Artificial Neural Networks

By: Bijan Moaveni

Email: *b_moaveni@iust.ac.ir* http://webpages.iust.ac.ir/b_moaveni/

Programs of the Course

- Aims of the Course
- Reference Books
- Preliminaries
- Evaluation

Aims of the Course

- 1. Discuss the fundamental techniques in Neural Networks.
- 2. Discuss the fundamental structures and its learning algorithms.
- 3. Introduce the new models of NNs and its applications.

Neural Network is an intelligent numerical computation method.

Learning Outcomes

- 1. Understand the relation between real brains and simple artificial neural network models.
- 2. Describe and explain the most common architectures and learning algorithms for Multi-Layer Perceptrons, Radial-Basis Function Networks and Kohonen Self-Organising Maps.
- 3. Explain the learning and generalization aspects of neural network systems.
- 4. Demonstrate an understanding of the implementation issues for common neural network systems.
- 5. Demonstrate an understanding of the practical considerations in applying neural networks to real classification, recognition, identification, approximation problems and control.

Course Evaluation

- 1. Course Projects 40%
- 2. Final Exam 50%
- 3. Conference Paper 10%

Reference Books

- Haykin S., Neural Networks: A Comprehensive Foundation., Prentice Hall, 1999.
- Hagan M.T., Dcmuth H.B. and Beale M., Neural Network Design, PWS Publishing Co., 1996.

Preliminaries

- 1. Matrices Algebra to Neural Network design and implementation.
- 2. MATLAB software for simulation. (NN toolbox is arbitrary).



Introduction

- 1. What are Neural Networks?
- 2. Why are Artificial Neural Networks Worth Noting and Studying?
- 3. What are Artificial Neural Networks used for?
- 4. Learning in Neural Networks
- 5. A Brief History of the Field
- 6. Artificial Neural Networks compared with Classical Symbolic A.I.
- 7. Some Current Artificial Neural Network
 ² Applications









Why are Artificial Neural Networks Worth Noting and Studying?

- 1. They are extremely powerful computational devices.
- 2. Parallel Processing makes them very efficient.
- 3. They can learn and generalize from training data so there is no need for enormous feats of programming.
- 4. They are particularly fault tolerant this is equivalent to the "graceful degradation" found in biological systems.
- 5. They are very noise tolerant so they can cope or deal with situations where normal symbolic (classic) systems would have difficulty.
- 6. In principle, they can do anything a symbolic or classic ⁷ system can do, and more.



Learning in Neural Networks

There are many forms of neural networks. Most operate by passing neural 'activations' through a network of connected neurons.

One of the most powerful features of neural networks is their ability to *learn* and *generalize* from a set of <u>training data</u>. They adapt the strengths/weights of the connections between neurons so that the final output activations are correct.

There are three broad types of learning:

- 1. Supervised Learning (i.e. learning with a teacher)
- 2. Reinforcement learning (i.e. learning with limited feedback)
- 3. Unsupervised learning (i.e. learning with no help)

There are most common learning algorithms for the most common types of neural networks.

9

A Brief History 1943 McCulloch and Pitts proposed the McCulloch-Pitts neuron • model **1949** Hebb published his book *The Organization of Behavior*, in which the Hebbian learning rule was proposed. **1958** Rosenblatt introduced the simple single layer networks now called Perceptrons. **1969** Minsky and Papert's book *Perceptrons* demonstrated the limitation of single layer perceptrons and almost the whole field went into hibernation. **1982** Hopfield published a series of papers on Hopfield networks. **1982** Kohonen developed the Self-Organising Maps that now bear his name. **1986** The Back-Propagation learning algorithm for Multi-Layer Perceptrons was rediscovered and the whole field took off again. **1990s** The sub-field of Radial Basis Function Networks is developed. **2000s** The power of Ensembles of Neural Networks and Support 10 Vector Machines becomes apparent.

A Brief History

1943 McCulloch and Pitts proposed the McCulloch-Pitts neuron model



Warren S. McCulloch (Nov., 16, 1898 – Sep., 24, 1969) American neurophysiologist and cybernetician

W. McCulloch and W. Pitts, 1943 "A Logical Calculus of the Ideas Immanent in Nervous Activity". In :Bulletin of Mathematical Biophysics Vol 5, pp 115-183.



























