

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205  
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



# Chapter 5

## Classical Artificial Neural Networks

Xuegong Zhang  
September 30, 2021



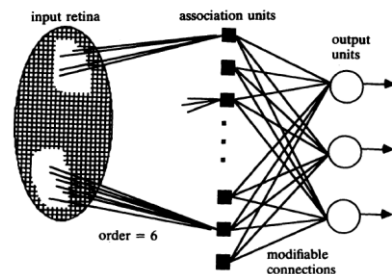
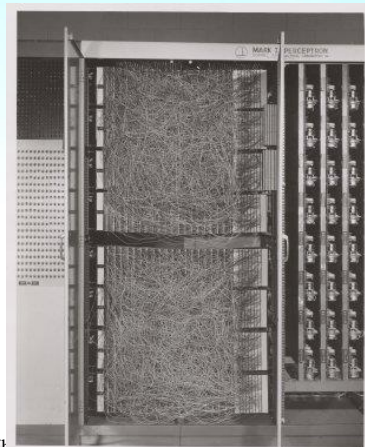
Xuegong Zhang

## Perceptron



Frank Rosenblatt, *The Perceptron – a perceiving and recognizing automaton*,  
Report 85-460-1, Cornell Aeronautical Laboratory, Jan. 1957

$$y = \text{sgn}\left(\sum_{i=1}^d w_i x_i + w_0\right)$$



M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

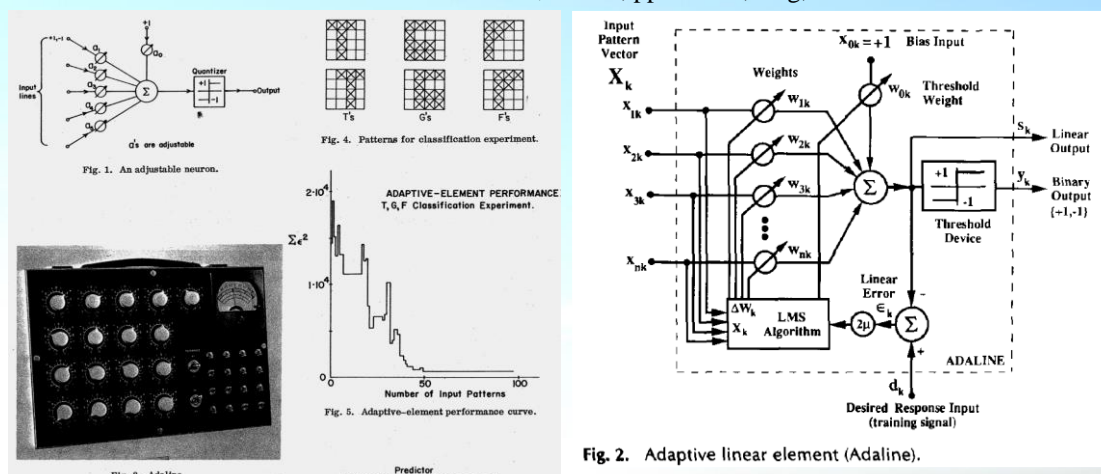
Xuegong Zi

<https://en.wikipedia.org/wiki/Perceptron>

2

# ADALINE

Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960



Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960

Widrow & Lehr, 30 years of adaptive neural networks: Perceptron, Madaline, and Backpropagation, *Proceedings of the IEEE*, 78(9): 1415-1442, 1990

Xuegong Zhang

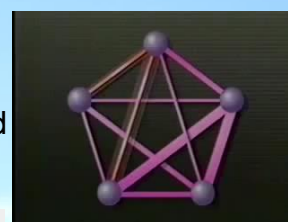
3

## An early prediction on perceptrons

- In a 1958 press conference organized by the US Navy, Rosenblatt made statements about the perceptron that caused a **heated controversy** among the fledgling AI community; based on Rosenblatt's statements, *The New York Times* reported:

**The Navy revealed the embryo of an electronic computer today that it expects will be able to walk, talk, see, write, reproduce itself and be conscious of its existence. Later perceptrons will be able to recognize people and call out their names and instantly translate speech in one language to speech and writing in another language, it was predicted.**<sup>17</sup>

