

# GENERIC MASTER THESIS TITLE

by

Geir Tore Ulvik

THESIS

for the degree of

MASTER OF SCIENCE



Faculty of Mathematics and Natural Sciences  
University of Oslo

November 2020



# Abstract

This is the abstract text.





To someone

This is a dedication to my cats.





# Acknowledgements

I acknowledge my acknowledgements.



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Theory</b>	<b>3</b>
2.1	Linear Regression . . . . .	4
2.2	Over- and underfitting . . . . .	5
2.3	The Bias-Variance Tradeoff . . . . .	6
2.4	Regularization . . . . .	7
2.5	Logistic Regression . . . . .	7
2.6	Gradient Descent . . . . .	7
2.6.1	Stochastic Gradient Descent . . . . .	7
2.7	Neural Networks . . . . .	8
2.7.1	Artificial Neurons . . . . .	8
2.7.2	Perceptron . . . . .	9
2.7.3	Backpropagation . . . . .	9
2.8	Assessing the Performance of Models . . . . .	9
2.8.1	Imbalanced Data in Classification . . . . .	9
2.8.2	Confusion Matrix . . . . .	10
2.8.3	Receiver Operating Characteristic . . . . .	10
2.9	Nuclear Science . . . . .	10
2.9.1	Shell Structure . . . . .	10
<b>3</b>	<b>Method</b>	<b>13</b>
3.1	Pretrained network . . . . .	14
<b>4</b>	<b>Results</b>	<b>15</b>
4.1	Classification . . . . .	16
4.1.1	Simulated data . . . . .	16
4.1.2	Experimental data . . . . .	16
4.2	Regression . . . . .	18
4.2.1	Simulated data . . . . .	18
<b>5</b>	<b>Discussion</b>	<b>21</b>

Contents

**6 Conclusion**

**23**

# **Chapter 1**

## **Introduction**

Start your chapter by writing something smart. Then go get coffee.

# Chapter 2

## Theory

## 2.1 Linear Regression

Suppose you have a data set  $\mathcal{L}$  consisting of the data  $\mathbf{X}_{\mathcal{L}} = \{(y_i, x_i), i = 0 \dots n-1\}$ . Each point is associated with a scalar target  $y_i$ , and a vector  $\mathbf{x}$  containing values for  $p$  input features. Assuming the target variable  $y_i$  is linear in the inputs, it can be written as a linear function of the features, given by

$$y_i = w_0x_{i,0} + w_1x_{i,1} + \dots + w_{p-1}x_{i,p-1} + \epsilon_i, \quad (2.1)$$

where  $\mathbf{w} = (w_0, w_1, \dots, w_{p-1})^T$  is a vector of length  $p$  containing unknown values, and  $\epsilon$  are the errors in our estimate. This gives us a system of linear equations, which can be written in matrix form as

$$\mathbf{y} = \mathbf{X}\mathbf{w} + \epsilon, \quad (2.2)$$

where

$$\mathbf{X} = \begin{bmatrix} x_{0,0} & x_{0,1} & x_{0,2} & \dots & x_{0,p-1} \\ x_{1,0} & x_{1,1} & x_{1,2} & \dots & x_{1,p-1} \\ x_{2,0} & x_{2,1} & x_{2,2} & \dots & x_{2,p-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ x_{n-1,0} & x_{n-1,1} & x_{n-1,2} & \dots & x_{n-1,p-1} \end{bmatrix} \quad (2.3)$$

The unknown values  $\mathbf{w}$  are commonly referred to as *weights* in machine learning literature. To find the best possible weights  $\mathbf{w}$  we want a suitable quantity to optimize - a **cost function**,  $\mathcal{C}$  (also referred to as an **objective function**). An example of such a function is the squared error - or the Euclidian vector norm, defined as

$$L_2(\mathbf{x}) = \|\mathbf{x}\|_2 = \left( \sum x_i^2 \right)^{\frac{1}{2}}. \quad (2.4)$$

From this we define the cost function

$$\mathcal{C} = \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2. \quad (2.5)$$

