

Artificial Neural Networks and Information Theory

Colin Fyfe,
Department of Computing and Information Systems,
The University of Paisley.
Room H242
Phone 848 3305.

Edition 1.2,
2000

Contents

1	Introduction	7
1.1	Objectives	7
1.2	Intelligence	7
1.3	Artificial Neural Networks	8
1.4	Biological and Silicon Neurons	8
1.5	Learning in Artificial Neural Networks	9
1.6	Typical Problem Areas	11
1.7	A short history of ANNs	12
1.8	A First Tutorial	12
1.8.1	Worked Example	12
1.8.2	Exercises	14
2	Information Theory and Statistics	17
2.1	Probability	17
2.2	Statistics	18
2.2.1	The Gaussian or Normal Distribution	20
2.3	Quantification of Information	20
2.3.1	Logarithms	20
2.3.2	Information	22
2.3.3	Entropy	23
2.3.4	Joint and Conditional Entropy	25
2.3.5	Mutual Information	26
2.3.6	The Kullback Leibler Distance	27
2.3.7	Entropy from a Continuous Distribution	27
2.3.8	Entropy and the Gaussian Distribution	28
2.4	Information Theory and the Neuron	29
2.4.1	One Neuron with Noise on the Output	29
2.4.2	Noise on the Inputs	31
2.4.3	More than one output neuron	32
2.5	Principal Component Analysis	33
2.6	A Silly Example	34
2.7	Exercise	36
3	Hebbian Learning	39
3.1	Simple Hebbian Learning	39
3.1.1	Stability of the Simple Hebbian Rule	40
3.2	Weight Decay in Hebbian Learning	41
3.3	Principal Components and Weight Decay	42
3.4	Oja's One Neuron Model	43
3.4.1	Derivation of Oja's One Neuron Model	43
3.5	Recent PCA Models	44
3.5.1	Oja's Subspace Algorithm	44

3.5.2	Oja's Weighted Subspace Algorithm	46
3.5.3	Sanger's Generalized Hebbian Algorithm	47
3.6	The InfoMax Principle in Linsker's Model	47
3.7	Regression	48
3.7.1	Minor Components Analysis	50
3.8	Your Practical Work	51
3.8.1	Annealing of Learning Rate	51
3.8.2	The Data	51
4	Anti-Hebbian Learning	53
4.1	The Novelty Filter	54
4.2	Földiák's First Model	55
4.2.1	An Example	56
4.2.2	Földiák's Second Model	56
4.3	Rubner and Schulten's Model	58
4.4	The Negative Feedback Model	58
4.4.1	Biological Interneurons	59
4.4.2	Extensions to the Interneuron Network	59
4.5	The Anti-Hebbian Synapse as Energy Minimiser	64
5	Objective Function Methods	67
5.1	The Adaline	67
5.2	The Backpropagation Network	68
5.2.1	The Backpropagation Algorithm	69
5.2.2	The XOR problem	69
5.2.3	Backpropagation and PCA	70
5.3	Using Principal Components	71
5.3.1	Preprocessing	71
5.3.2	Principal Components Pruning	72
5.4	Cross-Entropy as the Objective Function	74
5.5	The I-Max Model	75
5.5.1	An Example	77
5.6	Contextual Information	78
5.6.1	Information Transfer and Learning Rules	80
5.6.2	Results	82
5.7	The Paisley Connection	82
5.7.1	Introduction	82
5.7.2	The Canonical Correlation Network	83
5.7.3	Experimental Results	85
5.7.4	Artificial Data	85
5.7.5	Real data	86
5.7.6	Random Dot Stereograms	86
5.7.7	More than two data sets	88
5.7.8	Non-linear Correlations	88
5.8	Speeding up Error Descent	90
5.8.1	Mathematical Background	92
5.8.2	QuickProp	93
5.8.3	Line Search	93
5.8.4	Conjugate Gradients	94
5.9	Least Mean Square Error Reconstruction	95
5.10	Conclusion	96

