

# Towards impressive titles

Tobias Axelsson

29 mars 2018

# Acknowledgements

I am a student blalsadf

## **Abstract**

asdasd

## **Sammanfattning**

asdadasd

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Background . . . . .	3
1.2	Objective . . . . .	4
1.3	Delimitations . . . . .	4
1.4	Thesis structure . . . . .	4
<b>2</b>	<b>Literature Review</b>	<b>5</b>
2.1	Machine learning basics . . . . .	5
2.2	Supervised learning . . . . .	6
2.2.1	Classification predictive modeling . . . . .	6
2.2.2	Regression predictive modeling . . . . .	7
2.3	Generalization . . . . .	8
2.3.1	Cross-validation . . . . .	8
2.3.2	Regularization . . . . .	9
2.4	Linear and non-linear models . . . . .	9
2.5	Linear classification algorithms . . . . .	9
2.5.1	Decision trees . . . . .	9
2.6	Linear regression algorithms . . . . .	9
2.6.1	Multiple linear regression . . . . .	9
2.7	Non-linear classification and regression algorithms . . . . .	9
2.7.1	Neural network . . . . .	9
2.8	Imbalanced data . . . . .	9
2.9	Multicollinearity . . . . .	9
2.10	Dataset split . . . . .	9
2.11	Data pre-processing . . . . .	9
<b>3</b>	<b>Method</b>	<b>10</b>
3.1	Overview . . . . .	10
3.2	Research approach . . . . .	10
3.3	Research strategy . . . . .	10
3.4	Tools . . . . .	10

<b>4</b>	<b>Implementation and results</b>	<b>11</b>
4.1	Data collection . . . . .	11
4.2	Neural network . . . . .	11
4.2.1	First iteration . . . . .	11
4.2.2	Second iteration . . . . .	11
<b>5</b>	<b>Analysis</b>	<b>12</b>
5.1	Neural network . . . . .	12
<b>6</b>	<b>Conclusions and recommendations</b>	<b>13</b>
6.1	Conclusions . . . . .	13
6.2	Recommendations . . . . .	13
<b>7</b>	<b>Discussion</b>	<b>14</b>
7.1	Thesis process . . . . .	14
7.2	Validity and reliability . . . . .	14
7.3	Future work . . . . .	14

# Chapter 1

## Introduction

*The chapter starts with a background describing why road condition monitoring is important and who Trafikverket are, how road condition data is collected today and why the technology behind it needs improvement. An objective for the project is defined followed by its delimitations. Lastly, a thesis structure is presented to simplify navigation through different parts of the project.*

### 1.1 Background

Living in cold areas of the world usually means work for individual people, municipalities and companies in trying to maintain a non-winter-like infrastructure. This of course, also involves winter road maintenance. Salting and plowing roads is an investment in not only saving lives, but also in lowering socio-economic costs: In two scenarios on a road with 2 cm snow and a daily traffic flow of 2000 vehicles, one with a salted and ploughed road taking four hours to drive, and another scenario on the same road without winter maintenance taking five hours to drive. The total socio-economic costs are 3.5% higher in the non-maintained road, mainly due to increased travel time and thus higher accident costs [1].

Despite the socio-economic savings in performing winter road maintenance, it still represents a notable economic cost. Trafikverket, the agency in charge of road state road maintenance in Sweden, reported that winter road maintenance were roughly 18% of the total road maintenance costs in 2013 [2]. Local contractors are hired to carry out the plowing and salting of state roads, with requirements on both ends regarding when to plow, which roads to prioritize etc. Trafikverket has over 800 Road Weather Information Systems (RWIS)(Fig. 1.1) distributed across state roads in Sweden which are used by contractors to carry out winter road maintenance work [3].



Figure 1.1: RWIS Station at sensor site Myggsjön [4].

	Operation	worst-case cost	time complexity	
Table 1.1 shows	Insert $x$ into $l_i$	2	$O(1)$	asdasd
	Update $count_i$	1	$O(1)$	

## 1.2 Objective

The objective is to determine if a road surface temperature sensor can be simulated with prediction models based on historic data from road weather information systems.

## 1.3 Delimitations

The project will focus on supervised learning algorithms. There are also unsupervised learning algorithms that are not covered for the following reasons:

- Performance of unsupervised learning algorithms can be difficult to evaluate [5].
- The data provided for this project by Trafikverket is labeled, which means it is suitable for supervised learning algorithms.

See section 2 for more information on machine learning theory.