M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996



<https://pmirla.github.io/2016/08/16/AI-Winter.html>

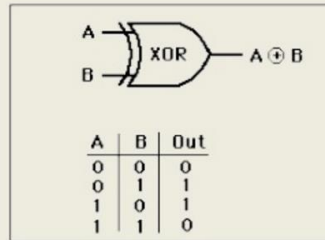
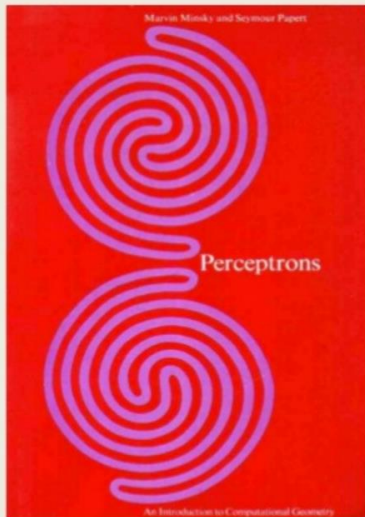


Xuegong Zhang

4

## Criticisms on Perceptrons

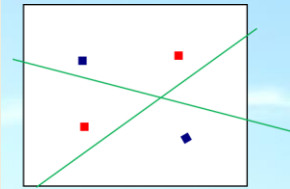
1969: Perceptrons can't do XOR!



<http://hyperphysics.phy-astr.gsu.edu/hbase/electronic/ietron/xor.gif>



Minsky & Papert



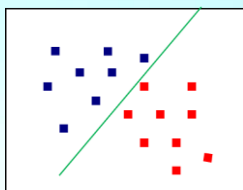
5

<http://www.i-programmer.info/images/stories/BabBag/AI/book.jpg>

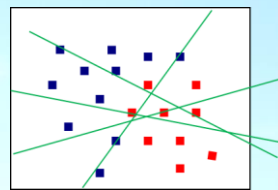
<https://c.constructingkids.files.wordpress.com/2013/05/minsky-papert-71-csolomon-x640.jpg>

<https://pmirla.github.io/2016/08/16/AI-Winter.html>

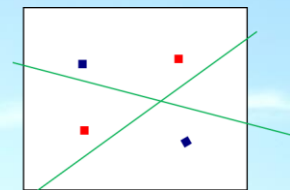
## Linearly non-separable cases



Linearly separable



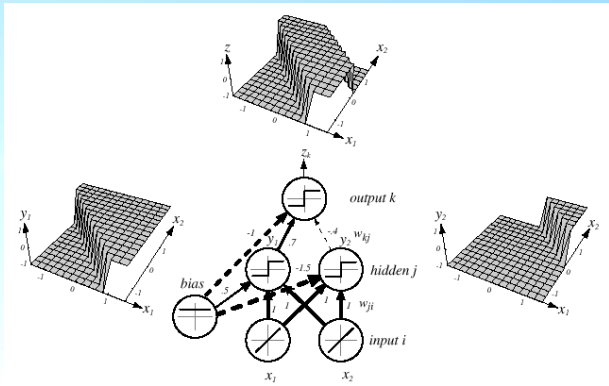
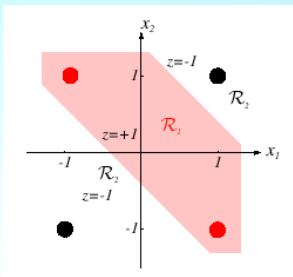
Not perfectly separable,  
but can achieve some  
reasonable separation  
with a linear classifier.



No reasonable linear  
solution or approximation.



