# IT/PC/B/T/411

# Machine Learning
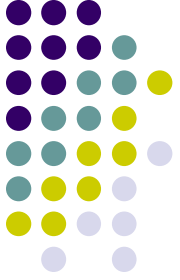
## Neural Networks

**Dr. Pawan Kumar Singh**

Department of Information Technology

Jadavpur University

pawankrsingh.cse@gmail.com

+91-6291555693

# Artificial Neural Networks

- Computational models inspired by the human brain:
  - Algorithms that try to mimic the brain.

  - Massively parallel, distributed system, made up of simple processing units (neurons)

  - Synaptic connection strengths among neurons are used to store the acquired knowledge.

  - Knowledge is acquired by the network from its environment through a learning process

# History

- late-1800's - Neural Networks appear as an analogy to biological systems

- 1960's and 70's – Simple neural networks appear
  - Fall out of favor because the perceptron is not effective by itself, and there were no good algorithms for multilayer nets

- 1986 – Backpropagation algorithm appears
  - Neural Networks have a resurgence in popularity
  - More computationally expensive

# Applications of ANNs

- ANNs have been widely used in various domains for:
  - Pattern recognition
  - Function approximation
  - Associative memory

*digit classificat*

*(+7 % uccurey)*

# Properties

- Inputs are flexible *Imp.*
  - any real values
  - Highly correlated or independent
- Target function may be discrete-valued, real-valued, or vectors of discrete or real values
  - Outputs are real numbers between 0 and 1
- Resistant to errors in the training data *Imp*
- Long training time *Imp*
- Fast evaluation *Imp*
- The function produced can be difficult for humans to interpret *Imp*

# When to consider neural networks

- Input is high-dimensional discrete or raw-valued
- Output is discrete or real-valued
- Output is a vector of values
- Possibly noisy data
- Form of target function is unknown
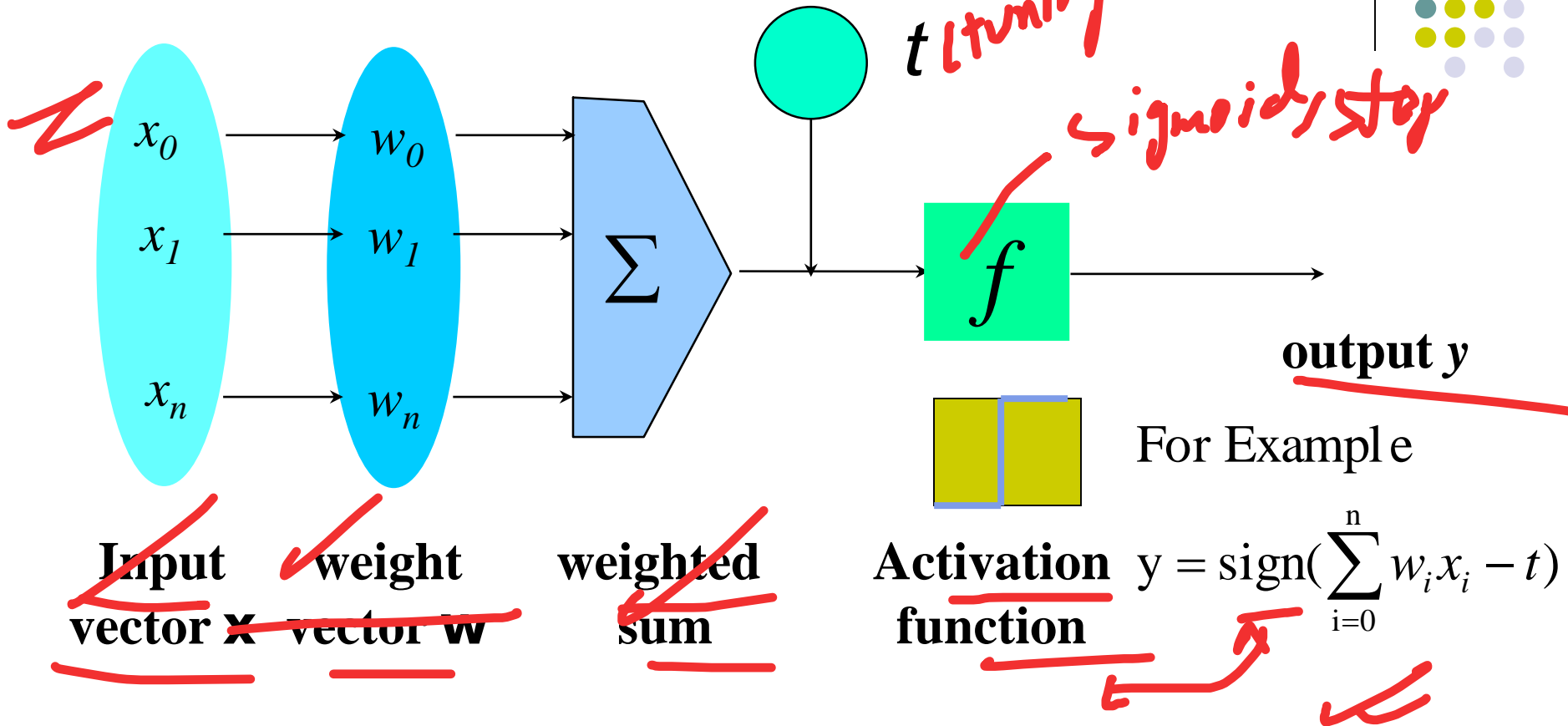- Human readability of the result is not important

