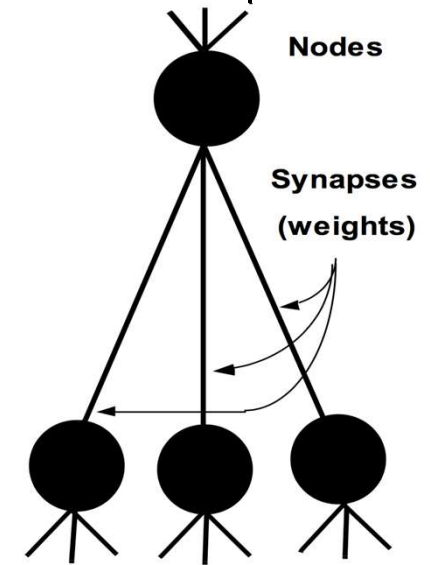
❖ **Consider humans:**
  - ❖ Neuron switching time
    ~ 0.001 second
  - ❖ Number of neurons
    ~ $10^{10}$
  - ❖ Connections per neuron
    ~ $10^{4-5}$
  - ❖ Scene recognition time
    ~ 0.1 second
  - ❖ 100 inference steps doesn't seem like enough

    → much parallel computation

❖ **Properties of artificial neural nets (ANN)**
  - ❖ Many neuron-like threshold switching units
  - ❖ Many weighted interconnections among units
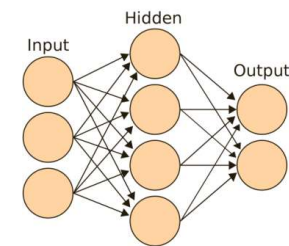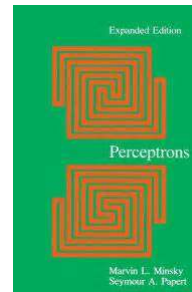  - ❖ Highly parallel, distributed processes



Dendrites

Axon

Synapses

Impulse

Nodes

Synapses
(weights)

# History of Neural Networks

❖ **The Beginnings** (1940s): M-P model, Hebbian learning theory,…
❖ **Golden Years** (1958~1969): The perceptron algorithm, Adaline,…
❖ **Winter** (1969~1980): Minsky's devastating criticism of perceptrons

❖ **Boom** (1980s): Hopfield net, Back Propagation algorithm
❖ **Winter** (1990s): Statistical learning theory, SVM
In the 1990's, many researchers abandoned neural because SVMs worked better, and there was no successful attempts to train deep networks.
❖ **The 3rd rise of NN** (2000-present): Deep learning

# Time Line

# Perceptions

Input units

Cough    Headache

weights

No disease    Pneumonia    Flu    Meningitis

Output units

$\Delta$ **rule**
*change weights to decrease the error*

$-$ $\dfrac{\text{what we got}}{\text{what we wanted}}$
*error*

# Perceptrons



**Output units**

**Input units**

Output of unit $j$:

$$o_j = 1/(1 + e^{-(a_j + \theta_j)})$$

$j$

Input to unit $j$: $a_j = \Sigma w_{ij} a_i$

Input to unit $i$: $a_i$

measured value of variable $i$

$i$

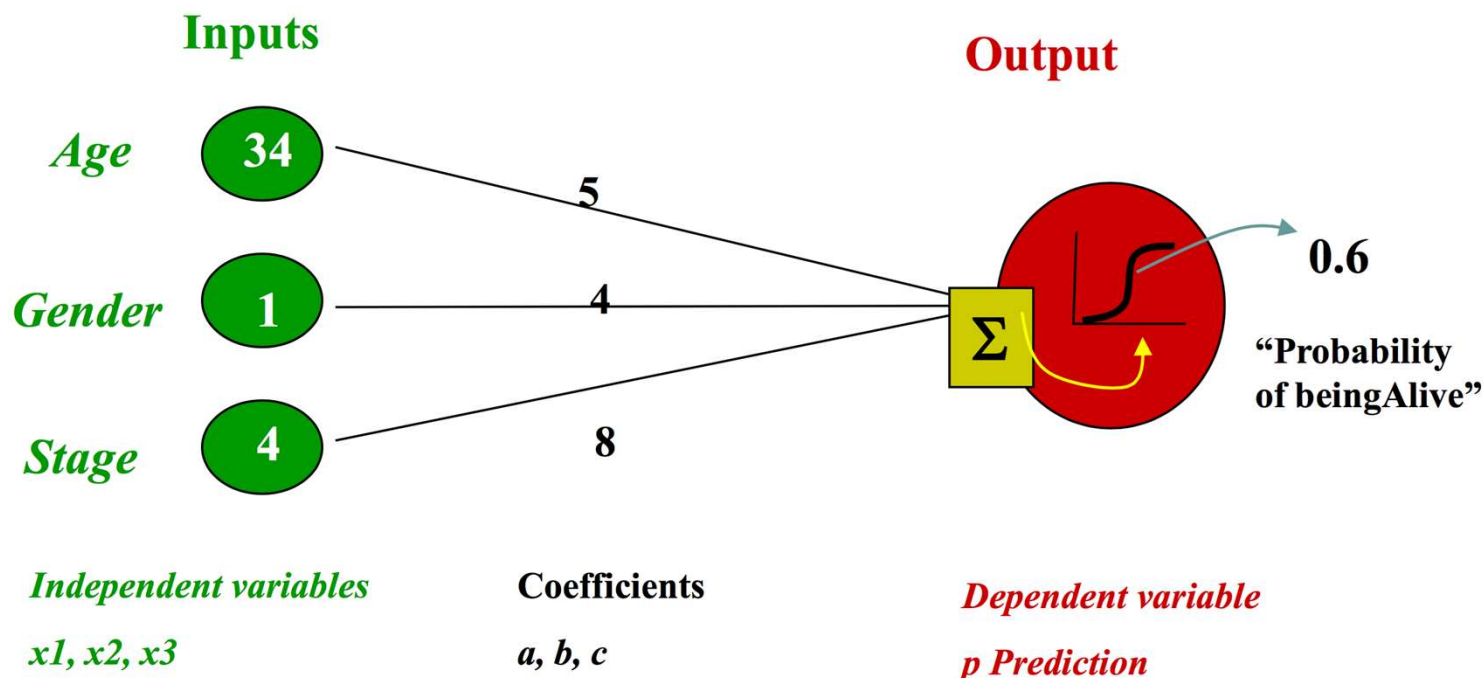# Jargon Pseudo-Correspondence
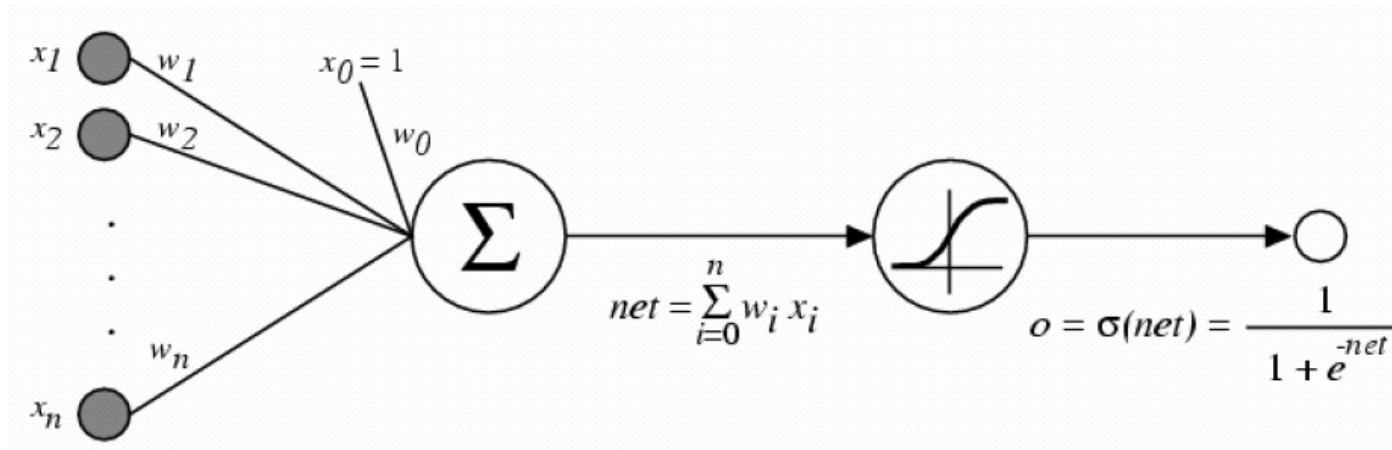
❖ Independent variable = input variable
❖ Dependent variable = output variable
❖ Coefficients = "weights"
❖ Estimates = "targets"