Implement nonlinear functions with multiple neurons



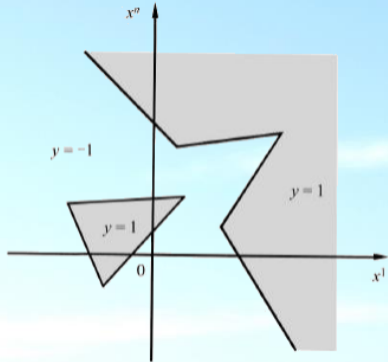
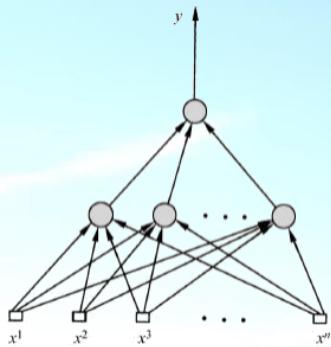
Duda, Hart & Stork, Pattern Classification (2nd edition), John Wiley & Sons, 2001

6

7



Early forms of multiple-layered perceptron (neural networks)

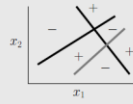


Xuegong Zhang

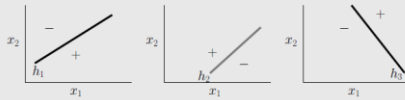
8

## Boolean decomposition of perceptrons

Consider a target function  $f$  whose '+' and '-' regions are illustrated below.

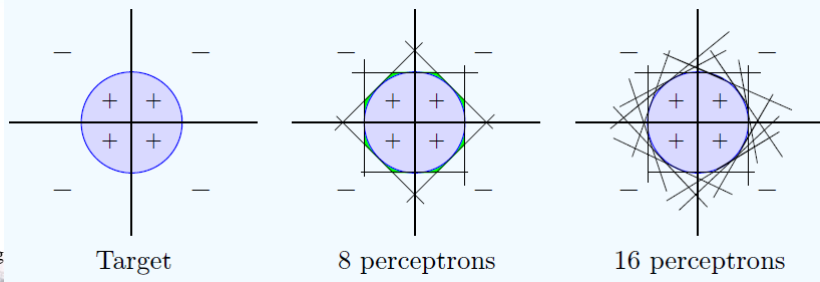
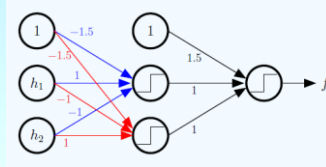
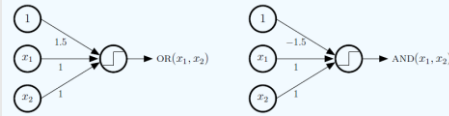


The target  $f$  has three perceptron components  $h_1, h_2, h_3$ :



Show that

$$f = \bar{h}_1 \bar{h}_2 h_3 + h_1 \bar{h}_2 h_3 + h_1 h_2 \bar{h}_3.$$



Xuegong Zhang

9

However, this field was closed  
for ~20 years (1960s-1980s)

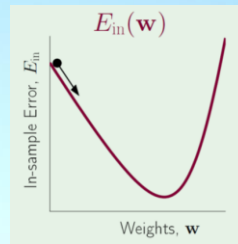
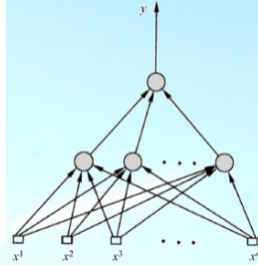
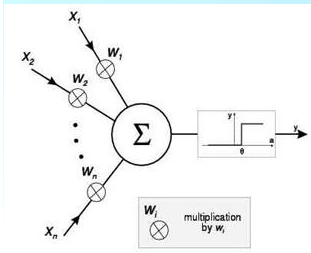
Xuegong Zhang

10





## Why didn't it work?



- The old trick for learning?
  - Cannot get gradient

Discriminant function

$$g_k(\mathbf{x}) \equiv y_k = \text{Sign} \left( \sum_j w_{jk} \text{Sign} \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$



Xuegong Zhang

11

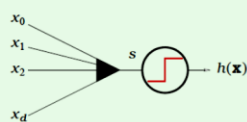
## Recall: Three major types of linear machines



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

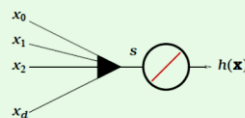
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



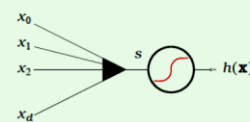
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

$$h(\mathbf{x}) = \theta(s)$$



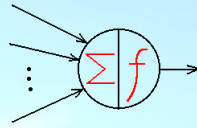
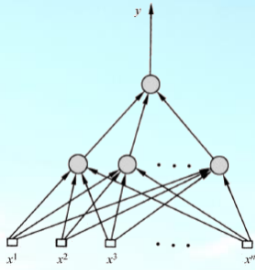
Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