## 1.4 Thesis structure

# Chapter 2

## Literature Review

*The chapter gives both general and specific information on theory used for this project. Mathematical statistics, regression and machine learning are covered in the first three sections, providing a general understanding of the field of study. Specific machine learning models are explained in the final three sections of the chapter.*

### 2.1 Machine learning basics

Machine learning as formally defined by Mitchell [6]: "A computer program is said to learn from experience  $E$  with respect to some class of tasks  $T$  and performance measure  $P$  if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ". This means that machine learning algorithms are used to solve a set of problems, measure its performance in doing so and ultimately improve in some way from previous experiences. For example, imagine a program designed to determine if a human face is in a photo or not. Since photos are taken at different distances, angles and faces have different characteristics such as eye color, skin color, distance between eyes and nose shape, implementing this "manually" may prove cumbersome. Instead of programming an algorithm to recognize faces, it can be programmed to *learn to recognize faces*. If the algorithm is allowed to analyze a dataset with thousands of photos of human faces, it could learn to distinguish a human face by recognizing parts of the face such as eyes, nose, mouth and where those parts are most likely placed to one another.

In essence, machine learning algorithms improve/learn in some way from analyzing a dataset. How they learn can be used to broadly categorize machine learning algorithms as either having supervised or unsupervised learning [7]. Supervised learning algorithms processes a labeled dataset while unsupervised learning tries to make sense of an unlabeled dataset [5]. This project does not concern algorithms related to unsupervised learning, as motivated in section 1.3.

## 2.2 Supervised learning

In supervised learning, the learner (algorithm) receives a dataset of labeled observations which is used to predict correct values for unseen data[5]. A database table storing weather-related data could for example have thousands of records (observations) where data in each record belong to certain columns (features) such as wind speed  $w_s$ , wind direction  $w_d$  and time  $t$ . The goal of supervised learning is to build a mapping function (model)

$$y = f_{map}(x) \quad (2.1)$$

such that when new input data is used,  $f_{map}$  is able to predict a correct output value [8]. The model is built from a dataset which is typically split into three parts [9]:

- Training dataset: Used to fit the model.
- Validation dataset: Used to give an unbiased evaluation of a model built from the training dataset and potentially update its hyperparameters. Hyperparameters are model parameters that are used in some learning algorithms. They are usually fixed before the training process begins [10].
- Test dataset: Gives an unbiased evaluation of the final model.

How the dataset "should" be split is brought up in section 2.10.

Supervised learning can be thought of as having a teacher supervising the algorithm. The correct answers are in the training data and the algorithm learns from being corrected by the teacher [8]. Going back to the forementioned example of the weather station to give a brief example of how a supervised machine learning algorithm works: Suppose a training, validation and test dataset is provided and one wishes to predict wind speed  $y = w_s$  based on wind direction and time  $x = [x_1, x_2] = [w_d, t]$ . During the training process, a supervised learning algorithm goes through the training dataset to build a model, as seen in Eq. 2.1, and possibly updated when validated against the validation dataset. Suppose the supervised learning algorithm used is multiple linear regression (see section 2.6.1) and a model is built from the training process:

$$w_s = f_{map}([w_d, t]) = \beta_0 + \beta_1 w_d + \beta_2 t = 4 + 0.2w_d + 1.7t \quad (2.2)$$

The model can then be tested with the test dataset to see how it performs on unseen data.

Estimation of continuous output variables, such as wind speed in the example presented above, is a regression problem. In supervised learning there are also algorithms associated with the problem of classification; how to categorize data [11].

### 2.2.1 Classification predictive modeling

In a classification problem, the computer is asked to place a new observation into one of  $k$  categories (labels),  $k \geq 2$  [7]. The problem of categorizing new email as spam or

not spam is an example of a classification problem. Google claims that their machine learning models can detect spam and phishing messages with 99.9% accuracy in their widely used Gmail application [9].

Another example of a classification problem, one that may well be the first that machine learning novices encounter, is classification of the Iris flower dataset. The dataset consists of 50 observations with four features: length and width of the sepals and petals, in centimeters. Based on this information, the problem is to classify to which of the following labels each observation belongs to [12]:

- Setosa
- Versicolour
- Virginica

How the classification is carried out depends on the algorithm used to build the model. These kind of algorithms are commonly known as classifiers. There are several classifiers that can be used for the Iris dataset, but their performance in doing so may differ. Performance of classifiers are typically measured in accuracy, which is the amount of correct predictions divided by the number of observations in the test dataset [7].

$$\text{accuracy} = \frac{\#\text{correct predictions}}{\#\text{observations}} \quad (2.3)$$

### 2.2.2 Regression predictive modeling

In contrast to classification problems, such as categorizing incoming email as spam or not spam, regression problems are about predicting continuous quantities. Regression models can have either real-valued or discrete input variables [11]. The model in eq. 2.2 is an example of a regression model since the goal is to predict a numerical value for wind speed. The problem could be translated into a classification problem by, for example stating that for given numerical intervals, the wind speed is categorized as being low, medium or high. This kind of conversion is known as discretization but even if the conversion is useful, it can result in surprising and/or poor performance [11]. This is why both classification and regression modeling are covered in this project because they are useful for different kind of predictions.

