

Final Course Assignment

Mark Niehues, Stefaan Hessmann
Mathematical Aspects in Machine Learning

July 17, 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). 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. There were some features that could be used for supervised training in a more or less straight forward way. First, the sex-feature which separates the passengers in male and female. We can conclude from fig. 1c that 'Sex' seems to be a meaningful feature for the learning machine as females had a much higher probability to survive.

Secondly, the wealth of a passenger expressed in the passengers Class and the fare he paid for the ticket has positive influence on his chance to survive. The fare feature was categorized into three groups of similar survival rates to permit overfitting (fig. 5). Surprisingly the port of embarkation also correlates with the survival rate as can be seen in figure 1a. Passengers that joined the Titanic at Cherbourg had higher survival chances than people who embarked at other harbours. The other features were engineered in a slightly more advanced way as described below.

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 %

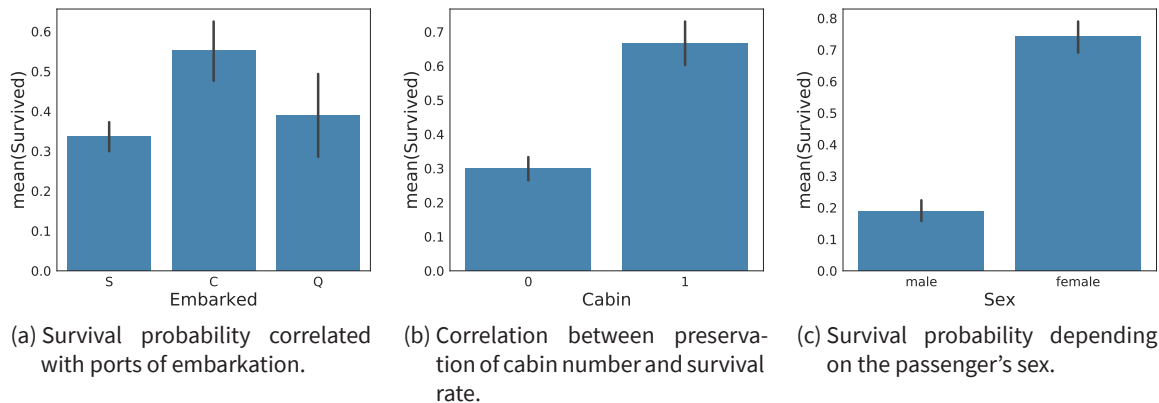


Figure 1: Relation between several features and the survival rate.

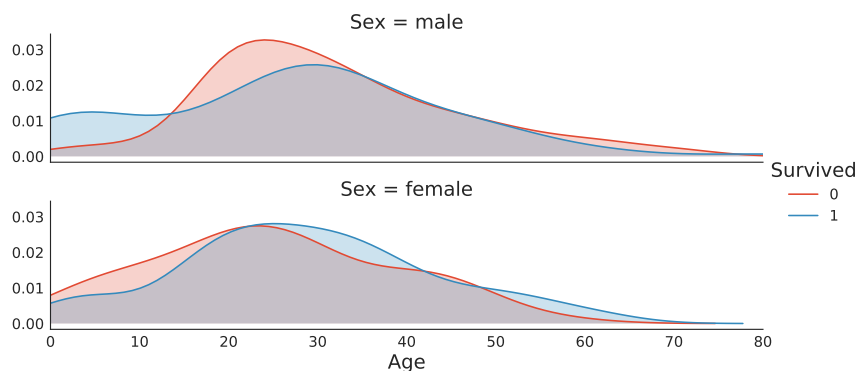


Figure 2: Survival probability depending on the passenger's age and sex.

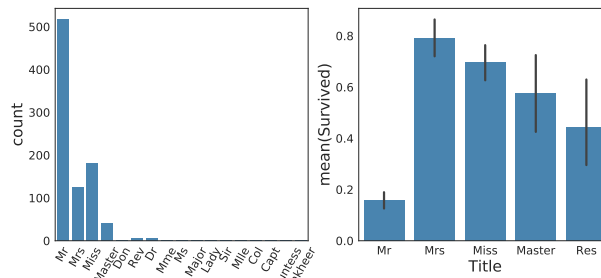


Figure 3: Total counts of the different titles (left) and survival probability in dependency of the title (right).