Human Nervous System

- The real structure of the human nervous corresponding to last block-diagram.
- It contains the neurons to transfer the signal form the receptors to brain and vice-versa to the effectors.









Components of Biological Neuron

4. *Dendrites* are fibres which come from the cell body and provide the receptive zones that receive activation from other neurons.

5. *Axons* are fibres acting as transmission lines that send activation to other neurons.

6. The junctions that allow signal transmission between the axons and dendrites are called *synapses*. The process of transmission is by diffusion of chemicals called *neurotransmitters* across the synaptic cleft.



Level of Brain Organization

There is a hierarchy of interwoven levels of organization:

- 1. Molecules and Ions
- 2. Synapses
- 3. Neuronal microcircuits
- 4. Dendrite trees
- 5. Neurons
- 6. Local circuits
- 7. Inter-regional circuits
- 8. Central nervous system

The ANNs we study in this module are crude approximations to levels 5 and 6.







The McCulloch and Pitts Neuron Analysis

• Note that the McCulloch-Pitts neuron is an extremely simplified model of real biological neurons. Some of its missing features include: non-binary inputs and outputs, non-linear summation, smooth thresholding, stochasticity, and temporal information processing.

• Nevertheless, McCulloch-Pitts neurons are computationally very powerful. One can show that assemblies of such neurons are capable of universal computation.









Implementing Logic Gates with M-P Neurons

According to the McCulloch-Pitts Neuron properties we can use it to implement the basic logic gates.

Not		
in	out	
1	0	
0	1	

And		
In_1	in ₂	out
1	1	1
1	0	0
0	1	0
0	0	0

OR			
in ₂	out		
1	1		
0	1		
1	1		
0	0		
	OR in ₂ 1 0 1 0		

5

What should we do to implement or realize a logic gate, Not/AND/OR, by N.N.?





























1st Mini Project

- 1. By using the perceptron learning rule generate a N.N. to represent a NOT gate.
- 2. By using the perceptron learning rule generate a N.N. to represent a AND gate.
- 3. By using the perceptron learning rule generate a N.N. to represent a OR gate.
- 4. Please show that the generalized error converge to constant value after a learning process.
- 5. Please test the above N.N.s by testing data?
- 6. Please check the above N. N.s with data which added to noise.
- 7. Repeat the learning process for above N.N.s in both with and without bias.
- 8. Please plot the updated weights.

20





1









3
















Also, the above goal can be obtained by minimizing the following cost function.

$$E = \sum_{k}^{p} \frac{1}{2} e(k)^2$$



Learning Rules

So, the <u>N.N. can be optimized</u> by <u>minimizing</u> the corresponding cost function with respect to the <u>synaptic weights</u> of network.

According to above explanation, <u>Widrow and Hoff in 1960</u> proposed a new method to update the weights based on *delta rule*.

$$\Delta w_i(k) = \eta e(k) x_i(k)$$

$$w_i(k+1) = w_i(k) + \eta e(k)x_i(k)$$

Learning Rules

Hebbian Learning rule:

Hebb's postulate of learning is the oldest and most famous of all learning rules.

His theory can be rephrased as a two-part as follows:

1. If two neurons on either side of a synapse (connection) are activated simultaneously (i.e. synchronously), then the strength of that synapse is selectively increased.

2. If two neurons on either side of a synapse are activated asynchronously, then that synapse is selectively weakened or eliminated.

8





Back-Propagation Algorithm

We look for a simple method of training in which the weights are updated on a *pattern-by-pattern* basis (online method). The adjustments to the weights are made in accordance with the respective <u>errors</u> computed for *each* pattern presented to the network.

The arithmetic average of these individual weight changes over the training set is therefore an *estimate* of the true change that would result from modifying the weights based on minimizing the cost function *E* over the entire training set.

$$E = \frac{1}{2} \sum_{j \in O.L.} e_j^2(n)$$



Back-Propagation Algorithm

$\Delta w_{ji}(n) = -\eta$	$n \frac{\partial E(n)}{\partial E(n)}$
	$\sqrt[n]{\partial w_{ji}(n)}$

•In words, gradient method could be thought of as a ball rolling down from a hill: the ball will roll down and finally stop at the valley.