**Examples:**

- Speech phoneme recognition
- Image classification

Digit Classification

- Financial prediction

Bank Credit Card churn

# A Neuron (= a perceptron)

*tuning parameters*

*sigmoid, step*

$x_0$

$x_1$

$x_n$

$w_0$

$w_1$

$w_n$

$\Sigma$

$t$

$f$

**output $y$**

For Example

**Input vector x**   **weight vector w**   **weighted sum**   **Activation function**   $y = \text{sign}(\sum_{i=0}^{n} w_i x_i - t)$
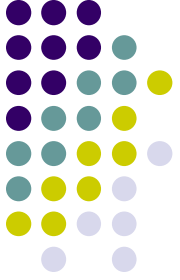
*dot product*

$\sigma(w^T x + b)$

- The *n*-dimensional input vector **x** is mapped into variable y by means of the scalar product and a nonlinear function mapping
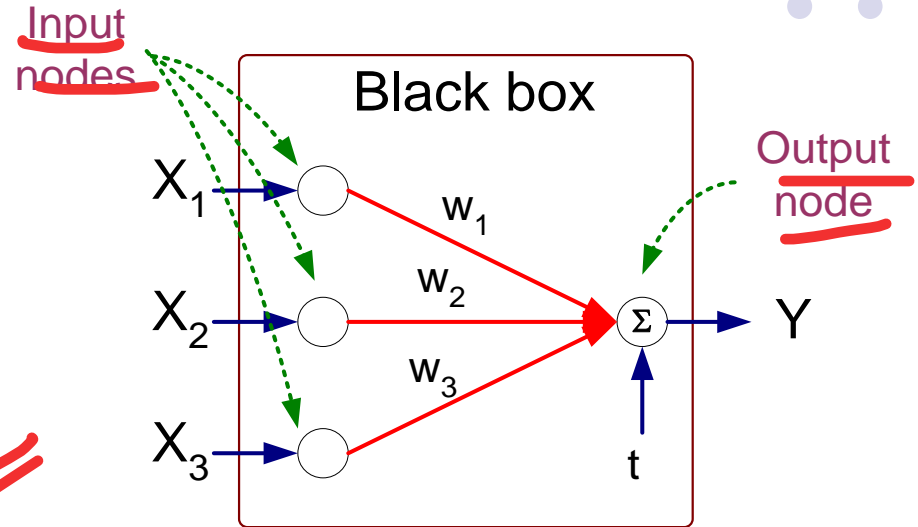
# Perceptron

- Basic unit in a neural network

- Linear separator

- Parts
  - N inputs, $x_1$ ... $x_n$
  - Weights for each input, $w_1$ ... $w_n$
  - A bias input $x_0$ (constant) and associated weight $w_0$
  - Weighted sum of inputs, $y = w_0x_0 + w_1x_1 + ... + w_nx_n$
  - A threshold function or activation function,
    - i.e 1 if $y > t$, -1 if $y <= t$

# Artificial Neural Networks (ANN)

- Model is an assembly of inter-connected nodes and weighted links

- Output node sums up each of its input value according to the weights of its links

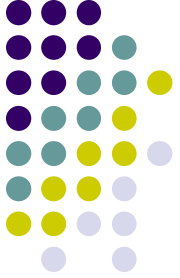- Compare output node against some threshold t

**Input nodes**

**Black box**

$X_1$

$w_1$

$X_2$ $w_2$ $\Sigma$ Y

$w_3$

$X_3$

t

**Output node**

**Perceptron Model**

$$Y = I(\sum_i w_i x_i - t)$$ or

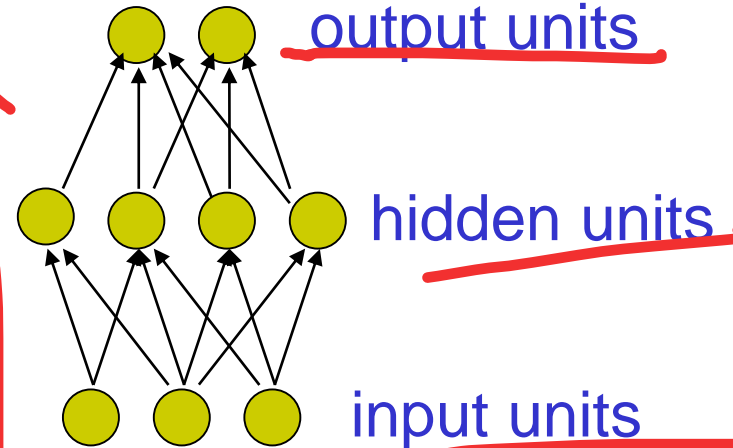$$Y = sign(\sum_i w_i x_i - t)$$

# Types of connectivity

**IMP**

- Feedforward networks
  - These compute a series of transformations
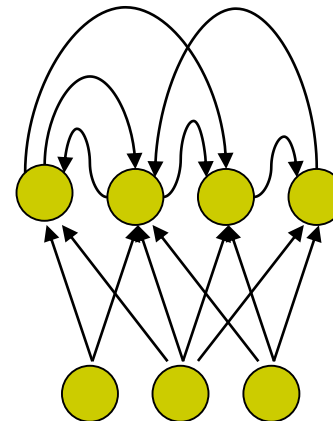  - Typically, the first layer is the input and the last layer is the output.
- Recurrent networks
  - These have directed cycles in their connection graph. They can have complicated dynamics.
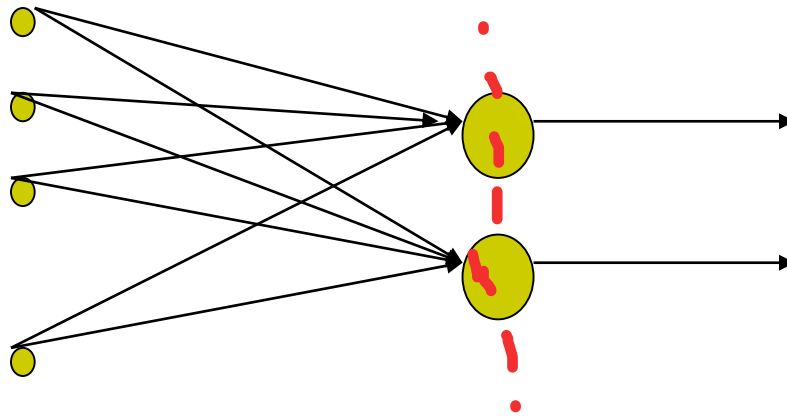  - More biologically realistic.