Performance of regression models can be measured by computing the mean squared error (MSE) of the model on the test dataset.

$$MSE_{test} = \frac{1}{n} \sum_i^n (y'_{test_i} - y_{test_i})^2 \quad (2.4)$$

where  $y'_{test_i}$  are predictions on the test and  $y_{test_i}$  are actual values [7]. It's a measurement of how close each prediction was to its corresponding target value on average.

## 2.3 Generalization

During the training process in supervised learning, a model is typically built based on its training data, and updated in order to reduce its training error. But the fundamental goal of machine learning is to generalize beyond observations in the training dataset since it's unlikely that the same exact observations are found again on unseen data [13]. This is why test datasets are used to measure performance of regression and classification models as seen in sections 2.2.2 and 2.2.1. Both training error, how well a performs on its training data, and generalization error, how well a model performs on unseen data, need to be considered in machine learning [7].

The terminology used to explain how well machine learning models learn and generalizes to new data is overfitting and underfitting [11] which are two central challenges in machine learning [7].

- Overfitting: Random fluctuations and statistical noise is learnt to the extent that it affects the model's ability to generalize [11]. Instead of learning the data trend in the training data, the model "memorizes" it [14]. Non-linear (see section ??) and non-parametric (see section ??) models are prone to overfitting as they are flexible in choosing mathematical functions to fit its training data [11].
- Underfitting: A model that performs poorly on both its training data and on generalization. If a good performance metric is used, underfitting is easy to detect [11].

The goal then, is to select a model that is somewhere between underfitting and overfitting. Underfitting is typically remedied by choosing alternative models, but the most common problem in applied machine learning is how to avoid overfitting [11]. According to Davide [14] the mere awareness of the issue of overfitting along with two powerful tools: cross-validation and regularization, can be enough to overcome the problem.

### 2.3.1 Cross-validation

There are several cross-validation techniques out of which the most commonly used is k-fold cross validation: Once a randomly shuffled dataset is split into a training and test dataset, the training dataset is further split into  $k$  folds. One fold is used as a validation dataset and the remainder for training. The idea is to iterate this process  $k$  times so that every fold has been used once as a validation set, and ultimately average the performance over  $k$  iterations . Using a value of  $k = 10$  is a common choice in practice and in which case it is called 10-fold cross-validation [14]. Furthermore, using  $k = 10$  seems to be optimal when it comes to optimizing run-time for the test, limiting bias (underfitting) and variance (overfitting) [15].

There is a variant of this technique called stratified k-fold cross-validation which is mostly used in classification problems. It can also be applied to regression problems [16] but the results from Breiman and Spector [17] show no improvement from using this technique for regression problems. In stratified k-fold cross-validation the folds are

created in such a way that each fold contains similar proportions of predictor labels as the full dataset. For example, think of the classification problem of classifying email as spam or not spam. If the ratio of spam/not spam is 20%/80% in the original dataset, then the same proportion is attempted to be maintained in each of the  $k$  folds. This technique tends to generate less bias and variance when compared to regular k-fold cross validation [18].

### 2.3.2 Regularization

The second method of overcoming overfitting is regularization. Regularization aims to limit freedom of trained models by adding penalties to its parameters. Regularization can be used in any parametric (see section ??) machine learning algorithm . Three different types of regularization is covered in the following subsections.

2.3.2.1  $L_1$  regularization (Lasso)

2.3.2.2  $L_2$  regularization (Ridge)

2.3.2.3  $L_1/L_2$  regularization (Elastic net)

## 2.4 Linear and non-linear models

### 2.5 Linear classification algorithms

2.5.1 Decision trees

### 2.6 Linear regression algorithms

2.6.1 Multiple linear regression

### 2.7 Non-linear classification and regression algorithms

2.7.1 Neural network

### 2.8 Imbalanced data

### 2.9 Multicollinearity

### 2.10 Dataset split

### 2.11 Data pre-processing

# **Chapter 3**

## **Method**

*The chapter covers strategies and methods used to achieve the objective of the project. Reasons for each choice of method or strategy are motivated and described in the sections, which are ordered chronologically.*

### **3.1 Overview**

### **3.2 Research approach**

### **3.3 Research strategy**

### **3.4 Tools**

## Chapter 4

# Implementation and results

Describe the process of collecting data, training and implementing machine learning algorithms with different methods.

### 4.1 Data collection

### 4.2 Neural network

#### 4.2.1 First iteration

#### 4.2.2 Second iteration

## Chapter 5

# Analysis

Analyze data from the implementation with respect to the objective of the study.

