

Introduction to the Math of Neural Networks (Beta-1)

Introduction to the Math of Neural Networks (Beta-1)

Jeff Heaton

Heaton Research, Inc.
St. Louis, MO, USA

Do not make illegal copies of this ebook

Title	Introduction to the Math of Neural Networks (Beta-1)
Author	Jeff Heaton
Published	November 01, 2011
Copyright	Copyright 2011 by Heaton Research, Inc., All Rights Reserved.
ISBN	978-1-60439-033-9
Price	9.99 USD
File Created	Sat Oct 15 17:44:58 CDT 2011

This eBook is copyrighted material, and public distribution is prohibited. If you did not receive this ebook from Heaton Research (<http://www.heatonresearch.com>), or an authorized bookseller, please contact Heaton Research, Inc. to purchase a licensed copy. DRM free copies of our books can be purchased from:

<http://www.heatonresearch.com/book>

If you did purchase this book, thankyou! Your purchase of this books supports the Encog Machine Learning Framework. <http://www.encog.org/>

Publisher: Heaton Research, Inc
Introduction to the Math of Neural Networks
October, 2011
Author: Jeff Heaton
Editor: WordsRU.com
Cover Art: Carrie Spear
ISBN's for all Editions:
978-1-60439-033-9, Paperback
978-1-60439-034-6, PDF
978-1-60439-035-3, Nook
978-1-60439-036-0, Kindle

Copyright ©2011 by Heaton Research Inc., 1734 Clarkson Rd. #107, Chesterfield, MO 63017-4976. World rights reserved. The author(s) created reusable code in this publication expressly for reuse by readers. Heaton Research, Inc. grants readers permission to reuse the code found in this publication or downloaded from our website so long as (author(s)) are attributed in any application containing the reusable code and the source code itself is never redistributed, posted online by electronic transmission, sold or commercially exploited as a stand-alone product. Aside from this specific exception concerning reusable code, no part of this publication may be stored in a retrieval system, transmitted, or reproduced in any way, including, but not limited to photo copy, photograph, magnetic, or other record, without prior agreement and written permission of the publisher.

Heaton Research, Encog, the Encog Logo and the Heaton Research logo are all trademarks of Heaton Research, Inc., in the United States and/or other countries.

TRADEMARKS: Heaton Research has attempted throughout this book to distinguish proprietary trademarks from descriptive terms by following the capitalization style used by the manufacturer.

The author and publisher have made their best efforts to prepare this book, so the content is based upon the final release of software whenever possible. Portions of the manuscript may be based upon pre-release versions supplied by software manufacturer(s). The author and the publisher make no representation or warranties of any kind with regard to the completeness or accuracy of the contents herein and accept no liability of any kind including but not limited to performance, merchantability, fitness for any particular purpose, or any

losses or damages of any kind caused or alleged to be caused directly or indirectly from this book.

SOFTWARE LICENSE AGREEMENT: TERMS AND CONDITIONS

The media and/or any online materials accompanying this book that are available now or in the future contain programs and/or text files (the “Software”) to be used in connection with the book. Heaton Research, Inc. hereby grants to you a license to use and distribute software programs that make use of the compiled binary form of this book’s source code. You may not redistribute the source code contained in this book, without the written permission of Heaton Research, Inc. Your purchase, acceptance, or use of the Software will constitute your acceptance of such terms.

The Software compilation is the property of Heaton Research, Inc. unless otherwise indicated and is protected by copyright to Heaton Research, Inc. or other copyright owner(s) as indicated in the media files (the “Owner(s)”). You are hereby granted a license to use and distribute the Software for your personal, noncommercial use only. You may not reproduce, sell, distribute, publish, circulate, or commercially exploit the Software, or any portion thereof, without the written consent of Heaton Research, Inc. and the specific copyright owner(s) of any component software included on this media.

In the event that the Software or components include specific license requirements or end-user agreements, statements of condition, disclaimers, limitations or warranties (“End-User License”), those End-User Licenses supersede the terms and conditions herein as to that particular Software component. Your purchase, acceptance, or use of the Software will constitute your acceptance of such End-User Licenses.

By purchase, use or acceptance of the Software you further agree to comply with all export laws and regulations of the United States as such laws and regulations may exist from time to time.

SOFTWARE SUPPORT

Components of the supplemental Software and any offers associated with them may be supported by the specific Owner(s) of that material but they are not supported by Heaton Research, Inc.. Information regarding any available support may be obtained from the Owner(s) using the information provided in the appropriate README files or listed elsewhere on the media.

Should the manufacturer(s) or other Owner(s) cease to offer support or decline to honor any offer, Heaton Research, Inc. bears no responsibility. This notice concerning support for the Software is provided for your information only. Heaton Research, Inc. is not the agent or principal of the Owner(s), and Heaton Research, Inc. is in no way responsible for providing any support for the Software, nor is it liable or responsible for any support provided, or not provided, by the Owner(s).

WARRANTY

Heaton Research, Inc. warrants the enclosed media to be free of physical defects for a period of ninety (90) days after purchase. The Software is not available from Heaton Research, Inc. in any other form or media than that enclosed herein or posted to www.heatonresearch.com. If you discover a defect in the media during this warranty period, you may obtain a replacement of identical format at no charge by sending the defective media, postage prepaid, with proof of purchase to:

Heaton Research, Inc.
Customer Support Department
1734 Clarkson Rd #107
Chesterfield, MO 63017-4976
Web: www.heatonresearch.com
E-Mail: support@heatonresearch.com

DISCLAIMER

Heaton Research, Inc. makes no warranty or representation, either expressed or implied, with respect to the Software or its contents, quality, performance, merchantability, or fitness for a particular purpose. In no event will Heaton Research, Inc., its distributors, or dealers be liable to you or any other party for direct, indirect, special, incidental, consequential, or other damages arising out of the use of or inability to use the Software or its contents even if advised of the possibility of such damage. In the event that the Software includes an online update feature, Heaton Research, Inc. further disclaims any obligation to provide this feature for any specific duration other than the initial posting.

The exclusion of implied warranties is not permitted by some states. Therefore, the above exclusion may not apply to you. This warranty provides you with specific legal rights; there

may be other rights that you may have that vary from state to state. The pricing of the book with the Software by Heaton Research, Inc. reflects the allocation of risk and limitations on liability contained in this agreement of Terms and Conditions.

SHAREWARE DISTRIBUTION

This Software may use various programs and libraries that are distributed as shareware. Copyright laws apply to both shareware and ordinary commercial software, and the copyright Owner(s) retains all rights. If you try a shareware program and continue using it, you are expected to register it. Individual programs differ on details of trial periods, registration, and payment. Please observe the requirements stated in appropriate files.

Contents

Introduction	xvii
0.1 Other Resources	xix
0.2 Structure of this Book	xix
1 Neural Network Activation	1
1.1 Understanding the Summation Operator	1
1.2 Calculating a Neural Network	2
1.3 Activation Functions	5
1.4 Bias Neurons	7
1.5 Chapter Summary	9
2 Error Calculation Methods	13
2.1 The Error Function	14
2.2 Calculating Global Error	15
2.3 Other Error Calculation Methods	16
2.3.1 Sum of Squares Error	16
2.3.2 Root Mean Square Error	16
2.4 Chapter Summary	17
3 Derivatives	19
3.1 Calculating the Slope of a Line	19

3.2	What is a Derivative?	21
3.3	Using Partial Derivatives	25
3.4	Using the Chain Rule	26
3.5	Chapter Summary	27
4	Backpropagation	29
4.1	Understanding Gradients	29
4.1.1	What is a Gradient	30
4.1.2	Calculating Gradients	32
4.1.3	Calculating the Node Deltas	33
4.1.4	Calculating the Individual Gradients	35
4.2	Applying Back Propagation	36
4.2.1	Batch and Online Training	37
4.2.2	Backpropagation Weight Update	37
4.3	Chapter Summary	40
5	RPROP	43
5.1	RPROP Arguments	44
5.2	Data Structures	46
5.3	Understanding RPROP	47
5.3.1	Determine Sign Change of Gradient	47
5.3.2	Calculate Weight Change	48
5.3.3	Modify Update Values	49
5.4	RPROP Update Examples	49
5.4.1	Training Iteration #1	50
5.4.2	Training Iteration #2	51
5.4.3	Training Iteration #8	53
5.5	Chapter Summary	54
6	Weight Initialization	57

6.1	Looking at the Weights	58
6.2	Ranged Randomization	59
6.3	Using Nguyen-Widrow	60
6.3.1	Performance of Nguyen-Widrow	60
6.3.2	Implementing Nguyen-Widrow	61
6.3.3	Nguyen-Widrow in Action	62
6.4	Chapter Summary	64
7	LMA Training	67
7.1	Calculation of the Hessian	69
7.2	LMA with Multiple Outputs	71
7.3	Overview of the LMA Process	72
7.4	Chapter Summary	73
8	Self-Organizing Maps	75
8.1	SOM Structure	76
8.1.1	Best Matching Unit	78
8.2	Training a SOM	79
8.2.1	SOM Training Example	79
8.2.2	Training the SOM Example	81
8.2.3	BMU Calculation Example	81
8.2.4	Example Neighborhood Functions	82
8.2.5	Example Weight Update	84
8.3	SOM Error Calculation	86
8.4	Chapter Summary	87
9	Normalization	89
9.1	Simple Normalization and Denormalization	90
9.1.1	Reciprocal Normalization	90
9.1.2	Reciprocal Denormalization	90

- 9.2 Range Normalization and Denormalization 91
 - 9.2.1 Range Normalization 91
 - 9.2.2 Range Denormalization 91
- 9.3 Chapter Summary 92

Introduction

- Math Needed for Neural Networks
- Other Resources
- Prerequisites

If you have read other books by me you will know that I try to shield the reader from the mathematics behind AI. Often you do not need to know the exact math that is used to train a neural network or perform a cluster operation. You simply want the result.

This is very much the idea of the Encog project. Encog is an advanced machine learning framework that allows you to perform many advanced operations such as neural networks, genetic algorithms, support vector machines, simulated annealing, and other machine learning methods. You are allowed to use these advanced techniques without the need to know what is happening behind the scenes.

However, sometimes you really do want to know what is going on behind the scenes. You do want to know that math that is involved. In this book you will learn what happens, behind the scenes, with a neural network. You will also be exposed to the math. I will present the material in mathematical terms.

There are already many neural network books that at first glance would appear as a math text. This is not what I seek to produce here. There are already several very good books that achieve a pure mathematical introduction to neural networks. My goal is to produce a mathematically based neural network book that targets someone with perhaps only a college algebra and computer programming background. These are the only two prerequisites for this book. Actually, there is a third prerequisite, but I will get to that in a moment.

Neural networks overlap several bodies of mathematics. Neural network goals, such as classification, regression and clustering come from statistics. The gradient descent that goes

into backpropagation, and other training methods, requires knowledge of Calculus. Advanced training, such as Levenberg Marquardt, require both Calculus and Matrix Mathematics.

To read nearly any academic level neural network, or machine learning, targeted book you will need some knowledge of Algebra, Calculus, Statistics and Matrix Mathematics. However, the reality is only need a relatively small amount of knowledge from each of these areas. The goal of this book, is to teach you enough math to understand neural networks and their training. You will understand exactly how a neural network functions, and should be able to implement your own in any computer language you are familiar with.

As areas of mathematics are needed, I will provide an introductory-level tutorial on the math. I only assume that you know basic algebra to start out with. This book will discuss such mathematical concepts as derivatives, partial derivatives, matrix transformation, gradient descent and more.

If you have not done this sort of math in awhile, I plan for this book to be a good refresher. If you have never done this sort of math, then this book could serve as a good introduction. If you are very familiar with math, you can still learn neural networks from this book. However, you may want to skip some of the sections that seem too basic.

This book is not about Encog. Nor is it about how to program in any particular programming language. I assume you will likely apply these principles to programming languages. If you want examples of how I apply the principles in this book there is Encog. This book is really more about the algorithms and mathematics behind neural networks.

I did say there was one other prerequisite to this book. That is, other than basic algebra and programming knowledge in any language. That is knowledge of what a neural network is, and how it is used. You should already be familiar with what neural networks are and how they are used. If you do not yet know how to use a neural network, you may want to start here.

<http://goo.gl/8ESx>

The above article, that I wrote, provides a brief crash course on what neural networks are. You may want to look at some of the Encog examples, as well. You can find more information about Encog at the following URL.

<http://www.heatonresearch.com/encog/>

If neural networks are cars, then this book is a mechanics guide. If I am going to teach you to repair and build cars, I have two basic assumptions in order of importance. First is that you've actually seen a car, and know what one is used for. The second assumption is that you can actually drive a car. If neither of these are true, then why do you care about

learning the internals of how a car works? The same is true of neural networks.

0.1 Other Resources

There are many other resources on the internet that will be very useful as you read through this book. They are listed in this section.

First, is the Khan Academy. This is a collection of YouTube videos that demonstrate many areas of mathematics. If you need additional review on any mathematical concept in this book, there is likely a video on the Khan Academy that covers it.

<http://www.khanacademy.org/>

Second is the Neural Network FAQ. This text-only resource has a great deal of information for neural networks.

<http://www.faqs.org/faqs/ai-faq/neural-nets/>

The Encog wiki has a fair amount of general information on machine learning. This information is not necessarily tied to Encog. There are articles in the Encog wiki that will be helpful as you complete this book.

http://www.heatonresearch.com/wiki/Main_Page

Finally, the Encog forums are a place where AI and neural networks can be discussed. These forums are fairly active and you will likely receive an answer from myself or one of the community members at the forum.

<http://www.heatonresearch.com/forum>

These resources should be helpful as you progress through this book.

0.2 Structure of this Book

The first chapter, “Neural Network Activation”, shows how the output from a neural network is calculated. Before you can see how to train and evaluate a neural network you must understand how a neural network produces its output.

The second chapter, named “Error Calculation”, demonstrates how to evaluate the output from a neural network. Neural networks begin with random weights. Training adjusts these weights to produce meaningful output.

The third chapter, “Understanding Derivatives”, focuses entirely on a very important Calculus topic. Derivatives, and partial derivatives, are used by several neural network training methods. This chapter will introduce you to those aspects of derivatives that are needed for this book.

Chapter 4, “Training with Backpropagation”, shows you how to apply knowledge from Chapter three towards training a neural network. Backpropagation is one of the oldest training techniques for neural networks. There newer, and much superior, training methods available. However, understanding backpropagation provides a very important foundation for RPROP, QPROP and LMA.

Chapter 5, “Faster Training with RPROP”, introduces resilient propagation (RPROP) which builds upon backpropagation to provide much quicker training times.

Chapter 6, “Weight Initialization”, shows how neural networks are given their initial random weights. Some sets of random weights perform better than others. This chapter looks at several, less than random, weight initialization methods.

Chapter 7, “LMA Training”, introduces the Levenberg Marquardt Algorithm (LMA). LMA is the most mathematically intense training method in this book. LMA sometimes offers very rapid training for a neural network.

Chapter 8, “Self Organizing Maps” shows how to create a clustering neural network. The SOM can be used to group data. The structure of the SOM is similar to the feedforward neural networks seen in this book.

Chapter 9, “Normalization” shows how numbers are normalized for neural networks. Neural networks typically require that input and output numbers be in the range of 0 to 1, or -1 to 1. This chapter shows how to transform numbers into that range.

Chapter 1

Neural Network Activation

- Summation
- Calculating Activation
- Activation Functions
- Bias Neurons

In this chapter you will see how to calculate the output for a feedforward neural network. Most neural networks are in some way based on the feedforward neural network. Seeing how this simple neural network is calculated will form the foundation for understanding training, and other more complex features of neural networks.

Several mathematical terms will be introduced in this chapter. You will be shown summation notation and simple mathematical formula notation. We will begin with a review of the summation operator.

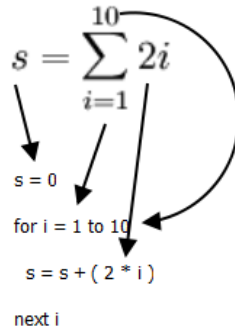
1.1 Understanding the Summation Operator

In this section, we will take a quick look at the summation operator. The summation operator, represented by the capital Greek letter sigma can be seen in Equation 1.1.

$$s = \sum_{i=1}^{10} 2i \tag{1.1}$$

The above equation is a summation. If you are unfamiliar with sigma notation, it is essentially the same thing as a programming **for** loop. Figure 1.1 shows Equation 1.1 reduced to pseudocode.

Figure 1.1: Summation Operator to Code

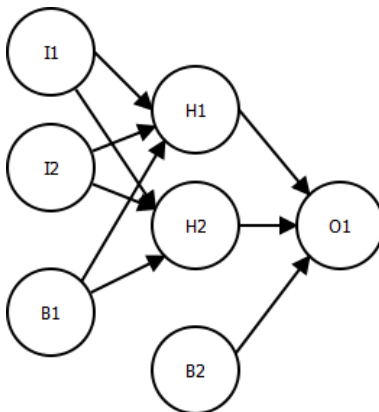


As you can see, the summation operator is very similar to a **for** loop. The information just below the sigma symbol specifies the starting value and the indexing variable. The information above the sigma specifies the limit of the loop. The information to the right of sigma specifies the value that is being summed.