6	Identifying Independent Sources	97
6.1	The Trouble with Principal Components	97
6.1.1	Independent Codes	97
6.1.2	A Typical Experiment	99
6.2	Competitive Learning	99
6.2.1	Simple Competitive Learning	100
6.3	Anti-Hebbian and Competitive Learning	100
6.3.1	Sparse Coding	100
6.3.2	Földiák's Sixth Model	102
6.3.3	Implementation Details	102
6.3.4	Results	103
6.4	Competitive Hebbian Learning	103
6.5	Multiple Cause Models	105
6.5.1	Saund's Model	105
6.5.2	Dayan and Zemel	107
6.6	Predictability Minimisation	107
6.7	Mixtures of Experts	108
6.7.1	An Example	109
6.8	The Paisley Connection	110
6.8.1	Non negative Weights and Factor Analysis	110
6.8.2	Non negative Outputs	111
6.8.3	Additive Noise	112
6.8.4	Topographical Ordering and Modular Noise.	113
6.8.5	Results	113
6.8.6	Additive Noise	115
6.9	Probabilistic Models	118
6.9.1	Mixtures of Gaussians	118
6.9.2	A Logistic Belief Network	121
6.9.3	The Helmholtz Machine and the EM Algorithm	121
6.9.4	The Wake-Sleep algorithm	122
6.10	Conclusion	123
7	Independent Component Analysis	129
7.1	A Restatement of the Problem	129
7.2	Jutten and Herault	131
7.2.1	An Example Separation	132
7.2.2	Learning the weights	133
7.3	Non-linear PCA	133
7.3.1	Simulations and Discussion	134
7.4	Information Maximisation	134
7.4.1	The Learning Algorithm	136
7.5	The Paisley Dimension	138
7.5.1	Example	140
7.6	Penalised Minimum Reconstruction Error	140
7.6.1	The Least Mean Square Error Network	140
7.6.2	Adding Competition	142
7.7	Conclusion	143
8	Learning	145
8.1	The Bias-Variance Dilemma	145
8.1.1	Decomposition of the Error	146
8.1.2	General Statistical Derivation	147
8.1.3	An Example	147
8.2	The VC Dimension	149

8.3	PAC Learning	151
8.3.1	Examples	151
8.4	Regularisation	152
8.4.1	Regression	152
8.4.2	Weight Decay	153
8.4.3	Eigenvector Analysis of Weight Decay	154
8.5	Radial Basis Functions	155
8.5.1	RBF and MLP as Approximators	156
8.5.2	Comparison with MLPs	159
8.6	Learning in Radial Basis Functions	160
8.6.1	Fixed Centers selected at Random	160
8.6.2	Self-organised Selection of Centres	162
8.6.3	Supervised Selection of Centres	162
8.7	Cross-Validation	162
8.8	Support Vector Machines	163
8.8.1	Classification	163
8.8.2	A Case Study	164
8.8.3	Simulations	165
8.9	Conclusion	166
9	Unsupervised Learning using Kernel Methods	171
9.1	Introduction	171
9.2	Kernel PCA	171
9.2.1	The Linear Kernel	171
9.2.2	Non linear Kernels	172
9.2.3	The Curse of Kernel Dimensionality	173
9.3	Kernel Principal Factor Analysis	174
9.3.1	Principal Factor Analysis	175
9.3.2	The Varimax Rotation	175
9.3.3	Kernel Principal Factor Analysis	175
9.3.4	Simulations	180
9.4	Kernel Exploratory Projection Pursuit	180
9.4.1	Exploratory Projection Pursuit	180
9.4.2	Kernel Exploratory Projection Pursuit	180
9.4.3	Simulations	181
9.5	Canonical Correlation Analysis	181
9.5.1	Kernel Canonical Correlation Analysis	184
9.5.2	Simulations	186
9.6	Conclusion	189
A	Linear Algebra	191
A.1	Vectors	191
A.1.1	Same direction vectors	191
A.1.2	Addition of vectors	191
A.1.3	Length of a vector	192
A.1.4	The Scalar Product of 2 Vectors	192
A.1.5	The direction between 2 vectors	192
A.1.6	Linear Dependence	192
A.1.7	Neural Networks	192
A.2	Matrices	193
A.2.1	Transpose	193
A.2.2	Addition	193
A.2.3	Multiplication by a scalar	193
A.2.4	Multiplication of Matrices	193