Feature: Age

The survival distributions in function of the passengers' ages are pointed out in figure 9 in the Appendix. The plot shows that children under twelve years were most likely to survive, whereas older children and young adults until about 30 years had bad chances.

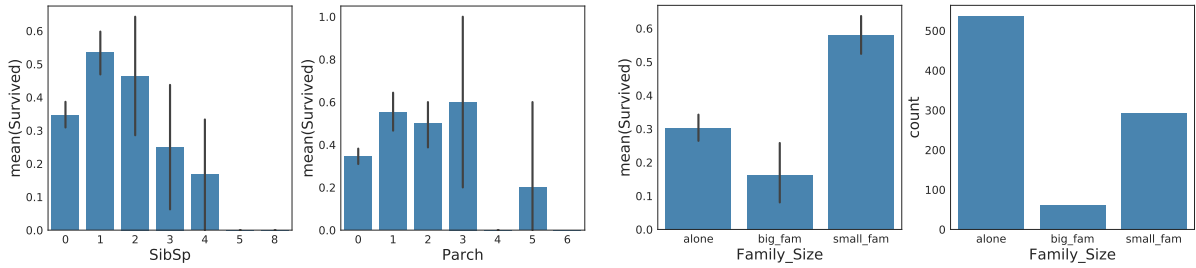
If we categorize the survival distribution depending on the age additionally by the sex feature as shown in figure 9, it turns out that young male passengers were likely to survive while males between about 12 and 30 years were unlikely to survive. This effect is inverted for females. From this follows that age is a feature that can have influence on the predictions of our learning machine, especially if it is combined with the sex feature.

Feature: Name and Title

The given dataset also contains a list of the passenger-names that have been involved in the titanic accident. At first sight a name does not appear to be a useful feature for our survival predictions, but the name-list contains also a persons title. Figure 3 shows the total number of occurrences for all titles that have been found. The titles 'Mr.', 'Master', 'Mrs.' and 'Miss' are relatively common, whereas the other titles occur only infrequently. Therefore all other titles are grouped into a group named 'Res'. The figure also shows the chances of survival for people with different titles on the right-handed plot. The plot points out that a persons title has an impact on the survival odds. It is obvious that there is a correlation between the title feature and the sex and the age feature.

Feature: Family Constellation

The dataset contains two categories that handle family constellations. 'SibSp' contains the number of siblings and spouses a passenger was travelling with and 'Parch' contains the number of parents and children aboard. Figure 4a points out the influence of family constellations on the survival rate. The errorbars imply that the number of datasets with more than two siblings and spouses or more than two parents and children aboard is too sparse for reliable predictions. For this reason the two features are grouped to a new feature called 'Family' for every point in the dataset. This new feature holds the total number of a passengers relatives aboard which is grouped again into three categories containing people travelling alone, small families with one to three relatives aboard and big families with more than three relatives aboard. The classification into these groups was due to the similar survival rate of the number of relatives inside a class. The impact of the family feature is displayed in figure 4b. The errorbars of the plot for the new feature are considerably smaller than before. In this way we can drop the 'SibSp' and the 'Parch' feature and replace it by the new family feature.



(a) Survival probability in dependency of the number of siblings / spouses (left) and parents / children (right) aboard. (b) Total number of datasets (left) survival rate (right) for the family feature.

Figure 4: Merging the SibSP and Parch feature

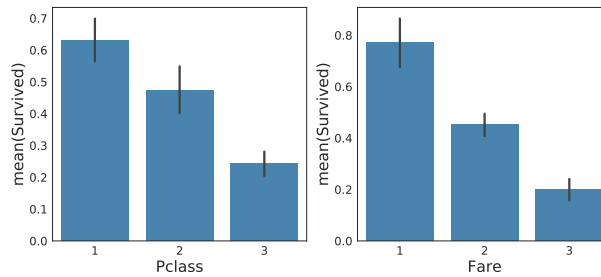


Figure 5: Correlation between Pclass (left), Fare (right) and survival rate.

