

Final Course Assignment

Mark Niehues, Stefaan Hessmann
Mathematical Aspects in Machine Learning

July 14, 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 [**Kaggle2017**] 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 ???. 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 ??? 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.

Table 1: Features and their amount of missing data.

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

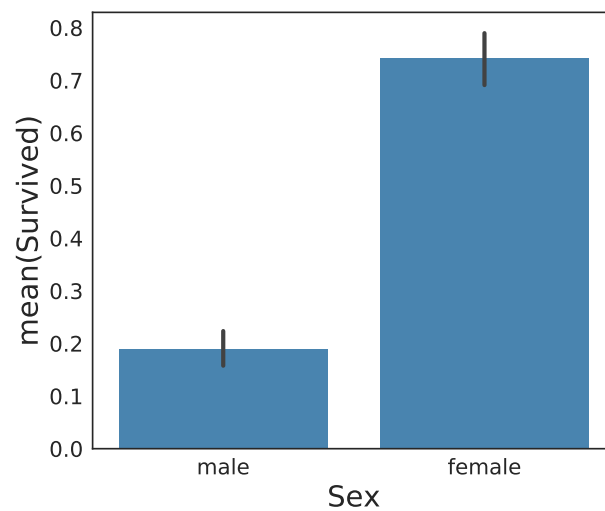


Figure 1: Survival probability depending on the passenger's sex.

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 [platt] 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.

As the mathematical background of the SVMs has been explained in the lecture and might be considered a standard solution for machine learning, the following introduction focuses on the main equations.

The initial problem is a linear separable dataset with the labels $y_i \in \{-1, 1\}$. The classifier that the SVM is supposed to compute will have the form

$$f(x) = \langle \omega, x \rangle + b \quad (1)$$

Now suppose we have a separating hyperplane and w is perpendicular. The main task of the SVM is to maximise the closest perpendicular distance between the hyperplane and the two classes. This is down by the following constraints

$$f(x) \geq 1 \text{ for } y_i = +1 \quad (2)$$

$$f(x) \leq -1 \text{ for } y_i = -1 \quad (3)$$

Consequently do points that lie on the hyperplane satisfy $f(x) = 0$.

From these set of equations follows that the minimal distance from the hyperplane to one of the datapoints is $d = \frac{1}{|\omega|}$ which shall be maximized. Introducing an additional factor that allows but penalizes non separable noise and reformulating the problem with Lagrange multipliers (the α_i) we get the following problem:

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m y^{(i)} y^{(j)} \alpha_i \langle x^{(i)}, x^{(j)} \rangle \quad (4)$$

$$\text{subject to } 0 \leq \alpha_i \leq C, i = 1, \dots, m \quad (5)$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0 \quad (6)$$

For the presented problem the Kuhn Tucker conditions define the α_i that represent an optimal solution. The KKT conditions are

$$\alpha_i = 0 \implies y^{(i)} (\langle \omega, x^{(i)} \rangle + b) \geq 1 \quad (7)$$

$$\alpha_i = C \implies y^{(i)} (\langle \omega, x^{(i)} \rangle + b) \leq -1 \quad (8)$$

$$0 \leq \alpha_i \leq C \implies y^{(i)} (\langle \omega, x^{(i)} \rangle + b) = -1 \quad (9)$$

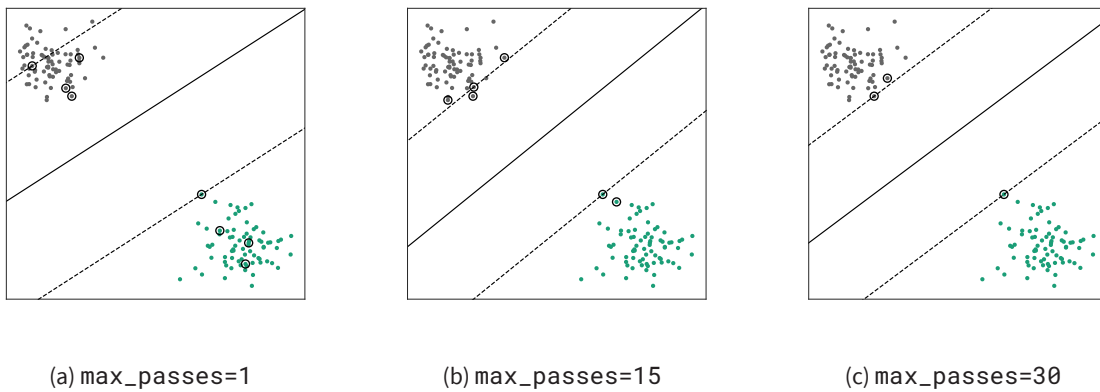


Figure 2: Result of our SVM depending on number of loops without any changing α that terminate the algorithm.

Appendix

A Pseudo Code of the SMO algorithm

B Code

```

2 import numpy as np

5 class Kernels:
6     """
7     Class that holds different Kernels
8     """
9     def __init__(self, gamma):
10         self.gamma = gamma
11         self.kernels = {
12             "rbf" : self.kernel_rbf,
13             "linear" : self.kernel_lin}

15     def get_kernel(self, kernel_name):
16         return self.kernels[kernel_name]

18     def kernel_lin(self, x, y):
19         """
20         Linear kernel
21         """
22         return x.dot(y)

24     def kernel_rbf(self, x, y):
25         """
26         RBF Kernel
27         """
28         d = x - y
29         return np.exp(-np.dot(d, d) * self.gamma)

```

Input: C : regularization parameter tol : numerical tolerance max_passes : max # of times to iterate over α '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 α '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 η by (14).
 - **if** ($\eta \geq 0$)
 - **continue** to next i .
 - Compute and clip new value for α_j using (12) and (15).
 - **if** ($|\alpha_j - \alpha_j^{(old)}| < 10^{-5}$)
 - **continue** to next i .
 - Determine value for α_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**