

VICTORIA UNIVERSITY OF WELLINGTON
Te Whare Wānanga o te Ūpoko o te Ika a Māui



School of Engineering and Computer Science
Te Kura Mātai Pūkaha, Pūrorohiko

PO Box 600
Wellington
New Zealand

Tel: +64 4 463 5341
Fax: +64 4 463 5045
Internet: office@ecs.vuw.ac.nz

Logical Neural Networks: Opening the black box

Daniel Thomas Braithwaite

Supervisor: Marcus Frean

Submitted in partial fulfilment of the requirements for
Bachelor of Science with Honours in Computer Science.

Abstract

Neural Networks (NN's) have been shown to be universal function approximators and have been employed to solve a number of non-trivial problems. The downside to these networks is that they are a black box, difficult to interpret the trained NN's weights. Previous research has shown that Logical Neural Networks applied to the MINST data set obtains comparable performance to standard NN's and a more understandable learned representation. We aim to build on this by providing evidence that builds a foundation for LLN's.

1. Introduction

Neural Networks (NN's) perform exceptionally well over a wide range of different problems. However once trained these networks become a black box, near impossible for a human to interpret the learned representation and understand patterns in the features identified by the NN.

Logical Neural Networks (LNN's) are NN's with activation functions which resemble logic gates. The idea being that if we are able to think about the output of neuron as a logical combination of inputs and weights then perhaps we can understand the learned structure of the network.

LNN's, have previously been shown to provide a simpler representation while still performing well [1] (in the context of the MINST data set). Despite the promise of LNN's they have not been deeply explored.

2. The Problem

LNN's are an exciting concept because of their interpretability, this combined with their reasonable performance (when compared to standard NN's) makes them a promising solution to the usual "black boxiness" of NN's. We suggest this problem still hasn't been solved as we have yet to develop a foundation for LNN's. First we ask, can it be demonstrated through experimentation that like standard NN's, LNN's are also universal function approximators? Before we can hope to answer this key question we also must ask whether we have the right parametrization of logical gates as these form the foundation of an LNN.

Regardless of our answer to the first question a next logical step for our investigation is to ask, what types of problems is an LNN best suited to? Which ones can it not only perform well on but also provide a interpretable learned representation. This will help build the foundation of LNN's, giving us a better understanding of how to configure them for different problems, i.e. network layout and which logical neurons to use.

3. Proposed Solution

To begin this project we take a step back to ensure we have the right approach. A network built out of neurons that resemble logic gates we would expect should be able to learn any given boolean formula. This would show we have good paramaterisations for our logical gates and be a fundamental step towards showing LNN's are universal function approximators. First we must address a key issue with our existing setup, namely our choice of "gates" is not a functionally complete set. So far LNN's have been shown to perform with a reasonable accuracy on the MINST data set using layers consisting of only Noisy-OR and Noisy-AND neurons. Consider the problem of representing some known boolean expression using discrete logic circuit, but we only have access to AND and OR gates. These two gates combined does not make a functionally complete set. We propose to implement a number of different sets of functionally complete neurons, such as NAND and NOR. If an LNN consisted of multiple layers of NAND's or NOR's we would expect to find groups of gates learning to act like some of our other gates. One issue we might encounter here is that

there is a set of weights which allows a given network to represent a given boolean expression but the network is unable to learn these weights. The gates we aim to parameterise and implement are AND, OR, NOT, NAND, NOR.

The first 4 weeks of this project would be taken to develop logical gate parameterisations, using the bibliography of the previous research and other resources. By the end of this we will have produced a bibliography and our gate parameterisations.

We would like to demonstrate that we can use our logical neurons to learn any given boolean expression (if we have an arbitrary number of neurons). We will design experiments by generating some boolean expressions and then simply feeding the expression into a LNN. It would be interesting to also feed these expressions into a standard NN for a comparison in performance. We might find that in practice that some boolean expressions can't be learned, in this case we will investigate and provide justification for this outcome. During these experiments we can also identify what training data will yield a LNN which generalizes effectively i.e. can we fix an input variable to be only 0 while training and obtain a network which generalizes to the case where that variable is 1? We expect to spend 7 weeks investigating and running experiments at the conclusion of which we will have evidence allowing us to provide an answer to our question.

Next we address a fundamental question, are LNN's universal function approximators? We design experiments in a similar fashion as before except this time rather than being boolean expressions we will generate integer and real valued functions, similarly comparing performance to standard NN. We allocate 7 weeks for this section, at the conclusion of which we will have evidence to provide an answer to the question posed above.

Now that we will have explored the foundation of LNN's, using what we have discovered we can use our LNN to solve standard benchmark problems and compare against performance of standard NN's. We would use K-Fold Cross Validation, which is commonly used for comparisons like this. It would then be useful to examine the trained LNN's and see if we can interpret how the network has learned to solve the problem. This is a logical place to return to as this really is the motivation behind LNN's. We allocate 5 weeks for this.

An overview of the proposed solution follows

- Produce bibliography and parameterisation of gates (est: 5 weeks)
- Demonstrating LNN's are/are not able to learn any known boolean function (est: 7 weeks)
- Demonstrating LNN's are/are not universal function approximators (est: 7 weeks)
- Evaluate LNN's on benchmark problems (est: 5 weeks)

This only allocates 24 weeks, the left over time is to account for slippages.

4. Evaluating your Solution

A solution to this problem should build the foundation to LNN's. Giving strong evidence through experiments that LNN's can or can't learn any known boolean function and strong evidence showing that LNN's are or are not universal function approximators. This would require an in depth investigation into various parameterisations of logic gates. Following

this it should identify where LNN's are most powerful and where a standard NN should be used instead. Finally a comparison of accuracy and interpretability of LLN's against standard NN's over various benchmark problems.

5. Resource Requirements

I do not for see any requirements I don't already have access to.

Bibliography

- [1] LAURENSEN, J. Learning logical activations in neural networks, 2016.