2.2 Missing Data

As already mentioned before, the dataset lacks some data for the features 'Cabin', 'Embarked' and 'Age'. In order to use these features for our learning machine it is necessary to deal with the missing data. The 'Cabin' data is mostly missing and it is not likely to find a reasonable way to estimate the data. By dividing the data into the categories 'preserved' and 'lost' as shown in figure 1b it turns out that whether the data is missing or not has an effect on the survival rate. For this reason we just reclassify the feature and this already solves the missing data problem for the 'Cabin' feature. More than 99 % of the 'Embarked' dataset is preserved so that missing data can be discarded. Only about 80 % of the 'Age' data is preserved. Therefore it is worth the effort to estimate the remaining 20 %. As described in section 2.1 the survival rate differs for males and females in three different age categories. Additionally we can see that the age feature correlates with a passengers cabin class (see figure 10) and it is obvious that a passengers title is correlated to an age. Based on this information we group all people by the titles 'Miss', 'Mrs', 'Master' and 'Mr + Res' and furthermore by their cabin class. This leads to twelve different categories. For every category we calculate the mean age and match it to the passengers with an unknown age.

This way it is possible to use more than 99 % of our dataset without dropping a feature.

2.3 Preparing the Dataset

Before the dataset can be used for the training of a learning machine, all non numeric features need to be transformed into numerics. If the feature can be arranged in a meaningful order all string variables are assigned to a number. This was done for the 'Family' feature. The other features ('Sex', 'Embarked', 'Title') were transformed to a binary dataset, because they could not be arranged in a meaningful way. Afterwards the features 'SibSp', 'Parch', 'Name', 'PassengerID' and 'Ticket' were dropped, because they are not useful for our learning machine. At last we scale all features on a scale between 0 and 1, so that all features have the same influence on the predictions.

3 Implementation of an easy SMO Algorithm

3.1 Brief Introduction

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 mainly based on [3].

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

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

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

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

To deal with linearly non separable data, the scalar products can be replaced by kernel functions $kernel(x_i, x_j)$.

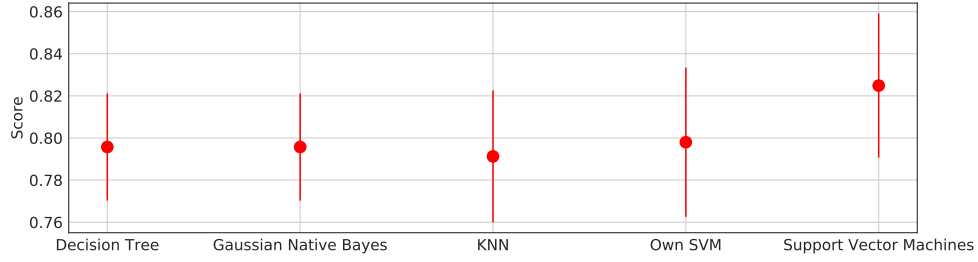


Figure 6: Comparison of different ML Algorithms and their score using the Titanic training set and K-Fold validation.

3.2 Description of the Implementation

Instead of trying to maximize the whole set of α the SMO algorithm exploits that the maximum will be reached when pairs α_i, α_j fulfil the KKT conditions (while it needs to be at least a pair, since the conditions imply linearity of two α values). Thus the SMO algorithm selects two α parameters (that do not meet the KKT conditions) and optimizes them. Afterwards the threshold gets adjusted according to the new values.

A big part of the actual publication from Platt deals with the heuristic of how two choose the α_i and α_j since this is a critical factor for the pace of its convergence as the number of possible pairs in a setup with m features is $m(m - 1)$. Accordingly the amount of time it takes to find the *critical* values is decisive for the algorithms performance.

