

Aarhus University

Decision Support Systems

## PROJECT REPORT

*Author:*

*Supervisor:*

Lasse Lildholdt  
(201507170)  
Stinus Skovgaard  
(201507170)  
Daniel Tøttrup  
(201507170)  
Johan Vasegaard  
(201507170)  
Frederik Madsen  
(201507170)

Christian Fischer Pedersen

16. marts 2020



# Indhold

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Simple linear regression / Multiple linear regression</b>	<b>2</b>
<b>3</b>	<b>Logisitic regression / Linear discriminant analysis</b>	<b>3</b>
<b>4</b>	<b>Cross validation / Bootstrap</b>	<b>4</b>
<b>5</b>	<b>Subset selection</b>	<b>5</b>
<b>6</b>	<b>Shrinkage methods / DImension reduction methods</b>	<b>6</b>

# Figurer

4.1	Logistic regression model illustration . . . . .	4
-----	--	---

## INTRODUCTION

In this project report, the reader will be presented with problem solutions for the course Decision support systems. Throughout the solution the reader will achieve knowledge on several different subject within the main area. Each topic will be presented with the theory along with solutions for appropriate exercise to validate the presented theory.

The main topics which will be handled in this report will be as follows:

- Simple linear regression / Multiple linear regression
- Logisitic regression / Linear discriminant analysis
- Cross validation / Bootstrap
- Subset selection
- Shrinkage methods / DImension reduction methods

# SIMPLE LINEAR REGRESSION / MULTIPLE LINEAR REGRESSION

The simple Linear Regression approach is a quick and simple method for fitting a line through a 2-dimensional dataset. It is assumed that there is an approximately linear relationship between the two dimensions. This can be written mathematically as:

$$Y \approx \beta_0 + \beta_1 * X \quad (2.1)$$

eq. (2.1) can also be seen as "Regressing Y onto X". As an example the dataset Advertising.csv contains sales of a certain product and advertisement money spent on certain media platforms. X represents TV advertising and Y represents sales. It is possible to regress sales onto TV.

In order to do this, we need to calculate the constants  $\beta_0$  and  $\beta_1$  which represent the intercept and slope terms in the linear model.

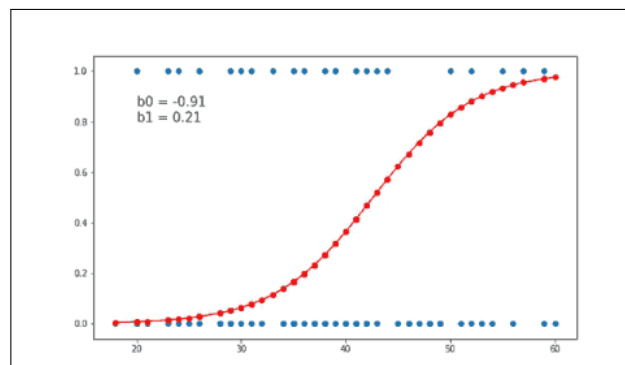
# **LOGISITIC REGRESSION / LINEAR DISCRIMINANT ANALYSIS**

## CROSS VALIDATION / BOOTSTRAP

When we are working with different machine learning models, we need two different data sets. The first set is often referred to as the training data set, and is used to create the model and define the curve seen in 4.1. The other dataset is referred to as the test set and is used to check if the model performs well. The data between the two sets can not be identical because we want to evaluate how the model performs on new data sets.

Because we often don't have well defined test sets available, we need some methods to gain test data. Here we have several different options. Some methods will not provide new data, but instead use the training data to estimate the test error (the performance). Another approach, which is the one discussed here, is a method where you hold out some part of the data, in the training process and use the hold out data as a test data set.

When this is done with several different splits (different hold out sets) and we evaluate our model using the approach, we are achieving cross validation. The size of the hold out set can vary. If we choose the heaviest computational split, we only contain one sample in the hold out set. This is called leave-one-out-cross-validation LOOCV.



**Figure 4.1:** Logistic regression model illustration

When we have gained several different splits from using cross validation, we need to evaluate how our model performs in these different testing scenarios. We often use mean least square to evaluate our model. The formula for calculating this for our model with cross validation can be seen in 4.1.

$$CV_{(k)} = \sum_{k=1}^K \frac{n_k}{n} MSE_k \quad (4.1)$$



## SUBSET SELECTION

# **SHRINKAGE METHODS / DIMENSION REDUCTION METHODS**