Logistic Regression Model (the sigmoid unit)

**Inputs**

*Age* 34

*Gender* 1

*Stage* 4

5

4

8

Σ

**Output**

0.6

"Probability of beingAlive"

*Independent variables*

*x1, x2, x3*

Coefficients

*a, b, c*

*Dependent variable*

*p Prediction*

# The perceptron learning algorithm



❖ Recall the nice property of sigmoid function $\dfrac{d\sigma}{dt} = \sigma(1-\sigma)$

❖ Consider regression problem f:X→Y , for scalar Y: $y = f(x) + \epsilon$

❖ Let's maximize the conditional data likelihood

$$\overrightarrow{w} = \arg\max_{\overrightarrow{w}} \ln \prod_i P(y_i|x_i; \overrightarrow{w})$$

$$\overrightarrow{w} = \arg\min_{\overrightarrow{w}} \sum_i \frac{1}{2}(y_i - \hat{f}(x_i; \overrightarrow{w}))^2$$

# Gradient Descent



Gradient

$$\nabla E[\vec{w}] \equiv \left[ \frac{\partial E}{\partial w_0}, \frac{\partial E}{\partial w_1}, \cdots \frac{\partial E}{\partial w_n} \right]$$

Training rule:

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

i.e.,

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

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_i$ = weigh i

$$\frac{\partial E[\vec{w}]}{\partial w_j} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum (t_d - o_d)^2$$

$$=$$

# Gradient Descent

$$\frac{\partial E_D[\vec{w}])}{\partial w_j} = \frac{\partial}{\partial w_i} \frac{1}{2} \sum_d (t_d - o_d)^2$$

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

$$= \sum_d (t_d - o_d)\left(-\frac{\partial o_d}{\partial w_i}\right)$$

$$= -\sum_d (t_d - o_d) \frac{\partial o_d}{\partial net_i} \frac{\partial net_d}{\partial w_i}$$

$$= -\sum_d (t_d - o_d) o_d (1 - o_d) x_d^i$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_i$ = weigh i

**Incremental mode:**
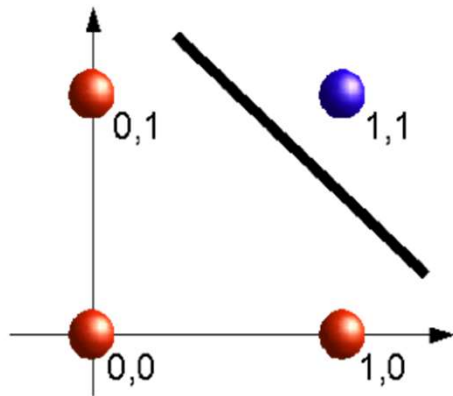Do until converge:
❖ For each training example **d** in **D**
1. compute gradient $\nabla \mathbf{E}_d[\mathbf{w}]$
**2.** $\vec{w} = \vec{w} - \eta \nabla E_d[\vec{w}]$
where
$$\nabla E_d[\vec{w}] = -(t_d - o_d)o_d(1 - o_d)\vec{x}_d$$