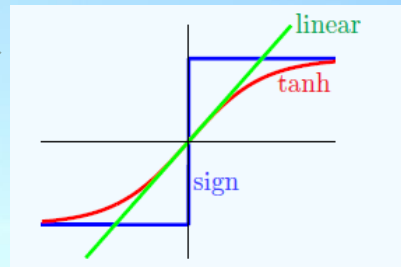
12



# Solution?



$f$ : activation function



$$g_k(\mathbf{x}) \equiv y_k = \text{Sign} \left( \sum_j w_{jk} \text{Sign} \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$

Homework:

Show that nonlinearity cannot be achieved if we use linear activation function for all nodes in a multiple layer perceptron model.



Sec 6.2 Problem 1, Duda et al, Pattern Classification, p. 335

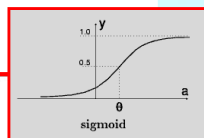
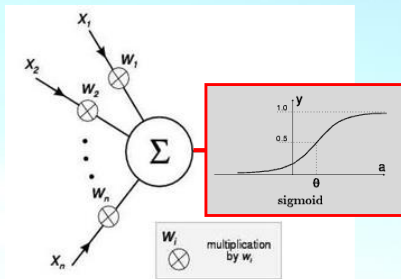
Xuegong Zhang

13

## MLP (Multi-layer Perceptron)

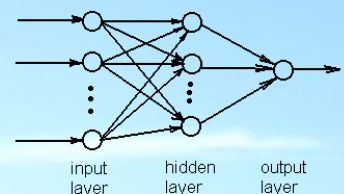


--- The most popular Artificial Neural Network (NN)



$$f(\alpha) = \frac{1}{1 + e^{-(\alpha - \theta)}}$$

Sigmoid function



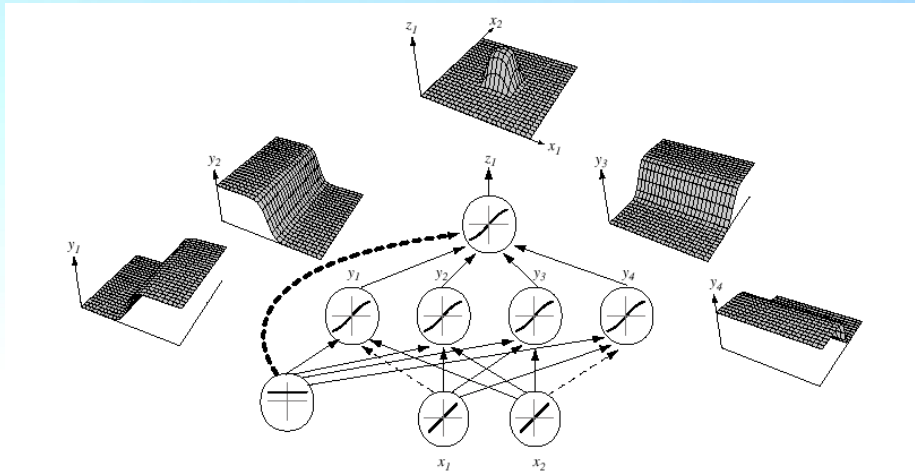
a 3-layer MLP

$$g_k(\mathbf{x}) \equiv y_k = f \left( \sum_j w_{jk} f \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$

Xuegong Zhang

14

## Nonlinearity with multiple layers



**FIGURE 6.2.** A 2-4-1 network (with bias) along with the response functions at different units; each hidden output unit has sigmoidal activation function  $f(\cdot)$ . In the case shown, the hidden unit outputs are paired in opposition thereby producing a “bump” at the output unit. Given a sufficiently large number of hidden units, any continuous function from input to output can be approximated arbitrarily well by such a network. From: Richard O. Duda, Peter E. Hart, and David G. Stork, *Pattern Classification*. Copyright © 2001 by John Wiley & Sons, Inc.