A.2.5 Identity	194
A.2.6 Inverse	194
A.3 Eigenvalues and Eigenvectors	194
B Calculus	195
B.1 Introduction	195
B.1.1 Partial Derivatives	195
B.1.2 Second Derivatives	196
C Backpropagation Derivation	197
C.1 The XOR problem	198

Chapter 1

Introduction

1.1 Objectives

After this chapter, you should

1. understand the basic building blocks of artificial neural networks (ANNs)
2. understand the two modes of operation in ANNs
3. understand the importance of learning in ANNs
4. be able to use a simple rule to create learning in an ANN
5. begin to understand the importance of linearly inseparable problems
6. know some of the problems on which ANNs have been used.

1.2 Intelligence

This course comprises an advanced course to those new information processing simulations which are intended to emulate the information processors which we find in biology.

Traditional artificial intelligence is based on high-level symbol processing i.e. logic programming, expert systems, semantic nets etc all rely on there being in existence some high level representation of knowledge which can be manipulated by using some type of formal syntax manipulation scheme - the rules of a grammar. Such approaches have proved to be very successful in emulating human prowess in a number of fields e.g.

- we now have software which can play chess at Grand Master level
- we can match professional expertise in medicine or the law using expert systems
- we have software which can create mathematical proofs for solving complex mathematical problems.

Yet there are still areas of human expertise which we are unable to mimic using software e.g. our machines have difficulty reliably reading human handwriting, recognising human faces or exhibiting common sense. Notice how low-level the last list seems compared to the list of achievements: it has been said that the difficult things have proved easy to program whereas the easy things have proved difficult.

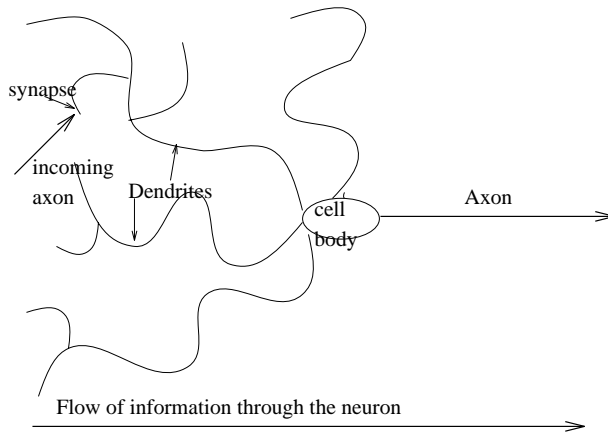


Figure 1.1: A simplified neuron

1.3 Artificial Neural Networks

Now tasks such as those discussed above seemingly require no great human expertise - young children are adept at many of these tasks. This explains the underlying presumption of creating artificial neural networks (ANNs): that the expertise which we show in this area is due to the nature of the hardware on which our brains run. Therefore if we are to emulate biological proficiencies in these areas we must base our machines on hardware (or simulations of such hardware) which seems to be a silicon equivalent to that found within our heads.

First we should be clear about what the attractive properties of human neural information processing are. They may be described as :

- Biological information processing is robust and fault-tolerant: early on in life¹, we have our greatest number of neurons yet though we daily lose many thousands of neurons we continue to function for many years without an associated deterioration in our capabilities
- Biological information processors are flexible: we do not require to be reprogrammed when we go into a new environment; we adapt to the new environment, i.e. we learn.
- We can handle fuzzy, probabilistic, noisy and inconsistent data in a way that is possible with computer programs but only with a great deal of sophisticated programming and then only when the context of such data has been analysed in detail. Contrast this with our innate ability to handle uncertainty.
- The machine which is performing these functions is highly parallel, small, compact and dissipates little power.

We shall therefore begin our investigation of these properties with a look at the biological machine we shall be emulating.