**Batch mode:**
Do until converge:
1. compute gradient $\nabla \mathbf{E}_D[\mathbf{w}]$
$$\vec{w} = \vec{w} - \eta \nabla E_D[\vec{w}]$$

# What decision surface does a perceptron define?



NAND

| X | Y | Z(color) |
|---|---|----------|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |



$\theta = 0.5$

$f(x_1w_1 + x_2w_2) = y$
$f(0w_1 + 0w_2) = 1$
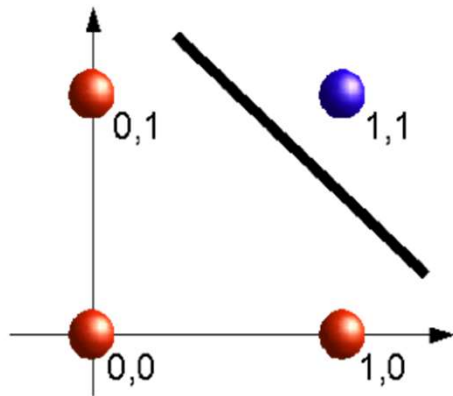$f(0w_1 + 1w_2) = 1$
$f(1w_1 + 0w_2) = 1$
$f(1w_1 + 1w_2) = 0$

$$f(a) = \begin{cases} 1, \text{ for } a > \theta \\ 0, \text{ for } a \leq \theta \end{cases}$$
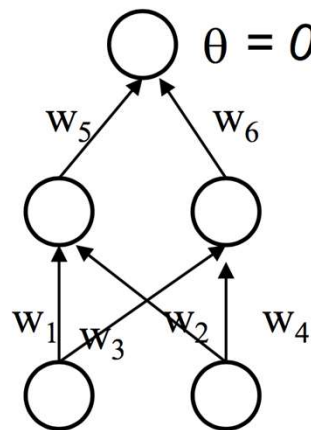
some possible values for $w_1$ and $w_2$

| $w_1$ | $w_2$ |
|-------|-------|
| 0.20 | 0.35 |
| 0.20 | 0.40 |
| 0.25 | 0.30 |
| 0.40 | 0.20 |

# What decision surface does a perceptron define?

| X | Y | Z(color) |
|---|---|----------|
| 0 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 0 |

XOR

$\theta = 0.5$ for all units

$w_5$  $w_6$

$w_1$  $w_3$  $w_2$  $w_4$

$$f(a) = \begin{cases} 1, \text{ for } a > \theta \\ 0, \text{ for } a \leq \theta \end{cases}$$

a possible set of values for $(w_1, w_2, w_3, w_4, w_5, w_6)$:
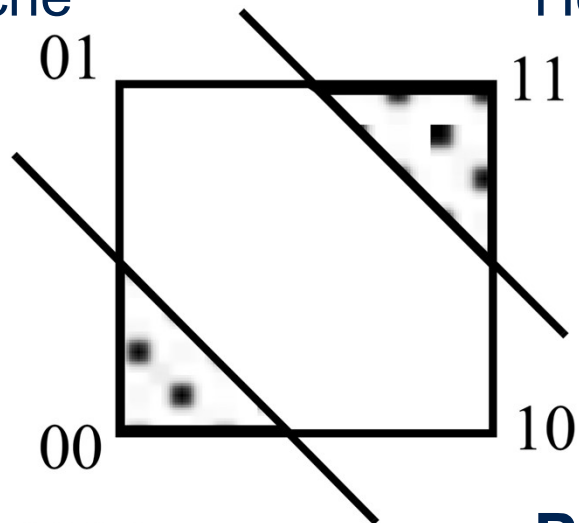(0.6,-0.6,-0.7,0.8,1,1)

# Non Linear Separation

**Meningitis**
No cough
Headache

**Flu**
Cough
Headache



No treatment
Treatment

**No disease**
No cough
No headache

**Pneumonia**
Cough
No headache

# Neural Network Model

# Combined logistic models



**Inputs**

*Age* 34

*Gender* 2

*Stage* 4

.6

.1

.7

.5

.8

Σ

**Output**

0.6

"Probability of beingAlive"