feedForward Networks

output units

hidden units

input units

Recurrent networks

# Different Network Topologies

- Single layer feed-forward networks
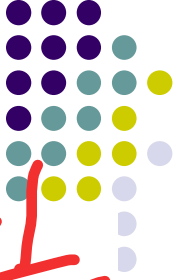  - Input layer projecting into the output layer



Single layer
network

Input
layer

Output
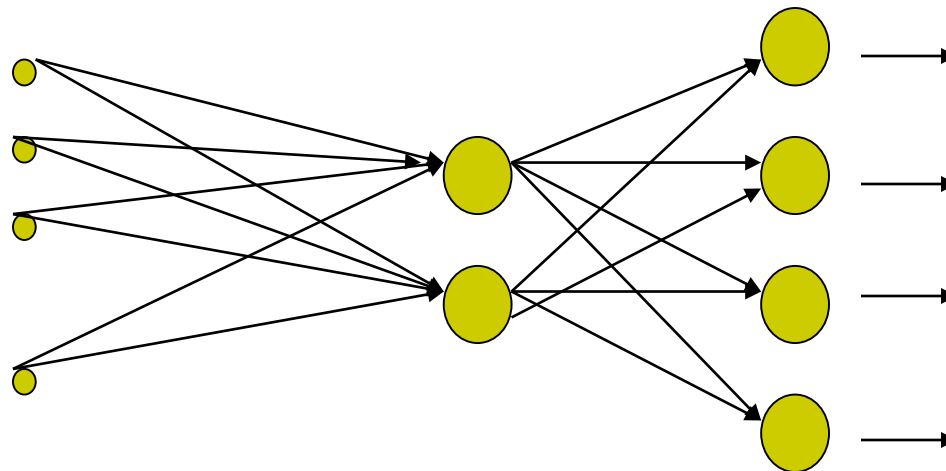layer

Classification
network.

# Different Network Topologies

**Multi-layer feed-forward networks** /1 *generate*

One or more hidden layers. Input projects only from previous layers onto a layer.

b/ *multi class classification*

2-layer or
1-hidden layer
*fully connected*
network

Input
layer

Hidden
layer

Output
layer

# Different Network Topologies

- Multi-layer feed-forward networks

$3 \times 3$

$3 \times 2$

$2 \times 1$

$x_1$

$x_2$

$x_3$

$h_\Theta(x)$

Layer 1  Layer 2  Layer 3  Layer 4

Input layer

Hidden layers

Output layer

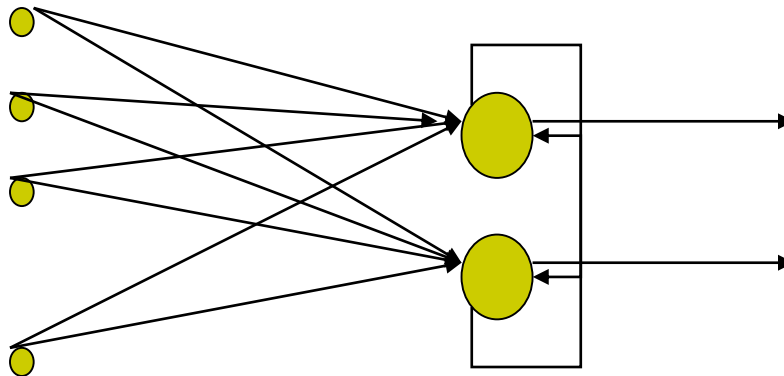$6 + 3 + 6 + 2 + 2 + 1 = 20$ training parameters

# Different Network Topologies

- Recurrent networks

  - A network with feedback, where some of its inputs are connected to some of its outputs (discrete time).

Recurrent network

Input layer

Output layer

# Algorithm for learning ANN

- Initialize the weights ($w_0$, $w_1$, ..., $w_k$)

- Adjust the weights in such a way that the output of ANN is consistent with class labels of training examples
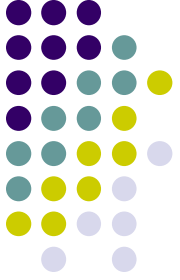
  - Error function:

