



UNIVERSITY OF EXETER

# Can statistics help us to understand deep learning?

Johannes Smit

May 2019

# Abstract

All theses/dissertations must include an abstract of approximately 300 words bound in with each copy and placed so as to follow the title page.

# Acknowledgements

Hi mum.

# Contents

<b>Abstract</b>	<b>i</b>
<b>Acknowledgements</b>	<b>ii</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Literature Review</b>	<b>2</b>
<b>3 Results</b>	<b>3</b>
3.1 Training the neural network . . . . .	4
3.1.1 Ordering of the datapoints . . . . .	4
3.2 Regression . . . . .	5
3.2.1 Linear regression . . . . .	5
3.2.2 Stepwise regression . . . . .	5
3.2.3 LASSO . . . . .	6
3.3 Gaussian process regression . . . . .	6
<b>4 Conclusion</b>	<b>7</b>
<b>A Code</b>	<b>8</b>

# Chapter 1

## Introduction

This is chapter one.

- What is machine learning?
- Machine learning is very complex to human understanding
- Machine learning is widely used in many high stakes scenarios — give examples
- It is important to be able to look inside the ‘black box’ of machine learning — explain how this would be useful in each of the given examples
- We can use statistical methods to try and recover the information in a ML algorithm — give some idea how

# Chapter 2

## Literature Review

This is chapter two.

- Machine Learning
  - Neural Networks
  - Deep Learning
- Stepwise Regression
- LASSO
- Gaussian Processes
- Previous papers on this subject

# Chapter 3

## Results

Plan:

- Introduce a simple example function
- Possibly also use a more complex example?
- Training a neural network
  - Keras/Tensorflow
  - Choosing the right size NN
  - Order of the input points matters
- Fitting a linear regression
- Using stepwise regression
- Using LASSO
- Using GPs

To demonstrate the possibility of using statistical methods to understand the process of deep learning, we use a simple function  $f(x) = x + 5 \sin(x)$ . 256 evenly spread datapoints were taken from this function, and noise following  $\epsilon \sim \mathcal{N}(0, 0.1)$  was applied.

## 3.1 Training the neural network

This function was learnt by a neural network with 8 layers, each with 10 neurons and the  $\tanh(\cdot)$  activation function, except for the final layer which used a linear activation function. The result of this learning is seen in Figure 3.1.



Figure 3.1: The output of the neural network.

### 3.1.1 Ordering of the datapoints

Due to the form of stochastic gradient descent used to train the model, if the order of the datapoints is not randomised, then the optimisation algorithm can more easily get stuck in a local minimum. An example of this is seen in Figure 3.2, where the same neural network has been trained on the example data.





Figure 3.2: The output of the same neural network trained on the same data but ordered differently.

## 3.2 Regression

One method of opening the black box of machine learning is to use regression with many different predictors.

### 3.2.1 Linear regression

In this case,  $x$ ,  $x^2$ ,  $\sin(x)$ ,  $\sin(2x)$ ,  $\sin(x/2)$ ,  $\cos(x)$ ,  $\cos(2x)$  and  $\cos(x/2)$  were used. This produced a complex model with many coefficients close to zero.

### 3.2.2 Stepwise regression

It is then possible to use stepwise regression to reduce the number of parameters in this linear model.

While the relevant estimators  $x$  and  $\sin(x)$  were identified and their coefficients fairly accurately estimated, a few other variables were also identified

as significant. This method is only likely to work with a perfect or near perfect fit with no noise, which is unrealistic for real applications.

### **3.2.3 LASSO**

Alternatively, we can use the Least Absolute Shrinkage and Selection Operator(LASSO) method to select the significant variables from the full model. The LASSO method constrains the sum of the absolute values of the model parameters, regularising the least influential parameters to zero. We then use k-fold cross validation to find the optimal hyperparameter  $\lambda$ .

## **3.3 Gaussian process regression**

# Chapter 4

## Conclusion

Plan:

- How well have each of the attempts worked?
- What could be improved?
- How useful would more research on this topic be?
- What should future research on this topic focus on?

# Appendix A

## Code

Appendix here.