*Independent variables*

**Weights**

**Hidden Layer**

**Weights**

*Dependent variable*

*Prediction*

# Combined logistic models

**Inputs**

*Age*    34

*Gender*    2

*Stage*    4

.2

.3

.2

.5

.8

Σ

**Output**

0.6

"Probability of beingAlive"

*Independent variables*    **Weights**    **Hidden Layer**    **Weights**    *Dependent variable*

*Prediction*

# Combined logistic models



**Inputs**

*Age* 34

*Gender* 1

*Stage* 4

.6
.2
.1
.3
.7
.2

.5

.8

Σ

**Output**

0.6

"Probability of beingAlive"

*Independent variables*

**Weights**

**Hidden Layer**

**Weights**

*Dependent variable*

*Prediction*

# Not really,
# no target for hidden units...



**Age** 34 .6 .2 Σ .4 .5 Σ 0.6

**Gender** 2 .1 .3 .2 "Probability of beingAlive"

**Stage** 4 .7 .2 .8

**Independent variables**   **Weights**   **Hidden Layer**   **Weights**   **Dependent variable**

**Prediction**

# Perceptrons

$$\vec{w} := \vec{w} + \eta \sum_d (t_d - o_d) o_d (1 - o_d) \vec{x}_d$$

**Input units**

Cough    Headache

*weights*

No disease    Pneumonia    Flu    Meningitis

**Output units**

△ **rule**
*change weights to decrease the error*

−  $\dfrac{\text{what we got}}{\text{what we wanted}}$
*error*

# Hidden Units and Backpropagation

# Back Propagation Algorithm

❖ Initialize all weights to small random numbers

Until convergence, Do

1. Input the training example to the network and compute the network outputs

2. For each output unit $k$

$$\delta_k = o_k^{(2)}(1 - o_k^{(2)})(o_k^{(2)} - t)$$

⬇

$$\frac{\partial E}{\partial net_k}, \text{net}_k \text{ is the input of output unit k}$$

$$net_k = \sum_j w_{jk}^{(2)} o_j^{(1)}$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_i$ = weigh i