1.2 Calculating a Neural Network

We will begin by looking at how a neural network calculates its output. You should already know the structure of a neural network from the book's introduction. Consider a neural network such as the one in Figure 1.2.

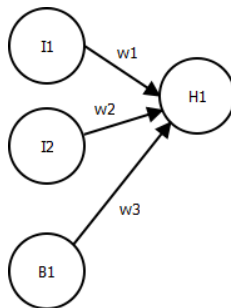
Figure 1.2: A Simple Neural Network



This neural network has one output neuron. As a result, it will have one output value. To calculate the value of this output neuron (**O1**), we must calculate the activation for each of the inputs into **O1**. The inputs that feed into **O1** are **H1**, **H2** and **B2**. The activation for **B2** is simply 1.0, because it is a bias neuron. However, **H1** and **H2** must be calculated independently. To calculate **H1** and **H2**, the activations of **I1**, **I2** and **B1** must be considered. Though **H1** and **H2** share the same inputs, they will not calculate to the same activation. This is because they have different weights. The weights are represented by lines in the above diagram.

First we must look at how one activation calculation is done. This same activation calculation can then be applied to the other activation calculations. We will examine how **H1** is calculated. Figure 1.3 shows only the inputs to **H1**.

Figure 1.3: Calculating H1's Activation



We will now see how to calculate H1. This relatively simple equation is shown in Equation 1.2.

$$h_1 = A\left(\sum_{c=1}^n (i_c * w_c)\right) \quad (1.2)$$

To understand Equation 1.2 we first examine the variables that go into it. For the above equation we have three input values, given by the variable **i**. The three input values are input values of **I1**, **I2** and **B1**. **I1** and **I2** are simply the input values that the neural network was provided to compute the output. **B1** is always 1, because it is the bias neuron.

There are also three weight values considered **w1**, **w2** and **w3**. These are the weighted connections between **H1** and the previous layer. Therefore, the variables to this equation are:

```

i [1] = first input value to the neural network
i [2] = second input value to neural network
i [3] = 1
w [1] = weight from I1 to H1
w [2] = weight from I2 to H1
w [3] = weight from B1 to H1
n = 3, the number of connections

```

Though the bias neuron is not really part of the input array, a one is always placed into the input array for the bias neuron. Treating the bias as a forward-only neuron makes the calculation much easier.

To understand Equation 1.2 we will consider it as pseudocode.

```

double w[3] // the weights
double i[3] // the input values
double sum = 0; // the sum
// perform the summation (sigma)
for c = 0 to 2
    sum = sum + ( w[c] * i[c] )
next
// apply the activation function
sum = A(sum)

```

Here we sum up each of the inputs times its respective weight. Finally this sum is passed to an activation function. Activation functions are a very important concept in neural network programming. In the next section we will examine activation functions.

1.3 Activation Functions

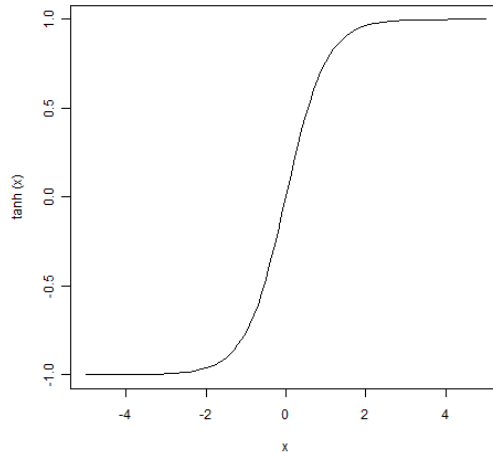
Activation functions are used very commonly in neural networks. Activation functions serve several important functions for a neural network. The primary reason to use an activation function is to introduce non-linearity to the neural network. Without this non-linearity a neural network could do little to learn non-linear functions. The output that we expect neural networks to learn is rarely linear.

The two most common activation functions are the sigmoid and hyperbolic tangent activation function. The hyperbolic tangent activation function is the more common of these two, as has a number range from -1 to 1, compared to the sigmoid function which is only from 0 to 1.

$$f(x) = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (1.3)$$

The hyperbolic tangent function is actually a trigonometric function. However, our use for it has nothing to do with trigonometry. This function was chosen for the shape of its graph. You can see a graph of the hyperbolic tangent function in Figure 1.4.

Figure 1.4: The Hyperbolic Tangent Function



Notice the range is from -1 to 1? This gives it a much wider range of numbers it can

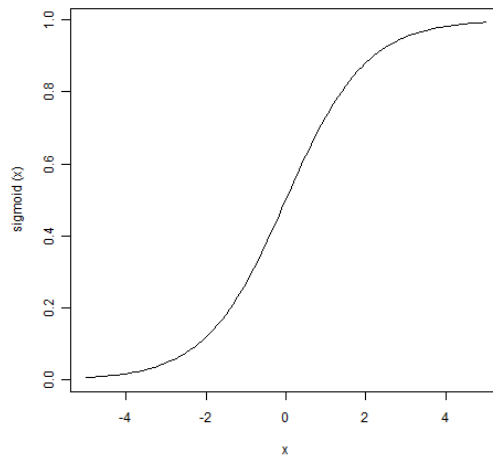
accept. Also notice how values beyond -1 to 1 are quickly scaled? This provides a consistent range of numbers for the network.

Now we will look at the sigmoid function. You can see this in Equation 1.4.

$$f(x) = \frac{1}{1 + e^{-x}} \quad (1.4)$$

The sigmoid function is also called the logistic function. Typically it does not perform as well as the hyperbolic tangent function. However, if you have all positive values in the training data, it can perform well. The graph for the sigmoid function is shown in Figure 1.5.

Figure 1.5: The Sigmoid Function

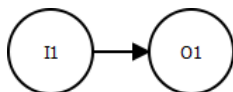


As you can see it scales numbers to 1.0. It also has a range that only includes positive numbers. It is less general purpose than hyperbolic tangent, but it can be useful. For example, the XOR operator, as discussed in the introduction, only has positive numbers. The sigmoid function outperforms the hyperbolic tangent function.

1.4 Bias Neurons

You may be wondering why bias values are even needed? Bias values allow a neural network to output a value of zero even when the input is near one. Adding a bias allows the output of the activation function to be shifted to the left or right on the x-axis. To see this, consider a simple neural network where a single input neuron **I1** is directly connected to an output neuron **O1**. The network shown in Figure 1.6 has no bias.

Figure 1.6: A Bias-less Connection

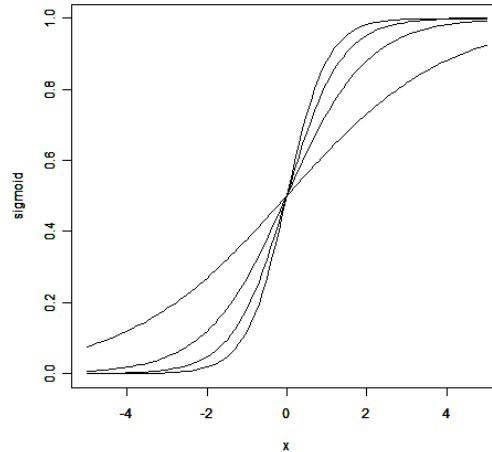


This network's output is computed by multiplying the input (\mathbf{x}) by the weight (\mathbf{w}). The result is then passing an Activation Function. In this case, we are using the Sigmoid Activation Function.

Consider the output of the sigmoid function for the following four weights.

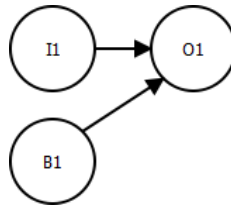
```
sigmoid(0.5*x)
sigmoid(1.0*x)
sigmoid(1.5*x)
sigmoid(2.0*x)
```

Given the above weights, the output of the sigmoid will be as seen in Figure 1.7.

Figure 1.7: Adjusting Weights

Changing the weight w alters the "steepness" of the sigmoid function. This allows the neural network to learn patterns. However what if wanted the network to output 0 when x is a value other than 0, such as 3? Only changing the steepness of the sigmoid will not accomplish this. You must be able to shift the entire curve to the right.

That is the purpose of bias. Adding a bias neuron causes the neural network to appear as Figure 1.8.

Figure 1.8: A Biased Connection

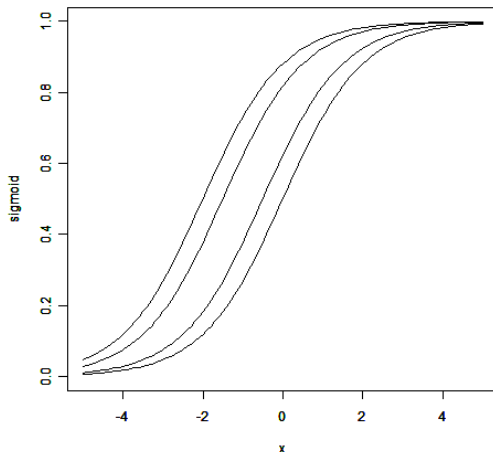
Now we calculate with the bias neuron present. We will calculate for several bias weights.

```
sigmoid(1*x + 1*1)
sigmoid(1*x + 0.5*1)
sigmoid(1*x + 1.5*1)
```

```
sigmoid(1*x + 2*1)
```

This produces the following plot, seen in Figure 1.9.

Figure 1.9: Adjusting Bias



As you can see, the entire curve now shifts.

1.5 Chapter Summary

This chapter demonstrated how a feedforward neural network calculates output. The output of a neural network is determined by calculating each successive layer, after the input layer. The final output of the neural network eventually reaches the output layer.

Neural networks make use of activation functions. An activation function provides non-linearity to the neural network. Because most of the data that a neural network seeks to learn is non-linear the activation functions must be non-linear. An activation function is applied after the weights and activations have been multiplied.

Most neural networks have bias neurons. Bias is an important concept for neural networks. Bias neurons are added to every non-output layer of the neural network. Bias neurons are different than ordinary neurons in two very important ways. First, the output from a

bias neuron is always one. Second, a bias neuron has no inbound connections. The constant value of one allows the a layer to respond with non-zero values even when the input to the layer is zero. This can be very important for certain data sets.

The neural networks will output values determined by the weights of the connections. These weights are usually set to random initial values. Training is the process where these random weights are adjusted to produce meaningful results. We need a way for the neural network to measure the effectiveness of the neural network. This measure is called error calculation. Error calculation is discussed in the next chapter.

Chapter 2

Error Calculation Methods

- Understanding Error Calculation
- The Error Function
- Error Calculation Methods
- How is the Error Used

In this chapter we will see how to calculate errors for a neural network. When performing supervised training a neural network's actual output must be compared against the ideal output specified in the training data. The difference between actual and ideal output is the error of the neural network.

Error calculation occurs at two levels. First there is the local error. This is the difference between the actual output of one individual neuron and the ideal output that was expected. The local error is calculated using an error function.

The local errors are aggregated together to form a global error. The global error is the measurement of how well a neural network performs to the entire training set. There are several different means by which a global error can be calculated. The global error calculation methods discussed in this chapter are listed here.

- Sum of Squares Error (ESS)
- Mean Square Error (MSE)

- Root Mean Square (RMS)

Usually you will simply use MSE. MSE is the most common means of calculating errors for a neural network. Later in the book we will use ESS. The Levenberg Marquardt Algorithm (LMA), which will be covered in Chapter 8, requires ESS. Lastly, RMS can be useful in certain situations.

2.1 The Error Function

We will start with the local error. The local error comes from the error function. The error function is fed the actual and ideal outputs for a single output neuron. The error function then produces a number that represents the error of that output neuron. Training methods will see to minimize this error.

This book will cover two error functions. The first is the simple linear error function that is most commonly used. The second is the arctangent error function that is introduced by the Quick Propagation training method. Arctangent error functions and Quick Propagation will be discussed in Chapter 6, “Quick Propagation”. This chapter will focus on the standard linear error function. The formula for the linear error function can be seen in Equation 2.1.

$$E = (i - a) \tag{2.1}$$

The linear error function is very simple. The error is simply the difference between the ideal (**i**) and actual (**a**) outputs from the neural network. The only requirement of the error function is that it produce an error that you would like to minimize.

To see an example of this, consider a neural network output neuron that produced 0.9 when it should have produced 0.8. The error for this neural network would be the difference between 0.8 and 0.9, which is -0.1.

In some cases you may not provide an ideal output to the neural network and still use supervised training. In this case you would write an error function that somehow evaluates the output of the neural network for the given input. This evaluation error function would need to assign some sort of a score to the neural network. A higher number would indicate less desirable output; a lower number would indicate more desirable output. The training process would attempt to minimize this score.

2.2 Calculating Global Error

Now that we have seen how to calculate the local error, we will move on to global error. Because MSE error calculation is the most common, we will begin here. You can see the equation that is used to calculate MSE in Equation 2.2.

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n E^2 \quad (2.2)$$

As you can see from the above equation, it makes use of the local error (E) that we defined in the last section. Each local error is squared and summed. The resulting sum is then divided by the total number of cases. In this way the MSE error is similar to a traditional average, except that each local error is squared. The squaring negates the effect of some errors being positive and others being negative. This is because a positive number squared is a positive number, just as a negative number square is also a positive number. If you are unfamiliar with the summation operator, shown as a capital Greek letter sigma, refer to Chapter 1.

The MSE error is typically written as a percentage. The goal is to decrease this error percentage as training progresses. To see how this is used consider the following program output.

```

Beginning training...
Iteration #1 Error:51.023786% Target Error: 1.000000%
Iteration #2 Error:49.659291% Target Error: 1.000000%
Iteration #3 Error:43.140471% Target Error: 1.000000%
Iteration #4 Error:29.820891% Target Error: 1.000000%
Iteration #5 Error:29.457086% Target Error: 1.000000%
Iteration #6 Error:19.421585% Target Error: 1.000000%
Iteration #7 Error:2.160925% Target Error: 1.000000%
Iteration #8 Error:0.432104% Target Error: 1.000000%
Input=0.0000,0.0000, Actual=0.0091, Ideal=0.0000
Input=1.0000,0.0000, Actual=0.9793, Ideal=1.0000
Input=0.0000,1.0000, Actual=0.9472, Ideal=1.0000
Input=1.0000,1.0000, Actual=0.0731, Ideal=0.0000
Machine Learning Type: feedforward
Machine Learning Architecture: ?:B->SIGMOID->4:B->SIGMOID->?
Training Method: lma
Training Args:

```

The above shows a program learning the XOR operator. Notice how the MSE error drops each iteration? Finally, but iteration eight the error is below one percent, and training stops.

2.3 Other Error Calculation Methods

Though MSE is the most common method of calculating global error, it is not the only method. In this section we will look at two other global error calculation methods.

2.3.1 Sum of Squares Error

The sum of squares method (ESS) uses a similar formula as the MSE error method. However, ESS does not divide by the number of elements. As a result, the ESS is not a percent. It is simply a number that is larger depending on how severe the error is. Equation 2.3 shows the MSE error formula.

$$\text{ESS} = \frac{1}{2} \sum_p E^2 \quad (2.3)$$

As you can see above, the sum is not divided by the number of elements. Rather, the sum is simply divided in half. This results in an error that is not a percent, but rather a total of the errors. Squaring the errors eliminates the effect of positive and negative errors.

Some training methods require that you use ESS. The Levenberg Marquardt Algorithm (LMA) requires that the error calculation method be ESS. LMA will be covered in Chapter 8.

2.3.2 Root Mean Square Error

The Root Mean Square (RMS) error method is very similar to the MSE method previously discussed. The primary difference is that the square root is taken of the sum. You can see the RMS formula in Equation 2.4.

$$\text{RMS} = \sqrt{\frac{1}{n} \sum_{i=1}^n E^2} \quad (2.4)$$

Root mean square error will always be higher than MSE. The following output shows the calculated error for all three error calculation methods. All three cases used the same actual and ideal values.

Trying from -1.00 to 1.00

```
Actual:[-0.36,0.07,0.55,0.05,-0.37,0.34,-0.72,-0.10,-0.41,-0.32]
Ideal:[-0.37,0.06,0.51,0.06,-0.36,0.35,-0.67,-0.09,-0.43,-0.33]
Error (ESS): 0.00312453
Error (MSE): 0.062491%
Error (RMS): 2.499810%
```

RMS is not used very often for neural network error calculation. RMS was originally created in the field of electrical engineering. I have not used RMS a great deal. Many research papers involving RMS show it being used for waveform analysis.

2.4 Chapter Summary

Neural networks start with random values for weights. These networks are then trained until a set of weights is found that provides output from the neural network that closely matches the ideal values from the training data. For training to progress a means is needed to evaluate the degree to which the actual output from the neural network matches the ideal output expected of the neural network.