  $$E = \sum_i \left[ Y_i - f(w_i, X_i) \right]^2$$

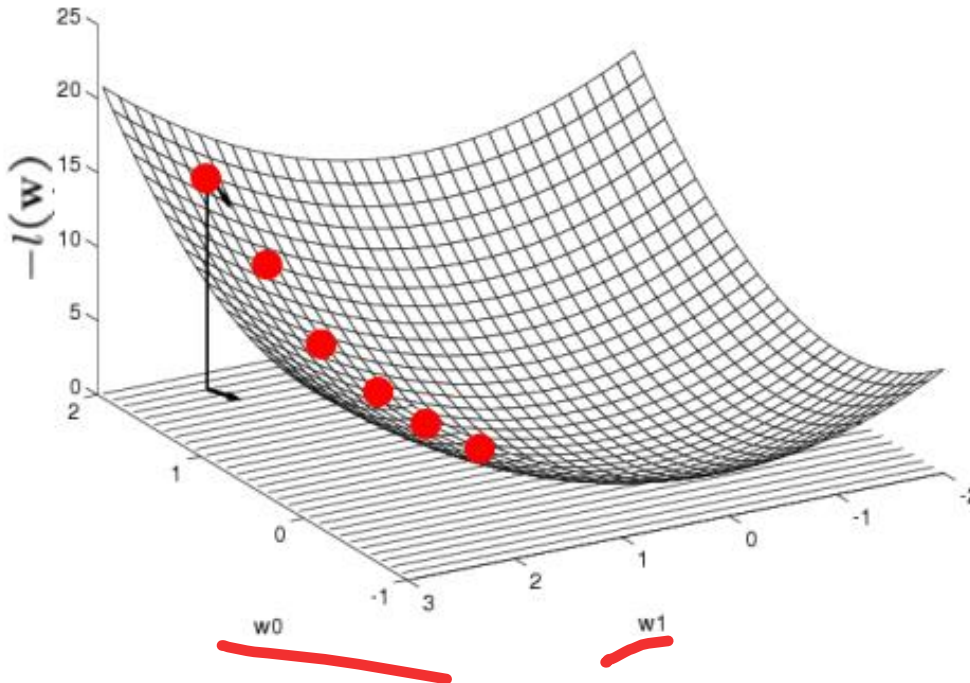  - Find the weights w's that minimize the above error function

    - e.g., gradient descent, backpropagation algorithm

# Optimizing concave/convex function

- Maximum of a concave function = minimum of a convex function

**Gradient ascent (concave) / Gradient descent (convex)**

**Gradient:**

$$\nabla_{\mathbf{w}} l(\mathbf{w}) = [\frac{\partial l(\mathbf{w})}{\partial w_0}, \ldots, \frac{\partial l(\mathbf{w})}{\partial w_d}]'$$
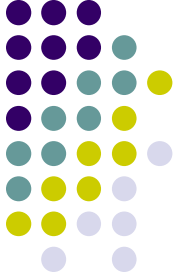
**Update rule:**    Learning rate, $\eta > 0$

$$\triangle \mathbf{w} = \eta \nabla_{\mathbf{w}} l(\mathbf{w})$$

$$w_i^{(t+1)} \leftarrow w_i^{(t)} + \eta \frac{\partial l(\mathbf{w})}{\partial w_i}\Big|_t$$

Gradient ascent rule

# GRADIENT DESCENT

Suppose we have a scalar function

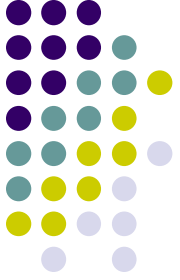$$f(w) : \Re \to \Re$$

We want to find a local minimum.

Assume our current weight is $w$

GRADIENT DESCENT RULE:

$$w \leftarrow w - \eta \frac{\partial}{\partial w} f(w)$$

$\eta$ is called the LEARNING RATE. A small positive number, e.g. $\eta = 0.05$

# Gradient Descent in "m" Dimensions

Given $\mathrm{f}(\mathbf{w}) : \Re^m \to \Re$

$$\nabla \mathrm{f}(w) = \begin{pmatrix} \dfrac{\partial}{\partial w_1} \mathrm{f}(w) \\ \vdots \\ \dfrac{\partial}{\partial w_m} \mathrm{f}(w) \end{pmatrix}$$ points in direction of steepest ascent.

$\left| \nabla \mathrm{f}(w) \right|$ is the gradient in that direction

GRADIENT DESCENT RULE: $\quad \mathrm{w} \leftarrow \mathrm{w} - \eta \nabla \mathrm{f}(\mathrm{w})$

Equivalently

$$w_j \leftarrow w_j - \eta \frac{\partial}{\partial w_j} \mathrm{f}(w)$$ ....where $w_j$ is the $j$th weight

"just like a linear feedback system"

# Linear Perceptrons

They are multivariate linear models:

$$\text{Out}(\boldsymbol{x}) = \boldsymbol{w}^{\mathrm{T}}\boldsymbol{x}$$

And "training" consists of minimizing sum-of-squared residuals by gradient descent.

$$E = \sum_{k} \left( \text{Out}(\mathbf{x}_k) - y_k \right)^2$$

$$= \sum_{k} \left( \mathbf{w}^{\mathrm{T}}\mathbf{x}_k - y_k \right)^2$$

QUESTION:  Derive the perceptron training rule.

# Linear Perceptron Training Rule

$$E = \sum_{k=1}^{R} (y_k - \mathbf{w}^T \mathbf{x}_k)^2$$

Gradient descent tells us we should update **w** thusly if we wish to minimize $E$:

$$w_j \leftarrow w_j - \eta \frac{\partial E}{\partial w_j}$$

So what's $\dfrac{\partial E}{\partial w_j}$?

# Linear Perceptron Training Rule
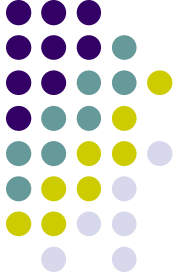
$$E = \sum_{k=1}^{R} (y_k - \mathbf{w}^T \mathbf{x}_k)^2$$

Gradient descent tells us we should update **w** thusly if we wish to minimize $E$:

$$w_j \leftarrow w_j - \eta \frac{\partial E}{\partial w_j}$$

So what's $\dfrac{\partial E}{\partial w_j}$ ?

$$\frac{\partial E}{\partial w_j} = \sum_{k=1}^{R} \frac{\partial}{\partial w_j} (y_k - \mathbf{w}^T \mathbf{x}_k)^2$$

$$= \sum_{k=1}^{R} 2(y_k - \mathbf{w}^T \mathbf{x}_k) \frac{\partial}{\partial w_j} (y_k - \mathbf{w}^T \mathbf{x}_k)$$

$$= -2 \sum_{k=1}^{R} \delta_k \frac{\partial}{\partial w_j} \mathbf{w}^T \mathbf{x}_k$$