# Back Propagation Algorithm

❖ Initialize all weights to small random numbers

Until convergence, Do

1. Input the training example to the network and compute the network outputs

2. For each output unit $k$

$$\delta_k = o_k^{(2)}(1 - o_k^{(2)})(o_k^{(2)} - t)$$

3. For each hidden unit $h$

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

$\dfrac{\partial E}{\partial net_h}$, net$_h$ is the input of hidden unit h
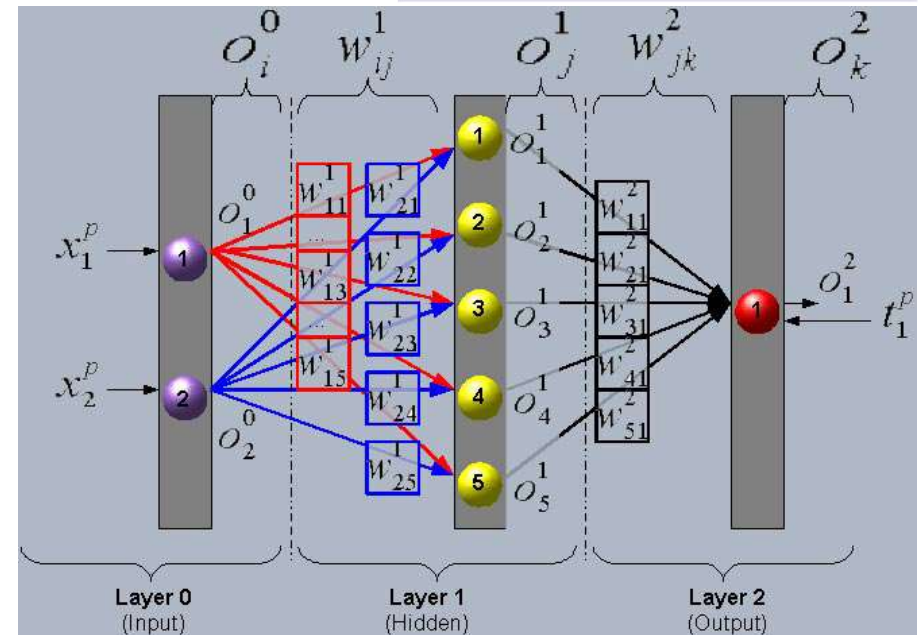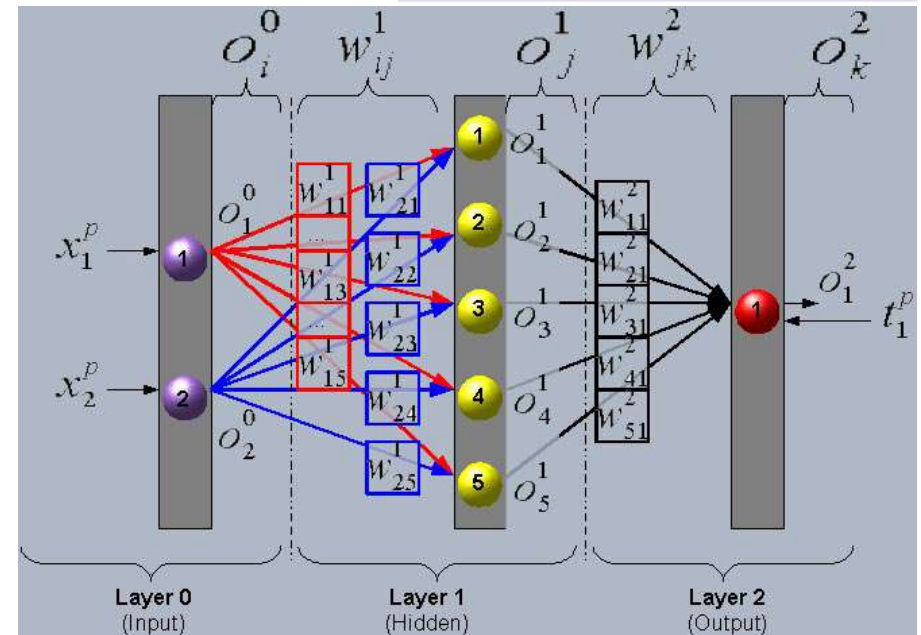
$$net_h = \sum_i w_{ij}^{(1)} O_i^0$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_i$ = weigh i