Back-Propagation Algorithm

15

 $\Delta w_{ji}(n) = \eta \underbrace{e_j(n)\varphi'(v_j(n))}_{\delta_j(n)} y_i(n) = \eta \delta_j(n) y_i(n)$

 $\delta_j(n)$: Local Gradient











Artificial Neural Networks Lecture 7 Some Notes on Back-Propagation

Learning Rate

The smaller we make the learning-rate parameter η , the smaller will the changes to the synaptic weights in the network be from one iteration to the next and the smoother will be the trajectory in weight space.

If, on the other hand, we make the learning-rate parameter η too large so as to speed up the rate of learning, the resulting large changes in the synaptic weights assume such a form that the network may become unstable (i.e., oscillatory).

Solution: A simple method of increasing the rate of learning and yet avoiding the danger of instability is to <u>modify the delta rule</u> by including a *momentum* term, as shown by' (Rumelhart et al., 1986a)

 $\Delta w_{ii}(n) = \eta \delta_i(n) y_i(n) + \alpha \Delta w_{ii}(n-1)$











Sequential Mode and Batch Mode

From an "on-line" operational point of view, the pattern mode of training is preferred over the batch mode, because it requires *less local storage* for each synaptic connection.

Moreover, given that the patterns are presented to the network in a random manner, the use of pattern-by-pattern updating of weights makes the search in weight space *stochastic* in nature, which, in turn, makes it less likely for the back-propagation algorithm to be trapped in a local minimum.

On the other hand, the use of batch mode of training provides a more accurate estimate of the gradient vector.

* So, the training process can be started with batch mode and then it can be changed to sequential mode.

Stopping Criteria The back-propagation algorithm is considered to have converged when the <u>Euclidean norm of the gradient vector</u> reaches a sufficiently small gradient threshold. The drawback of this convergence criterion is that, for successful trials, learning times may be long. The back-propagation algorithm is considered to have converged when the absolute rate of change in the average squared error (ΔE_{av}) per epoch is sufficiently small. Typically, the rate of change in the average squared error is considered to be small enough if it lies in the range of 0.1 to 1 percent per epoch; sometimes, a value as small as 0.01 percent per epoch is used.

Stopping Criteria

• Another useful criterion for convergence is as follows. After each learning iteration, the network is tested for its generalization performance. The learning process is stopped when the generalization performance is adequate.

Initializing in Back-Propagation

In Lee et al. (1991), a formula for the *probability of premature saturation* in back-propagation learning has been derived for the <u>batch mode</u> of updating, and it has been verified using computer simulation. The essence (core) of this formula may be summarized as follows:

1. Incorrect saturation is avoided by choosing the initial values of the synaptic weights and threshold levels **of** the network to be uniformly distributed inside a *small* range of values.

2. Incorrect saturation is less likely to occur *when the number of hidden neurons is maintained low*, consistent with a satisfactory operation of the network.

3. Incorrect saturation rarely occurs when the neurons of the network operate in their *linear regions*.

Note: For pattern-by-pattern updating, computer simulation results show similar trends to the batch mode of operation referred to herein

Heuristics for making the Back-Propagation Algorithm Perform Better

1. A M.L.P trained with the back-propagation algorithm may, in general, learn faster (in terms of the number of training iterations required) when the <u>asymmetric</u> sigmoidal activation function are used in neuron model. than when it is non-symmetric.

Asymmetric function:

$$\varphi(-v) = -\varphi(v)$$

Heuristics for making the Back-Propagation Algorithm Perform Better

2. It is important that the **desired values** are chosen within the range of the sigmoid activation functions.

Otherwise, the back-propagation algorithm tends to drive the free parameters of the network to infinity, and thereby slow down the learning process by orders of magnitude.



Heuristics for making the Back-Propagation Algorithm Perform Better

3. The initialization of the synaptic weights and threshold levels of the network should be *uniformly distributed* inside a *small range*. The reason for making the range small is to reduce the likelihood of the neurons in the network saturating and producing small error gradients.

However, the range should not be made too small, as it can cause the error gradients to be very small and the learning therefore to be initially very slow.

Heuristics for making the Back-Propagation Algorithm Perform Better

4. All neurons in the multilayer Perceptron should desirably learn at the same rate.

Typically, the <u>last layers</u> tend to <u>have larger local gradients</u> than the layers at the front end of the network. Hence, the learning-rate parameter η should be assigned a *smaller value* in the *last layers* than the front layers.