*...where...*

$$\delta_k = y_k - \mathbf{w}^T \mathbf{x}_k$$

$$= -2 \sum_{k=1}^{R} \delta_k \frac{\partial}{\partial w_j} \sum_{i=1}^{m} w_i x_{ki}$$

$$= -2 \sum_{k=1}^{R} \delta_k x_{kj}$$

# Linear Perceptron Training Rule

$$E = \sum_{k=1}^{R} (y_k - \mathbf{w}^T \mathbf{x}_k)^2$$

Gradient descent tells us we should update **w** thusly if we wish to minimize $E$:

$$w_j \leftarrow w_j - \eta \frac{\partial E}{\partial w_j}$$

*...where...*

$$\frac{\partial E}{\partial w_j} = -2 \sum_{k=1}^{R} \delta_k x_{kj}$$

$$\boxed{w_j \leftarrow w_j + 2\eta \sum_{k=1}^{R} \delta_k x_{kj}}$$

We frequently neglect the 2 (meaning we halve the learning rate)

# The "Batch" perceptron algorithm

1) Randomly initialize weights $w_1$ $w_2$ ... $w_m$

2) Get your dataset (append 1's to the inputs if you don't want to go through the origin).

3) for $i$ = 1 to R

$$\delta_i := y_i - \mathbf{W}^{\mathrm{T}}\mathbf{x}_i$$

4) for $j$ = 1 to m

$$w_j \leftarrow w_j + \eta \sum_{i=1}^{R} \delta_i x_{ij}$$

5) if $\sum \delta_i^2$ stops improving then stop.  Else loop back to 3.   *epoch* .

$$\delta_i \leftarrow y_i - \mathrm{w}^{\mathrm{T}}\mathrm{x}_i$$

$$w_j \leftarrow w_j + \eta\delta_i x_{ij}$$
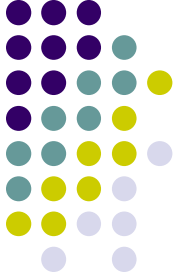
**A RULE KNOWN BY MANY NAMES**

**The LMS Rule**

**The Widrow Hoff rule**

**The delta rule**

*Classical conditioning*

**The adaline rule**

# Perceptrons and Boolean Functions

- Can learn any disjunction of literals

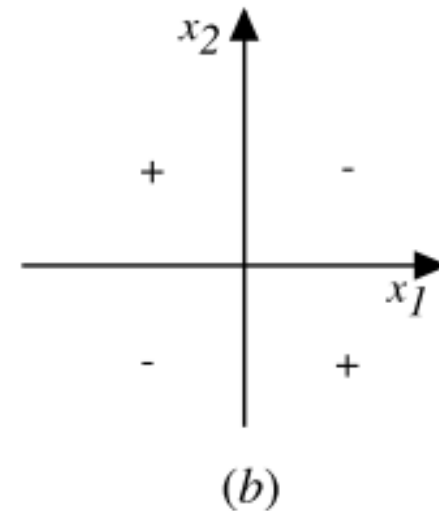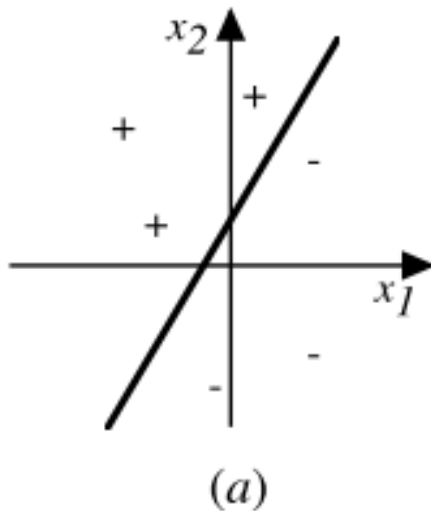  e.g. $x_1 \wedge \sim x_2 \wedge \sim x_3 \wedge x_4 \wedge x_5$
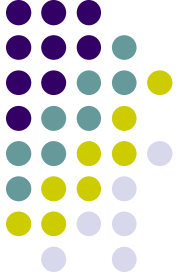
- Can learn majority function

  $$f(x_1, x_2 \ldots x_n) = \begin{cases} 1 \text{ if } n/2 \ x_i\text{'s or more are} = 1 \\ 0 \text{ if less than } n/2 \ x_i\text{'s are} = 1 \end{cases}$$

- What about the exclusive or function?

  $$f(x_1, x_2) = x_1 \ \forall \ x_2 =$$
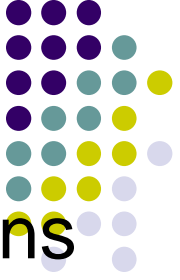  $$(x_1 \wedge \sim x_2) \vee (\sim x_1 \wedge x_2)$$

# Decision surface of a perceptron

$(a)$

$(b)$

- Decision surface is a hyperplane
  - Can capture linearly separable classes
- Non-linearly separable
  - Use a network of them

# Multilayer Networks

The class of functions representable by perceptrons is limited

$$\text{Out}(x) = g\left(\mathbf{w}^{\text{T}}\mathbf{x}\right) = g\left(\sum_j w_j x_j\right)$$

*Use a wider representation !*

$$\text{Out}(x) = g\left(\sum_j W_j g\left(\sum_k w_{jk} x_{jk}\right)\right)$$

rows

This is a nonlinear function
Of a linear combination
Of non linear functions
Of linear combinations of inputs

# A 1-HIDDEN LAYER NET

$N_{INPUTS} = 2$      (INS)