In machine learning, it is most common to cast the optimization as a minimization problem ("minimize the cost"). Our task is then to find an approximation

$$\hat{\mathbf{y}} = \mathbf{X}\mathbf{w} \quad (2.6)$$

which minimizes this cost function. To find the minimum we need a differentiation. To simplify that process, we rewrite the cost function on matrix form

$$\begin{aligned} \mathcal{C} &= \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2, \\ \mathcal{C} &= (\hat{\mathbf{y}} - \mathbf{X}\mathbf{w})^T (\hat{\mathbf{y}} - \mathbf{X}\mathbf{w}). \end{aligned}$$



To minimize we take the derivative with respect to the weights  $w$ , and find the minima by setting the derivative equal to zero

$$\nabla_w \mathcal{C} = \nabla_w (\hat{y} - Xw)^T (\hat{y} - Xw), \quad (2.7)$$

$$= -2X^T \hat{y} + 2X^T Xw, \quad (2.8)$$

$$0 = -2X^T \hat{y} + 2X^T Xw, \quad (2.9)$$

$$X^T \hat{y} = X^T Xw, \quad (2.10)$$

$$w = (X^T X)^{-1} X^T \hat{y}. \quad (2.11)$$

This requires matrix  $X^T X$  to be invertible to get the solution [1].

The residuals  $\epsilon$  are given by

$$\epsilon = y - \hat{y} = y - Xw,$$

and with

$$X^T (y - Xw) = 0,$$

we have

$$X^T \epsilon = X^T (y - Xw) = 0,$$

meaning that the solution for  $w$  is the one which minimizes the residuals. This method of regression is known as Ordinary Least Squares.

## 2.2 Over- and underfitting

In machine learning, when fitting a model to a data set the goal is nearly always to predict values or classify samples from regions of data the model has not seen. This is not a simple task, especially taking into consideration that data is rarely, if ever, noiseless. When extrapolating to unseen regions we must take steps to ensure the model complexity is appropriate - we want it to fit the signal, not the noise. First off - what do the terms "overfit" and "underfit" mean? An overfit model will typically perform well during the fitting procedure, but when presented with data outside the fitted region its performance decreases considerably. An underfit model lacks the expressive power to capture core signal variations in the data. Mehta et. al [2] demonstrates this concept through polynomial regression.

In machine learning literature and practice, you will encounter the concept of splitting the available data into two - training data and test data. We fit, or 'train' the model on the training data, and then assess the performance of the model on the test data. This practice lets us evaluate whether the model is overfitting to unseen data by comparing performance on the training data and test data.

## 2.3 The Bias-Variance Tradeoff

Considering the same dataset  $\mathcal{L}$  consisting of the data, let us assume that the true data is generated from a noisy model

$$\mathbf{y} = f(\mathbf{x}) + \epsilon,$$

where  $\epsilon$  is normally distributed with mean zero and standard deviation  $\sigma^2$ .

In our derivation of the ordinary least squares method we defined an approximation (2.6) to the function  $f$  in terms of the weights  $\mathbf{w}$  and the input matrix  $\mathbf{X}$  which together define our model, that is  $\hat{\mathbf{y}} = \mathbf{X}\mathbf{w}$ .

Thereafter we found the parameters  $\mathbf{w}$  by optimizing the means squared error via the so-called cost function

$$\mathcal{C}(\mathbf{X}, \mathbf{w}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2 = \mathbb{E} \left[ (\mathbf{y} - \hat{\mathbf{y}})^2 \right].$$

We can rewrite this as

$$\mathbb{E} \left[ (\mathbf{y} - \hat{\mathbf{y}})^2 \right] = \frac{1}{n} \sum_i (f_i - \mathbb{E} [\hat{\mathbf{y}}])^2 + \frac{1}{n} \sum_i (\hat{y}_i - \mathbb{E} [\hat{\mathbf{y}}])^2 + \sigma^2.$$

The three terms represent the square of the bias of the learning method, which can be thought of as the error caused by the simplifying assumptions built into the method. The second term represents the variance of the chosen model and finally the last terms is variance of the error  $\epsilon$ .

