

# Final Course Assignment

Mark Niehues, Stefaan Hessmann  
Mathematical Aspects in Machine Learning

13. Juli 2017

## 1 Introduction

In the past course we dealt with the broad mathematical foundations of machine learning. To get an idea of what the consequences of those mathematical theorems and approaches are and to get in touch with the standard Python tools, we have evaluated an comparatively easy data science example found on [kaggle.com](https://www.kaggle.com). Since this was our first machine learning project, we decided to deal with an rather simple problem. The example dataset [1] consists of the historic passenger records of the disastrous Titanic maiden voyage in 1912. The goal in this Challenge was to predict if a passenger survived the accident based on informations like for example age, sex and payed ticket fare. Therefore it is in terms of machine learning a **classification problem** with two classes: survived and not survived.

Inspired by sample solutions from the website, we first took a deeper look on the dataset and tried to select the most significant influences by reviewing the statistical properties of the dataset. In the following we implemented an naive Sequential Minimal Optimization (SMO) algorithm and ran a few tests with them in order to finally compare the results with other machine learning algorithms.

## 2 Applying Machine Learning Methods on the Titanic Disaster

### 2.1 Dataset

The given dataset consists of a CSV-file containing data of 891 passengers. The dataset contains an ID for every passenger, a label if the passenger has survived the disaster and the features that are described in table 1. It can be noticed that some of the features are incomplete.

After loading the dataset, it is necessary to process the data for our learning machine. Therefore the different features will be investigated to select meaningful features and the missing data needs to be handled.

### 2.2 Feature: Sex

The sex-feature divides the passengers into the categories 'female' and 'male'. Figure 1 shows the probability of survival for male and female passengers. It is obvious that females had a much higher probability to survive than the male passengers. 'Sex' seems to be a useful feature for the learning machine.

Tabelle 1: Features and their amount of missing data.

<b>PassengerId</b>	Unique ID for every passenger	0.0 %
<b>Survived</b>	Survived (1) or died (0)	0.0 %
<b>Pclass</b>	Passenger's class	0.0 %
<b>Name</b>	Passenger's name	0.0 %
<b>Sex</b>	Passenger's sex	0.0 %
<b>Age</b>	Passenger's age	<b>19.87 %</b>
<b>SibSp</b>	Number of siblings/spouses aboard	0.0 %
<b>Parch</b>	Number of parents/children aboard	0.0 %
<b>Ticket</b>	Ticket number	0.0 %
<b>Fare</b>	Ticket-price	0.0 %
<b>Cabin</b>	Number of the passenger's cabin	<b>77.10 %</b>
<b>Embarked</b>	Port of embarkation	<b>0.22 %</b>

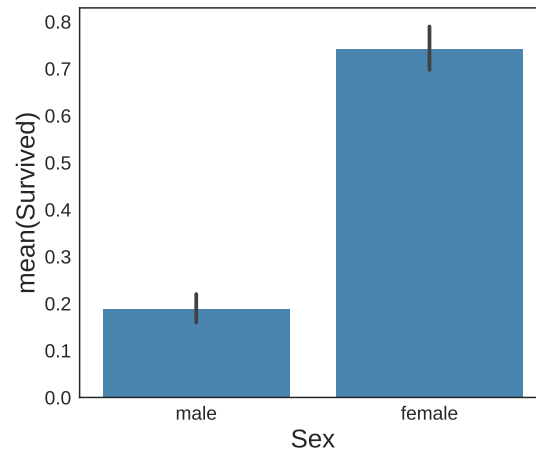


Abbildung 1: Survival probability depending on the passenger's sex.

## 2.3 Feature: Age

The impact of a passengers' age on their probability to survive categorized by their sex is pointed out

## 3 Implementation of an easy SMO Algorithm

To get a better understanding of what a Support Vector Machine does, we decided to implement one on our own using several publications. Most of them were based on the important paper of Platt [2] where he introduced a new approach for the calculation of the Support Vectors that improves the performance a lot. This algorithm is called Sequential Minimal Optimization. Performance was not the highest priority for us but instead understandability and the costs of implementation. Therefore we implemented a less complex version of the algorithm presented in Platt's paper.