$N_{HIDDEN} = 3$      (HID)    no - bias

$$v_1 = g\left(\sum_{k=1}^{N_{INS}} w_{1k} x_k\right)$$

$$v_2 = g\left(\sum_{k=1}^{N_{INS}} w_{2k} x_k\right)$$

$$v_3 = g\left(\sum_{k=1}^{N_{INS}} w_{3k} x_k\right)$$

$$Out = g\left(\sum_{k=1}^{N_{HID}} W_k v_k\right)$$

$x_1$

$x_2$

$w_{11}$

$w_{21}$

$w_{31}$

$w_{12}$

$w_{22}$

$w_{32}$

$W_1$

$w_2$

$w_3$

# OTHER NEURAL NETS



2-Hidden layers + Constant Term

"JUMP" CONNECTIONS



$$\text{Out} = g\left(\sum_{k=1}^{N_{INS}} w_{0k} x_k + \sum_{k=1}^{N_{HID}} W_k v_k\right)$$

# Multi-layer Networks

- Linear units inappropriate
  - No more expressive than a single layer
- „Introduce non-linearity
  - Threshold not differentiable
- „Use sigmoid function

# The Sigmoid

$$g(h) = \frac{1}{1 + \exp(-h)}$$



Now we choose  **w**  to minimize

$$\sum_{i=1}^{R} \left[ y_i - \mathrm{Out}(\mathrm{x}_i) \right]^2 = \sum_{i=1}^{R} \left[ y_i - g(\mathrm{w}^\mathrm{T} \mathrm{x}_i) \right]^2$$

# Sigmoid Unit



$x_1$ $w_1$ $x_0 = 1$
$x_2$ $w_2$ $w_0$
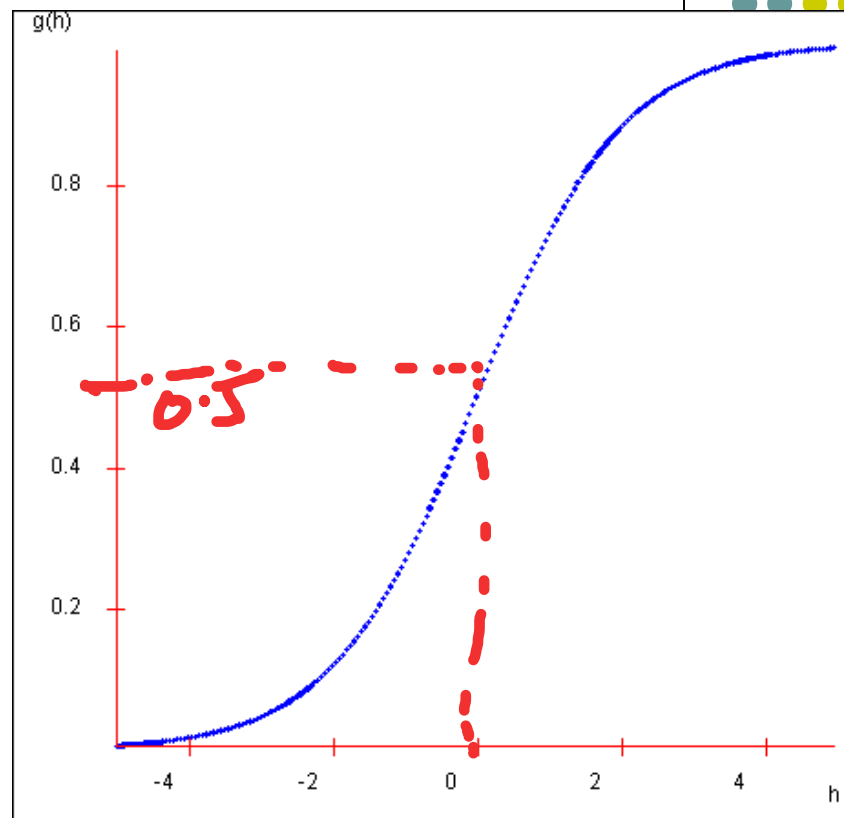$w_d$
$x_d$

$$net = \sum_{i=0}^{d} w_i x_i$$

$$o = \sigma(net) = \frac{1}{1 + e^{-net}}$$

*About bit* (handwritten annotation)

$\sigma(x)$ is the sigmoid function /activation function (also linear, threshold)

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

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

We can derive gradient decent rules to train

- One sigmoid unit
- *Multilayer networks* of sigmoid units →
  Backpropagation

# Backpropagation

$$\text{Out}(x) = g\left(\sum_j W_j g\left(\sum_k w_{jk} x_k\right)\right)$$

Find a set of weights $\{W_j\}, \{w_{jk}\}$

to minimize

$$\sum_i \left(y_i - \text{Out}(x_i)\right)^2$$

by gradient descent.

That's it!

That's the backpropagation algorithm.

# Backpropagation

- Iteratively process a set of training tuples & compare the network's prediction with the actual known target value

- For each training tuple, the weights are modified to **minimize the mean squared error** between the network's prediction and the actual target value

- Modifications are made in the "**backwards**" direction: from the output layer, through each hidden layer down to the first hidden layer, hence "**backpropagation**"

- Steps
  - Initialize weights (to small random #s) and biases in the network
  - Propagate the inputs forward (by applying activation function)
  - Backpropagate the error (by updating weights and biases)
  - Terminating condition (when error is very small, etc.)