* Neurons with **many inputs** should have a **smaller learning-rate** parameter than neurons with few inputs.

Heuristics for making the Back-Propagation Algorithm Perform Better

5. For <u>on-line operation</u>, <u>pattern-by-pattern updating</u> rather than batch updating should be used for weight adjustments.

For pattern-classification problems involving a large and redundant database, pattern-by-pattern updating tends to be orders of magnitude faster than batch updating.

6. The *order* in which the training examples are *presented to the network* should be *randomized* (shuffled) from one epoch to the next. This form of randomization is critical for <u>improving the speed of convergence</u>.

16

Heuristics for making the Back-**Propagation Algorithm Perform Better**

7. Learning-rate:

In previous lectures and projects we studied the important effect of learningrate in back-propagation learning algorithm. Here, some new methods to improve the learning-rate value is introduced.

 $\eta(n)$

 $0.1\eta_{0}$

0.01₁₀

Standard LMS algo

Conventional learning rate: $\eta = \eta_0$

In each iteration the learning rate value decreases (stochastic approximation):

 $\eta(n) = \frac{\eta_0}{n}$

 $1 + \frac{n}{2}$ τ



(Where, τ is search time constant)



















A new method to tune the learning-rate

*Delta-bar-Delta*This method is applicable to learning rates in MLP and F.MLP.

$$\eta(k) = \begin{cases} \eta(k-1) + \alpha & \delta(k-1)\delta(k) > 0\\ b\eta(k-1) & \delta(k-1)\delta(k) < 0\\ 0 & Otherwise \end{cases}$$
$$10^{-4} \le \alpha \le 10^{-1}$$
$$0.5 \le b \le 0.9$$

Artificial Neural Networks

Lecture 9

Some Applications of Neural Networks (1) (Function Approximation)

1

2



Many <u>computational models</u> can be described as <u>functions mapping</u> some <u>numerical input vectors</u> to <u>numerical outputs</u>. The outputs corresponding to some input vectors may be known from training data, but we may not know the mathematical function describing the actual process that generates the outputs from the input vectors.

Function approximation is the task of learning or constructing a function that generates approximately the same outputs from input vectors <u>as the process being modeled</u>, based on available training data.

$$x \longrightarrow y=f(x) \longrightarrow y$$



























3rd Mini Project

By using of an arbitrary neural network (MLP) approximate the function which is presented in example 1.

1st Part of Final Project

By using of an arbitrary neural network (MLP) approximate the function which is presented in example 2.

In this project You can use all hints which are introduced in previous lectures but, you should explain their effects (score: 2 points)

Artificial Neural Networks

Lecture 10

Some Applications of Neural Networks (2) (System Identification)

1





System Identification

Two facts make the MLP a powerful tool for approximating the functions or identifying the systems:

Multilayer feedforward neural networks are universal approximators:

It was proved by Cybenko (1989) and Hornik et al. (1989) that any continuous mapping over a compact domain can be approximated as accurately as necessary by a feedforward neural network with one hidden layer.

The back propagation algorithm:

This algorithm which performs stochastic gradient descent, provides an effective method to train a feedforward neural network to approximate a given continuous function over a compact domain.







State Space model for Identification

In this structure, the states of the N.N. model provide an approximation or estimation to the states of the system.

A natural *performance criterion* for the model would be the sum of the squares of the errors between the system and the model outputs:

$$E(k) = \frac{1}{2} \sum_{k} \|y(k) - \hat{y}(k)\|^{2} = \sum_{k} \|e(k)\|^{2}$$

Dynamic Back Propagation:

$$w_{h} \in NN_{h} \longrightarrow \Delta w_{h} = -\eta_{h} \frac{dE(k)}{dw_{h}(k)}$$

$$w_{f} \in NN_{f} \longrightarrow \Delta w_{f} = -\eta_{f} \frac{dE(k)}{dw_{f}(k)} = -\eta_{f} \sum_{j=1}^{n} \frac{\partial E(k)}{\partial z_{j}(k)} \cdot \frac{dz_{j}(k)}{dw_{f}(k)}$$

$$\frac{dz_{j}(k)}{dw_{f}} = \sum_{l=1}^{n} \frac{\partial z_{j}(k)}{\partial z_{l}(k-1)} \cdot \frac{\partial z_{l}(k-1)}{\partial w_{f}} + \frac{\partial z_{j}(k)}{\partial w_{f}} = 8$$


Input-Output model for Identification

Clearly, choosing the <u>state space models</u> for identification requires the use of <u>dynamic back propagation</u>, which is computationally a very **intensive procedure**. At the same time, to avoid instabilities while training, one needs to use small learning rate to adjust the parameters, and this in turn results in **long convergence times**.