1.4 Biological and Silicon Neurons

A simplified neuron is shown in Figure 1.1. Information is received by the neuron at synapses on its dendrites. Each synapse represents the junction of an incoming axon from another neuron with a dendrite of the neuron represented in the figure; an electro-chemical transmission occurs at the synapse which allows information to be transmitted from one neuron to the next. The information is then transmitted along the dendrites till it reaches the cell body where a summation of the electrical impulses reaching the body takes place and some function of this sum is performed. If

¹ Actually several weeks before birth

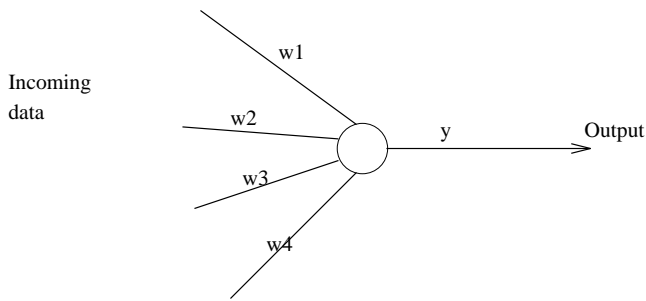


Figure 1.2: The artificial neuron. The weights model the synaptic efficiencies. Some form of processing not specified in the diagram will take place within the cell body.

this function is greater than a particular threshold the neuron will fire: this means that it will send a signal (in the form of a wave of ionisation) along its axon in order to communicate with other neurons. In this way, information is passed from one part of the network of neurons to another. It is crucial to recognise that synapses are thought to have different efficiencies and that these efficiencies change during the neuron's lifetime. We will return to this feature when we discuss learning.

We generally model the biological neuron as shown in Figure 1.2. The inputs are represented by the input vector \mathbf{x} and the synapses' efficiencies are modelled by a weight vector \mathbf{w} . Therefore the single output value of this neuron is given by

$$y = f\left(\sum_i w_i x_i\right) = f(\mathbf{w} \cdot \mathbf{x}) = f(\mathbf{w}^T \mathbf{x}) \quad (1.1)$$

You will meet all 3 ways of representing the operation of summing the weighted inputs. Sometimes $f()$ will be the identity function i.e. $f(\mathbf{x}) = \mathbf{x}$. Notice that if the weight between two neurons is positive, the input neuron's effect may be described as excitatory; if the weight between two neurons is negative, the input neuron's effect may be described as inhibitory.

Consider again Figure 1.2. Let $w_1 = 1$, $w_2 = 2$, $w_3 = -3$ and $w_4 = 3$ and let the activation function, $f()$, be the Heaviside (threshold) function such that

$$f(t) = \begin{cases} 1 & \text{if } t > 0 \\ 0 & \text{if } t \leq 0 \end{cases} \quad (1.2)$$

Now if the input vector $\mathbf{x} = (x_1, x_2, x_3, x_4) = (1, 2, 1, 2)$ then the activation of the neuron is $\mathbf{w} \cdot \mathbf{x} = \sum_j w_j x_j = 1*1 + 2*2 + 1*(-3) + 2*3 = 8$ and so $y = f(8) = 1$. However if the input vector is $(3, 1, 2, 0)$, then the activation is $3*1 + 1*2 + 2*(-3) + 0*3 = -1$ and so $y = f(-1) = 0$.

Therefore we can see that the single neuron is an extremely simple processing unit. The power of neural networks is believed to come from the accumulated power of adding many of these simple processing units together - i.e. we throw lots of simple and robust power at a problem. Again we may be thought to be emulating nature, as the typical human has several hundred billion neurons. We will often imagine the neurons to be acting in concert in layers such as in Figure 1.3.

In this figure, we have a set of inputs (the input vector, \mathbf{x}) entering the network from the left-hand side and being propagated through the network via the weights till the activation reaches the output layer. The middle layer is known as the hidden layer as it is invisible from outwith the net: we may not affect its activation directly.