**function** BACK-PROP-LEARNING(*examples*, *network*) **returns** a neural network
  **inputs**: *examples*, a set of examples, each with input vector **x** and output vector **y**
      *network*, a multilayer network with $L$ layers, weights $w_{i,j}$, activation function $g$
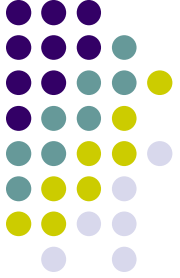  **local variables**: $\Delta$, a vector of errors, indexed by network node

  **repeat**
    **for each** weight $w_{i,j}$ in *network* **do**
      $w_{i,j} \leftarrow$ a small random number
    **for each** example $(\mathbf{x}, \mathbf{y})$ **in** *examples* **do**
      / $\star$ *Propagate the inputs forward to compute the outputs* $\star$ /
      **for each** node $i$ in the input layer **do**
        $a_i \leftarrow x_i$
      **for** $\ell = 2$ **to** $L$ **do**
        **for each** node $j$ in layer $\ell$ **do**
          $in_j \leftarrow \sum_i w_{i,j}\, a_i$
          $a_j \leftarrow g(in_j)$
      / $\star$ *Propagate deltas backward from output layer to input layer* $\star$ /
      **for each** node $j$ in the output layer **do**
        $\Delta[j] \leftarrow g'(in_j) \times (y_j - a_j)$
      **for** $\ell = L - 1$ **to** $1$ **do**
        **for each** node $i$ in layer $\ell$ **do**
          $\Delta[i] \leftarrow g'(in_i) \sum_j w_{i,j}\, \Delta[j]$
      / $\star$ *Update every weight in network using deltas* $\star$ /
      **for each** weight $w_{i,j}$ in *network* **do**
        $w_{i,j} \leftarrow w_{i,j} + \alpha \times a_i \times \Delta[j]$
  **until** some stopping criterion is satisfied
  **return** *network*
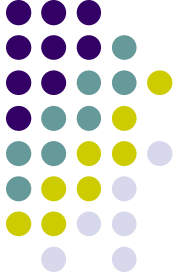
# How A Multi-Layer Neural Network Works?

- The **inputs** to the network correspond to the attributes measured for each training tuple

- Inputs are fed simultaneously into the units making up the **input layer**

- They are then weighted and fed simultaneously to a **hidden layer**

- The number of hidden layers is arbitrary, although usually only one

- The weighted outputs of the last hidden layer are input to units making up the **output layer**, which emits the network's prediction

- The network is **feed-forward** in that none of the weights cycles back to an input unit or to an output unit of a previous layer

- From a statistical point of view, networks perform **nonlinear regression**: Given enough hidden units and enough training samples, they can closely approximate any function

# Defining a Network Topology

- First decide the **network topology:** # of units in the *input layer*, # of *hidden layers* (if > 1), # of units in *each hidden layer*, and # of units in the *output layer*

- Normalizing the input values for each attribute measured in the training tuples to [0.0—1.0]

- One **input** unit per domain value, each initialized to 0

- **Output**, if for classification and more than two classes, one output unit per class is used

- Once a network has been trained and its accuracy is **unacceptable**, repeat the training process with a *different network topology* or a *different set of initial weights*

# Backpropagation and Interpretability

- Efficiency of backpropagation: Each **epoch** (one interation through the training set) takes O(|D| * *w*), with |D| tuples and *w* weights, but # of epochs can be exponential to n, the number of inputs, in the worst case

- **Rule extraction from networks:** network pruning
  - Simplify the network structure by removing weighted links that have the least effect on the trained network
  - Then perform link, unit, or activation value clustering
  - The set of input and activation values are studied to derive rules describing the relationship between the input and hidden unit layers

- **Sensitivity analysis:** assess the impact that a given input variable has on a network output.  The knowledge gained from this analysis can be represented in rules
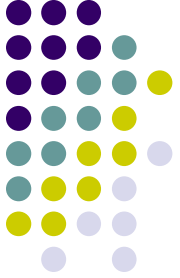
# Neural Network as a Classifier

- Weakness
  - Long training time
  - Require a number of parameters typically best determined empirically, e.g., the network topology or "structure."
  - Poor interpretability: Difficult to interpret the symbolic meaning behind the learned weights and of "hidden units" in the network
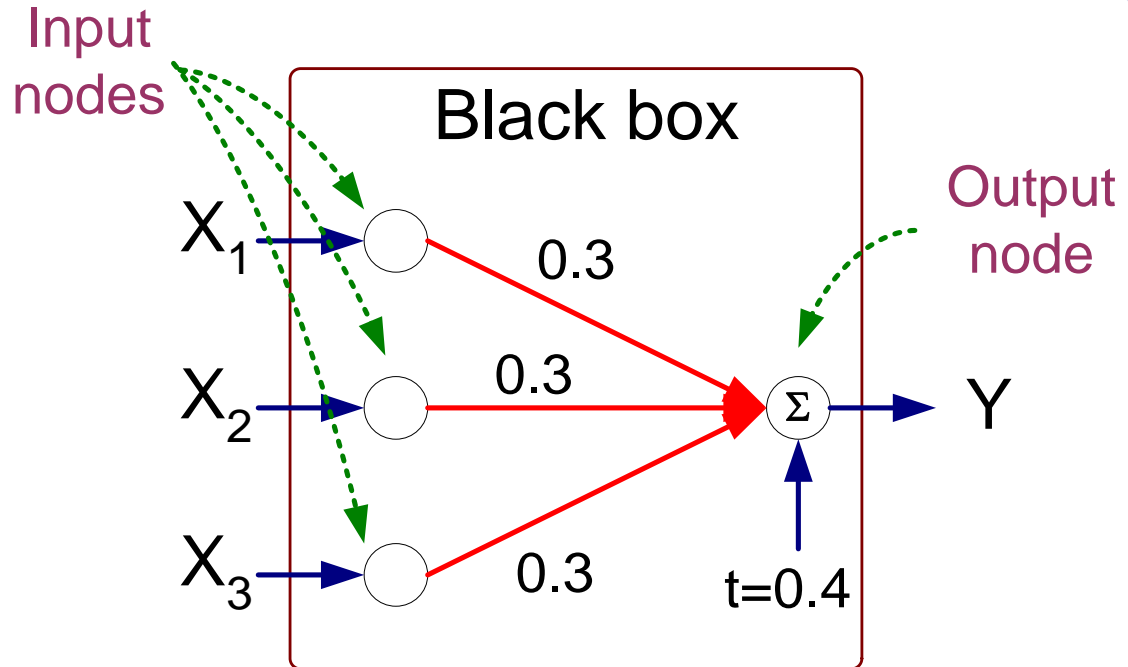
- Strength
  - High tolerance to noisy data
  - Ability to classify untrained patterns
  - Well-suited for continuous-valued inputs and outputs
  - Successful on a wide array of real-world data
  - Algorithms are inherently parallel
  - Techniques have recently been developed for the extraction of rules from trained neural networks
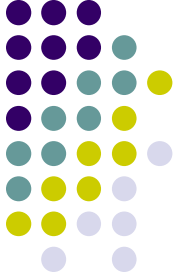
# Artificial Neural Networks (ANN)