Xuegong Zhang

Duda, Hart &amp; Stork, Pattern Classification (2nd edition), John Wiley &amp; Sons, 2001

15

## The Kolmogorov theorem



Any continuous function  $g(\mathbf{x})$  defined on the unit hypercube  $I^n$  ( $I = [0,1]$  and  $n \geq 2$ ) can be represented in the form

$$g(\mathbf{x}) = \sum_{j=1}^{2n+1} \Xi_j \left( \sum_{i=1}^d \Psi_{ij}(x_i) \right)$$

for properly chosen functions  $\Xi_j$  and  $\Psi_{ij}$ .

Xuegong Zhang

Duda, Hart &amp; Stork, Pattern Classification, John Wiley &amp; Sons, 2001, p.287

16



## The representative power of MLP



$$g_k(\mathbf{x}) \equiv y_k = f \left( \sum_j w_{jk} f \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$



### Conclusion:

MLP is capable of “implementing” any continuous function from input to output, given sufficient number of hidden units, proper activation function and weights.

Xuegong Zhang

Duda, Hart & Stork, Pattern Classification, John Wiley & Sons, 2001, p.287

17

## 1-minute break



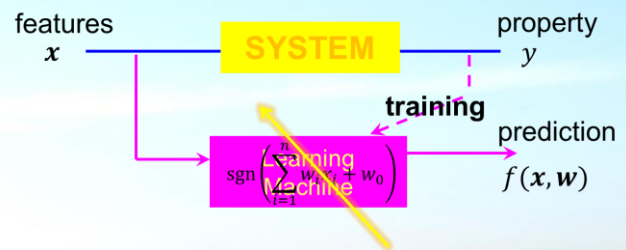
Xuegong Zhang

18

## The basic elements of supervised learning



- It needs a teacher.
  - We (people) **design** it (features, the model) and **train** it.
- We need materials to train it. / It needs materials to learn from.
  - Training data
- We need to tell what is the goal of the learning.
  - **Objective function**
- We need to tell it how to learn.
  - **Learning algorithm**



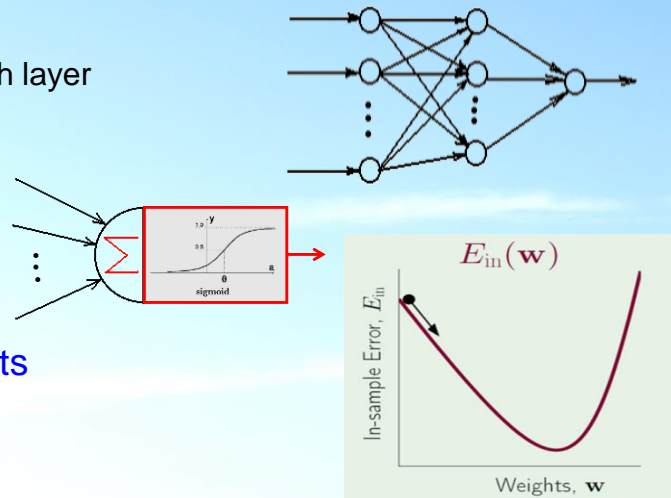
Xuegong Zhang

19

## How to design a MLP?



- The NN structure:
  - # of layers and # of nodes at each layer
  - Input: features to use
  - Output: real values
    - Coding according to the task
- Choose the activation function
  - Sigmoid
- Algorithm for training the weights
  - Gradient Descent



Xuegong Zhang

20

### Basic concepts of ML: Perceptron

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min J_P(\boldsymbol{\alpha}) = \sum_{y_j \in Y^k} (-\boldsymbol{\alpha}^T \mathbf{y}_j)$
  - We need to tell it how to learn.
    - Learning algorithm:  $\boldsymbol{\alpha}(k+1) = \boldsymbol{\alpha}(k) - \rho_k \nabla J = \boldsymbol{\alpha}(k) + \rho_k \sum_{y_j \in Y^k} \mathbf{y}_j$

### Basic concepts of ML: Logistic Regression

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j})$
  - We need to tell it how to learn.
    - Learning algorithm:  $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

### Basic concepts of ML: Linear Regression

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in R$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min E = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2$
  - We need to tell it how to learn.
    - Learning algorithm:  $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

### Multilayer Perceptron ?

Xuegong Zhang

21

### Basic concepts of ML: MLP

- How can we make a learning machine?