1.5 Learning in Artificial Neural Networks

There are two modes in artificial neural networks:

1. activation transfer mode when activation is transmitted throughout the network

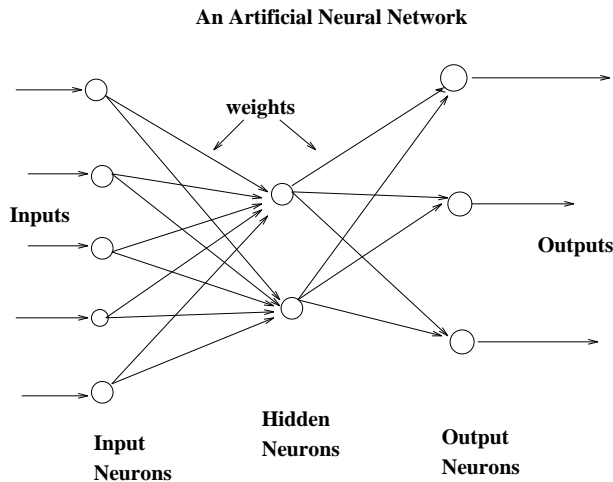


Figure 1.3: A typical artificial neural network consisting of 3 layers of neurons and 2 connecting layers of weights

2. learning mode when the network organises usually on the basis of the most recent activation transfer.

We will now consider the learning mode.

We stated that neural networks need not be programmed when they encounter novel environments. Yet their behaviour changes in order to adapt to the new environment. Such behavioural changes are due to changes in the weights in the network. We call the changes in weights in a neural network learning. The changes in weights in an artificial neural network are intended to model the changing synaptic efficiencies in real neural networks: it is believed that our learning is due to changes in the efficiency with which synapses pass information between neurons.

There are 3 main types of learning in a neural network:

Supervised learning: with this type of learning, we provide the network with input data and the correct answer i.e. what output we wish to receive given that input data. The input data is propagated forward through the network till activation reaches the output neurons. We can then compare the answer which the network has calculated with that which we wished to get. If the answers agree, we need make no change to the network; if, however, the answer which the network is giving is different from that which we wished then we adjust the weights to ensure that the network is more likely to give the correct answer in future if it is again presented with the same (or similar) input data. This weight adjustment scheme is known as supervised learning or learning with a teacher. The tutorial at the end of this chapter gives an example of supervised learning.

Unsupervised learning: with this type of learning, we only provide the network with the input data. The network is required to self-organise (i.e. to teach itself) depending on some structure in the input data. Typically this structure may be some form of redundancy in the input data or clusters in the data.

Reinforcement learning: is a half-way house between these two: we provide the network with the input data and propagate the activation forward but only tell the network if it has produced a right or a wrong answer. If it has produced a wrong answer some adjustment of the weights is done so that a right answer is more likely in future presentations of that particular piece of input data.

For many problems, the interesting facet of learning is not just that the input patterns may be learned/classified/identified precisely but that this learning has the capacity to generalise. That

is, while learning will take place on a set of *training patterns* an important property of the learning is that the network can generalise its results on a set of *test patterns* which it has not seen during learning. One of the important consequences here is that there is a danger of overlearning a set of training patterns so that new patterns which are not part of the training set are not properly classified.

For much of this course we will concentrate on unsupervised learning; the major exceptions occur in Chapters 5 and 8 in which we will use supervised learning methods.

1.6 Typical Problem Areas

The number of application areas in which artificial neural networks are used is growing daily. Here we simply produce a few representative types of problems on which neural networks have been used

Pattern completion: ANNs can be trained on sets of visual patterns represented by pixel values. If subsequently, a part of an individual pattern (or a noisy pattern) is presented to the network, we can allow the network's activation to propagate through the network till it converges to the original (memorised) visual pattern. The network is acting like a content-addressable memory. Typically such networks have a recurrent (feedback as opposed to a feedforward) aspect to their activation passing. You will sometimes see this described as a network's *topology*.

Classification: An early example of this type of network was trained to differentiate between male and female faces. It is actually very difficult to create an algorithm to do so yet an ANN has been shown to have near-human capacity to do so.