This chapter began by introducing the concept of local and global error. Local error is the error used to measure the difference between the actual and ideal output of an individual output neuron. This error is calculated using an error function. Error functions are only used to calculate local error.

Global error is the total error of the neural network across all output neurons and training set elements. Three different techniques were presented in this chapter for the calculation of global error. Mean Square Error (MSE) is the most commonly used. Sum of Squared Errors (ESS) is used by some training methods to calculate error. Root Mean Square (RMS) can be used to calculate the error for certain applications. RMS was created for the field of electrical engineering for waveform analysis.

The next chapter will introduce a mathematical concept known as derivatives. Derivatives come from Calculus and will be used to analyze the error functions and adjust the weights to minimize this error. In this book we will learn about several propagation training techniques. All propagation training techniques use derivatives to calculate update values for the weights of the neural network.

Chapter 3

Derivatives

- Slope of a Line
- What is a Derivative
- Partial Derivatives
- Chain Rule

In this chapter we will look at the mathematical concept of a derivative. Derivatives are used in many aspects of neural network training. A basic understanding of derivatives will be useful for proper understanding of the next few chapters that will focus on training a neural network.

The concept of a derivative is central to an understanding of Calculus. The topic of derivatives is very large and could easily consume several chapters. I am only going to explain those aspects of differentiation that are important to the understanding of neural network training. If you are already familiar with differentiation you can safely skim, or skip, this chapter.

3.1 Calculating the Slope of a Line

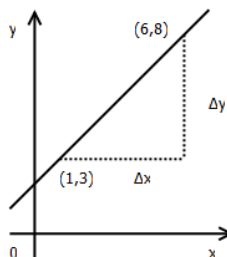
The slope of a line is a numerical quality of a line that tells you the direction and steepness of a line. In this section we will see how to calculate the slope of a straight line. In the next section we will see how to calculate the slope of a curved line at a single point.

The slope of a line is defined as the “rise” over the “run”, or the change in y over the change in x . The slope of a line can be written in the form of Equation 3.1.

$$m = \frac{\Delta y}{\Delta x} = \frac{y_2 - y_1}{x_2 - x_1} \quad (3.1)$$

This can be visualized graphically in Figure 3.1.

Figure 3.1: Slope of a Line

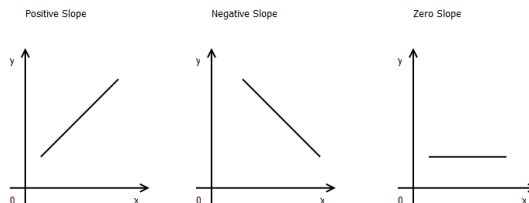


We could easily calculate the slope of the above line using Equation 3.1. Filling in the numbers for the two points we have on the line produces the following.

$$(8-3) / (6-1) = 1$$

The slope of this line is one. This is a positive slope. When a line has a positive slope it goes up left to right. When a line has a negative slope it goes down left to right. When a line is horizontal the slope is 0 and when the line is vertical the slope is undefined. Figure 3.2 several slopes for comparison.

Figure 3.2: Several Slopes



A straight line, such as seen above, can be written in slope-intercept form. Equation 3.2 shows the slope intercept form of an equation.

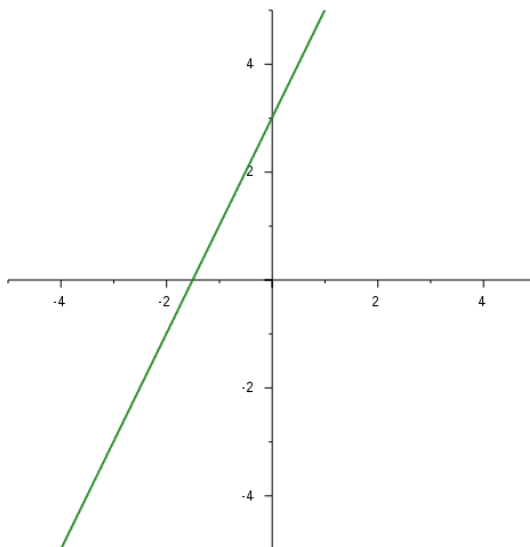
$$y = mx + b \quad (3.2)$$

Where \mathbf{m} is the slope of the line and \mathbf{b} is the y-intercept, which is the y-coordinate of the point where the line crosses the y axis. To see this in action, consider the chart of the following equation.

$$f(x) = 2x + 3$$

This equation can be seen graphically in Figure 3.3.

Figure 3.3: The Graph of $2x+3$



As you can see from the above diagram, the line intercepts the y-axis at 3.

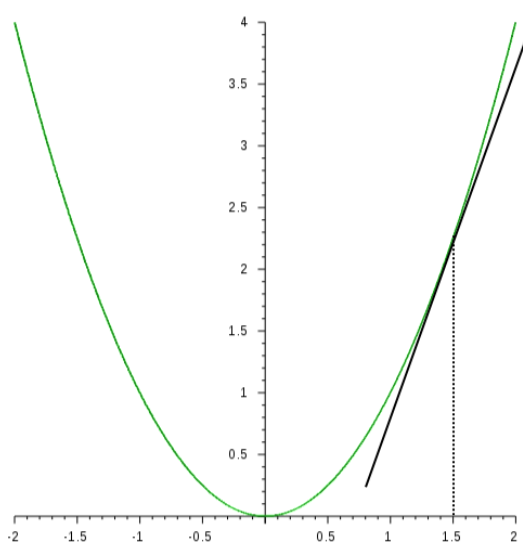
3.2 What is a Derivative?

In the last section we saw that we can easily calculate the slope of any straight line. It is very rare that we will work with a simple straight line. Usually, we will be faced with the sort of curves that we saw in the last chapter while examining the activation functions of neural networks. A derivative allows us to take the derivative of a line at one point. Consider the following simple equation.

$$y = x^2 \quad (3.3)$$

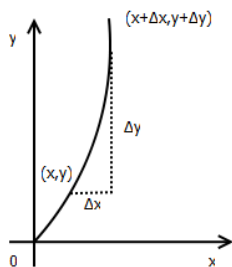
You can see equation 3.3 graphed in Figure 3.4.

Figure 3.4: Graph of x Squared



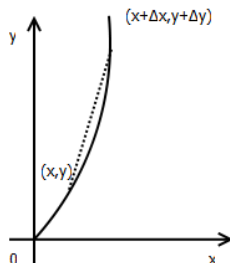
In the above chart we would like to obtain the derivative at 1.5. The chart of x squared is given by the u-shaped line. The slope at 1.5 is given by the straight line that just barely touches the “u-shaped” line at 1.5. This straight line is called a tangent line. If we take the derivative of Equation 3.2, we are left with an equation that will provide us with the slope of Equation 3.2 at any point x . It is relatively easy to derive such an equation. To see how this would be done consider Figure 3.5.

Figure 3.5: Calculate Slope at X



Here we are given a point at (x,y) . This is the point that we would like to find the derivative for. However, we need two points to calculate a slope. Therefore we create a second point that is equal to x and y plus **delta-x** and **delta-y**. You can see this imaginary line in Figure 3.6.

Figure 3.6: Slope of a Secant Line



The imaginary line above is called a secant line. The slope of this secant line is close to the slope at (x,y) , but it is not the exact number. As **delta-x** and **delta-y** become closer to zero, the slope of the secant line becomes closer to the instantaneous slope at (x,y) . We can use this fact to write an equation that is the derivative of Equation 3.2.

Before we look specially at Equation 3.2 we will look at the general case of how to find a derivative for any function $f(x)$. This formula uses a constant **h** that defines a second point by adding **h** to **x**. The smaller that **h** becomes, the more accurate of a value we are given for the slope of the line at **x**. Equation 3.4 shows the slope of the secant line between **x** and **h**.

$$m = \frac{f(x+h) - f(x)}{(x+h) - x} \quad (3.4)$$

The derivative is equal to the above equation as **h** approaches zero. This is shown in Equation 3.5.

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{(x+h) - x} \quad (3.5)$$

You should notice two things about the above formula. Notice the apostrophe above **f**? This designates the function as a derivative, and is pronounced "f-prime". The second is the word **lim**. This designates a limit. The arrow at the bottom specifies "as h approaches zero".

When taking a limit, sometimes the limit is undefined at the value it approaches. Therefore the limit is the value either at, or close to, the value the limit is approaching. In many cases the limit can be determined simply by solving the formula with the approached value substituted for \mathbf{x} .

The above formula can be simplified by removing redundant \mathbf{x} terms in the denominator. This results in Equation 3.6 which is the definition of a derivative as a slope.

$$f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \quad (3.6)$$

Now that we have a formula for a derivative, we will see a simple application in the next section. We will now use this equation to finally determine the derivative for Equation 3.2. Equation 3.2 is simply the function of \mathbf{x} squared.

Using the formula from the last section it is easy to take the derivative of a formula such as \mathbf{x} squared. Equation 3.6 shows Equation 3.5 modified to use \mathbf{x} squared in place of $\mathbf{f}(\mathbf{x})$.

$$f'(x) = \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} \quad (3.7)$$

Equation 3.7 can be expanded as follows. This is using a simple algebraic rule that allows us to expand the term $(\mathbf{x}+\mathbf{h})$ squared. This results in Equation 3.8.

$$f'(x) = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} \quad (3.8)$$

This allows us to remove redundant \mathbf{x} terms in the numerator.

$$f'(x) = \lim_{h \rightarrow 0} \frac{2xh + h^2}{h} \quad (3.9)$$

We can also cancel out the \mathbf{h} terms in the numerator with the same term in the denominator. This leaves us with Equation 3.10.

$$f'(x) = \lim_{h \rightarrow 0} 2x + h \quad (3.10)$$

We can now evaluate the limit at zero. This produces the final general formula for the derivative shown in Equation 3.11.

$$f'(x) = 2x \quad (3.11)$$

The above equation would be found in the front cover of many Calculus text books. Simple derivative formulas like this are useful for converting common equations into derivative form. Calculus text books usually have these derivative formulas listed in a table. Using this table, more complex derivatives can be obtained. I will not review how to obtain the derivative any arbitrary function. Generally, when I want to take the derivative of an arbitrary function I use a program called R. R can be obtained from the following URL.

<http://www.r-project.org/>

The following R command could be used to find the derivative of x squared.

```
D( expression(x^2), "x" )
```

For a brief tutorial on R, visit the following URL.

http://www.heatonresearch.com/wiki/Brief_R_Tutorial

3.3 Using Partial Derivatives

So far we have only seen "total derivatives". A partial derivative of a function of several variables is the derivative of the function with respect to one of those variables. All other variables will be held constant. This differs from total derivative, in which all variables are allowed to vary.

The partial derivative of a function f with respect to the variable z is variously denoted by the following forms.

$$f'_z, f_z, \partial_z f, \text{ or } \frac{\partial f}{\partial z} \quad (3.12)$$

The form that will be used in this book is shown here.

$$\frac{\partial f}{\partial z} = \dots \quad (3.13)$$

To see an example of partial derivatives consider the function f that has more than one variable.

$$z = f(x, y) = x^2 + xy + y^2. \quad (3.14)$$

The derivative of \mathbf{z} , with respect to \mathbf{x} is given as follows. The variable \mathbf{y} is treated as a constant.

$$\frac{\partial z}{\partial x} = 2x + y \quad (3.15)$$

Partial derivatives are an important concept for neural networks. We will typically take the partial derivative of the error of a neural network with respect to each of the weights. This will be covered in greater detail in the next chapter.

3.4 Using the Chain Rule

There are many different rules in Calculus to allow you to take derivatives manually. We just saw an example of the power rule. This rule states that given the following equation.

$$f(x) = x^n \quad (3.16)$$

The derivative of $f(x)$ will be as follows.

$$f'(x) = nx^{n-1} \quad (3.17)$$

This allows you to quickly take the derivative of any power. There are many other derivative rules. They are very useful to know. However, if you do not wish to learn manual differentiation you can generally get by using a program such as R.

However, there is one more rule that is very useful to know. This rule is called the chain rule. The chain rule deals with composite functions. A composite function is nothing more than when one function takes the results of a second function as input. This may sound complex, but programmers make use of composite functions all the time. In Java, the following is an example of a composite function call.

```
System.out.println( Math.pow(3,2) );
```

This is a composite function because we take the result of the function **pow** and feed it to **println**.

Mathematically we write this as follows. Consider if we had functions **f** and **g**. If we wished to pass the value of 5 to **f**, and then pass the result of **f** onto **g**, we would use the following expression.

$$(f \circ g)(5) \tag{3.18}$$

The chain rule of calculus gives us a relatively easy way to calculate the composite of two functions. The chain rule is given in Equation 3.12.

$$(f \circ g)'(t) = f'(g(t))g'(t) \tag{3.19}$$

The chain rule will be very valuable to calculate the derivative across an entire neural network. A neural network is very much a composite function. In a typical three layer neural network the output of the input layer flows to the hidden layer and finally to the output layer.

3.5 Chapter Summary

In this chapter we took a look at derivatives. Derivatives are a core concept in Calculus. A derivative is defined as the slope of a curved line for one individual value of \mathbf{x} . The derivative can also be thought of as the instantaneous rate of change at the point \mathbf{x} . Derivatives can be calculated manually, or using a software package such as R.

Derivatives are very important for neural network training. The derivatives of activation functions are used to calculate the error gradient with respect to individual weights. Various training algorithms make use of these gradients to determine how to best update the neural network weights.

In the next chapter we will look at backpropagation. Backpropagation is a training algorithm that adjusts the weights of neural networks to produce more desirable output from the neural network. Backpropagation works by calculating the partial derivative of the error function with respect to each of the weights.

Chapter 4

Backpropagation

- Understanding Gradients
- Calculating Gradients
- Understanding Backpropagation
- Momentum and Learning Rate

So far we have only looked at how to calculate the output from a neural network. The output from the neural network is a result of applying the input to the neural network across the weights of several layers. In this chapter we will see how these weights are adjusted to produce outputs that are closer to the desired output.

This process is called training. Training is an iterative process. To make use of training, you perform multiple training iterations. It is hoped that the training iterations will lower the global error of the neural network. Global and local error were discussed in Chapter 2.

4.1 Understanding Gradients

The first step is to calculate the gradients of the neural network. The gradients are used to calculate the slope, or gradient, of the error function for a particular weight. This allows the training method to know to either increase or decrease the weight. There are number of different training methods that make use of gradients. These training methods are called propagation training. This book will discuss the following propagation training methods.

- Backpropagation
- Resilient Propagation
- Quick Propagation

This chapter will focus on using the gradients to train the neural network using back propagation. The next few chapters will focus on the other propagation methods.

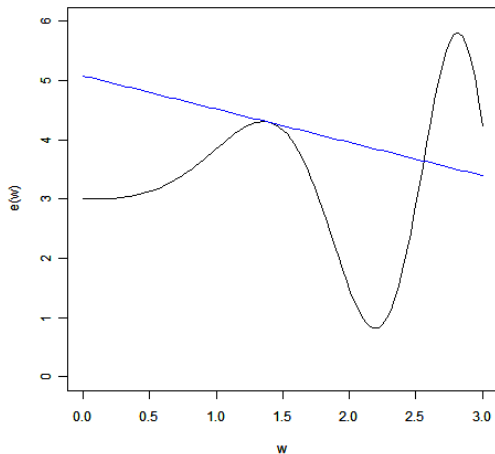
4.1.1 What is a Gradient

We will begin by looking at what a gradient is. Basically training is a search. You are searching for the set of weights that will cause the neural network to have the lowest global error for a training set. If we had an infinite amount of computation resources we would simply try every possible combination of weights and see which provided the absolute best global error.

Because we do not have unlimited computing resources we will have to take some sort of short cut. A shortcut is essentially all that neural network training methods really are. Each training method is a clever way of finding an optimal set of weights without doing an impossible exhaustive search.

Consider a chart that shows the global error of a neural network for each possible weight. This graph might look something like Figure 4.1.

Figure 4.1: Gradient



Looking at this chart you can easily see the optimal weight. The optimal weight is the location where the line has the lowest y value. The problem is, we do not get to see the entire graph. This would be the exhaustive search previously mentioned. We only see the error for the current value of the weight. However, we can determine the slope of the error curve at a particular weight. In the above chart we see the slope of the error curve at 1.5. The slope is given by the straight line that just barely touches the error curve at 1.5. This is this slope is the gradient. In this case the slope is -0.5622 .

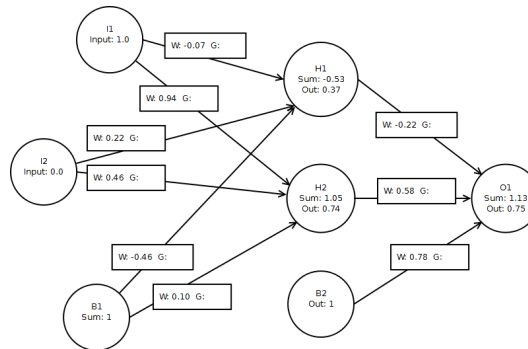
The gradient is the instantaneous slope of the error function at the specified weight. This is the same definition for the derivative that we learned in Chapter 3. The gradient is given by the derivative of the error curve at that point. This line tells us something about how steep the error function is at the given weight.