Input-Output Model of plant:

Consider the difference Equation corresponding to a typical linear plant:

$$y(k) = \sum_{i=1}^{n} a_{i} y(k-i) + \sum_{j=1}^{n-1} b_{j} u(k-j)$$





Theorem Let S_{nl} be the nonlinear system, and $S_{linearized}$ its linearization around the equilibrium point. If $S_{linearized}$ is observable, then S_{nl} is locally strongly observable. *Furthermore, locally,* S_{nl} *can be realized by an inputoutput model.*

Observability Matrix $\varphi_o = \begin{bmatrix} c \\ cA \\ cA \end{bmatrix}$

Input-Output model for Identification

Neural Network Implementation:

If strong observability conditions are known (or assumed) to be satisfied in the system's region of operation with n state variables, then the identification procedure using a feedforward neural network is quite straightforward.

At each instant of time, the **inputs** to the network consisting of the **system's past** *n* **input** values and **past** *n* **output values** (all together 2n), are fed into the neural network.

The network's output is compared with the next observation of the system's output to yield the error

 $e(k+1) = y(k+1) - \tilde{y}(k+1) = y(k+1) - \tilde{h}[Y_{l}(k-n+1), U_{l}(k-n+1)]$

The weights of the network are then adjusted using static back propagation to minimize the sum of the squared error.







2nd Part of Final Project

By using of an arbitrary neural network (MLP) identify the discrete nonlinear plant which is presented in example 2 (Score: 1 points).

- By using a test signal, show that the N.N. identifier perform a appropriate input-output model of plant.
- By using of the PRBS signal, repeat the identifying procedure and compare the results.

The material of this lecture is based on:

Omid Omidvar and David L. Elliott, **Neural Systems for Control**, Academic Press; 1st edition (1997).

Artificial Neural Networks

Lecture 11

Some Applications of Neural Networks (3) (Control)































3rd Part of Final Project

In this project, you should find a practical plant in papers and by using of NN controllers provide a suitable closed loop control performance. (Score: 2 points)

- In this project you can use of any NN controllers structure which are presented in this lecture.
- In this project you can use of any NN controllers which are introduced in papers and text books (score: +1 point).



Artificial Neural Networks

Lecture 12

Recurrent Neural Networks

Recurrent Neural Networks

The **conventional feedforward neural networks** can be used to approximate *any* **spatiality finite function**. That is, for functions which have a *fixed* input space there is always a way of encoding these functions as neural networks.

For example in function approximation, we can use the automatic learning techniques such as backpropagation to find the weights of the network if sufficient samples from the function is available.

Recurrent neural networks are fundamentally different from feedforward architectures in the sense that they not only operate on an <u>input space</u> but also on an <u>internal *state* space</u>.

These are proposed to learn sequential or time varying patterns.



























Radial Basis Function

Radial functions are a special class of function. Their characteristic feature is that their **response** decreases (or increases) monotonically with distance from a central point.

$$y_i = \varphi(\|x - c_i\|)$$

A typical radial function is the Gaussian which in the case of a scalar input is























Comparison of RBF Networks and MLP [1]

Radial-basis function (RBF) networks and multilayer perceptrons are examples of nonlinear layered feedforward networks. They are both universal approximators.

However, these two networks differ from each other in several important respects, as:

- 1. An RBF network (in its most basic form) has a single hidden layer, whereas an MLP may have one or more hidden layers.
- Typically, the computation nodes of an MLP, be they located in a <u>hidden or output</u> layer, share a common neuron model. On the other hand, the computation nodes in the hidden layer of an RBF network are quite different and serve a different purpose from those in the output layer of the network.