For a derivation of this equation, we refer to the work of Mehta et. al [2]. The authors also provide a list of universal lessons to keep in mind, which we quote here:

- Fitting is not predicting. Fitting existing data well is fundamentally different from making predictions about new data.
- Using a complex model can result in overfitting. Increasing a model's complexity (i.e number of fitting parameters) will usually yield better results on the training data. However when the training data size is small and the data are noisy, this results in *overfitting* and can substantially degrade the predictive performance of the model.
- For complex datasets and small training sets, simple models can be better at prediction complex models due to the bias-variance tradeoff. It takes less data to train a simple model than a complex one. Therefore, even though the correct model is guaranteed to have better predictive performance for an infinite amount of training data (less bias), the training errors stemming from finite-size sampling (variance) can cause simpler models to outperform the more complex model when sampling is limited.

## 2.4 Regularization

With the computing resources available today, increasing model complexity to deal with underfitting is usually a simple task. However, this computational freedom has led to overfitting being the common challenge to overcome.

## 2.5 Logistic Regression

Differently to linear regression, classification problems are concerned with outcomes taking the form of discrete variables. For a specific physical problem, we'd like to identify its state, say whether it is an ordered or disordered system. (cite?) One of the most basic examples of a classifier algorithm is logistic regression.

For a logistic regressor to improve it needs a way to track its performance. This is the purpose of a cost function. Essentially, the cost function says something about how wrong the model is in classifying the input. The objective in machine learning, and logistic regression, is then to minimize this error.

The cost function used in this project is called the **cross-entropy**, or the 'negative log likelihood', and takes the form

$$\mathcal{C}(w) = - \sum_{i=1}^n (y_i(w_0 + w_1 x_i) - \log(1 + \exp(w_0 + w_1 x_i))) \quad (2.12)$$

## 2.6 Gradient Descent

Minimizing the cost function is done using Gradient Descent. The gist of it is that to optimize the weights or coefficients, and biases to minimize the cost function, we can change their values to

$$\frac{\partial \mathcal{C}(w)}{\partial w} = -X^T (\hat{y} - \hat{p}) \quad (2.13)$$

### 2.6.1 Stochastic Gradient Descent

The stochastic gradient descent method address some of the shortcomings of the normal gradient descent method. The gradient descent method is for instance sensitive to the choice of learning rate (cite?).

The underlying idea of stochastic gradient descent comes from observing that the cost function we want to minimize, almost always can be written as a sum over  $n$  data points. (cite?). Which gives

$$C(w) = \sum_{i=1}^n c_i(\mathbf{x}_i w) \quad (2.14)$$

(cite?)

This means that we also can find the gradient as a sum over  $i$  gradients as follows:

$$\Delta_w C(w) = \sum_i^n \Delta_w c_i(\mathbf{x}_i w) \quad (2.15)$$

(cite?)

Randomness is included by only taking the gradient on a subset of data.

## 2.7 Neural Networks

Artificial neural networks (ANN) are computational systems that can learn to perform tasks by considering examples, generally without being programmed with any task-specific rules. It is supposed to mimic a biological system, wherein neurons interact by sending signals in the form of mathematical functions between layers. All layers can contain an arbitrary number of neurons, and each connection is represented by a weight variable.

In natural science, ANNs have already found numerous applications. In statistical physics, they have been applied to detect phase transitions in 2D Ising and Potts models, lattice gauge theories, and different phases of polymers, or solving the Navier-Stokes equation in weather forecasting. Deep learning has also found interesting applications in quantum physics. Various quantum phase transitions can be detected and studied using ANNs, topological phases, and even non-equilibrium many-body localization. Representing quantum states as ANNs quantum state tomography are among some of the impressive achievements to reveal the potential of DNNs to facilitate the study of quantum systems.

In quantum information theory, it has been shown that one can perform gate decompositions with the help of neural.

The applications are not limited to the natural sciences. There is a plethora of applications in essentially all disciplines, from the humanities to life science and medicine.