# Back Propagation Algorithm

❖ Initialize all weights to small random numbers

  Until convergence, Do

  1. Input the training example to the network and compute the network outputs

  2. For each output unit $k$

$$\delta_k = o_k^{(2)}(1 - o_k^{(2)})(o_k^{(2)} - t)$$

  3. For each hidden unit $h$

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

  4. Update each network weight $w_{i,j}$

$$\frac{\partial E}{\partial w_{i,j}^{(1)}} = \delta_j x_i, \qquad \frac{\partial E}{\partial w_{j,k}^{(2)}} = \delta_k o_j,$$
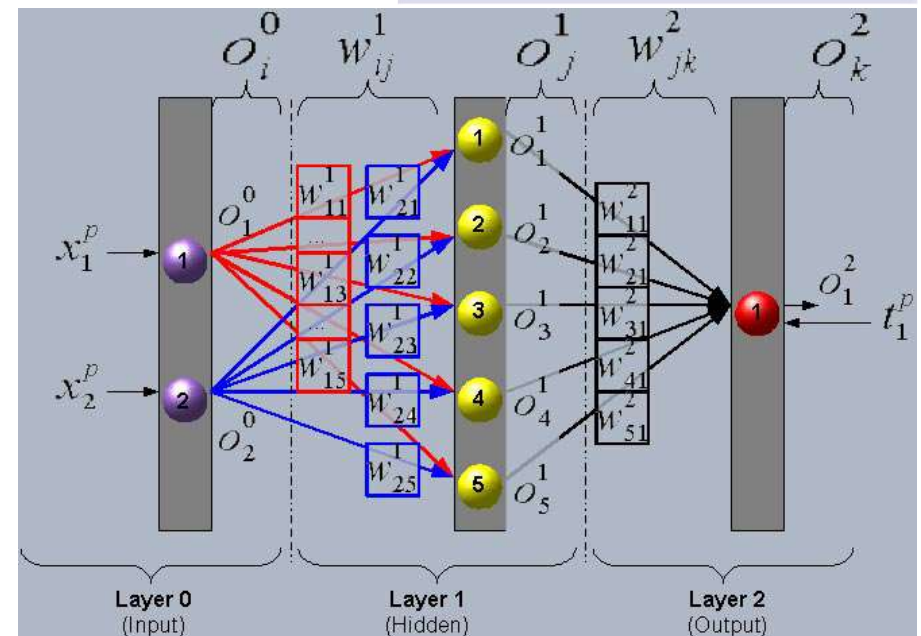
$$w_{i,j} := w_{i,j} + \Delta w_{i,j}$$

$x_d$ = input

$t_d$ = target output

$o_d$ = observed unit output

$w_i$ = weigh i



multiply δ (the unit at the output end of the weight) by the value for the unit at the input end of the weight

# More on Backpropatation

❖ It is doing 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 run multiple times)
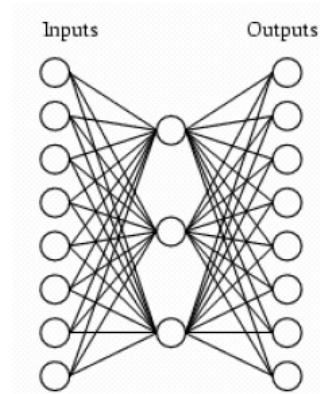❖ Often include weight *momentum* α

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

❖ Minimizes error over *training* examples
  ❖ Will it generalize well to subsequent testing examples?
❖ Training can take thousands of iterations,→very slow!
❖ Using network after training is very fast

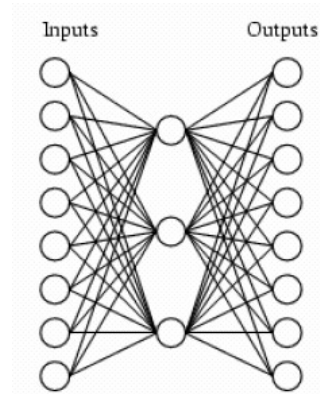# Learning Hidden Layer Representation

❖ A network:



❖ A target function:

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

❖ Can this be learned?

# Learning Hidden Layer Representation

❖ A network:



Inputs                Outputs

❖ Learned hidden layer representation:

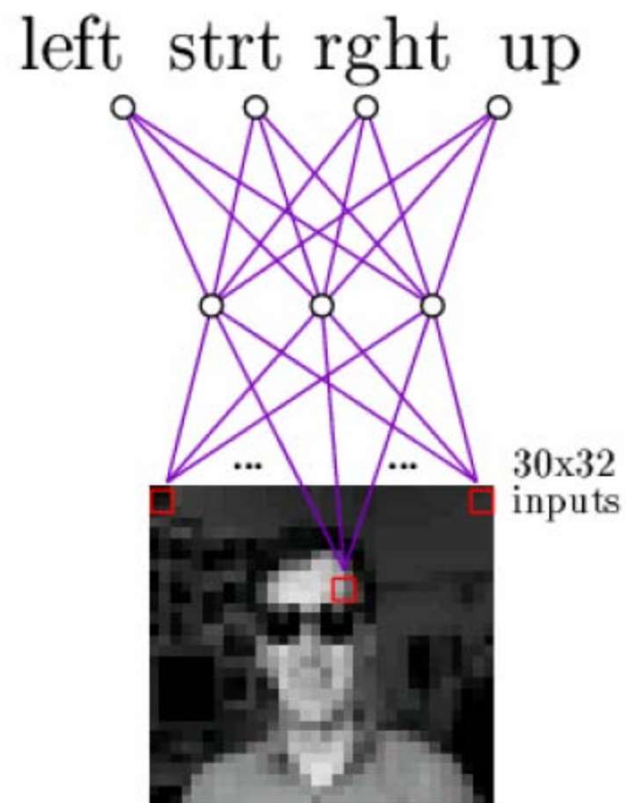| Input | Hidden Values | | | Output |
|---|---|---|---|---|
| 10000000 → | .89 | .04 | .08 → | 10000000 |
| 01000000 → | .01 | .11 | .88 → | 01000000 |
| 00100000 → | .01 | .97 | .27 → | 00100000 |
| 00010000 → | .99 | .97 | .71 → | 00010000 |
| 00001000 → | .03 | .05 | .02 → | 00001000 |
| 00000100 → | .22 | .99 | .99 → | 00000100 |
| 00000010 → | .80 | .01 | .98 → | 00000010 |
| 00000001 → | .60 | .94 | .01 → | 00000001 |