## 4 Appendix

Listing 1: Hello World

```
0 # Copyright (C) 2017 Mark Niehues, Stefaan Hessmann
1 #
2 # This program is free software: you can redistribute it and/or modify
3 # it under the terms of the GNU General Public License as published by
4 # the Free Software Foundation, either version 3 of the License, or
5 # (at your option) any later version.
6 #
7 # This program is distributed in the hope that it will be useful,
8 # but WITHOUT ANY WARRANTY; without even the implied warranty of
9 # MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the
10 # GNU General Public License for more details.
11 #
12 # You should have received a copy of the GNU General Public License
13 #
14
15
16 import numpy as np
17
18
19 class Kernels:
20     """
21     Class that holds different Kernels
22     """
23     def __init__(self, gamma):
24         self.gamma = gamma
25         self.kernels = {
26             "rbf" : self.kernel_rbf,
27             "linear": self.kernel_lin}
28
29     def get_kernel(self, kernel_name):
30         return self.kernels[kernel_name]
31
32     def kernel_lin(self, x, y):
33         """
34         Linear kernel
35         """
36         return x.dot(y)
37
38     def kernel_rbf(self, x, y):
39         """
40         RBF Kernel
```

**Input:**

$C$ : regularization parameter  
 $tol$ : numerical tolerance  
 $max\_passes$ : max # of times to iterate over  $\alpha$ 's without changing  
 $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$ : training data

**Output:**

$\alpha \in \mathbb{R}^m$ : Lagrange multipliers for solution  
 $b \in \mathbb{R}$ : threshold for solution

- Initialize  $\alpha_i = 0, \forall i, \quad b = 0$ .
- Initialize  $passes = 0$ .
- **while** ( $passes < max\_passes$ )
  - $num\_changed\_alphas = 0$ .
  - **for**  $i = 1, \dots, m$ ,
    - Calculate  $E_i = f(x^{(i)}) - y^{(i)}$  using (2).
    - **if** ( $(y^{(i)} E_i < -tol \ \&\& \ \alpha_i < C) \ || \ (y^{(i)} E_i > tol \ \&\& \ \alpha_i > 0)$ )
      - Select  $j \neq i$  randomly.
      - Calculate  $E_j = f(x^{(j)}) - y^{(j)}$  using (2).
      - Save old  $\alpha$ 's:  $\alpha_i^{(old)} = \alpha_i, \alpha_j^{(old)} = \alpha_j$ .
      - Compute  $L$  and  $H$  by (10) or (11).
      - **if** ( $L == H$ )
        - **continue** to next  $i$ .
      - Compute  $\eta$  by (14).
      - **if** ( $\eta >= 0$ )
        - **continue** to next  $i$ .
      - Compute and clip new value for  $\alpha_j$  using (12) and (15).
      - **if** ( $|\alpha_j - \alpha_j^{(old)}| < 10^{-5}$ )
        - **continue** to next  $i$ .
      - Determine value for  $\alpha_i$  using (16).
      - Compute  $b_1$  and  $b_2$  using (17) and (18) respectively.
      - Compute  $b$  by (19).
      - $num\_changed\_alphas := num\_changed\_alphas + 1$ .
    - **end if**
  - **end for**
  - **if** ( $num\_changed\_alphas == 0$ )
    - $passes := passes + 1$
  - **else**
    - $passes := 0$
- **end while**

Abbildung 2: Pseudo Code of the SMO algorithm. Taken from: [3]

```
41         """
42         d = x - y
43         return np.exp(-np.dot(d, d) * self.gamma)
```

## Literatur

- [1] Kaggle. *Titanic: Machine Learning from Disaster*. 13. Juli 2017. URL: <https://www.kaggle.com/c/titanic>.
- [2] John Platt. „Sequential minimal optimization: A fast algorithm for training support vector machines“. In: (1998).
- [3] Unknown. *The Simplified SMO Algorithm*. 13. Juli 2017. URL: <http://cs229.stanford.edu/materials/smo.pdf>.