### 2.7.1 Artificial Neurons

The field of artificial neural networks has a long history of development, and is closely connected with the advancement of computer science and computers in general. A model of artificial neurons was first developed by McCulloch

and Pitts in 1943 [3] to study signal processing in the brain and has later been refined by others. The general idea is to mimic neural networks in the human brain, which is composed of billions of neurons that communicate with each other by sending electrical signals. Each neuron accumulates its incoming signals, which must exceed an activation threshold to yield an output. If the threshold is not overcome, the neuron remains inactive, i.e. has zero output.

This behaviour has inspired a simple mathematical model for an artificial neuron.

$$y = f\left(\sum_{i=1}^n w_i x_i\right) = f(u) \text{labelartificialNeuron} \quad (2.16)$$

!et Here, the output  $y$  of the neuron is the value of its activation function, which have as input a weighted sum of signals  $x_1, \dots, x_n$  received by  $n$  other neurons.

### 2.7.2 Perceptron

#### Multilayer Perceptron

### 2.7.3 Backpropagation

## 2.8 Assessing the Performance of Models

If classification accuracy is not enough to gauge whether a model is performing well, or well in the desired way, alternative way to measure performance must be explored. For cases of imbalanced data there are a few widely used methods that reveal information about the model that the simple accuracy metric can't.

### 2.8.1 Imbalanced Data in Classification

The physical world is rarely balanced.

A common challenge in classification is imbalanced data, in which a large amount of the labeled data belongs to just one or a few of the classes. For binary classification, if 90% of the data belongs to one of the classes, then the classifier is likely to end up placing every single input in that class, as it will bring its accuracy to 90%. Technically, this accuracy is correct, but it's not very useful since the decision isn't at all affected by the features of the input. Accuracy alone isn't a good enough measure of performance to reveal this.

## 2.8.2 Confusion Matrix

A confusion matrix is an  $n$  by  $n$  matrix containing correct classifications on the diagonal, and false positives and negatives in the off-diagonal elements. An example of such a matrix could be the following table: In the table above

	True Cat	True Dog	True Rabbit
Predicted Cat	<b>5</b>	2	0
Predicted Dog	3	<b>3</b>	2
Predicted Rabbit	0	1	<b>11</b>

**Table 2.1:** Confusion matrix for an example classification where the classes are Cat, Dog and Rabbit. Correct classifications in bold.

(2.1), the diagonal elements  $i = j$  are the correct classifications, while the other elements correspond to cases where the model predicted class  $i$  but should've predicted class  $j$ . The confusion matrix thus gives information about false positives and false negatives, in addition to classification accuracy. This is very useful in cases where for example false positives can be readily ignored or filtered later, but false negatives may have severe consequences. An example of this could be detection of cancer, in which a false positive can be ruled out from further testing, while a false negative may lead to a patient being sent home when actually needing help. For a more in-depth look at confusion matrices see [4].

## 2.8.3 Receiver Operating Characteristic

The Receiver Operating Characteristic (ROC) is a widely used measure of a classifier's performance. The performance is measured as the effect of the true positive rate (TPR) and the false positive rate (FPR) as a function of thresholding the positive class. To evaluate the ROC curve for a model, traditionally the Area Under the Curve (AUC) is used, which ranges from 0 (an ideal "opposite" classifier) to 1.0 (an ideal classifier) with 0.5 indicating a random choice classifier. For a thorough explanation of ROC curves and the underlying concepts, see [4].

## 2.9 Nuclear Science

### 2.9.1 Shell Structure

- Comprehensive and predictive model of atomic nuclei
  - Evolving structure of atomic nuclei as a function of protons and neutrons from first principles

- Understanding the origin of the elements
  - Explosive nucleosynthesis
- Use of atomic nuclei to test fundamental symmetries
- Search for new applications of isotopes and solutions to societal problems

### Nomenclature

A nucleus  $Y$  has  $Z$  protons and  $N$  neutrons with a mass of  $A = Z + N$ . This is written as  ${}^A_ZY_N$ . For a given nucleus there may be several