Used with a training technique, this can provide some insight into how the weight should be adjusted for a lower error. Now that we have seen what a gradient is, we will see how to actually calculate a gradient in the next section.

4.1.2 Calculating Gradients

We will now look at how to actually calculate the gradient. We will calculate an individual gradient for each weight. A weight is a connection between two neurons. I will show you the equations, as well as how to apply them to an actual neural network with real numbers. The neural network that we will use is shown in Figure 4.2.

Figure 4.2: An XOR Network



The neural network above is a typical three layer feedforward network like we have seen before. The circles indicate neurons. The lines connecting the circles are the weights. The rectangles in the middle of the connections give the weight for each connection.

The problem that we must now face is how to calculate the partial derivative for the output of each neuron. This can be accomplished using a method based on the chain rule of Calculus. The chain rule was discussed in Chapter 3. We will begin with one training set element. For Figure 4.2 we are providing an input of $[1,0]$ and expecting an output of $[1]$. You can see the input being applied on the above figure as the first input neuron has an input value of 1.0 and the second input neuron has an input value of 0.0.

This input feeds through the network and eventually produces an output. The exact process by which the output and sums are calculated was covered in Chapter 1. Backpropagation has both a forward and backward pass. When forward pass occurred when the output of the neural network was calculated. We will calculate the gradients for only this item in the training set. Other items in the training set will have different gradients. How the gradients for each individual training set element are combined will be discussed later in this chapter when we discuss "Batch and Online Training".

We are now ready to calculate the gradients. There are several steps involved in calculating the gradients for each weight. These steps are summarized here.

- Calculate the error, based on the ideal of the training set
- Calculate the node delta for the output neurons
- Calculate the node delta for the interior neurons
- Calculate individual gradients

These steps will be covered in the following sections.

4.1.3 Calculating the Node Deltas

The first step is to calculate a constant value for every node, or neuron, in the neural network. We will start with the output nodes and work our way backwards through the neural network. This is where the term backpropagation comes from. We initially calculate the errors for the output neurons and propagate these errors backwards through the neural network.

The value that we will calculate for each node is called the node delta. The term layer delta is also sometimes used to describe this value as well. Layer delta describes the fact that these deltas are calculated one layer at a time. The method for calculating the node deltas differs depending on if you are calculating for an output or interior node. The output neurons are obviously, all output nodes. The hidden and input neurons are the interior nodes. The equation to calculate the node delta is provided in Equation 4.1.

$$\delta_i = \begin{cases} -E f'_i & , \text{ output nodes} \\ f'_i \sum_k w_{ki} \delta_k & , \text{ interior nodes} \end{cases} \quad (4.1)$$

We will calculate the node delta for all hidden and non-bias neurons. There is no need to calculate the node delta for the input and bias neurons. The node delta can easily be calculated for input and bias neurons using the above equation. However, these values are not needed for the gradient calculation. As you will soon see, gradient calculation for a weight only looks at the neuron that the weight is connected to. Bias and input neurons are only the beginning point for a connection. They are never the end point.

We will begin by using the formula for the output neurons. You will notice that the formula uses a value **E**. This is the error for this output neuron. You can see how to calculate **E** from Equation 4.2.

$$E = (a - i) \quad (4.2)$$

You may recall a similar equation as Equation 4.2 from Chapter 3. This is the error function. Here we subtract the ideal from the actual. For the example neural network provided earlier this becomes the following.

$$E = 0.75 - 1.00 = -0.25$$

Now that we have **E**, we can calculate the node delta for the first (and only) output node. Filling in Equation 4.1 we get the following.

$$-(-0.25) * dA(1.1254) = 0.185 * 0.25 = 0.05$$

The value of 0.05 will be used for the node delta of output neuron. In the above equation, **dA** represents the derivative of the activation function. For this example, we are using a sigmoid activation function. The sigmoid activation function is shown in Equation 4.3.

$$s(x) = \frac{1}{1 + e^{-x}} \quad (4.3)$$

The derivative of the sigmoid function is shown in Equation 4.4.

$$f'(x) = s(x) * (1.0 - s(x)) \quad (4.4)$$

Now that the node delta has been calculated for the output neuron, we should calculate it for the interior neurons as well. The equation to calculate the node delta for interior neurons was provided in Equation 4.1. Applying this for the first hidden neuron we have the following.

$$dA(\text{sum of H1}) * (\text{O1 Node Delta} * \text{Weight of H1} \rightarrow \text{O1})$$

You will notice that Equation 4.1 called for the summing up a number of items based on the number of inbound connections to output neuron one. Because there is only one inbound connection to output neuron one there is only one value to sum. This value is the product of the output neuron one node delta and the weight between hidden neuron one and output neuron one.

Filling in actual values for the above expression, we are left with the following.

$$dA(-0.53) * (0.05 * -0.22) = -0.0025$$

The value -0.0025 is the node delta for the first hidden neuron. Calculating the second hidden neuron follows exactly the same form as above. The second neuron would be computed as follows.

$$dA(\text{sum of H2}) * (\text{O1 Node Delta} * \text{Weight of H2} \rightarrow \text{O1})$$

Plugging in actual numbers, we have the following.

$$dA(1.05) * (0.05 * 0.58) = 0.0055$$

The value of 0.0055 is the node delta for the second hidden neuron.

As previously explained, there is no reason to calculate the node delta for either the bias neurons or the input neurons. We now have every layer delta needed to calculate a gradient for each weight in the neural network. Calculation of the individual gradients will be discussed in the next section.

4.1.4 Calculating the Individual Gradients

We can now calculate the individual gradients. Unlike the layer deltas there is only one equation used to calculate the actual gradient. A gradient is calculated using Equation 4.5.

$$\frac{\partial E}{\partial w_{(ik)}} = \delta_k \cdot o_i \quad (4.5)$$

The above equation calculates the partial derivative of the error (**E**) with respect to each individual weight. The partial derivatives are the gradients. Partial derivatives were discussed in Chapter 3. To determine an individual gradient multiply the node delta for the target neuron by the weight from the source neuron. In the above equation **k** represents the target neuron and **i** represents the source neuron.

To calculate the gradient for the weight from **H1** to **O1** the following values would be used.

$$\text{output}(h1) * \text{nodeDelta}(o1) \\ (0.37 * 0.05) = 0.01677$$

It is important to note that in the above equation we are multiplying by the output of hidden 1, not the sum. When dealing directly with a derivative you should supply the sum. Otherwise, you would be indirectly applying the activation function twice. In the above equation we are not dealing directly with the derivative, so we use the regular node output. The node output has already had the activation function applied.

Once the gradients are calculated the individual positions of the weights no longer matter. We can simply think of the weights and gradients as single dimensional arrays. The individual training methods that we will look at will treat all weights and gradients equal. It does not matter if a weight is from an input neuron or an output neuron. It is only important the correct weight be used with the correct gradient. The ordering of this weight and gradient array is arbitrary. However, Encog uses the following order for the above neural network.

Weight/Gradient 0:	Hidden 1	->	Output 1
Weight/Gradient 1:	Hidden 2	->	Output 1
Weight/Gradient 2:	Bias 2	->	Output 1
Weight/Gradient 3:	Input 1	->	Hidden 1
Weight/Gradient 4:	Input 2	->	Hidden 1
Weight/Gradient 5:	Bias 1	->	Hidden 1
Weight/Gradient 6:	Input 1	->	Hidden 2
Weight/Gradient 7:	Input 2	->	Hidden 2
Weight/Gradient 8:	Bias 1	->	Hidden 2

The various learning algorithms will make use of the weight and gradient arrays. It is also important to note that these are two separate arrays. There is a weight array, as well as a gradient array. Both arrays will be exactly the same length. Training a neural network is nothing more than adjusting the weights to provide desirable output. In this chapter we will see how backpropagation uses the gradient array to modify the weight array.

4.2 Applying Back Propagation

Backpropagation is a simple training method that uses the calculated gradients of a neural network to adjust the weights of the neural network. As these weights are adjusted the neural network should produce more desirable output. The global error of the neural network should fall as it is trained. Before we can examine the backpropagation weight update process we must look at two different ways that the gradients can be calculated.

4.2.1 Batch and Online Training

We have already seen how to calculate the gradients for an individual training set element. Earlier in this chapter we saw how we could calculate the gradients for a case where the neural network was given an input of $[1,0]$ and an output of $[1]$ was expected. This works fine for a single training set element. However, most training sets have many elements. There are two different ways to handle multiple training set elements. These two approaches are called online and batch training.

Online training implies that you modify the weights after every training set element. Using the gradients that you obtained for the first training set element you calculate and apply a change to the weights. Training progresses to the next training set element and also calculates an update to the neural network. This training progresses until every training set element has been used. At this point one iteration, or epoch, of training has completed.

Batch training also makes use of every training set element. However, the weights are not updated for every training set element. Rather the gradients for each training set element are summed. Once every training set element has been used, the neural network weights can be updated. At this point the iteration is considered complete.

Sometimes a batch size will be set. For example, you might have a training set size of 10,000 elements. You might choose to update the weights of the neural network every 1,000 elements. This would cause the neural network weights to be updated 10 times during the training iteration.

Online training was the original method used for backpropagation. However, online training is inefficient, as the neural network must be constantly updated. Additionally, online training is very difficult to implement in a multi-threaded manor that will take advantage of multi-core processors. Because of these reasons batch training is generally preferable to online training.

4.2.2 Backpropagation Weight Update

We are now ready to update the weights. As previously mentioned, we will treat the weights and gradients as a single dimensional array. Give these two arrays we are ready to calculate the weight update for an iteration of backpropagation training. The formula to update the weights for backpropagation is shown in Equation 4.6.

$$\Delta w_{(t)} = \epsilon \frac{\partial E}{\partial w_{(t)}} + \alpha \Delta w_{(t-1)} \quad (4.6)$$

The above equation calculates the change in weight for each element in the weight array. You will also notice that the above equation calls for the weight change from the previous iteration. These values must be kept in another array.

The above equation calculates the weight delta to be the product of the gradient and the learning rate (represented by epsilon). Additionally the product of the previous weight change and the moment value (represented by alpha) is added. The learning rate and momentum are two parameters that must be provided to the backpropagation algorithm. Choosing values for learning rate and momentum is very important to the performance of the training. Unfortunately, the process for determining learning rate and momentum is mostly trial and error.

The learning rate scales the gradient and can slow down or speed up learning. A learning rate below zero will slow down learning. For example, a learning rate of 0.5 would decrease every gradient by 50%. A learning rate above 1.0 would speed up training. In reality the learning rate is almost always below zero.

Choosing too high of a learning rate will cause your neural network to fail to converge. A neural network that is failing to converge will generally have a high global error that simply bounces around, rather than converging to a low value. Choosing too low of a learning rate will cause the neural network to take a great deal of time to converge.

Like the learning rate, the momentum is also a scaling factor. Momentum determines what percent of the previous iteration's weight change should be applied to this iteration. Momentum is optional. If you do not want to use momentum, just specify a value of zero.

Momentum is a technique that was added to backpropagation to help the training find its way out of local minima. Local minima are low points on the error graph that are not the true global minimum. Backpropagation has a tendency to find its way into a local minimum and not find its way back out again. This will cause the training to converge to a higher undesirable error. Momentum gives the neural network some force in its current direction and may allow it to force through a local minimum.

We are now ready to plug values into Equation 4.6 and calculate a weight delta. We will calculate a weight change for the first iteration using the neural network we previously used in this chapter. So far we have only calculated gradients for one training set element. There are still four other training set elements to calculate for. We will sum all four gradients

together to apply batch training. The batch gradient is calculated by summing the individual gradients, which are listed here.

```
0.01677795762852397
-0.05554301180824532
0.021940533165555356
-0.05861906411780882
```

Earlier in the chapter we calculated the first gradient, listed above. If you would like to see the calculation for the others, this information can be found at the following URL.

http://www.heatonresearch.com/wiki/Back_Propagation

Summing all of these gradients produces the following batch update gradient.

```
-0.07544358513197481
```

For this neural network we will use a learning rate of 0.7 and a momentum of 0.3. These are just arbitrary values. However, they do work well for training an XOR neural network. Plugging these values into Equation 4.6 results in the following.

```
delta = (0.7 * -0.0754) + (0.3 * 0.0) = -0.052810509592382364
```

This is the first training iteration, so the previous delta value is 0.0. The momentum has no effect on the first iteration. This delta value will be added to the weight to alter the neural network for the first training iteration. All of the other weights in this neural network will be updated in the same way, according to their calculated gradient.

This first training iteration will lower the neural network's global error slightly. Additional training iterations will further lower the error. The following program output shows the convergence of this neural network.

```
Epoch #1 Error:0.3100155809627523
Epoch #2 Error:0.2909988918032235
Epoch #3 Error:0.2712902750837602
Epoch #4 Error:0.2583119003843881
Epoch #5 Error:0.2523050561276289
Epoch #6 Error:0.2502986971902545
Epoch #7 Error:0.2498182295192154
Epoch #8 Error:0.24974245650541688
Epoch #9 Error:0.24973458893806627
Epoch #10 Error:0.24972923906975902
...
Epoch #578 Error:0.010002702374503777
```

```
Epoch #579 Error:0.009947830890527089
Neural Network Results:
1.0,0.0, actual=0.9040102333814147,ideal=1.0
0.0,0.0, actual=0.09892634022671229,ideal=0.0
0.0,1.0, actual=0.904020682439766,ideal=1.0
1.0,1.0, actual=0.10659032105865764,ideal=0.0
```

Each iteration, or epoch, decreases the error. Once the error drops below one percent, the training stops. You can also see the output from the neural network for the XOR data. The answers are not exactly correct, however, it is very clear that that the two training cases that should be 1.0 are much closer to 1.0 than the others.

4.3 Chapter Summary

In this chapter you were introduced to backpropagation. Backpropagation is one of the oldest, and most commonly used, training algorithms available for neural networks. Backpropagation works by calculating a gradient value for each weight in the network. Many other training methods also make use of these gradient values.

There is one gradient for each weight in the neural network. Calculation of the gradients is a step by step process. The first step in calculating the gradients is to calculate the error for each of the outputs of the neural network. This error is for one training set element. The gradients for all training set elements may be batched together later in the process.

Once the error for the output layer has been calculated values can be calculated for each of the output neurons. These values are called the node deltas for each of the output neurons. We must calculate the node deltas for the output layer first. We calculate the node deltas for each layer of the neural network working out way backwards to the input layer. This is why this technique is called backpropagation.

Once the node deltas have been calculated it is very easy to calculate the gradients. At the end you will have all of the gradients for one training set element. If you are using online training you will now use these gradients to apply a change to the weights of the neural network. If you are using batch training, you will sum the gradients from each of the training set elements into a single set of gradients for the entire training set.

Backpropagation must be provided a learning rate and momentum. Both of these are configuration items that will have an important effect on the training speed of your neural network. Learning rate specifies how fast the weights should be updated. Too high of a

learning rate will cause a network to become unstable. Too low of a learning rate will cause the neural network to take too long to train. Momentum allows the neural network to escape local minima. Local minima are low points in the error graph that are not the true global minimum.

Choosing values for the learning rate and momentum can be tricky. Often it is just a matter of trial and error. The Resilient Propagation training method (RPROP) requires no parameters to be set. Further, RPROP often trains much faster than backpropagation. RPROP will be covered in the next chapter.

Chapter 5

RPROP

- Understanding Error Calculation
- The Error Function
- Error Calculation Methods
- How is the Error Used

In the last chapter we looked at backpropagation. Backpropagation is one of the oldest, and most popular, methods for training a neural network. Unfortunately, backpropagation is also one of the slowest methods for training a neural network. In this chapter we will take a look at a faster training technique called resilient propagation, or RPROP.

RPROP works very similarly to backpropagation. Both backpropagation and RPROP must first calculate the gradients for the weights of the neural network. The way in which the gradients are used is how backpropagation and RPROP differ.

In a simple XOR example typically takes backpropagation over 500 iterations to converge to a solution with an error rate of below one percent. It will usually take RPROP around 30 to 100 iterations to accomplish the same thing. This great performance increase is one reason that RPROP is a very popular training algorithm.

Another feature in the popularity of RPROP is that there are no necessary training parameters to the RPROP algorithm. When you make use of backpropagation you must specify the learning rate and momentum. These two parameters can have a huge impact on

the effectiveness of your training. RPROP does include a few training parameters, however, they can almost always be left at their default settings.

There are several variants of the RPROP protocol. The following table lists some of the variants.

- RPROP+
- RPROP-
- iRPROP+
- iRPROP-

This book will focus on classic RPROP, as described in a 1994 paper by Martin Reidmiller entitled "Rprop - Description and Implementation Details". The other four variants described above are relatively minor adaptations to classic RPROP. For more information about all variants of RPROP visit the following URL.

