

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

1	Introduction	1
2	Theory	3
2.1	Supervised Learning	4
2.1.1	Linear Regression	4
2.1.2	Classification	5
2.1.3	Cost Functions	6
2.1.4	Gradient Descent	6
2.2	Neural Networks	7
2.2.1	Perceptron	7
2.2.2	Backpropagation	7
2.3	Assessing the Performance of Models	7
2.3.1	Imbalanced Data in Classification	7
2.3.2	Confusion Matrix	7
2.3.3	Receiver Operating Characteristic	8
3	Results	9
4	Discussion	11
5	Conclusion	13

Chapter 1

Introduction

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

Chapter 2

Theory

2.1 Supervised Learning

Prediction and classification. From standard regression methods for predicting a continuous variable, to classifying with logistic regression and neural networks.

2.1.1 Linear Regression

Ordinary Least Squares

Suppose you have a data set consisting of N data points. Each point is associated with a scalar target y_i , and a vector 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

$$y_i = \beta_0 x_{i,0} + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1} + \varepsilon_i, \quad (2.1)$$

where $\beta = (\beta_0, \beta_1, \dots, \beta_{p-1})^T$ is a vector of length p containing unknown values, and ε are the variables, random or otherwise, that influence the target values, but are unaccounted for in the features x . The above equation (2.1) can be written on vector form as

$$y_i = x_i^T \beta + \varepsilon_i, \quad (2.2)$$

and further, in matrix notation

$$\hat{y} = \mathbf{X}\hat{\beta} + \hat{\varepsilon}. \quad (2.3)$$

In order to find the best possible coefficients β we want suitable quantity to optimize - a **cost function** \mathcal{C} . One example is the average squared error, which we want to minimize. This function is defined as

$$\mathcal{C}(\beta) = \frac{1}{N} \sum_{i=0}^{N-1} (y_i - f(x_i))^2 \quad (2.4)$$

This can be rewritten in matrix notation

$$\mathcal{C}(\beta) = (\hat{y} - \mathbf{X}\beta)^T (\hat{y} - \mathbf{X}\beta) \quad (2.5)$$

As we want to minimize the residual sum of squares, we require

$$\frac{\partial \mathcal{C}(\hat{\beta})}{\partial \hat{\beta}} = 0 = \mathbf{X}^T (\hat{y} - \mathbf{X}\hat{\beta}) \quad (2.6)$$

This can be rewritten as

$$\hat{X}^T \hat{y} = \hat{X}^T \hat{X} \hat{\beta} \quad (2.7)$$

We then assume that the matrix $\hat{X}^T \hat{X}$ is invertible to get the solution. [James et al.(2000)Jam

Ridge Regression

OLS will not give a solution if the design matrix \hat{X} is singular or near singular, as it depends on $\hat{X}^T \hat{X}$ being invertible. (cite?)

Ridge regression is known as a shrinkage method. It shrinks the regression coefficients $\hat{\beta}$ by imposing a penalty on their size (cite?). We are in other words adding a diagonal component to the matrix to invert. From OLS we therefore make the following change

$$\hat{X}^T \hat{X} \rightarrow \hat{X}^T \hat{X} + \lambda \hat{I} \quad (2.8)$$

where I is the identity matrix. (cite?). This gives

$$RSS(\lambda) = (\hat{y} - \hat{X}\hat{\beta})^T (\hat{y} - \hat{X}\hat{\beta}) - \lambda \hat{\beta}^T \hat{\beta} \quad (2.9)$$

and

$$\hat{\beta}^{ridge} = (\hat{X}^T \hat{X} + \lambda \hat{I})^{-1} \hat{X}^T \hat{y} \quad (2.10)$$

$\lambda > 0$ is a complexity parameter that controls the amount of shrinkage. The larger the value of λ the greater amount of shrinkage. [James et al.(2000)]James, Witten, Hastie, and

This can also be written in a non-vector format which makes explicit the size constraint on the parameters. It also makes it easier to see the difference from Lasso regression:

$$\hat{\beta}^{ridge} = \operatorname{argmin}_{\beta} \sum_{i=1}^N \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2 \quad (2.11)$$

subject to

$$\sum_{j=1}^p \beta_j^2 \leq t \quad (2.12)$$

There is a one to one correspondence between λ and t . [James et al.(2000)]James, Witten, Hastie, and

2.1.2 Classification

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.

2.1.3 Cost Functions

In order for a logistic regressor to improve it needs a way to track how it's performing. 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

$$\mathcal{C}(\hat{\beta}) = - \sum_{i=1}^n (y_i(\beta_0 + \beta_1 x_i) - \log(1 + \exp(\beta_0 + \beta_1 x_i))) \quad (2.13)$$

2.1.4 Gradient Descent

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

$$\frac{\partial \mathcal{C}(\hat{\beta})}{\partial \hat{\beta}} = -\hat{X}^T (\hat{y} - \hat{p}) \quad (2.14)$$

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

$$C(\beta) = \sum_{i=1}^n c_i(\mathbf{x}_i \beta) \quad (2.15)$$

(cite?)

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

$$\Delta_{\beta} C(\beta) = \sum_i^n \Delta_{\beta} c_i(\mathbf{x}_i \beta) \quad (2.16)$$

(cite?)

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

2.2 Neural Networks

2.2.1 Perceptron

Multilayer Perceptron

2.2.2 Backpropagation

2.3 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.

2.3.1 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.

2.3.2 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	5	2	0
Predicted Dog	3	3	2
Predicted Rabbit	0	1	11

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 [Fawcett(2006)].

2.3.3 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 [Fawcett(2006)].

Chapter 3

Results

Chapter 4

Discussion

Chapter 5

Conclusion

Bibliography

[James et al.(2000)]James, Witten, Hastie, and Tibshirani] 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.

[Fawcett(2006)] 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.