Comparison of RBF Networks and MLP [1]

- 3. The hidden layer of an RBF network is nonlinear, whereas the output layer is linear. On the other hand, the hidden and output layers of an MLP used as a classifier are usually all nonlinear; however, when the MLP is used to solve nonlinear regression problems, a linear layer for the output is usually the preferred choice.
- 4. The argument of the activation function of each hidden unit in an RBF network computes the *Euclidean norm (distance)* between the input vector and the center of that unit. On the other hand, the activation function of each hidden unit in an MLP computes the *inner product* of the input vector and the synaptic weight vector of that unit.

Comparison of RBF Networks and MLP [1]

5. MLPs construct *global* approximations to nonlinear input-output mapping. Consequently, they are capable of generalization in regions of the input space where little or no training data are available.

On the other hand, RBF networks using exponentially decaying localized nonlinearities (e.g., Gaussian functions) construct *local* approximations to nonlinear input-output mapping, with the result that these networks are capable of fast learning and reduced sensitivity to the order of presentation of training data.

15

<text><text><equation-block><text><text><text>



Artificial Neural Networks

Lecture 14

Hopfield Neural Network







Hopfield Neural Network

The energy *E* for the whole network can be determined from energy function as the following equation:

$$E = -\frac{1}{2}\sum_{i}\sum_{j}w_{ij}y_{i}y_{j} - \sum_{i}x_{i}y_{i} + \sum_{i}y_{i}\theta_{i}$$

So:
$$\Delta E_i = -\left(\sum_j w_{ij} y_j + x_i - \theta_i\right) \Delta y_i$$

 Δy_i is *positive* when the terms in brackets is *positive*; and Δy_i becomes *negative* in the other case. Therefore the energy increment for the whole network ΔE will always decrease however the input changes.

5

Hopfield Neural Network

So, the following two statements can be introduced:

- 1. The energy function *E* is a Lyapunov function.
- 2. The HNN is a stable in accordance with Lyapunov's Theorem.

The ability to minimize the energy function in a very short convergence time makes the HN described above be very useful in solving the problems with solutions obtained through minimizing a cost function.

Therefore, this cost function can be rewritten into the form of the energy function as *E* if the synaptic weights w_{ij} and the external input x_i can be determined in advance.

Hopfield Neural Network

Hopfield networks can be implemented to operate in two modes:

- *Synchronous mode* of training Hopfield networks means that all neurons fire at the same time.

- *Asynchronous mode* of training Hopfield networks means that the neurons fire at random.

7

Example: Consider a Hopfield network with three neurons

 $W = \begin{bmatrix} 0 & -0.4 & 0.2 \\ -0.4 & 0 & 0.5 \\ 0.2 & 0.5 & 0 \end{bmatrix}$

Let the state of the network be: $y(0) = [1, 1, 0]^T$.

Hopfield Neural Network Example: 0 -0.4 0.2 0 0.5 -0.4W =y(0) = 10 0.2 0.5 0 **Synchronous** Asynchronous mode mode $y(0) = [1,1,0]^T$ $y_{I}(1) = [0,1,0]^{T}$ $y(1) = [0,1,0]^T$ $y(1) = [1,0,0]^T$ *y*(1)=[1,1,1] $y_2(1) = [0,0,0]^T$ $y=y_3(1)=[0,0,0]^T$

State table of Hopfield N.N.					
A Hopfield net with n neurons has 2^n possible states, assuming that each neuron output produces two values 0 and 1.					
The state table for the above example Hopfield network with 3 neurons is given below.					
Init.	state if	state if	state if		
state	N1 fires	N2 fires	N3 fires		•
000	100	000	<u>000</u>		
001	101	011	000		
010	010	<u>000</u>	011		
011	011	011	011		
100	100	100	101		
101	101	111	101		
110	<i>010</i>	100	111		
111	011	111	111		
	-			9]

Hopfield N.N. as BAM

Hopfield networks are used as **content-addressable memory** or **Bidirectional Associative Memory (BAM).** The content-addressable memory is such a device that returns a pattern when given a noisy or incomplete version of it.

In this sense a content-addressable memory is *error-correcting* as it can override provided inconsistent information.

The discrete Hopfield network as a memory device operates in two phases: *storage phase* and *retrieval phase*.