Optimisation: It is notoriously difficult to find algorithms for solving optimisation problems. A famous optimisation problem is the Travelling Salesman Problem in which a salesman must travel to each of a number of cities, visiting each one once and only once in an optimal (i.e. least distance or least cost) route. There are several types of neural networks which have been shown to converge to 'good-enough' solutions to this problem i.e. solutions which may not be globally optimal but can be shown to be close to the global optimum for any given set of parameters.

Feature detection: An early example of this is the phoneme producing feature map of Kohonen: the network is provided with a set of inputs and must learn to pronounce the words; in doing so, it must identify a set of features which are important in phoneme production.

Data compression: There are many ANNs which have been shown to be capable of representing input data in a compressed format losing as little of the information as possible; for example, in image compression we may show that a great deal of the information given by a pixel to pixel representation of the data is redundant and a more compact representation of the image can be found by ANNs.

Approximation: Given examples of an input to output mapping, a neural network can be trained to approximate the mapping so that a future input will give approximately the correct answer i.e. the answer which the mapping should give.

Association: We may associate a particular input with a particular output so that given the same (or similar) output again, the network will give the same (or a similar) output again.

Prediction: This task may be stated as: given a set of previous examples from a time series, such as a set of closing prices for the FTSE, to predict the next (future) sample.

Control: For example to control the movement of a robot arm (or truck, or any non-linear process) to learn what inputs (actions) will have the correct outputs (results).

1.7 A short history of ANNs

The history of ANNs started as long ago as 1943 when McCulloch and Pitts showed that simple neuron-like building blocks were capable of performing all possible logic operations. At that time too, Von Neumann and Turing discussed interesting aspects of the statistical and robust nature of brain-like information processing, but it was only in the 1950s that actual hardware implementations of such networks began to be produced. The most enthusiastic proponent of the new learning machines was Frank Rosenblatt who invented the *perceptron* machine, which was able to perform simple pattern classification.

However, it became apparent that the new learning machines were incapable of solving certain problems and in 1969 Minsky and Papert wrote a definitive treatise, 'Perceptrons', which clearly demonstrated that the networks of that time had limitations which could not be transcended. The core of the argument against networks of that time may be found in their inability to solve XOR problems (see Chapter 5). Enthusiasm for ANNs decreased till the mid '80s when first John Hopfield, a physicist, analysed a particular class of ANNs and proved that they had powerful pattern completion properties and then in 1986 the subject really took off when Rumelhart, McClelland and the PDP Group rediscovered powerful learning rules which transcended the limitations discussed by Minsky and Papert.

1.8 A First Tutorial

This tutorial will emphasise learning in neural nets. Recall that learning is believed to happen at the synapses (the meeting points between neurons) and that we model synaptic efficiency with weights.

You are going to hand simulate a simple neural net (see Figure 1.4) performing classification according to the AND (see table) OR and XOR rules. We will use a network with three input neurons - the two required for the input values and a third neuron known as the bias neuron. The bias always fires 1 (i.e. is constant).

You will work in groups of 3 - one person in charge of selecting the input, one in charge of calculating the feedforward value of the neuron, and one person in charge of calculating the change in weights.

To begin with, the group selects random (between 0 and 1) values for the three weights shown in the figure.

1. The INPUTER selects randomly from the set of patterns shown in the table and places the cards in the appropriate places on the table.
2. The FEEDFORWARDER multiplies the weights by the input patterns to calculate the output.

$$y = \sum_{i=0}^2 w_i x_i \quad (1.3)$$

Then $y = 1$ if $y > 0$, $y = -1$ if $y < 0$.

3. The WEIGHTCHANGER changes the weights *when the value of y is not the same as the target y_T* according to the formula

$$w_i = w_i + \eta * (y_T - y) * x_i \quad (1.4)$$

Steps (1)-(3) are repeated as often as necessary.

1.8.1 Worked Example

Let us have initial values $w_0 = 0.5$, $w_1 = 0.3$, $w_2 = 0.7$ and let η be 0.1.