<http://www.heatonresearch.com/wiki/RPROP>

In the next sections we will see how the RPROP algorithm is implemented.

5.1 RPROP Arguments

As previously mentioned, one advantage of RPROP over backpropagation is that no training arguments need to be provided in order for RPROP to be used. That is not to say that there are no configuration settings for RPROP. The configuration settings for RPROP do not usually need to be changed from their defaults. However, if you really want to, there are several configuration settings that you can set for RPROP training. These configuration settings are listed here.

- Initial Update Values
- Maximum Step

As you will see in the next section RPROP keeps an array of update values for the weights. This determines how large of a change will be made to each weight. This is something like the

learning rate in backpropagation, only much better. There is an update value for every weight in the neural network. This allows the update values to be fine tuned to each individual weight, as training progresses. Some backpropagation algorithms will vary the learning rate and momentum as learning progresses. The RPROP approach is much better, because it does not have a single learning rate for the entire neural network, as backpropagation does.

These update values must start out somewhere. The "initial update values" argument defines this. By default this argument is set to a value of 0.1. Generally this default should never be changed. One possible exception to this is an already trained neural network. If the neural network is already trained then some of the initial update values are going to be too strong for the neural network. The neural network will regress for many iterations before it is able to improve. An already trained neural network may benefit from a much smaller initial update.

Another approach for an already trained neural network is to have saved the update values once training stops. This will allow you to resume training without the initial spike in training errors that you would normally see when resuming resilient propagation training. This method will only work if you are resuming resilient propagation. If you were previously training the neural network with a different training algorithm, then you will not have an array of update values to restore from.

As training progresses the gradients will be used to adjust the updates up and down. The "maximum step" argument defines the maximum upward step size that can be taken over the update values. The default value for the max step argument is 50. It is unlikely that you would need to change the value of this argument.

In addition to arguments defined for RPROP, there are also constants. These constants are simply values that are kept by RPROP during processing. These values are never changed.

- Delta Minimum (1e-6)
- Negative Eta (0.5)
- Positive Eta (1.2)
- Zero Tolerance (1e-16)

Positive and negative eta will be described in the next sections. The zero tolerance defines how close a number should be to zero before that number is equal to zero. In computer

programming it is typically bad practice to directly compare a floating point number to zero. This is because the number would have to be exactly equal to zero.

Delta minimum specifies the minimum value that one of the update values can go to. This is important, because if an update value were to go to zero it would never be able to increase beyond zero.

5.2 Data Structures

There are several data structures that must be kept in memory while RPROP training is performed. These structures are all arrays of floating point numbers. They are summarized here.

- Current Weight Change Values
- Last Weight Change Values
- Current Gradient Values
- Previous Gradient Values
- Current Update Values

The current update values are kept to hold the current update values for the training. If you wish to be able to resume training at some point this is the array that must be stored. There is one update value per weight. These update values cannot go below the minimum delta constant. Likewise, these update values cannot exceed the maximum step argument.

The last weight delta value must also be tracked. Backpropagation kept this value for momentum. RPROP uses this value in a different way than backpropagation. We will see how this array is used in the next section. The current weight change is kept long enough to change the weights and then be copied to the previous weight change.

The current and previous gradients are needed. RPROP is particularly interested when the sign changes from the current gradient to the previous gradient. This indicates that an action must be taken with regard to the update values. This is covered in the next section.

5.3 Understanding RPROP

In the last few sections the arguments, constants and data structures necessary for RPROP were covered. In this section, we will see exactly how to run through an iteration of RPROP. In the next section, we will apply real numbers to RPROP and see how training iterations progress for an XOR training. We will train exactly the same network that we did with backpropagation. This will give us a good idea of the difference in performance of backpropagation compared to RPROP.

While discussing backpropagation two weight update methods were mentioned: online and batch. RPROP does not support online training. All weight updates used with RPROP will be performed in batch mode. Because of this, each iteration of RPROP will receive gradients that are the sum of the individual gradients of each training set. This is consistent with using backpropagation in batch mode.

There are three distinct steps in an iteration of an RPROP iteration. They are covered in the next three sections.

5.3.1 Determine Sign Change of Gradient

At this point we should have the gradients. These gradients nearly exactly the same as the gradients calculated by the backpropagation algorithm. The only difference is that RPROP uses a gradient that is the inverse of the backpropagation gradient. This is easy enough to shift. Simply place a negative operator in front of every backpropagation gradient. Because the same process is used for both RPROP and backpropagation gradients, this process will not be repeated here. To learn how to calculate a gradient, refer to Chapter 4.

The first step is to compare the gradient of the current iteration to the gradient of the previous iteration. If there is no previous iteration then we can assume that the previous gradient was zero.

To determine if the gradient sign has changed we will use the *sgn* function. The *sgn* function is defined in Equation 5.1.

$$\text{sgn}(x) = \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases} \quad (5.1)$$

The sign function returns the sign of the number provided. If *x* is less than zero the result is

-1. If x is greater than zero then the result is 1. If x is equal to zero, then the result is zero. I usually implement the sign function to use a tolerance for zero, since it is nearly impossible for floating point operations to hit zero precisely on a computer.

To determine if the gradient has changed sign Equation 5.2 is used.

$$c = \frac{\partial E^{(t)}}{\partial w_{ij}} \cdot \frac{\partial E^{(t-1)}}{\partial w_{ij}} \quad (5.2)$$

Equation 5.2 will result in a constant \mathbf{c} . This value is evaluated to be either close to zero, negative or positive. A negative value for \mathbf{c} indicates that the sign has changed. A positive value indicates no change in sign for the gradient. A value near zero indicates that there was either a very small change in sign or nearly a change in sign. To see all three of these outcomes, consider the following situations.

-1 * 1 = -1 (negative, changed from negative to positive)
 1 * 1 = 1 (positive, no change in sign)
 1.0 * 0.000001 = 0.000001 (near zero, almost changed signs, but not quite)

Now that we have calculated the constant \mathbf{c} , which gives some indication of sign change, we can calculate the weight change. This is covered in the next section.

5.3.2 Calculate Weight Change

Now that we have change in sign of the gradient, we see what is done in each of the three cases mentioned in the previous section. These three cases are summarized in Equation 5.3.

$$\Delta w_{ij}^{(t)} = \begin{cases} -\Delta_{ij}^{(t)} & , \text{ if } c > 0 \\ +\Delta_{ij}^{(t)} & , \text{ if } c < 0 \\ 0 & , \text{ otherwise} \end{cases} \quad (5.3)$$

This equation calculates the actual weight change for each iteration. If the value of \mathbf{c} is positive, then the weight change will be equal to the negative of the weight update value. Similarly, if the value of \mathbf{c} is negative the weight change will be equal to the positive of the weight update value. Finally, if the value of \mathbf{c} is near zero, then there will be no weight change.

5.3.3 Modify Update Values

The weight update values, seen in the previous section, are used each iteration to update the weights of the neural network. There is a separate weight update value for every weight in the neural network. This works much better than a single learning rate for the entire neural network, as was done with backpropagation. These weight update values are modified during each training iteration, as seen in Equation 5.4.

$$\Delta_{ij}^{(t)} = \begin{cases} \eta^+ \cdot \Delta_{ij}^{(t-1)} & , \text{ if } c > 0 \\ \eta^- \cdot \Delta_{ij}^{(t-1)} & , \text{ if } c < 0 \\ \Delta_{ij}^{(t-1)} & , \text{ otherwise} \end{cases} \quad (5.4)$$

The weight update values are modified in a way very similar to how the weights themselves are modified. Just like the weights, the weight update values are changed based on the value c , previously calculated.

If the value of c is positive, then the weight update value will be multiplied by the value of positive eta. Similarly, if the value of c is negative the weight update value will be multiplied by negative eta. Finally, if the value of c is near zero, then there will be no change to the weight update value.

5.4 RPROP Update Examples

We will now look at a few iterations of a RPROP train of a neural network. The critical stats on this neural network are listed here.

```
Layers: 3(input , hidden , output)
Input Neurons: 2
Hidden Neurons: 2
Output Neurons: 1
Activation Function: Sigmoid
Bias Neurons: Yes
Total weights: 9
```

The initial random weights for this neural network are given here.

```
Weight 0: H1->O1: -0.22791948943117624
Weight 1: H2->O1: 0.581714099641357
Weight 2: B2->O1: 0.7792991203673414
```

```

Weight 3: I1->H1: -0.06782947598673161
Weight 4: I2->H1: 0.22341077197888182
Weight 5: B1->H1: -0.4635107399577998
Weight 6: I1->H2: 0.9487814395569221
Weight 7: I2->H2: 0.46158711646254
Weight 8: B1->H2: 0.09750161997450091

```

The neural network will be trained for the XOR operator, just as was performed in Chapter 4. However, in this chapter RPROP will be used in place of backpropagation.

5.4.1 Training Iteration #1

For the first training iteration the weight update values are all set to 0.1. This is their default starting point. We begin by calculating the gradients of each of the weights. This is the same gradient calculation as was performed for backpropagation, so the calculation will not be covered here. The gradients are provided here.

```

Gradient 0: H1->O1: -0.07544358513197481
Gradient 1: H2->O1: -0.12346935587390481
Gradient 2: B2->O1: -0.18705713395934637
Gradient 3: I1->H1: 0.005292326734004241
Gradient 4: I2->H1: 0.0049016107791925246
Gradient 5: B1->H1: 0.010264148244428655
Gradient 6: I1->H2: -0.0068740307089347
Gradient 7: I2->H2: -0.005236293038814788
Gradient 8: B1->H2: -0.02103299467864286

```

Now that we have the gradients, we must calculate the weight change with the previous gradients. There are no previous gradients, so their values are all zero. The calculation of \mathbf{c} is shown here.

```

c 0: H1->O1: 0.0 * -0.07544358513197481 = 0
c 1: H2->O1: 0.0 * -0.12346935587390481 = 0
c 2: B2->O1: 0.0 * -0.18705713395934637 = 0
c 3: I1->H1: 0.0 * 0.005292326734004241 = 0
c 4: I2->H1: 0.0 * 0.0049016107791925246 = 0
c 5: B1->H1: 0.0 * 0.010264148244428655 = 0
c 6: I1->H2: 0.0 * -0.0068740307089347 = 0
c 7: I2->H2: 0.0 * -0.005236293038814788 = 0
c 8: B1->H2: 0.0 * -0.02103299467864286 = 0

```

From this we can determine each of the weight change values. The value of \mathbf{c} is zero in all cases. Because of this the weight change value is the negative of the weight update value. This results in a weight change of -0.01 for every weight.

```
Weight Change 0: -0.01
Weight Change 1: -0.01
Weight Change 2: -0.01
Weight Change 3: -0.01
Weight Change 4: -0.01
Weight Change 5: -0.01
Weight Change 6: -0.01
Weight Change 7: -0.01
Weight Change 8: -0.01
```

This results in the following weights.

```
Weight 0: H1->O1: -0.3279194894311762
Weight 1: H2->O1: 0.48171409964135703
Weight 2: B2->O1: 0.6792991203673414
Weight 3: I1->H1: 0.03217052401326839
Weight 4: I2->H1: 0.3234107719788818
Weight 5: B1->H1: -0.3635107399577998
Weight 6: I1->H2: 0.8487814395569221
Weight 7: I2->H2: 0.36158711646254804
Weight 8: B1->H2: -0.0024983800254990973
```

This ends the first iteration of the RPROP. Some implementations of RPROP will suppress any weight changes for the first training iteration to allow it to “initialize”. However skipping the weight update is generally unnecessary, as the weights are starting from random values anyhow.

The first iteration is now complete. The first iteration serves as little more than an initialization iteration.

5.4.2 Training Iteration #2

We begin the second iteration by again calculating the gradients of each of the weights. The gradients will have changed, because the underlying weights changed. We will use these new gradient values to calculate a new value of \mathbf{c} for each weight.

```
c 0: H1->O1: -0.0754435851317481 * -0.07564714780823276 = 1
c 1: H2->O1: -0.12346935587390481 * -0.10495082682420408 = 1
```

```

c 2: B2->O1: -0.18705713395934637 * -0.16712652502209419 = 1
c 3: I1->H1: 0.005292326734004241 * 0.007147520399328029 = 1
c 4: I2->H1: 0.0049016107791925246 * 0.00657604229900621 = 1
c 5: B1->H1: 0.010264148244428655 * 0.013445893781261988 = 1
c 6: I1->H2: -0.0068740307089347 * -0.006335334269910348 = 1
c 7: I2->H2: -0.005236293038814788 * -0.004772389693042953 = 1
c 8: B1->H2: -0.02103299467864286 * -0.01641086135590903 = 1

```

From this we can determine each of the weight change values. The value of c is one in all cases. This means that the sign of none of the gradients changed. Because of this the weight change value is the negative of the weight update value. This results in a weight change of -0.01 for every weight.

```

Weight Change 0: -0.01
Weight Change 1: -0.01
Weight Change 2: -0.01
Weight Change 3: -0.01
Weight Change 4: -0.01
Weight Change 5: -0.01
Weight Change 6: -0.01
Weight Change 7: -0.01
Weight Change 8: -0.01

```

This results in the following weights.

```

-0.4479194894311762
0.36171409964135703
0.5592991203673414
0.1521705240132684
0.4434107719788818
-0.24351073995779982
0.7287814395569221
0.24158711646254805
-0.12249838002549909

```

Everything continues to move in the same direction as iteration one. The update values will also move up from 0.1 to 0.12, which is the direction specified by positive η . This ends iteration 2.

As you can see, the training will continue in the present path until one of the gradients changes signs. The update values will continue to grow, and the weights will continue to change by a value based on the sign of their gradient and the weight update value. This continues until iteration #8 when the update values will begin to take on different values.

5.4.3 Training Iteration #8

We begin the ninth iteration by again calculating the gradients of each of the weights. The weights have changed over the iterations that we skipped. This will cause the gradients to change. We will use the previous gradients and these new gradients to calculate a new value of \mathbf{c} for each weight.

```
c 0: H1->O1: -0.024815885942307825 * -0.009100949661316388 = 1
c 1: H2->O1: -0.023266017866306714 * -0.0091094801332018 = 1
c 2: B2->O1: -0.045150613856680816 * -0.01879394429346648 = 1
c 3: I1->H1: 7.967367771930835E-4 * -0.001991948196987281 = -1
c 4: I2->H1: 8.920300028002447E-4 * -0.0015368274160319773 = -1
c 5: B1->H1: 0.005021027721309641 * 0.0011945995248252954 = 1
c 6: I1->H2: -0.0010888731437352726 * -3.140593778629495E-4 = 1
c 7: I2->H2: -7.180424572079871E-4 * -1.4136606514778948E-4 = 1
c 8: B1->H2: -0.002215756957751366 * -7.264418447125096E-4 = 1
```

As you can see from above there are two weights that now have a -1 value for \mathbf{c} . This indicates that the gradient sign has changed. This means we have passed a local minima, and error is now climbing. We must do something to stop this increase in error for weights three and four.

Instead of multiplying the weight update value by positive eta, we will now use negative eta for weights three and four. This will scale the weight back from the value of 0.04319 to 0.021599. As previously discussed, the value for negative eta is a constant defined to be 0.5. At the end of this iteration, the weight update values are as follows.

```
Update 0: H1->O1: 0.05183999999999999
Update 1: H2->O1: 0.05183999999999999
Update 2: B2->O1: 0.05183999999999999
Update 3: I1->H1: 0.021599999999999998
Update 4: I2->H1: 0.021599999999999998
Update 5: B1->H1: 0.05183999999999999
Update 6: I1->H2: 0.05183999999999999
Update 7: I2->H2: 0.05183999999999999
Update 8: B1->H2: 0.05183999999999999
```

As you can see all weights have continued increasing except for weights three and four. These two weights have been scaled back, and will receive a much smaller update the next iteration. Additionally, the last gradient for each of these two weights is set to zero. This will prevent a modification of the weight update values in the next training iteration. This is accomplished, because \mathbf{c} will be zero in the next step because the previous gradient is zero.

This process will continue until the global error of the neural network falls to an acceptable level. You can see the complete training process here.

```
Epoch #1 Error:0.3100155809627523
Epoch #2 Error:0.2888003866116162
Epoch #3 Error:0.267380775814409
Epoch #4 Error:0.25242444534566344
Epoch #5 Error:0.25517114662144347
Epoch #6 Error:0.25242444534566344
Epoch #7 Error:0.2508797332883249
Epoch #8 Error:0.25242444534566344
Epoch #9 Error:0.2509114256134314
...
Epoch #127 Error:0.029681452838256468
Epoch #128 Error:0.026157454894821013
Epoch #129 Error:0.023541442841907054
Epoch #130 Error:0.0253591989944982
Epoch #131 Error:0.020825411676740083
Epoch #132 Error:0.01754524879617848
Epoch #133 Error:0.015171808565942009
Epoch #134 Error:0.012948657050164597
Epoch #135 Error:0.011092515418846417
Epoch #136 Error:0.009750156492866442
```

This same set of weights took 579 iterations with backpropagation. As you can see RPROP outperforms backpropagation by a considerable margin. With more advanced random weight generation, the number of iterations needed for RPROP can be brought down even further. This will be demonstrated in the next chapter.