- It needs a teacher.

→ The model:  $g(\mathbf{x}) = f(\sum_j w_{jk} f(\sum_i w_{ij} x_i + w_{j0}) + w_{k0})$

- We need materials to train it. / It needs materials to learn from.

→ Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in R$

- We need to tell what is the goal of the learning.

→ Objective function:  $\min E = \frac{1}{2} \sum_{j=1}^N (g(\mathbf{x}_j) - y_j)^2$

- We need to tell it how to learn.

→ Learning algorithm:  $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

$$\nabla E = \partial E / \partial \mathbf{w} ?$$



Xuegong Zhang

22

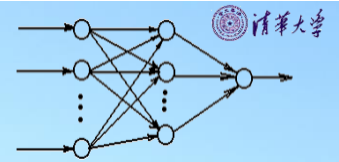
## The challenge

$$\min E = \frac{1}{2} \sum_{j=1}^N (g(\mathbf{x}_j) - y_j)^2$$

$$\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$$

$$g(\mathbf{x}) = f \left( \sum_j w_{jk} f \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$

$$\nabla E = \partial E / \partial \mathbf{w}$$

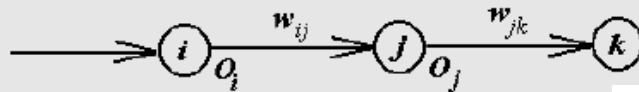


Xuegong Zhang

23

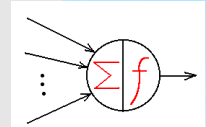
## Back-Propagation (BP) of the Error

LeCun, 1986; Rumelhart, Hinton & Williams, 1986; Parker, 1985



for node  $j$ ,  $net_j = \sum_i w_{ij} O_i$ ,  $O_j = f(net_j)$

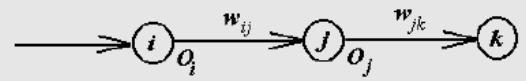
for output node:  $\hat{y}_j = O_j$  is the actual output,  $y_j$  is the desired output, error



$$g(\mathbf{x}) = f \left( \sum_j w_{jk} f \left( \sum_i w_{ij} x_i + w_{j0} \right) + w_{k0} \right)$$

Xuegong Zhang

24



error:

$$E = \frac{1}{2} \sum_j (y_j - \hat{y}_j)^2 \quad \text{Objective function}$$

gradients:

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ij}} = \delta_j O_i, \quad \delta_j \equiv \frac{\partial E}{\partial net_j}$$

Xuegong Zhang

25

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \frac{\partial net_j}{\partial w_{ij}} = \delta_j O_i, \quad \delta_j \equiv \frac{\partial E}{\partial net_j}$$

$$net_j = \sum_i w_{ij} O_i, \quad O_j = f(net_j)$$

for output nodes:  $O_j = \hat{y}_j, \quad \delta_j = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial net_j} = -(y_j - \hat{y}_j) f'(net_j)$

for hidden nodes (back-propagation):

$$\delta_j = \frac{\partial E}{\partial net_j} = \sum_k \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j} = \sum_k \delta_k w_{jk} f'(net_j)$$

The Chain Rule for Derivatives



Xuegong Zhang

26



for output nodes:  $O_j = \hat{y}_j$ ,  $\delta_j = \frac{\partial E}{\partial \hat{y}_j} \frac{\partial \hat{y}_j}{\partial net_j} = -(y_j - \hat{y}_j) f'(net_j)$

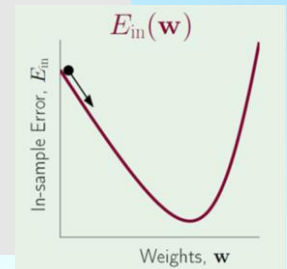
for hidden nodes (back-propagation):

$$\delta_j = \frac{\partial E}{\partial net_j} = \sum_k \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial O_j} \cdot \frac{\partial O_j}{\partial net_j} = \sum_k \delta_k w_{jk} f'(net_j)$$

adaption (at iteration  $t$ ):

$$w_{ij}(t+1) = w_{ij}(t) + \Delta w_{ij}(t), \quad \Delta w_{ij}(t) = -\eta \delta_j(t) O_i(t)$$

where:  $\eta$  - step length (learning rate)