Bias	first input	second input	target output
1	1	1	1
1	1	-1	-1
1	-1	1	-1
1	-1	-1	-1

Table 1.1: The values for the AND patterns

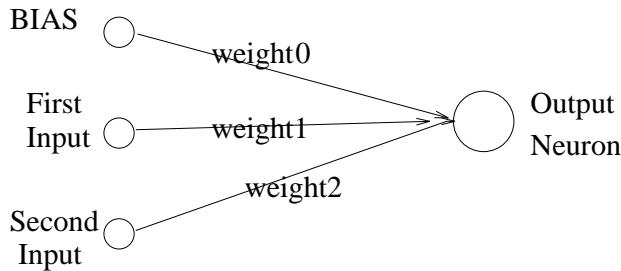


Figure 1.4: The Simple Neural Network

1. “Randomly” select pattern 1. The FEEDFORWARDER calculates $y = 0.5 + 0.3 + 0.7 = 1.5$. So $y=1$ which is the same as the target and so the WEIGHTCHANGER does nothing.
2. Imagine pattern 2 is chosen. The FEEDFORWARDER calculates $y = 0.5 + 0.3 - 0.7 = 0.1$. So $y=1$. Now $y_T = -1$ and so WEIGHTCHANGER calculates

$$\begin{aligned} w_0 &= w_0 + 0.1 * (-2) * 1 = 0.5 - 0.2 = 0.3 \\ w_1 &= w_1 + 0.1 * (-2) * 1 = 0.3 - 0.2 = 0.1 \\ w_2 &= w_2 + 0.1 * (-2) * (-1) = 0.7 + 0.2 = 0.9 \end{aligned}$$

3. Now pattern 3 is chosen. The FEEDFORWARDER calculates $y = 0.3 - 0.1 + 0.9 = 1.1$. So $y=1$ and $y_T = -1$ and so WEIGHTCHANGER calculates

$$\begin{aligned} w_0 &= w_0 + 0.1 * (-2) * 1 = 0.3 - 0.2 = 0.1 \\ w_1 &= w_1 + 0.1 * (-2) * (-1) = 0.1 + 0.2 = 0.3 \\ w_2 &= w_2 + 0.1 * (-2) * 1 = 0.9 - 0.2 = 0.7 \end{aligned}$$

4. Now pattern 4 is chosen. The FEEDFORWARDER calculates $y = 0.1 - 0.3 - 0.7 = -0.9$. So $y = -1$ and $y_T = -1$. Therefore the WEIGHTCHANGER does nothing.
5. Now select pattern 2. The FEEDFORWARDER calculates $y = 0.1 - 0.3 + 0.7 = 0.5$. Then $y=1$ but $y_T = -1$. WEIGHTCHANGER calculates

$$\begin{aligned} w_0 &= w_0 + 0.1 * (-2) * 1 = 0.1 - 0.2 = -0.1 \\ w_1 &= w_1 + 0.1 * (-2) * (-1) = 0.3 + 0.2 = 0.5 \\ w_2 &= w_2 + 0.1 * (-2) * 1 = 0.7 - 0.2 = 0.5 \end{aligned}$$

6. Now all of the patterns give the correct response (try it).

We can draw the solution found by using the weights as the parameters of the line $ax_1 + bx_2 + c = 0$. Using $a = w_1, b = w_2, c = w_0$ we get

$$0.5x_1 + 0.5x_2 - 0.1 = 0 \tag{1.5}$$

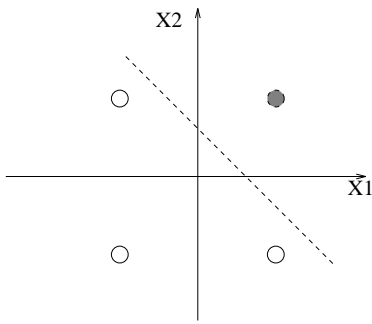


Figure 1.5: The line joining (0,0.2) and (0.2,0) cuts the 4 points into 2 sets correctly