5.5 Chapter Summary

In this chapter we saw the resilient propagation (RPROP) training method. This training method is much more efficient than backpropagation. Considerably fewer training iterations are necessary for RPROP compared to backpropagation.

RPROP works by keeping an array of update values used to modify the weights of the neural network. The gradients are not used to update the weights of the neural network directly. Rather, the gradients influence how the weight update values are changed for each iteration. Only the sign of the gradient is used to determine the direction to take the update values. This gives a considerable improvement over backpropagation.

So far we have randomized each neural network with purely random numbers. This does not always lead to the fastest training times. As we will see in the next chapter we can make modifications to the random numbers that the neural network is initialized to. These small changes will yield faster training times.

Chapter 6

Weight Initialization

- Ranged Random Weights
- Trained and Untrained Neural Networks
- Nguyen-Widrow Weight Init

Neural networks must start with their weights initialized to random numbers. These random weights provide the training algorithms with a starting point for the weights. If all of the weights of a neural network were set to zero, the neural network would never train to an acceptable error level. This is because a zeroed weight matrix would place the neural network into a local minimum from which it can never escape.

Often the weights of neural networks are simply initialized with random numbers between a specific range. The range -1 to +1 is very popular. These random weights will change as the neural network is trained to produce acceptable outputs. In the next section we will take a look at what trained and untrained neural networks look like.

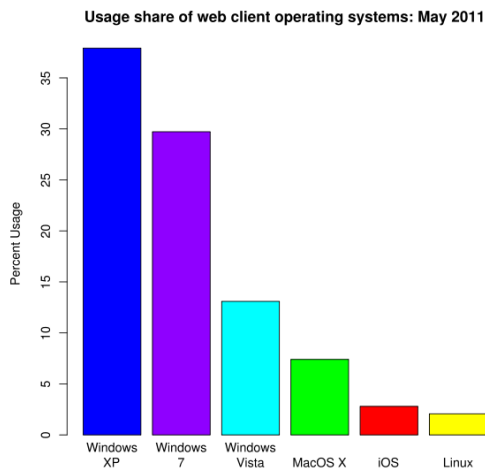
Later in this chapter we will see that these random weights can be modified to help the neural network to train faster. The Nguyen-Widrow weight initialization algorithm is a popular technique to adjust these starting weights. The Nguyen-Widrow weight initialization algorithm puts the weights into a position that is much more conducive to training. This means fewer training iterations to get the neural network to an acceptable error rate.

6.1 Looking at the Weights

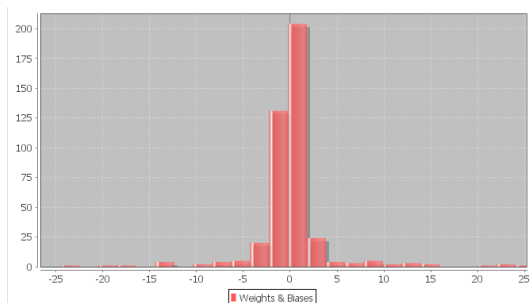
In previous chapters we've looked at the weights of a neural network as an array of numbers. You can't typically glance at a weight array and see any sort of meaningful pattern. However, if the weights are represented graphically patterns begin to emerge.

One common way to view the weights of a neural network is with a special type of chart called a histogram. A histogram is made up of vertical bars that count the number of occurrences in a population. You've probably seen histograms many times before. For example, Figure 6.1 shows the popularity of operating systems.

Figure 6.1: Histogram of OS Popularity (from Wikipedia)



The y-axis shows the number of occurrences of each of the groups in the x-axis. We can use a histogram to look at the weights of a neural network. You can typically tell a trained, from an untrained neural network by looking at this histogram. Figure 6.2 shows a trained neural network.

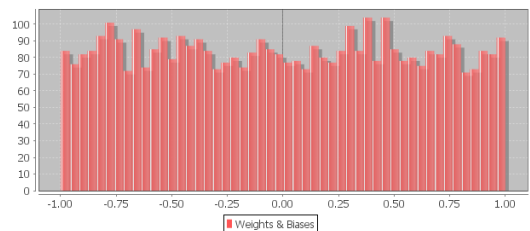
Figure 6.2: A Trained Neural Network

A neural network histogram is the same as concept as the operating system histogram shown earlier. The y-axis specifies how many weights fell into the ranges specified by the numbers on the x-axis. This allows you to see the distribution of the weights.

Most trained neural networks will look something like the above chart. Their weights will be very tightly clustered around zero. A trained neural network will typically look like a vary narrow gaussian curve.

6.2 Ranged Randomization

In the last section we saw what a trained neural network looks like in a weight histogram. Untrained neural networks can have a variety of appearances. How the weight histogram will look will be determined by the weight initialization method used.

Figure 6.3: A Ranged Randomization

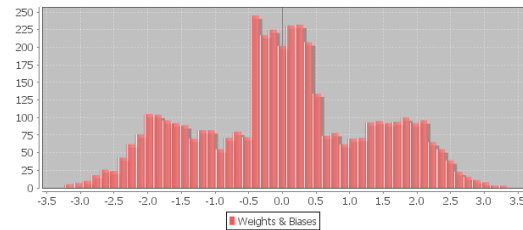
Range randomization produces a very simple looking chart. The more weights there are the flatter the top will be. This is because the random number generator should be giving

you an even distribution of numbers. If you are randomizing to the range of -1 to 1, you would expect to have approximately the same number of weights above zero as below.

6.3 Using Nguyen-Widrow

We will now look at the Nguyen-Widrow weight initialization method. The Nguyen-Widrow method starts out just like the range randomized method. Random values are chosen between -0.5 and +0.5. However, a special algorithm is employed to modify the weights. The histogram of a Nguyen-Widrow weight initialization looks like Figure 6.4.

Figure 6.4: The Nguyen-Widrow Initialization



As you can see, the Nguyen-Widrow initialization has a very distinctive pattern. There is a large distribution of weights between -0.5 and 0.5. However it then gradually rises and then rapidly falls off to around -3.0 and +3.0.

6.3.1 Performance of Nguyen-Widrow

You may be wondering how much of an advantage Nguyen-Widrow can have. The following shows the average number of training iterations needed to train a neural network initialized by range random and Nguyen-Widrow.

Average iterations needed (lower is better)
Range random: 502.86
Nguyen-Widrow: 454.88

As you can see from the above information, the Nguyen-Widrow outperforms the range randomizer.

6.3.2 Implementing Nguyen-Widrow

This technique was invented by Derrick Nguyen and Bernard Widrow. It was first introduced in their paper “Improving the learning speed of 2-layer neural networks by choosing initial values of the adaptive weights” in the Proceedings of the International Joint Conference on Neural Networks, 3:21-26, 1990.

To implement an Nguyen-Widrow randomization first initialize the neural network with random weight values in the range -0.5 to +0.5. This is exactly the same technique as was described earlier in this article for the ranged random numbers.

The Nguyen-Widrow randomization technique is efficient because it assigns each hidden neuron to a range of the problem. To do this we must map the input neurons to the hidden neurons. We calculate a value, called beta, that establishes these ranges. You can see the calculation of beta in Equation 6.1.

$$\beta = 0.7h^{\frac{1}{i}} \quad (6.1)$$

The variable \mathbf{h} represents the number of hidden neurons in the first hidden layer, whereas the variable \mathbf{i} represents the number of input neurons. We will calculate the weights taking each hidden neuron one at a time. For each hidden neuron we calculate the Euclidean norm for of all inputs to the current hidden neuron. This is done with Equation 6.2.

$$n = \sqrt{\sum_{i=0}^{i < w_{max}} w_i^2} \quad (6.2)$$

Beta will stay the same for every hidden neuron. However, the norm must be recalculated for each hidden neuron. Once the beta and norm values have been calculate, the random weights can be adjusted. The equation below shows how weights are adjusted using the previously calculated values.

$$w_{t+1} = \frac{\beta w_t}{n} \quad (6.3)$$

All inbound weights to the current hidden neuron are adjusting using the same norm. This same process is repeated for each hidden neuron.

You may have noticed that we are only specifying how to calculate the weights between the input layer and the first hidden layer. The Nguyen-Widrow method does not specify

how to calculate the weights between a hidden layer and the output. Likewise the Nguyen-Widrow method does not specify how to calculate the weights between multiple hidden layers. All weights outside of the first layer and first hidden layer are simply initialized to a value between -0.5 and +0.5.

6.3.3 Nguyen-Widrow in Action

We will now walk through the random weight initialization for a small neural network. We will look at a neural network that has two input neurons and a single output neuron. There is a single hidden layer with two neurons. Bias neurons are present on the input and hidden layers.

We begin by initializing the weights to random numbers in the range -0.5 to +0.5. The starting weights are shown here.

Weight 0: H1→O1:	0.23773012320107711
Weight 1: H2→O1:	0.2200753094723884
Weight 2: B2→O1:	0.12169691073037914
Weight 3: I1→H1:	0.5172524211645029
Weight 4: I2→H1:	-0.5258712726818855
Weight 5: B1→H1:	0.8891383322123643
Weight 6: I1→H2:	-0.007687742622070948
Weight 7: I2→H2:	-0.48985643968339754
Weight 8: B1→H2:	-0.6610227585583137

First we must calculate beta, which is given in Equation 6.1. This calculation is shown here.

$\text{Beta} = 0.7 * (\text{hiddenNeurons} ^ (1.0/\text{inputCount}))$
--

Filling in the variables we have.

$\text{Beta} = 0.7 * (2.0 ^ (1.0/2.0)) = 0.9899$
--

We are now ready to modify the weights. We will only modify weights three through eight. Weights zero through two are not covered by the Nguyen-Widrow algorithm, and are simply set to random values between -0.5 and 0.5.

The weights will be recalculated in two phases. First we will recalculate all of the weights from the bias and input neurons to hidden neuron one. To do this we must calculate the Euclidean norm, or magnitude, for hidden neuron one. From Equation 6.2, we have the following.

$$\text{Norm Hidden 1} = \sqrt{(\text{weight 3})^2 + (\text{weight 4})^2 + (\text{weight 5})^2}$$

Filling in the weights we have the following.

$$\text{Norm Hidden 1} = \sqrt{(0.5172^2) + (-0.5258)^2 + (0.8891^2)} = 1.155$$

We will now look at how to calculate the first weight.

$$\text{New Weight} = (\text{beta} * \text{Old Weight}) / \text{Norm Hidden 1}$$

Filling in values we have.

$$\text{New Weight} = (0.9899 * 0.5172) / 1.155 = 0.4432$$

All three weights feeding hidden the first neuron are transformed in this way.

$$\begin{aligned} 0.5172524211645029 &\rightarrow 0.44323151482681195 \text{ (as just seen)} \\ -0.5258712726818855 &\rightarrow -0.4506169739523903 \\ 0.8891383322123643 &\rightarrow 0.7618990530577658 \end{aligned}$$

We now repeat the same process for neuron two. The value for beta stays the same. However, the magnitude must be recalculated. The recalculated magnitude for neuron two is given here.

$$\text{Neuron Two Magnitude} = 0.8227815750355376$$

Using this we can now recalculate weights six through eight.

$$\begin{aligned} -0.007687742622070948 &\rightarrow -0.009249692928269318 \\ -0.48985643968339754 &\rightarrow -0.5893825884595142 \\ -0.6610227585583137 &\rightarrow -0.79532547274779 \end{aligned}$$

This results in the final values for all of the weights.

$$\begin{aligned} \text{Weight 0: H1} \rightarrow \text{O1: } &0.23773012320107711 \\ \text{Weight 1: H2} \rightarrow \text{O1: } &0.2200753094723884 \\ \text{Weight 2: B2} \rightarrow \text{O1: } &0.12169691073037914 \\ \text{Weight 3: I1} \rightarrow \text{H1: } &0.44323151482681195 \\ \text{Weight 4: I2} \rightarrow \text{H1: } &-0.4506169739523903 \\ \text{Weight 5: B1} \rightarrow \text{H1: } &0.7618990530577658 \\ \text{Weight 6: I1} \rightarrow \text{H2: } &-0.009249692928269318 \\ \text{Weight 7: I2} \rightarrow \text{H2: } &-0.5893825884595142 \\ \text{Weight 8: B1} \rightarrow \text{H2: } &-0.79532547274779 \end{aligned}$$

These weights, though still random, better utilize the hidden neurons. These are the weights that we will begin training with.

6.4 Chapter Summary

In this chapter we saw how the weights of a neural network are initialized. These weights must be set to random values. If the weights were all set to zero, the neural network would never train properly. There are several different ways that the weights of a neural network are commonly initialized.

The first is range randomization. This process initializes each weight of the neural network to a random value in a specific range. The range used is typically -1 to 1 or 0 to 1. Though range randomization can provide very good results for neural network training, there are better methods.

Nguyen-Widrow weight initialization can provide significantly better training times than simple range randomization. The Nguyen-Widrow randomization technique is efficient because it assigns each hidden neuron to a range of the problem.

So far all training presented has focused on attempting to minimize the error. In the next chapter we will look the Levenberg Marquardt (LMA) training algorithm. LMA can often train in fewer iterations than resilient propagation (RPROP).

Chapter 7

LMA Training

- Newton's Method
- Calculating the Hessian
- The Levenberg-Marquardt Algorithm
- Multiple Output Neurons

The Levenberg-Marquardt algorithm (LMA) is a very efficient training method for Neural Networks. In many cases LMA will outperform RPROP. Because of this it is an important training algorithm that should be considered by any neural network programmer.

LMA is a hybrid algorithm that is based on both Newton's Method and Gradient Descent (backpropagation). This allows LMA to have strengths of both. Gradient Descent is guaranteed to converge to a local minimum; however, it is quite slow. GNA is quite fast but often fails to converge. By using a damping factor to interpolate between the two, a hybrid method is created. To understand how this works we first examine Newton's method. Newton's Method is shown here.

$$W_{min} = W_0 - H^{-1}g \quad (7.1)$$

You will notice several variables shown above. The result of the above equation is deltas that can be applied to the weights of the neural network. The variable \mathbf{H} represents the Hessian. The will be covered in the next section. The variable g represents the gradients of the neural

network. You will also notice the -1 “exponent” on the variable \mathbf{H} . This specifies that we are doing a matrix decomposition of the variables \mathbf{H} and \mathbf{g} .

I could easily spend an entire chapter on matrix decomposition. I decided not to. We will simply treat matrix decomposition as a black box atomic operator. Much as if I tell you to take the square root of a number. I do not explain how to calculate a square root, I simply assume that you will use the square root function in your programming language of choice.

There is a common piece of code for matrix decomposition that is included in the JAMA package. This public domain code, which was adapted from a FORTRAN program, has been used in many mathematical computer applications. I will not cover how to actually perform matrix decomposition here. I will assume that you use a package such as JAMA, or another source.

There are several types of matrix decomposition. The decomposition that we are going to use is the LU decomposition. This decomposition requires a square matrix. This works well because the Hessian matrix has the same number of rows as columns. There is a row and column for every weight in the neural network. The LU decomposition takes the Hessian, which is a matrix of the second derivatives of the partial derivatives of the output of each of the weights. The LU decomposition solves the Hessian by the gradients. These gradients are the square of the error of each weight. These are the same gradients as we calculated in Chapter 5, except they are squared. Because the errors are squared we are required to use the sum of square error when dealing with LMA.

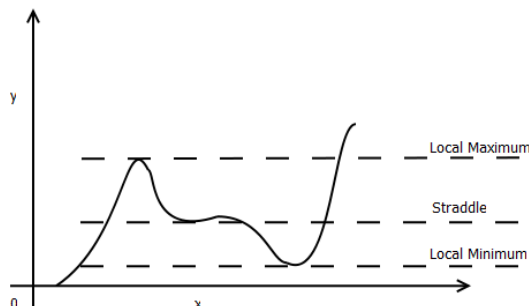
If you have never heard the term “second derivative” before, the second derivative is the derivative of the first derivative. Recall from Chapter 2 that the derivative of a function is the slope at any point. This slope points in the direction that the curve is approaching a local minimum. The second derivative is also a slope. It points in a direction to minimize the first derivative. The goal of Newton’s method, as well as the LMA algorithm is to reduce all of the gradients to zero.

Interestingly, this goal does not include the error. Newton’s Method and LMA are somewhat oblivious to the error. They simply seek to squash all of the gradients to zero. In reality they are not completely oblivious to the error, as the error is used to calculate the gradients.

Newton’s Method will converge the weights of a neural network to a local minimum, a local maximum or a straddle position. This is done by minimizing all of the gradients (first derivatives) to zero. The derivatives will all be zero at local minima, maxima or straddle

position. These three points are shown in Figure 7.1.