- Isotopes - nuclei with the same number of protons, but varying number of neutrons
  - ${}^{66}\text{Ni}$ ,  ${}^{67}\text{Ni}$ ,  ${}^{68}\text{Ni}$ ,  ${}^{69}\text{Ni}$ ,  ${}^{70}\text{Ni}$
- Isotones - nuclei with the same number of neutrons, but varying number of protons
  - ${}^{70}\text{Zn}$ ,  ${}^{69}\text{Cu}$ ,  ${}^{68}\text{Ni}$ ,  ${}^{67}\text{Co}$ ,  ${}^{66}\text{Fe}$
- Isobars - nuclei with the same number of nucleons  $A$ 
  - ${}^{68}\text{Fe}$ ,  ${}^{68}\text{Co}$ ,  ${}^{68}\text{Ni}$ ,  ${}^{68}\text{Cu}$ ,  ${}^{68}\text{Zn}$





# **Chapter 3**

## **Method**

## 3.1 Pretrained network

# **Chapter 4**

## **Results**

## 4.1 Classification

Our primary objective is to separate experimental data into two possible categories. Supervised learning typically leverages large amounts of labeled data to learn representations of the classes present in the inputs. A challenge in Physics is that experiments generate enormous amounts of data, but it is not labeled. Hand-labeling data is time-consuming and cost-ineffective. A possible solution to this challenge is to train models on simulated data. The idea being that simulated and experimental data are similar enough that the learned patterns from simulated data are, to some degree, transferrable to experimental data. The simulated data contains two classes of events - a single electron and a double electron.

### 4.1.1 Simulated data

In table 4.1 the performance of each model trained is reported using the f1-score. The model architecture for each model is described in (ref appendix). As a benchmark, we are including a state of the art pretrained network ([VGG HERE]) applied to the data using the approach described in 3.1. To give a broader picture of the performance we also include the calculated confusion matrix values in table 4.2. Even if f1-scores are similar, there may still be differences in either of the confusion matrix values.

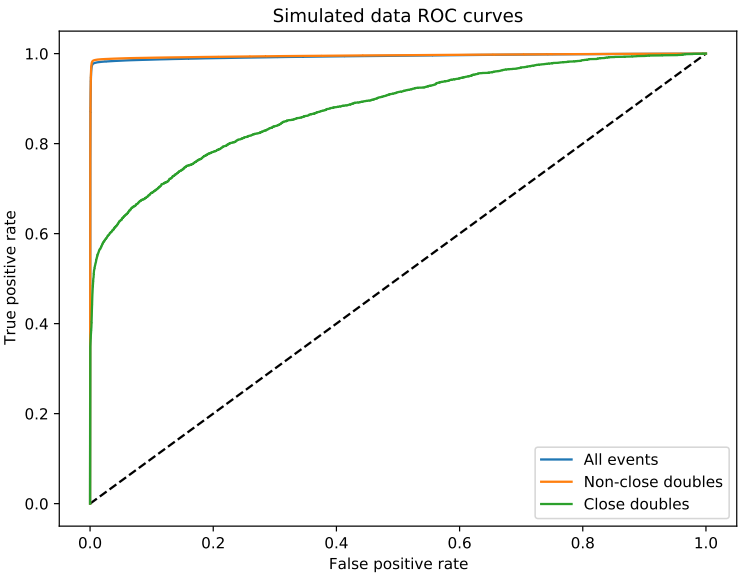
### 4.1.2 Experimental data

**Table 4.1:** Mean F1-scores for classification of simulated data using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

Logistic	Dense	Convolutional	Pretrained VGG16
0.734 $\pm 7.727 \times 10^{-3}$	0.907 $\pm 1.329 \times 10^{-2}$	0.959 $\pm 6.286 \times 10^{-3}$	0.894 $\pm 1.591 \times 10^{-2}$

**Table 4.2:** Mean confusion matrix values for classification of simulated data using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