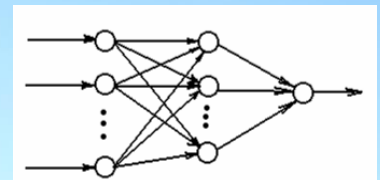
Xuegong Zhang

27

## The Back-Propagation (BP) Algorithm



- Purpose: to train MLP weights with data
- Goal: **to minimize the error** (MSE)
- Algorithm: **gradient-decreasing**
  - **forward**: compute the output with certain input  
input nodes >>> hidden nodes >>> output nodes
  - **compare**: discrepancy between output and desired
  - **backward**: propagate the error back through the network  
input nodes <<< hidden nodes <<< output nodes
  - **adaptation**: adapt the weights according to the errors



Xuegong Zhang

28



## Pseudo-codes for BP with sigmoid nodes



1°. Initialize weights (with small random values),  $t = 0$

2°. Apply a training sample  $x = [x_1, \dots, x_n]^T \in R^n$  with desired output  $D = [d_1, \dots, d_m]^T \in R^m$

3°. Forward calculation:  $Y = [y_1, \dots, y_m]^T \in R^m$ ,

$$y_l = f\left(\sum_{k=1}^{n_2} w_{kl} f\left(\sum_{j=1}^{n_1} w_{jk} f\left(\sum_{i=1}^n w_{ij} x_i\right)\right)\right), l = 1, \dots, m$$

4°. Adjust weights from the output layer. For layer  $l$ ,

$$w_{ij}^l(t+1) = w_{ij}^l(t) - \eta \delta_j^l x_i^{l-1}, j = 1, \dots, n_l, i = 1, \dots, n_{l-1}$$

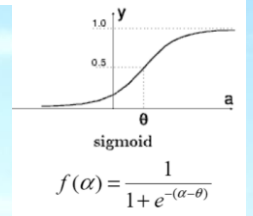
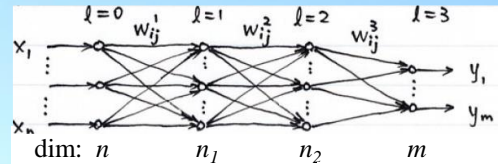
where for the output layer

$$\delta_j^l = -y_j(1 - y_j)(d_j - y_j), j = 1, \dots, m$$

and for hidden layers

$$\delta_j^l = x_j^l(1 - x_j^l) \sum_{k=1}^{n_{l+1}} \delta_k^{l+1} w_{jk}^{l+1}(t), j = 1, \dots, n_l$$

5°. Loop: if stop criterion not met, set  $t = t + 1$  and go to 2° with another sample.



$$f(\alpha) = 1/(1 + e^{-\alpha})$$

$$f'(\alpha) = f(\alpha)(1 - f(\alpha))$$

Xuegong Zhang

29

## Basic concepts of ML: MLP



• How can we make a learning machine?

– It needs a teacher.

→ The model:  $g(\mathbf{x}) = f\left(\sum_j w_{jk} f\left(\sum_i w_{ij} x_i + w_{j0}\right) + w_{k0}\right)$

– We need materials to train it. / It needs materials to learn from.

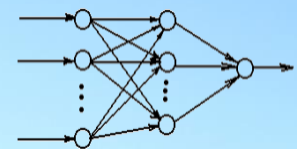
→ Training data:  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in R$

– We need to tell what is the goal of the learning.

→ Objective function:  $\min E = \frac{1}{2} \sum_{j=1}^N (g(\mathbf{x}_j) - y_j)^2$

– We need to tell it how to learn.