Figure 7.1: Local Minimum, Straddle and Local Maximum



It is up to the algorithm implementation to ensure that local maximum's and straddle points are filtered out. The above algorithm works by taking the matrix decomposition of the Hessian matrix and the gradients. The Hessian matrix is typically estimated. There are various means of doing this. If the Hessian is inaccurate this can greatly throw off Newton's Method. LMA enhances Newton's Algorithm to the following formula.

$$W_{min} = W_0 - (H + \lambda I)^{-1}g \quad (7.2)$$

Here we add a damping factor multiplied by an identity matrix. Lambda is the damping factor and "I" represents the identity matrix. An identity matrix is a square matrix with all zeros except for a NW line of ones). As lambda increases the hessian will be factored out of the above equation. As lambda decreases the hessian becomes more significant than gradient descent. This allows the training algorithm to interpolate between gradient descent and Newton's Method. Higher lambda favors gradient descent, lower lambda favors Newton.

A training iteration of LMA begins with a low lambda and increases it until a desirable outcome is produced. This can be seen in the following flowchart.

7.1 Calculation of the Hessian

The Hessian matrix is a square matrix with rows and columns equal to the number of weights in the neural network. Each cell in this matrix represents the second order derivative of the output of the neural network with respect to a given weight combination. The Hessian is shown in Equation 7.3.

$$H(e) = \begin{bmatrix} \frac{\partial^2 e}{\partial w_1^2} & \frac{\partial^2 e}{\partial w_1 \partial w_2} & \cdots & \frac{\partial^2 e}{\partial w_1 \partial w_n} \\ \frac{\partial^2 e}{\partial w_2 \partial w_1} & \frac{\partial^2 e}{\partial w_2^2} & \cdots & \frac{\partial^2 e}{\partial w_2 \partial w_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 e}{\partial w_n \partial w_1} & \frac{\partial^2 e}{\partial w_n \partial w_2} & \cdots & \frac{\partial^2 e}{\partial w_n^2} \end{bmatrix} \quad (7.3)$$

It is important to note that the Hessian is symmetrical about the diagonal. This can be used to enhance performance of the calculation. Calculation of the Hessian can be accomplished by calculating the gradients.

$$\frac{\partial E}{\partial w_{(i)}} = 2(y - t) \frac{\partial y}{\partial w_{(i)}} \quad (7.4)$$

The second [[derivative]] of above equation becomes an element of the Hessian matrix. This is calculated with the following formula.

$$\frac{\partial^2 E}{\partial w_i \partial w_j} = 2 \left(\frac{\partial y}{\partial w_i} \frac{\partial y}{\partial w_j} + (y - t) \frac{\partial^2 y}{\partial w_j \partial w_i} \right) \quad (7.5)$$

The above formula can easily be calculated if not for the second component. This component involves the second partial derivative is difficult to calculate. This component is actually not that important and can be dropped. For an individual training case it might be very important. However, the second component is multiplied by the error of that training case. We assume that the errors in a training set are independent and evenly distributed about zero. On an entire training set they should essentially cancel each other out. This is not a perfect assumption. However, we only seek to approximate the Hessian.

This results in the following equation.

$$\frac{\partial^2 E}{\partial w_i \partial w_j} = 2 \left(\frac{\partial y}{\partial w_i} \frac{\partial y}{\partial w_j} \right) \quad (7.6)$$

Using the above equation is, of course, only an approximation of the true Hessian. However, the simplification of the algorithm to calculate the second derivative is well worth the loss in accuracy. Any loss in accuracy will most likely be accounted for with an increase in lambda.

To calculate the Hessian and gradients we must calculate the partial first derivatives of the output of the neural network. Once we have these partial first derivatives the above equations allow us to easily calculate the Hessian and gradients.

Calculation of the first derivatives of the output of the neural network is very similar to the process that we performed to calculate the gradients for backpropagation. The primary difference is that the derivative of the output is taken this time. In standard backpropagation the derivative of the error function is taken. We will not review the entire backpropagation process here. Chapter 4 covers backpropagation and gradient calculation.

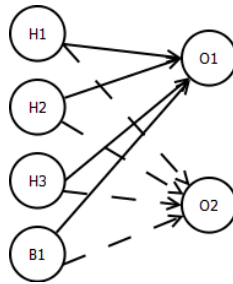
7.2 LMA with Multiple Outputs

Some implementations of the LMA algorithm only support a single output neuron. This is primarily because the LMA algorithm has its roots in mathematical function approximation. In mathematics functions typically only return a single value. As a result, there is not a great deal of information on support for multiple neuron support.

Support for multiple output neurons involves simply summing each cell of the Hessian as the additional output neurons are calculated. It is as if you calculated a separate Hessian matrix for each output neuron, and then summed the Hessian matrices together. This is the approach that Encog uses, and it leads to very fast convergence times.

One important aspect to consider with multiple outputs is that not every connection will be used. Depending on which output neuron you are currently calculating for, there will be unused connections for the other output neurons. It is very important that the partial derivative for each of these unused connections be set to zero when the other output neuron is being calculated for.

For example, consider a neural network that has two output neurons and three hidden neurons. Each of these two output neurons would have a total of four connections from the hidden layer. There would be three from the three hidden neurons, and a fourth for the bias neuron. This segment of the neural network would look like Figure 7.2.

Figure 7.2: Calculating Output Neuron 1

Here we are calculating output neuron one. Notice that there are also four connections for output neuron two? It is critical that the derivatives of the four connections to output neuron two calculate to zero when output neuron one is calculated. If this process is not followed, the Hessian will not work for multiple outputs.

7.3 Overview of the LMA Process

So far we have only seen the math behind LMA. LMA must be part of an algorithm for it to be effective. The LMA process can be summarized as the following steps.

1. Calculate the first derivative of output of the neural network with respect to every weight
2. Calculate the Hessian
3. Calculate the gradients of the error (ESS) with respect to every weight
4. Either set lambda to a low value (first iteration) or the lambda of the previous iteration
5. Save the weights of the neural network
6. Calculate delta weight based on the lambda, gradients and Hessian
7. Apply the deltas to the weights and evaluate error
8. If error has improved, end the iteration
9. If error has not improved increase lambda (up to a max lambda) restore the weights and go back to step 6

As you can see, the process for LMA revolves around setting the lambda value low and then slowly increasing it if the error rate does not improve. It is important to save the weights at each change in lambda so that they can be restored if the error does not improve.

7.4 Chapter Summary

In this chapter we saw how to use the Levenberg-Marquardt Algorithm (LMA) to train a neural network. The LMA algorithm is very efficient and will of the outperform backpropagation and RPROP. LMA works by attempting to minimize all of the error gradients to zero.

LMA is a hybrid algorithm. LMA works by using both Newton's Method, as well as regular gradient descent. Gradient descent is essentially the same thing as backpropagation. The beauty of LMA is that a damping factor, called lambda, is used to interpolate between gradient descent and Newton's Method. We start out favoring Newton's Method; however, if the error fails to improve we introduce more gradient descent learning.

So far we have only looked at the feedforward neural network. In the next two chapters we will introduce two new types of neural network. You will learn about the Self Organizing Map (SOM) and a Radial Basis Function neural network. Both types of neural network share a great deal with the feedforward neural networks you have seen thus far in this book.

Chapter 8

Self-Organizing Maps

- SOM Structure
- Training a SOM
- Neighborhood Functions
- SOM Error Calculation

So far this book has focused on the feedforward neural network. In this chapter we will look at a different type of neural network. This chapter will focus on the Self-Organizing Map (SOM). Though the SOM can be considered a type of feedforward neural network, it is used very differently. A SOM is used for unsupervised classification.

Self-Organizing maps start with random weights, just like the neural networks we've seen so far. However, the means by which these weights are organized to produce meaningful output is very different than a feedforward neural network. SOM's make use of unsupervised training.

Unsupervised training works very differently than the supervised training that we have been using thus far. Unsupervised training sets only specify the input to the neural network. There is no ideal output provided. Because of this the SOM learns to cluster, or map, your data into a specified number of classes.

The number of output neurons in the SOM defines how many output classes you would like to cluster the input data into. The number of input neurons defines the attributes about each of the training set items that you would like the neural network to cluster based on.

This is exactly the same as was true for feedforward neural networks. You must present data to the neural network through a fixed number of input neurons.

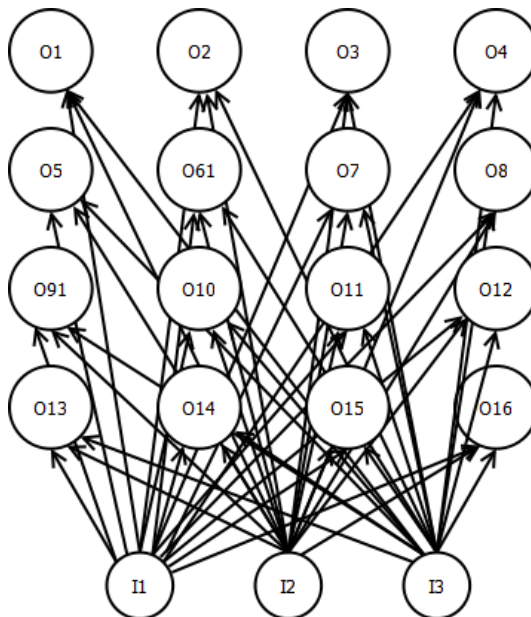
SOM's are only used for classification, which is they divide the input into classes. SOM's are not generally used for regression. Regression is where the neural network predicts one, or more, numeric values. With a SOM, you simply provide the total number of classes, and the SOM automatically your training data into these classes. Additionally, the SOM should be able to organize data it was never trained with into these classes.

8.1 SOM Structure

The SOM has a structure somewhat similar to the neural networks that we have seen already in this book. However, there are some very important differences. Some of the differences between SOM's and traditional feedforward neural networks are shown here.

- SOM networks have no bias neurons
- SOM networks have no hidden layers
- SOM networks have no activation functions
- SOM networks must have more than one output neuron
- Output neurons in a 1d, 2d, 3d, etc lattice

Consider a SOM that has six input neurons and three output neurons. This SOM is designed to classify data into three groups. The incoming data items that will be classified each have six values. This SOM can be seen in Figure 8.1.

Figure 8.1: A 2D SOM with 16 Outputs and 3 Inputs

As you can see from the above diagram there are no hidden layers and there are no bias neurons. You will also notice that the output neurons are arranged in a lattice. In this case the output neurons are arranged in a 2D grid. This grid will become very important when the neural network is trained. Neurons near each other on the grid will be trained together. The grid does not need to be 2D, grids can be 1D, 3D or of an even higher number of dimensions. The topology of this grid is defined by the neighborhood function that is used to train the neural network. Neighborhood functions will be described later in this chapter.

A SOM can be thought of as reducing a high number of dimensions to a lower number of dimensions. The higher number of dimensions is provided by the input neurons. One dimension for each input neuron. The lower number of dimensions is the configuration of the output neuron lattice. In the case of Figure 8.1, we are reducing three dimensions to two.

In Chapter 1 we saw how a feedforward neural network calculates its output. The method by which a SOM calculates its output is very different. The output of a SOM is the “winning” output neuron for a given input. This winning neuron is also called the Best Matching Unit, or BMU. We will see how to calculate the BMU in the next section.

8.1.1 Best Matching Unit

The best matching unit is the output neuron whose weights most closely matches the input being provided to the neural network. Consider the SOM shown in Figure 8.1. There are 16 output neurons and 3 input neurons. This means that each of the 16 output neurons has 3 weights, one from each of the input neurons. The input to the SOM would also be three numbers, as there are three input neurons. To determine the BMU we determine the output neuron whose three weights most closely match the three input values fed into the SOM.

It is easy enough to calculate the BMU. To do this, we loop over every output neuron and calculate the Euclidean distance between the output neuron's weights and the input values. Whatever output neuron has the lowest Euclidean distance is the BMU. The Euclidean distance is simply the distance between two multi-dimensional points. If you are dealing with two dimensions, the Euclidean distance is simply the length of a line drawn between the two points.

The Euclidean distance is used often in Machine Learning. It is a quick way to compare two arrays of numbers that have the same number of elements. Consider three arrays, named array **a**, array **b** and array **c**. The Euclidean distance between array **a** and array **b** is 10. The Euclidean distance between array **a** and array **c** is 20. In this case, the contents of array **a** more closely match array **b** than they do array **c**.

We will now look at how to actually calculate the Euclidean distance. Equation 8.1 shows this formula.

$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2} \quad (8.1)$$

The above equation shows us the Euclidean distance **d** between two arrays **p** and **q**. The above equation also states that **d(p,q)** is the same as **d(q,p)**. This simply states that the distance is the same no matter what end you start at. Calculation of the Euclidean distance is no more than summing the squares of the difference of each array element. Finally, the square root of this sum is taken. This square root is the Euclidean distance. The output neuron with the lowest Euclidean distance is the BMU.

8.2 Training a SOM

In the previous chapters we learned several methods for training a feedforward neural network. We learned about such techniques as backpropagation, RPROP and LMA. These are all supervised training methods. Supervised training methods work by adjusting the weights of a neural network to produce the correct output for a given input.

A supervised training method will not work for a SOM. SOM networks require unsupervised training. In this section we will learn to train a SOM with an unsupervised method. The training technique generally used for SOM networks is summarized in Equation 8.2.

$$W_v(t+1) = W_v(t) + \theta(v, t)\alpha(t)(D(t) - W_v(t)) \quad (8.2)$$

The above equation shows how the weights of a SOM neural network are updated as training progresses. The current training iteration is noted by the letter \mathbf{t} , and the next training iteration is noted by $\mathbf{t}+1$. Equation 8.2 allows us to see how weight \mathbf{v} is adjusted for the next training iteration.

The variable \mathbf{v} denotes that we are performing the same operation on every weight. The variable \mathbf{W} represents the weights. The symbol **theta** is a special function, called a neighborhood function. The variable \mathbf{D} represents the current training input to the SOM.

The symbol alpha denotes a learning rate. The learning rate changes for each iteration. This is why Equation 8.2 shows the learning rate with the symbol t , as the learning rate is attached to the iteration. The learning rate for a SOM is said to be monotonically decreasing. Monotonically decreasing simply means that the learning rate only falls, and never increases.

8.2.1 SOM Training Example

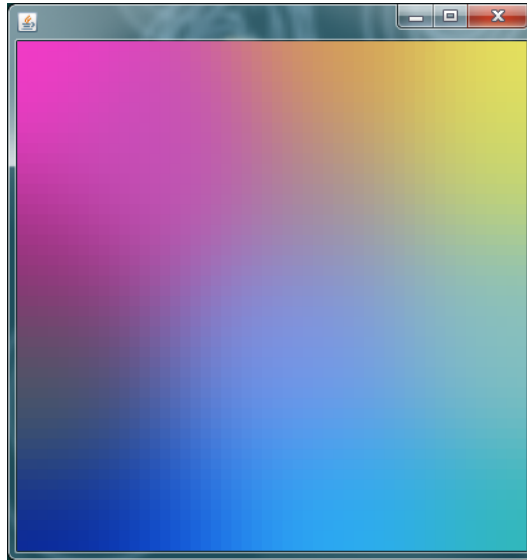
We will now see how to actually train a SOM. We will apply the equation presented in the previous section. However, we will approach it more from an algorithm perspective so we can see the actual learning strategy behind the equation. To see the SOM in action we need a very simple example. We will use a SOM very similar to the one that we saw in Figure 8.1. This SOM will attempt to match colors. However, instead of the 4x4 lattice we saw in Figure 8.1, we will have a lattice of 50x50. This results in a total of 2,500 output neurons.

The SOM will use its three input neurons to match colors to the 2,500 output neurons. The three input neurons will contain the red, blue and green components of the color that

is currently being submitted to the SOM. For training, we will generate 15 random colors. The SOM will learn to cluster these colors.

This sort of training is demonstrated in one of the Encog examples. You can see the output from this example program in Figure 8.2.

Figure 8.2: Mapping Colors



As you can see from the above figure, similar colors are clustered together. Additionally, there are 2,500 output neurons, and only 15 colors that were trained with. This network could potentially recognize up to 2,500 colors. The fact that we trained with only 15 colors means we have quite a few unutilized output neurons. These output neurons will learn to recognize colors that are close to the 15 colors that we trained with.

What you are actually seeing in Figure 8.2 are the weights of SOM network that has been trained. As you can see, even though the SOM was only trained to recognize 15 colors, it is able to recognize quite a few colors. Any new color provided to the SOM will be mapped to one of the 2,500 colors seen in the above image. The SOM can be trained to recognize more classes than it its provided training data. This is defiantly the case in Figure 8.2. The unused output neurons will end up learning to recognize data that falls between elements of the smaller training set.