# Expressive Capabilities of ANNs

❖ 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 arbitrary small
    error, by network with one hidden layer [Cybenko 1989; Hornik et al 1989]
  - ❖ Any function can be approximated to arbitrary accuracy by a network with two hidden layers [Cybenko 1988].
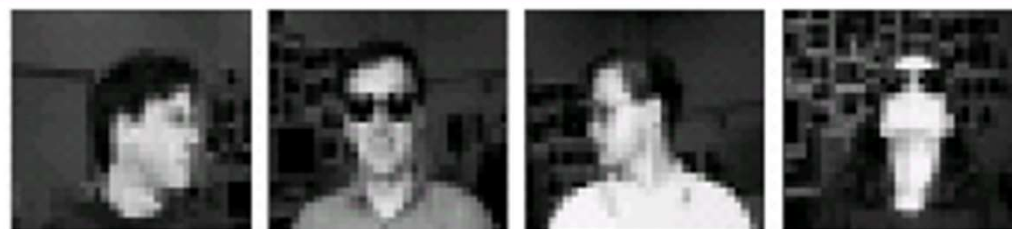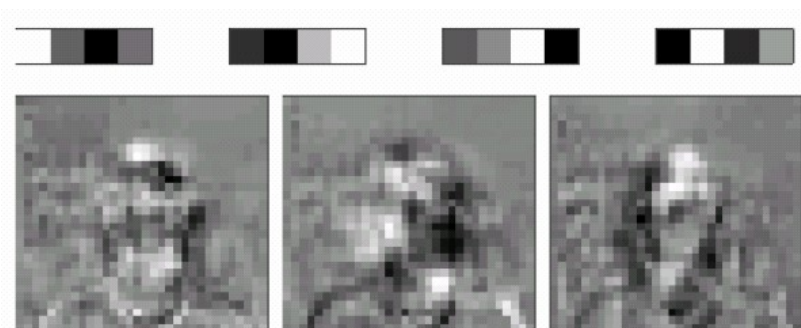
# Application: ANN for Face Reco.

❖ The model



left strt rght up

30x32 inputs

❖ The learned hidden unit weights



Typical input images

http://www.cs.cmu.edu/~tom/faces.html

# Artificial neural networks – what you should know

- ❖ Highly expressive non-linear functions
- ❖ Highly parallel network of logistic function units
- ❖ Minimizing sum of squared training errors
  - ❖ Gives MLE estimates of network weights if we assume zero mean Gaussian noise on output values
- ❖ Minimizing sum of sq errors plus weight squared(regularization)
  - ❖ MAP estimates assuming weight priors are zero mean Gaussian
- ❖ Gradient descent as training procedure
  - ❖ How to derive your own gradient descent procedure
- ❖ Discover useful representations at hidden units
- ❖ Local minima is greatest problem
- ❖ Overfitting, regularization, early stopping