→ Learning algorithm:  $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$  via the BP algorithm



Xuegong Zhang

30

### Basic concepts of ML: Perceptron

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in \{-1, 1\}$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min J_P(\alpha) = \sum_{y_j \in Y^k} (-\alpha^T y_j)$
  - We need to tell it how to learn.
    - Learning algorithm:  $\alpha(k+1) = \alpha(k) - \rho_k \nabla J = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$

### Basic concepts of ML: Logistic Regression

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $h(x) = \theta(w^T x)$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in \{-1, 1\}$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min E(w) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j w^T x_j})$
  - We need to tell it how to learn.
    - Learning algorithm:  $w(k+1) = w(k) - \rho_k \nabla E$

### Basic concepts of ML: Linear Regression

- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $f(x) = \sum_{i=0}^d w_i x_i = w^T x$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in R$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min E = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2$
  - We need to tell it how to learn.
    - Learning algorithm:  $w(k+1) = w(k) - \rho_k \nabla E$

### Basic concepts of ML: MLP

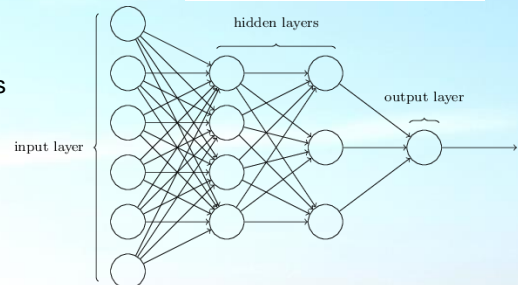
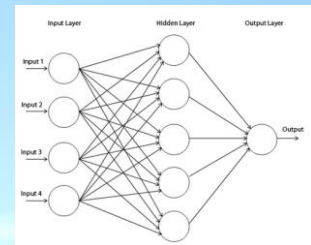
- How can we make a learning machine?
  - It needs a teacher.
    - The model:  $g(x) = f(\sum_j w_{jk} f(\sum_i w_{ij} x_i + w_{j0}) + w_{k0})$
  - We need materials to train it. / It needs materials to learn from.
    - Training data:  $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_j \in R^{d+1}, y_j \in R$
  - We need to tell what is the goal of the learning.
    - Objective function:  $\min E = \frac{1}{N} \sum_{j=1}^N (g(x_j) - y_j)^2$
  - We need to tell it how to learn.
    - Learning algorithm:  $w(k+1) = w(k) - \rho_k \nabla E$  via the BP algorithm

Xuegong Zhang

31

## What can MLP do?

- As a general nonlinear classifier
  - Coding class labels
    - Two-class: 0 vs. 1
  - Network structure
    - 3-layer MLP (with 1 hidden layer)
      - any convex decision regions
    - 4-layer MLP (with 2 hidden layers)
      - any connected/disconnected decision regions



Xuegong Zhang

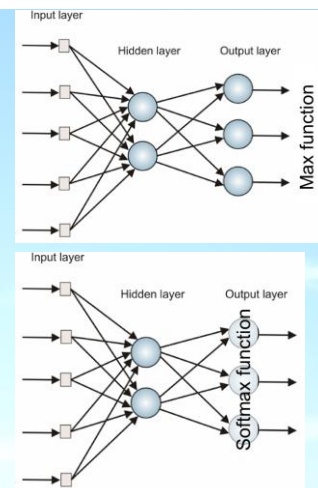
32

- Multi-class classification with MLP

- Multi-class:

- One-hot coding (1-of-C coding)  
e.g., 3 classes coded as: [1,0,0], [0,1,0], [0,0,1]
- Softmax output

$$y_j = \frac{e^{w_j \cdot o_h}}{\sum_{k=1}^K e^{w_k \cdot o_h}}, j = 1, \dots, K$$



Xuegong Zhang

33

## Homework



- Problems (Pr3)
  1. Prove that multilayer neural network with linear activation functions cannot achieve nonlinearity.
  2. Derive the BP algorithms for tanh() activation function, and for SoftMax output function.
- Computer exercises (Ex2)
  - Find a package of KNN
  - Describe its algorithm
  - Code your own MLP program with BP learning.
  - Experiment on the medical dataset
- Deadline:
  - Oct. 6 (Wednesday), 23:00
- Deadline:
  - Oct. 13 (Wednesday), 23:00



Xuegong Zhang

34