8.2.2 Training the SOM Example

We will now look at how the SOM network is trained for the colors example. To begin with all of the weights of the SOM network are randomized to values between -1 and +1. A training set is now generated for 15 random colors. Each of these 15 random colors will have three input values. Each of the three input values will be a random value between -1 and +1. For the red, green and blue values -1 represents that the color is totally off, and +1 represents that the color is totally on.

We will see how the SOM is trained for just one of these 15 colors. The same process would be used for the remaining 14 colors. Consider if we were training the SOM for the following random color.

-1,1,0.5

We will see how the SOM will be trained with this training elemtn.

8.2.3 BMU Calculation Example

The first step would be to compare this input against every output neuron in the SOM and find the Best Matching Unit (BMU). BMU was discussed earlier in this chapter. The BMU can be calculated by finding the smallest Euclidean distance in SOM. The random weights in the SOM are shown here.

Output Neuron 1: -0.2423, 0.4837, 0.8723
 Output Neuron 2: -0.5437, -0.8734, 0.2234
 Ë
 Output Neuron 2500: -0.1287, 0.9872, -0.8723

Of course, we are skipping quite a few of the output neurons. Normally, you would calculate the Euclidean distance for all 2,500 output neurons. Just calculating the Euclidean distance for the above three neurons should give you an idea of how this is done. Using Equation 8.1 we calculate the Euclidean distance between the input and neuron one.

$$\text{sqrt}((-0.2423 - -1)^2 + (0.4837 - 1)^2 + (0.8723 - 0.5)^2) = 0.9895$$

A similar process can be used to calculate neuron two.

$$\text{sqrt}((-0.5437 - -1)^2 + (-0.8734 - 1)^2 + (0.2234 - 0.5)^2) = 1.947$$

Similarly, we can also calculate neuron 2,500.

$$\text{sqrt}((-0.1287 - -1)^2 + (0.9872 - 1)^2 + (-0.8723 - 0.5)^2) = 1.6255$$

Now that we have calculated all of the Euclidean distances, we can determine the BMU. The BMU is neuron one. This is because the distance of 0.9895 is the lowest. Now that we have a BMU, we can update the weights.

8.2.4 Example Neighborhood Functions

We will now loop over every weight in the entire SOM and use Equation 8.2 to update them. The idea is that we will modify the BMU neuron to be more like the training input. However, in addition to modifying the BMU neuron to be more like the training input, we will also modify neurons in the neighborhood around the BMU to be more like the input as well. However, the further a neuron is from the BMU the less of an impact this weight change will have.

Determining the amount of change that will happen to a weight is the job of the neighborhood function. Any radial basis function (RBF) can be used as a neighborhood function. A radial basis function (RBF) is a real-valued function whose value depends only on the distance from the origin.

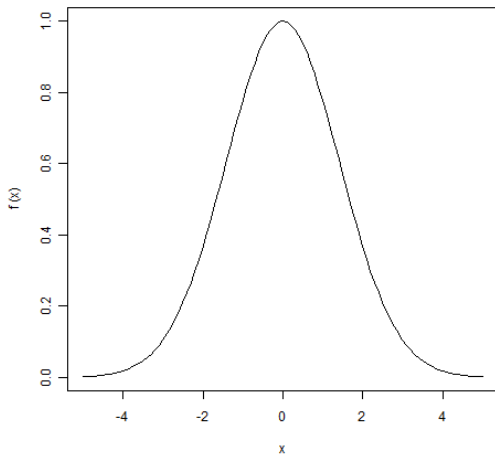
The Gaussian function is the most common choice for a neighborhood function. The Gaussian function is shown in Equation 8.3.

$$f(x_1, x_2, \dots, x_n) = e^{-v}, \quad (8.3)$$

This equation continues as follows.

$$v = \sum_{i=0}^{i < n} \frac{x_i - c^2}{2w^2} \quad (8.4)$$

You can see the Gaussian function graphed in Figure 8.3. Where \mathbf{n} is the number of dimensions, \mathbf{c} is the center, \mathbf{w} is the width of the Gaussian curve. The number of dimensions is equal to the number of input neurons. The width starts out at some fixed number and decreases as learning progresses. By the final training iteration the width should be one.

Figure 8.3: Gaussian Function Graphed

From the above figure, you can see that the Gaussian function is a radial basis function. The value only depends on the distance from the origin. That is to say that $\mathbf{f}(\mathbf{x})$ has the same value regardless of if \mathbf{x} is -1 or $+1$.

Looking at Figure 8.3, you can see how the Gaussian function scales the amount of training received by each neuron. The BMU would have zero distance from itself, so the BMU would receive full training of 1.0. As you move further away from the BMU, in either direction, the amount of training quickly falls off. As a neuron that was -4 or $+4$ from the BMU would receive hardly any training at all.

The Gaussian function is not the only function available for SOM training. Another common choice is the Ricker wavelet, or "Mexican hat" neighborhood function. This function is generally only known as the "Mexican Hat" function in the Americas, due to its resemblance to a "sombrero". In technical nomenclature this function is known as the Ricker wavelet, where it is frequently employed to model seismic data. The equation for the Mexican Hat function is shown in Equation 8.4.

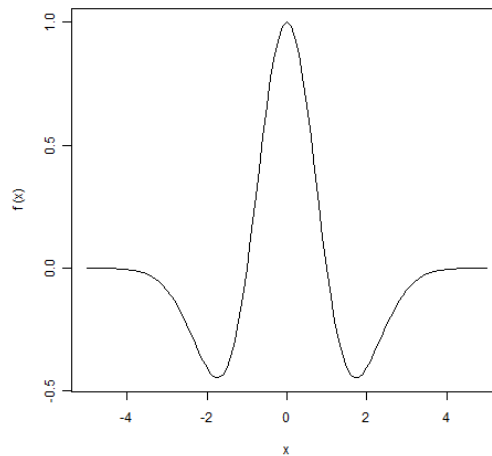
$$f(x_1, x_2, \dots, x_n) = (1 - v)e^{-0.5v}, \quad (8.5)$$

This equation continues with the following.

$$v = \sum_{i=0}^{i < n} (x - c)^2 \quad (8.6)$$

You will see the value of the Mexican Hat neighborhood function when you examine its graph in Figure 8.4.

Figure 8.4: Graph of the Mexican Hat



Just as before, the BMU is at the origin of \mathbf{x} . As you can see from the above chart, the Mexican Hat function actually punishes neurons just at the edge of the radius of the BMU. Their share of the learning will actually be less than zero. As you get even further from the BMU, the learning again returns to zero, and remains at zero. However, neurons just at the edge will have negative learning applied. This can cause the SOM to classes to better differentiate among each other.

8.2.5 Example Weight Update

Now that we have seen how to calculate the BMU and the neighborhood function, we are finally ready to calculate the actual weight update. The formula for the weight update was given in Equation 8.2. We will now calculate the components of this equation for an actual

weight update. Recall that the weights given earlier in this chapter. Neuron one had the following weights.

Output Neuron 1: $-0.2423, 0.4837, 0.8723$

Neuron one was the BMU when provided with the following input.

$-1, 1, 0.5$

As you can see, the BMU is somewhat similar to the input. During training, we now want to modify the BMU to be even more like the input. We will also modify the neighbors to receive some of this learning as well.

The “ultimate” learning is to simply copy the input into the weights. Then the Euclidean distance of the BMU will become zero and the BMU is perfectly trained for this input vector. However, we do not want to go that extreme. We simply want to move the weights in the direction of the input. This is where equation 8.2 comes in. We will use Equation 8.2 for every weight in the SOM, not just the BMU weights. For this example, we will start with the BMU. Calculating the update to the first weight, which is -1 , we have.

$$w = -0.2423 + (N * r * (-1\dot{U}(-0.2423))) \quad (8.7)$$

Before we calculate the above equation, we will take a look at what we are actually doing. Look at the term at the far right. We are taking the difference between the input and the weight. As I said before, the ultimate is just to assign the input to the weight. If we simply added the last term to the weight, the weight would be the same as the input vector.

We do not want to be this extreme. Therefore we scale this difference. First we scale it by the learning rate, which is the variable r . Then we also scale it by the result of the neighborhood function, which is N . For the BMU, N is the value 1.0 and has no effect. If the learning rate were 1.0, as well, then the input would actually be copied to the weight. However a learning rate is never above 1.0. Additionally, the learning rate typically decays as the learning iterations progress. Considering that we have a learning rate of 0.5, our weight update becomes the following.

$$w = -0.2423 + (1.0 * 0.5 * (-1 - -0.2423)) = -0.62115 \quad (8.8)$$

As you can see, the weight moved from -0.2423 to -0.62115 . This moved the weight closer to the input of -1 . We can perform a similar update for the other two weights that feed into the BMU. This can be seen here.

$$w = 0.4837 + (1.0 * 0.5 * (1 - 0.4837)) = 0.74185 \quad w = 0.8723 + (1.0 * 0.5 * (0.5 - 0.8723)) = 0.68615 \quad (8.9)$$

As you can see, both weights moved closer to the input values.

The neighborhood function is always 1.0 for the BMU. However, consider if we were to calculate the weight update for Neuron Two, which is not the BMU. We would need to calculate the neighborhood function, which was given in Equation 8.3. This assumes we are using a Gaussian neighborhood function. The Mexican Hat function could also be used.

The first step is to calculate \mathbf{v} . Here we use a width \mathbf{w} of 3.0. When using Gaussian for a neighborhood function the center \mathbf{c} is always 0.0.

$$v = ((x1 - c1)^2 / 2w^2) + ((x2 - c2)^2 / 2w^2) \quad (8.10)$$

We will plug in these values. The values $\mathbf{x1}$ and $\mathbf{x2}$ specify how far away the current neuron is from the BMU. The value for $\mathbf{x1}$ specifies the column distance, and the value of $\mathbf{x2}$ specifies the row distance. Because neuron two is on the same row, then $\mathbf{x2}$ will be zero. Neuron two is only one column forward of the BMU, so $\mathbf{x1}$ will be 1. This gives us the following calculation for v .

$$v = ((0 - 0) / (2 * 3)^2) + ((1 - 0) / (2 * 3)^2) = 0.0277 \quad (8.11)$$

Now that v has been calculated we can calculate the Gaussian.

$$\exp(-v) = 0.9726 \quad (8.12)$$

As you can see, a neuron so close to the BMU gets nearly all the training that the BMU received. Other than using the above value for the neighborhood function, the weight calculation for neuron two is the same as neuron one.

8.3 SOM Error Calculation

While training feedforward neural networks we would typically calculate an error number to indicate if training is successful. A SOM is a unsupervised Neural Network. Because it is unsupervised, the error cannot be calculated by normal means. In fact, because it is unsupervised, there really is no error at all.

If there is no known, ideal data then there is nothing to calculate an error against. It would be helpful to see some sort of number to indicate the progression of training.

Many neural network implementations do not even report an error for the SOM. Additionally, many articles and books about SOM's do not provide a way to calculate an error. As a result, there is no standard way to calculate an error for a SOM.

Yet, there is a way to report an error. The error for a SOM can be thought of as the worst (or longest) Euclidean distance of any of the best matching units. This is a number that should decrease as training progresses. This can give you an indication of when your SOM is no longer learning. However, it is only an indication. It is not a true error rate. It is also more of a "distance" than an error percent.

8.4 Chapter Summary

This chapter introduced a new neural network and training method. The Self Organizing Map automatically organizes data into classes. No ideal data is provided; the SOM simply groups the input data into similar classes. The number of grouping classes must remain constant.

Training a SOM requires you to determine the best matching unit (BMU). The BMU is the neuron whose weights most closely match the input that we are training. The weights of the BMU are adjusted to more closely match the input that we are training. Additionally, weights near the BMU are also adjusted.

In this next chapter we will look at another neural network architecture. In this chapter we used. This neural network type is called a RBF Neural Network. RBF networks are similar to feedforward neural networks in that they can have multiple output neurons that produce floating point numbers. RBF networks do not classify like SOM's do. Yet RBF networks make use of the same Radial Basis Functions that a SOM uses.

Chapter 9

Normalization

- SOM Structure
- Training a SOM
- Neighborhood Functions
- SOM Error Calculation

The input to neural networks in this book has so far been mostly in either the range of -1 to +1 or 0 to +1. The XOR operator, for example, had input of either 0 or 1. In Chapter 9 we normalized colors to a range between -1 and +1. Neural Networks generally require their input and output to be in either of these two ranges. The problem is that most real-world data is not in this range.

To account for this we must convert real-world data into one of these ranges before that data is fed to the neural network. Additionally, we must convert data coming out of the neural network back to its normal numeric range. The process of converting real-world data to a specific range is called normalization. This process is called normalization.

You should normalize both input and ideal data. Because you are normalizing the ideal data, your neural network will be trained to return normalized data from the output neurons. Generally, you will want to convert the normalized output back to real-world numbers. The reverse of normalization is called denormalization. Both normalization and denormalization will be covered in this chapter.

9.1 Simple Normalization and Denormalization

We will begin by looking at a very simple means of normalization. This method is called reciprocal normalization. This normalization method supports both normalization and denormalization. However, reciprocal normalization is limited in that you cannot specify the range to normalize into. Reciprocal normalization is always normalizing to a number in the range between -1 and 1.

9.1.1 Reciprocal Normalization

Reciprocal normalization is very easy to implement. It requires no analysis of the data. As you will see with range normalization the entire data set must be analyzed prior to normalization. Equation 9.1 shows how to use reciprocal normalization.

$$f(x) = \frac{1}{x} \quad (9.1)$$

To see Equation 9.1 in use consider normalizing the number five.

$$f(5.0) = 1.0/5.0 = 0.2$$

As you can see, the number five has been normalized to 0.2. This is useful, as 5.0 is outside of the range of 0 to 1.

9.1.2 Reciprocal Denormalization

You will also likely need to denormalize the output from a neural network. It is very easy to denormalize a number which has been normalized reciprocally. This can be done with Equation 9.2.

$$f(x) = x^{-1} \quad (9.2)$$

To see Equation 9.2 in use consider denormalizing the number 0.2.

$$f(0.2) = (0.2)^{-1} = 5.0$$

As you can see, we now completed a round trip. We normalized 5.0 to 0.2, and then denormalized 0.2 back to 5.0.

9.2 Range Normalization and Denormalization

Range normalization is more advanced than simple reciprocal normalization. With range normalization you are allowed to specify the range that you will normalize into. This allows you to more effectively utilize the range offered by your activation function. In general, if you are using the sigmoid activation function, you should use a normalization range of 0 to 1. If you are using the hyperbolic tangent activation function, you should use the range from -1 to 1.

Range normalization must know the high and low values for the data that it is to normalize. To do this the training set must typically be analyzed to obtain a high and low value. If you are going to use additional data, not contained in the training set, your data range must include this data as well. Once you choose the input data range data outside of this range will not produce good results. Because of this it may be beneficial to enlarge the analyzed range by some amount to account for the true range.

9.2.1 Range Normalization

Once you have obtained the high and low values of your data set you are ready to normalize. Equation 9.3 can be used to perform the range normalization.

$$f(x) = \frac{(x-d_L)(n_H-n_L)}{(d_H-d_L)} + n_L \quad (9.3)$$

The above equation uses several constants. The constant **d** represents the high and low of the data to be normalized. The constant **n** represents the high and low that we are to normalize into.

Consider if we were to normalize into the range of -1 to 1. The data range is 10 to 30. The number we would like to normalize is 12. Plugging in numbers we have.

$$(((12-10)*(1-(-1)))/(30-10))+(-1) = -0.8$$

As you can see, the number 12 has been normalized to -0.8.

9.2.2 Range Denormalization

We will now see how to denormalize a ranged number. Equation 9.3 will perform this denormalization.

$$f(x) = \frac{(d_L - d_H)x - (n_H \cdot d_L) + d_H \cdot n_L}{(n_L - n_H)} \quad (9.4)$$

Plugging in the numbers to denormalize -0.8 we are left with the following.

$$((10 - 30) * -0.8 - 1 * 10 + 30 * -1) / (-1 - 1) = 12$$

Again we have made a round trip. We normalized 12 to -0.8. We then denormalized -0.8 back to 12.

9.3 Chapter Summary

This chapter covered normalization. Normalization is the process where data is forced to conform to a specific range. The usual range is either -1 to +1, or 0 to 1. The range you will choose is usually dependant on the activation function you are using. This chapter covered two different types of normalization.

Reciprocal normalization is a very simple normalization technique. This normalization technique normalizes numbers to the range -1 to 1. Reciprocal normalization simply takes the reciprocal normalization simply divides the number to normalize by 1.

Range randomization is more complex. However, range randomization allows you to randomize to any range you like. Additionally, range randomization must know the range of the input data. While this does allow you to make use of the entire normalization range, it also means that the entire data set must be analyzed ahead of time.

This chapter completes this book. I hope you have learned more about the lower levels of neural networks; as well as some of the math behind them. If you would like to apply code to these techniques, you may find my books “Introduction to Neural Networks for Java” or “Introduction to Neural Networks for C#” useful. These books focus less on the mathematics and more on how to implement these techniques in Java or C#.