| $X_1$ | $X_2$ | $X_3$ | Y |
|-------|-------|-------|---|
| 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 |
| 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 |
| 0 | 0 | 0 | 0 |

Input nodes

Output node

Black box

$X_1$ → 0.3

$X_2$ → 0.3 → Σ → Y

$X_3$ → 0.3    t=0.4

$$Y = I(0.3X_1 + 0.3X_2 + 0.3X_3 - 0.4 > 0)$$

$$\text{where } I(z) = \begin{cases} 1 & \text{if } z \text{ is true} \\ 0 & \text{otherwise} \end{cases}$$

# Linear Perceptrons

They are multivariate linear models:

$$\text{Out}(\boldsymbol{x}) = \boldsymbol{w}^{\mathsf{T}}\boldsymbol{x}$$

And "training" consists of minimizing sum-of-squared residuals by gradient descent.

$$E = \sum_k \left(\text{Out}(x_k) - y_k\right)^2$$

$$= \sum_k \left(w^{\mathsf{T}}x_k - y_k\right)^2$$

QUESTION:  Derive the perceptron training rule.

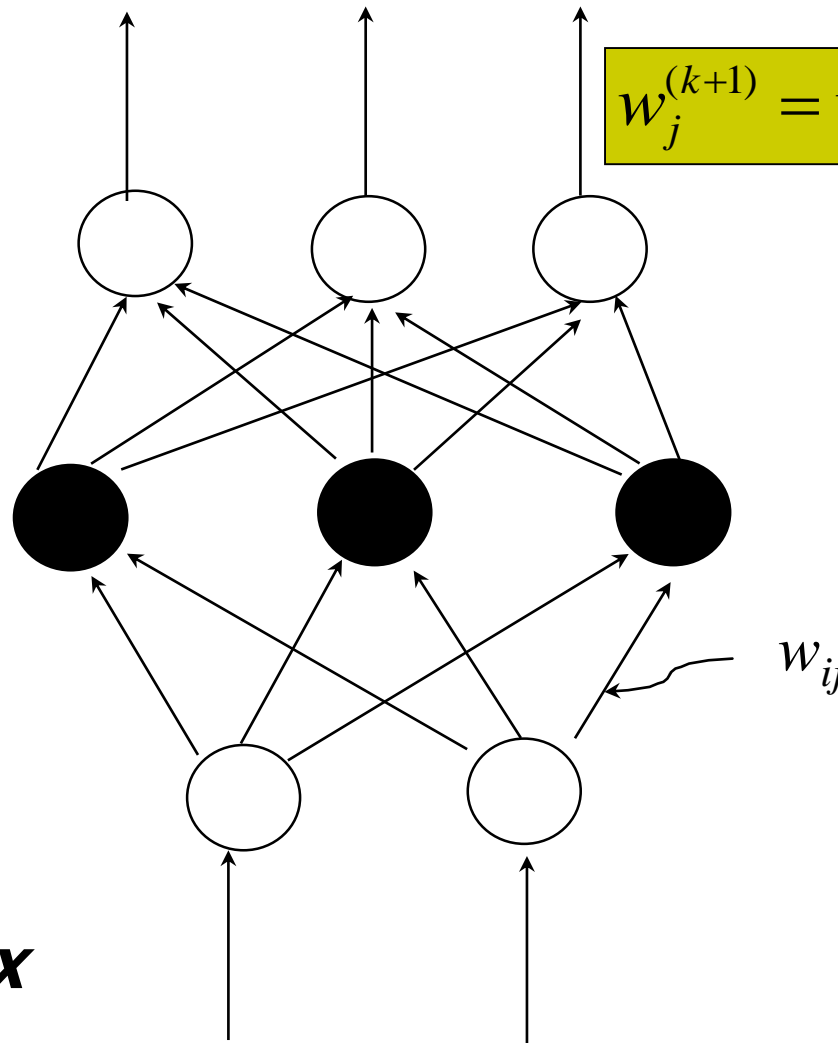# A Multi-Layer Feed-Forward Neural Network

**Output vector**

**Output layer**

$$w_j^{(k+1)} = w_j^{(k)} + \lambda(y_i - \hat{y}_i^{(k)})x_{ij}$$

**Hidden layer**

$$w_{ij}$$

**Input layer**

**Input vector: X**

# General Structure of ANN



x$_1$    x$_2$    x$_3$    x$_4$    x$_5$

Input Layer

Hidden Layer

Output Layer

y

Input          Neuron $i$          Output

I$_1$   w$_{i1}$

I$_2$   w$_{i2}$       $S_i$   Activation function $g(S_i)$   $O_i$        O$_i$

I$_3$   w$_{i3}$

threshold, t

Training ANN means learning the weights of the neurons