	TN	FP	FN	TP
Logistic	$1.28 \times 10^5$ $\pm 1.546 \times 10^4$	$6.16 \times 10^4$ $\pm 1.546 \times 10^4$	$4.39 \times 10^4$ $\pm 1.118 \times 10^4$	$1.46 \times 10^5$ $\pm 1.118 \times 10^4$
Dense	$1.72 \times 10^5$ $\pm 1.125 \times 10^4$	$1.85 \times 10^4$ $\pm 1.125 \times 10^4$	$1.73 \times 10^4$ $\pm 5.205 \times 10^3$	$1.73 \times 10^5$ $\pm 5.205 \times 10^3$
Convolutional	$1.89 \times 10^5$ $\pm 7.287 \times 10^2$	$1.39 \times 10^3$ $\pm 7.290 \times 10^2$	$1.37 \times 10^4$ $\pm 2.692 \times 10^3$	$1.76 \times 10^5$ $\pm 2.692 \times 10^3$
Pretrained VGG16	$1.79 \times 10^5$ $\pm 9.079 \times 10^3$	$1.15 \times 10^4$ $\pm 9.080 \times 10^3$	$2.71 \times 10^4$ $\pm 9.409 \times 10^3$	$1.63 \times 10^5$ $\pm 9.408 \times 10^3$



**Figure 4.1:** generic text

## 4.2 Regression

On data already classified, we attempt to predict the energy of the events and the position of origin for the electrons. Because there is a travel distance between the ejection site and the scintillator array, the positions aren't necessarily the locations of the highest-intensity pixels in the detector images.

### 4.2.1 Simulated data

Position - Single electron

Energy - Single electron

Position - Double electron

Energy - Double electron

**Table 4.3:** Mean R2-scores for regression of positions of origin, on single events in simulated data, using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

Linear	Dense	Convolutional	Pretrained VGG16
0.801 $\pm 3.664 \times 10^{-3}$	0.99 $\pm 8.185 \times 10^{-4}$	0.997 $\pm 2.401 \times 10^{-4}$	0.891 $\pm 6.864 \times 10^{-3}$

**Table 4.4:** Mean R2-scores for regression of energy values, on single events in simulated data, using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

Linear	Dense	Convolutional	Pretrained VGG16
0.926 $\pm 3.691 \times 10^{-2}$	0.931 $\pm 3.364 \times 10^{-2}$	0.936 $\pm 3.282 \times 10^{-2}$	0.935 $\pm 1.945 \times 10^{-2}$

**Table 4.5:** Mean R2-scores for regression of positions of origin, on double events in simulated data, using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

Linear	Dense	Convolutional	Pretrained VGG16
0.364 $\pm 5.796 \times 10^{-3}$	0.471 $\pm 1.809 \times 10^{-3}$	0.473 $\pm 2.394 \times 10^{-3}$	0.357 $\pm 1.079 \times 10^{-2}$

**Table 4.6:** Mean R2-scores for regression of energy values, on double events in simulated data, using multiple models. Error estimates are the standard deviation in results from k-fold cross-validation with  $K = 5$  folds.

Linear	Dense	Convolutional	Pretrained VGG16
0.415 $\pm 6.501 \times 10^{-2}$	0.416 $\pm 6.632 \times 10^{-2}$	0.428 $\pm 5.027 \times 10^{-2}$	0.404 $\pm 5.308 \times 10^{-2}$





# **Chapter 5**

## **Discussion**



## **Chapter 6**

## **Conclusion**



# Bibliography

- [1] G. James, D. Witten, T. Hastie, and R. Tibshirani. *An introduction to Statistical Learning*, volume 7. 2000. ISBN 978-1-4614-7137-0. doi: 10.1007/978-1-4614-7138-7.
- [2] P. Mehta, M. Bukov, C. H. Wang, et al. A high-bias, low-variance introduction to Machine Learning for physicists. *Physics Reports*, 810:1–124, 2019. ISSN 03701573. doi: 10.1016/j.physrep.2019.03.001.
- [3] W. S. McCulloch and W. Pitts. A logical calculus of the ideas immanent in nervous activity. *The Bulletin of Mathematical Biophysics*, 5(4):115–133, dec 1943. ISSN 00074985. doi: 10.1007/BF02478259. URL <https://link.springer.com/article/10.1007/BF02478259>.
- [4] T. Fawcett. An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8):861–874, 2006. ISSN 01678655. doi: 10.1016/j.patrec.2005.10.010.