### 5.1 Neural network

## Chapter 6

# Conclusions and recommendations

6.1 Conclusions

6.2 Recommendations

## Chapter 7

# Discussion

### 7.1 Thesis process

### 7.2 Validity and reliability

Validity and reliability of the conclusions. Needed?

### 7.3 Future work

# Bibliography

- [1] A. Arvidsson, “The Winter Model – A new way to calculate socio-economic costs depending on winter maintenance strategy”, *Cold Regions Science and Technology Volume 136, April 2017, Pages 30-36*, 2017.
- [2] Trafikverket, *Trafikverkets årsredovisning 2015*, 2016.
- [3] ——, (2017). Vinterväghållning, [Online]. Available: <https://www.trafikverket.se/resa-och-trafik/underhall-av-vag-och-jarnvag/Sa-skoter-vi-vagar/Vintervaghallning/> (visited on 02/07/2018).
- [4] Pelpet. (2010). Rwis Station at sensor site Myggjön, [Online]. Available: [https://commons.wikimedia.org/wiki/File:Rwis\\_station\\_Myggjon\\_01.JPG](https://commons.wikimedia.org/wiki/File:Rwis_station_Myggjon_01.JPG) (visited on 02/06/2018).
- [5] M. Mehryar, R. Afshin, and T. Ameet, *Foundations of machine learning*. Ser. Adaptive Computation and Machine Learning. MIT Press, Cambridge, MA, 2012, ISBN: 978-0-262-01825-8. [Online]. Available: <http://proxy.lib.ltu.se/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=msn&AN=MR3057769&lang=sv&site=eds-live&scope=site>.
- [6] T. Mitchell, *Machine Learning*. McGraw Hill, 1997.
- [7] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [8] J. Brownlee. (2016). Supervised and Unsupervised Machine Learning Algorithms, [Online]. Available: <https://machinelearningmastery.com/supervised-and-unsupervised-machine-learning-algorithms/> (visited on 02/28/2018).
- [9] F. Lardinois. (2017). Google says its machine learning tech now blocks 99.9% of Gmail spam and phishing messages, [Online]. Available: <https://techcrunch.com/2017/05/31/google-says-its-machine-learning-tech-now-blocks-99-9-of-gmail-spam-and-phishing-messages/> (visited on 03/19/2018).
- [10] X. Amatriain. (). What are hyperparameters in machine learning?, [Online]. Available: <https://www.quora.com/What-are-hyperparameters-in-machine-learning> (visited on 03/20/2018).
- [11] J. Brownlee. (). Difference Between Classification and Regression in Machine Learning, [Online]. Available: <https://machinelearningmastery.com/classification-versus-regression-in-machine-learning/> (visited on 03/22/2018).

- [12] R. Fisher. (). Iris Data Set, [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/iris> (visited on 03/19/2018).
- [13] P. Domingos, “A Few Useful Things to Know About Machine Learning.”, *Communications of the ACM*, vol. 55, no. 10, pp. 78 –87, 2012, ISSN: 00010782. [Online]. Available: <http://proxy.lib.ltu.se/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=buh&AN=82151052&lang=sv&site=eds-live&scope=site>.
- [14] D. Chicco, “Ten quick tips for machine learning in computational biology.”, *Bio-Data Mining*, vol. 10, pp. 1 –17, 2017, ISSN: 17560381. [Online]. Available: <http://proxy.lib.ltu.se/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=aph&AN=126676440&lang=sv&site=eds-live&scope=site>.
- [15] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and regression trees*. Ser. Wadsworth Statistics/Probability Series. Wadsworth Advanced Books and Software, Belmont, CA, 1984, ISBN: 0-534-98053-8. [Online]. Available: <http://proxy.lib.ltu.se/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=msn&AN=MR726392&lang=sv&site=eds-live&scope=site>.
- [16] S. Lowe. (). Stratified Validation Splits for Regression Problems, [Online]. Available: <http://scottclowe.com/2016-03-19-stratified-regression-partitions/> (visited on 03/29/2018).
- [17] “SUBMODEL SELECTION AND EVALUATION IN REGRESSION - THE X-RANDOM CASE.”, *INTERNATIONAL STATISTICAL REVIEW*, vol. 60, no. 3, pp. 291 –319, n.d. ISSN: 03067734. [Online]. Available: <http://proxy.lib.ltu.se/login?url=http://search.ebscohost.com/login.aspx?direct=true&db=eds-wsc&AN=A1992KC6230004&lang=sv&site=eds-live&scope=site>.
- [18] R. Kohavi, “A Study of Cross-Validation and Bootstrap for Accuracy Estimation and Model Selection”, Morgan Kaufmann, 1995, pp. 1137–1143.