During the *storage phase* the network learns the weights after presenting the training examples. The training examples for this case of automated learning are binary vectors, called also fundamental memories. The weights matrix is learned using the <u>Widrow-Hoff rule</u>. According to this rule when an input pattern is passed to the network and the estimated network output does not match the given target, the corresponding weights are modified by a small amount.

The difference from the single-layer perceptron is that no error is computed, rather the target is taken directly for weight updating.



Outer product Learning

Learning: Suppose that we wish to *store* a set of <u>*N*-dimensional</u> vectors (binary words), denoted by $\{\xi_{\mu}, \mu = 1, 2, ..., M\}$. We call these *M* vectors fundamental memories, representing the patterns to be memorized by the network.

The **<u>outer product</u>** learning rule, that is, the generalization of *Hebb*'s learning rule:

$$\mathbf{W} = \frac{1}{N} \left(\sum_{\mu=1}^{M} \boldsymbol{\xi}_{\mu} \boldsymbol{\xi}_{\mu}^{T} - M \mathbf{I} \right)$$

From these defining equations of the synaptic weights matrix, we note the following:

- The output of each neuron in the network is fed back to all other neurons.
- There is no self-feedback in the network (i.e., $w_{ii} = 0$).
- The weight matrix of the network is symmetric. (i.e., $W^T = W$)
Learning Algorithm

Initialization: Let the testing vector become initial state $\mathbf{x}(0)$

Repeat

-update asynchronously the components of the state $\mathbf{x}(t)$

$$u_{j}(n) = \sum_{i \neq j} w_{ij} y_{i}(n-1) + x_{j}(n) > 0 \implies y_{i}(n) = 1$$
$$u_{j}(n) = \sum_{i \neq j} w_{ij} y_{i}(n-1) + x_{j}(n) < 0 \implies y_{i}(n) = 0$$

-continue this updating until the state remains unchanged

until convergence

Generate output: return the stable state (fixed point) as a result. The network finally produces a time invariant state vector \mathbf{y} which satisfies the *stability condition*: $\mathbf{y} = \operatorname{sgn}(\mathbf{W}\mathbf{y} + \mathbf{b})$

13

Learning Algorithm

During the *retrieval phase* a testing vector called probe is presented to the network, which initiates computing the neuron outputs and developing the state.

After sending the training input to the recurrent network its output changes for a number of steps until reaching a stable state.

The selection of the <u>next neuron</u> to fire is <u>asynchronous</u>, while the modifications of the state are deterministic.

After the state evolves to a stable configuration, that is the state is not more updated, the network produces a solution.

This state solution can be envision as a fixed point of the dynamical network system. The solution is obtained after adaptive training.



















History of Genetic Algorithm

Genetic Algorithms (GAs) are adaptive random search algorithm premised on the <u>evolutionary ideas</u> of natural selection and genetic. The basic concept of GAs is designed to *simulate processes in natural system* necessary for evolution, specifically those that follow the principles first laid down by Charles Darwin of survival of the fittest.

As such they represent an intelligent exploitation of a random search within a defined search space to solve a problem.

Genetic algorithms originated from the studies of *cellular automata*, conducted by *John Holland* and his colleagues in 60s at the University of Michigan. Research in GAs remained largely theoretical until the mid-1980s, when The *First International Conference* on Genetic Algorithms was held at <u>The University of Illinois</u>.



Genetic Algorithm

GAs were introduced as a computational analogy of adaptive systems. They are *modeled loosely* on the principles of the evolution via natural selection, employing a population of individuals that undergo selection in the presence of variation-inducing operators such as **mutation** and **recombination** (**crossover**). A fitness function is used to evaluate individuals, and reproductive success varies with **fitness**.

The Algorithms can be summarized as:

- 1. Randomly generate an initial population M(0)
- 2. Compute and save the fitness u(*m*) for each individual *m* in the current population M(t).
- 3. Define selection probabilities p(*m*) for each individual *m* in M(t) so that p(*m*) is proportional to u(*m*)
- 4. Generate M(t+1) by probabilistically selecting individuals from M(t) to produce offspring via genetic operators (Crossover and Mutation)
- 5. Repeat step 2 until satisfying solution is obtained.





Gas are useful and efficient when • The search space is large, complex or poorly understood. • Domain knowledge is scarce or expert knowledge is difficult to encode to narrow the search space. • No mathematical analysis is available. • Traditional search methods fail.



