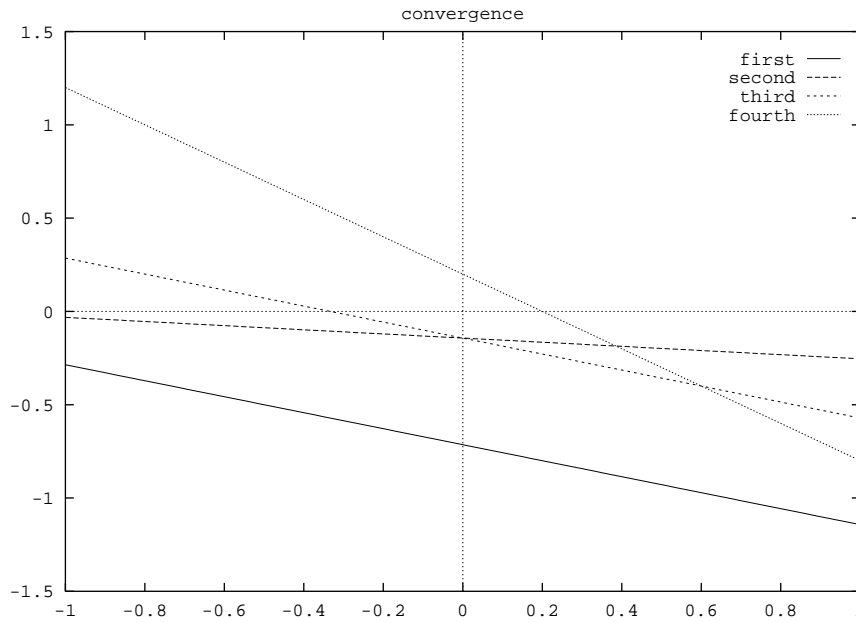


Figure 1.6: The iterative convergence of the network to a set of weights which can perform the correct mapping is shown diagrammatically here.

which we can draw by getting two points. If $x_1 = 0$, then $0.5x_2 = 0.1$ and so $x_2 = 0.2$ which gives us one point (0,0.2). Similarly we can find another point (0.2,0) which is drawn in Figure 1.5. Notice the importance of the BIAS weight: it allows a solution to be found which does not go through the origin; without a bias we would have to have a moving threshold.

We can see the convergence of $w_0 + w_1x_1 + w_2x_2 = 0$ in Figure 1.6. Notice that initially 2 patterns are correctly classified, very quickly a third is correctly classified and only on the fourth change are all 4 patterns correctly classified.

1.8.2 Exercises

1. Repeat the worked example with different initial values for the weights. (Objectives 3, 4).
2. Repeat for the OR patterns. (Objectives 3, 4).
3. Experiment with different initial values, learning rates. (Objectives 3, 4).
4. Now do the XOR problem. Don't spend too long on it - it's impossible. (Objective 5).

Nut	type A - 1	type A - 2	type A - 3	type B - 1	type B - 2	type B - 3
Length (cm)	2.2	1.5	0.6	2.3	1.3	0.3
Weight (g)	1.4	1.0	0.5	2.0	1.5	1.0

Table 1.2: The lengths and weights of six instances of two types of nuts

5. Draw the XOR coordinates and try to understand why it is impossible. (Objective 5).
6. Now we will attempt to train a network to classify the data shown in Table 1.2. Then we will train a network with the input vectors , \mathbf{x} , equal to
 - (1, 2.2, 1.4) with associated training output 1 (equal to class A)
 - (1, 1.5, 1.0) with associated training output 1
 - (1, 0.6, 0.5) with associated training output 1
 - (1, 2.3, 2.0) with associated training output -1 (equal to class B)
 - (1, 1.3,1.5) with associated training output -1
 - (1, 0.3,1.0) with associated training output -1
 Note that the initial 1 in each case corresponds to the bias term
7. Describe in your own words an Artificial Neural Network and its operation (Objectives 1, 2).
8. Describe some typical problems which people have used ANNs to solve. Attempt to describe what features of the problems have made them attractive to ANN solutions. (Objective 6).