Selection

In the last example:

The cumulative probability q_k for each chromosome is:

$q_1 = 0.111180$,	$q_2 = 0.208695$,	$q_3 = 0.262534$
$q_4 = 0.427611$,	$q_5 = 0.515668,$	$q_6 = 0.582475$
$q_7 = 0.683290$,	$q_8 = 0.794234$,	$q_9 = 0.942446$
$q_{10} = 1.000000$		

Now we are ready to spin the roulette wheel 10 (population size) times, and each time we select a chromosome. So, *r* sequence can be generated randomly:

0.301431	0.322062	0.766503	0.881893
0.350871	0.583392	0.177618	0.343242
0.032685	0.197577		



Crossover

One of the important GA operators which can help us to search the corresponding space is **Crossover**:

In crossover procedure there are two steps:

- 1. Define the crossover rate (p_c) to select the chromosomes for crossover.
- 2. Choose the crossover method (e.g. one-cut-point) and generate the new chromosomes.

In the last example: $p_c = 0.25$

0.6257	0.2668	0.2886	0.2951
$0.1632 < p_c$	0.5674	$0.0859 < p_c$	0.3928
0.7707	0.5486		

25

$\boldsymbol{v}_5' = [100110110100101101000000010111001]$

 $v_7' = [001110101110011000000010101001000]$

<section-header><section-header><section-header><section-header><text><text><equation-block><equation-block><equation-block><equation-block><equation-block>

Mutation

To prevent the GA of trapped in local minimum, Mutation operator is employed. In mutation procedure the following 2 steps are important.

- 1. Define the mutation rate (p_m) to select genes.
- 2. Generate "number of genes*population size" random numbers (r_m) and by comparing those with mutation rate choose the corresponding genes which satisfy the following equation to mutate $(0 \rightarrow 1 \text{ and } 1 \rightarrow 0)$.

 $r_m < p_m$

27

he last example:	p = 0.01		
	I m		
Random_num.	Bit_position	ChromNo.	Bit_No.
0.009857	105	4	6
0.003113	164	5	32
0.000946	199	7	1
0.001282	329	10	32















	G.A. solution	
Evaluation Here you ca chromosome	n see the corresponding fitness function for parent	
	$eval(v_1) = f(4.954222, 0.169225) = 10.731945$ $eval(v_2) = f(-4.806207, -1.630757) = 12.110259$ $eval(v_3) = f(4.672536, -1.867275) = 11.788221$ $eval(v_4) = f(1.897794, -0.196387) = 5.681900$ $eval(v_5) = f(-2.127598, 0.750603) = 6.757691$ $eval(v_6) = f(-3.832667, -0.959655) = 9.194728$ $eval(v_7) = f(-3.792383, 4.064608) = 11.795402$	
Now, we ge	$eval(v_8) = f(1.182745, -4.712821) = 11.559363$ $eval(v_9) = f(3.812220, -3.441115) = 12.279653$ $eval(v_{10}) = f(-4.515976, 4.539171) = 14.251764$ enerate a sequence of random numbers: 0.828211 0.199683 0.639149 0.629170 0.957427 0.149358 0.304788 0.058504 0.149693 0.326670	36



Matation	G.A. solution
Mutation:	
X	bit_pos chrom_num variable random_num 11 6 x ₁ 0.081393
offspring	$\longrightarrow v'_{5} = \{-4.068506, -0.959655\}$
The fitness v	alue for each offspring:
	• •
	$eval(v'_1) = f(-4.444387, -1.383817) = 11.927451$
	$eval(v'_1) = f(-4.444387, -1.383817) = 11.927451$ $eval(v'_2) = f(-4.194488, -1.206594) = 10.566867$
	$eval(v'_1) = f(-4.444387, -1.383817) = 11.927451$ $eval(v'_2) = f(-4.194488, -1.206594) = 10.566867$ $eval(v'_3) = f(-3.683262, -4.521950) = 13.449167$
	$eval(v'_1) = f(-4.444387, -1.383817) = 11.927451$ $eval(v'_2) = f(-4.194488, -1.206594) = 10.566867$ $eval(v'_3) = f(-3.683262, -4.521950) = 13.449167$ $eval(v'_4) = f(-1.311703, -3.631985) = 10.538330$