Nevertheless, in this assignment a very simple heuristic is implemented in order to keep the code simple and understandable: the pairs are just purely randomly selected.

Fig. 11 in the Appendix shows the pseudo code of our implementation. Basically the algorithm consists of an outer and inner loop. The inner loop iterates through the α_i and checks weather it violates the KKT conditions. If this is the case, randomly a second parameter α_j is selected and will be adjusted using that the optimal α_j is given by

$$\alpha'_j = \alpha_j - \frac{y^{(j)}(E_i - E_j)}{\eta} \quad (10)$$

where η can be interpreted as the second derivative of the generic function $W(\alpha)$ [3]. E_i is the current error on the sample x_i . After that the new parameter is cropped to boundaries that follow from equation 9 and 6. Then the opposing parameter is calculated exploiting the linearity the threshold can be updated by using the classifier function $f(x)$ and either one *alpha* that lays within $(0, C)$ or if this is true for both, their arithmetic mean.

The outer loop counts how often the inner loop fails to find a partner for optimization or the optimization yields to no significant changes (significance is defined by the user). The algorithm terminates when a certain, user defined number of passes is reached.

3.3 Comparison with SciKit SVM

The implementation does not make a lot use of Numpy's vectorization skills and therefore performs poor even considering that it is Python code. Still it reaches satisfying scores with the titanic train set in comparison to other SciKit algorithms as can be seen in fig. 6.

On the other hand and less surprisingly, the performance is quite poor compared to the SVM Module from the SciKit framework¹. One can deduce from figure 7 that the prefactor as well as the exponential behaviour is significantly worse.

¹<http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>

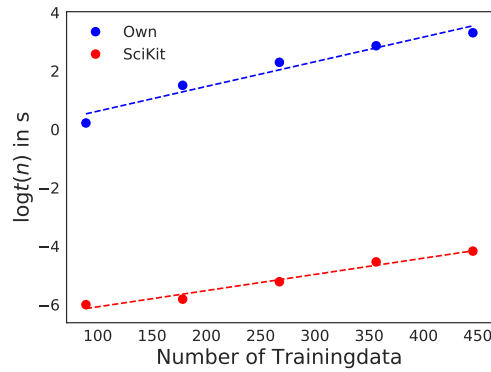


Figure 7: Comparison the self made implementation and the optimized SciKit implementation. The different intercepts ($\Delta \approx 6$) as well as the different slopes (factor 1.6) express the smaller prefactor and exponent of runtime of the SciKit's implementation.

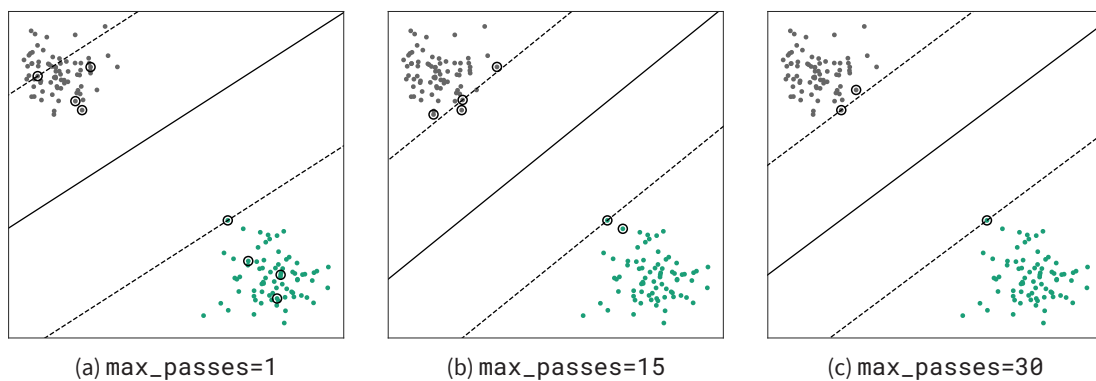


Figure 8: Results of our SVM on a linear separable test set depending on after how many changeless loops it terminates.

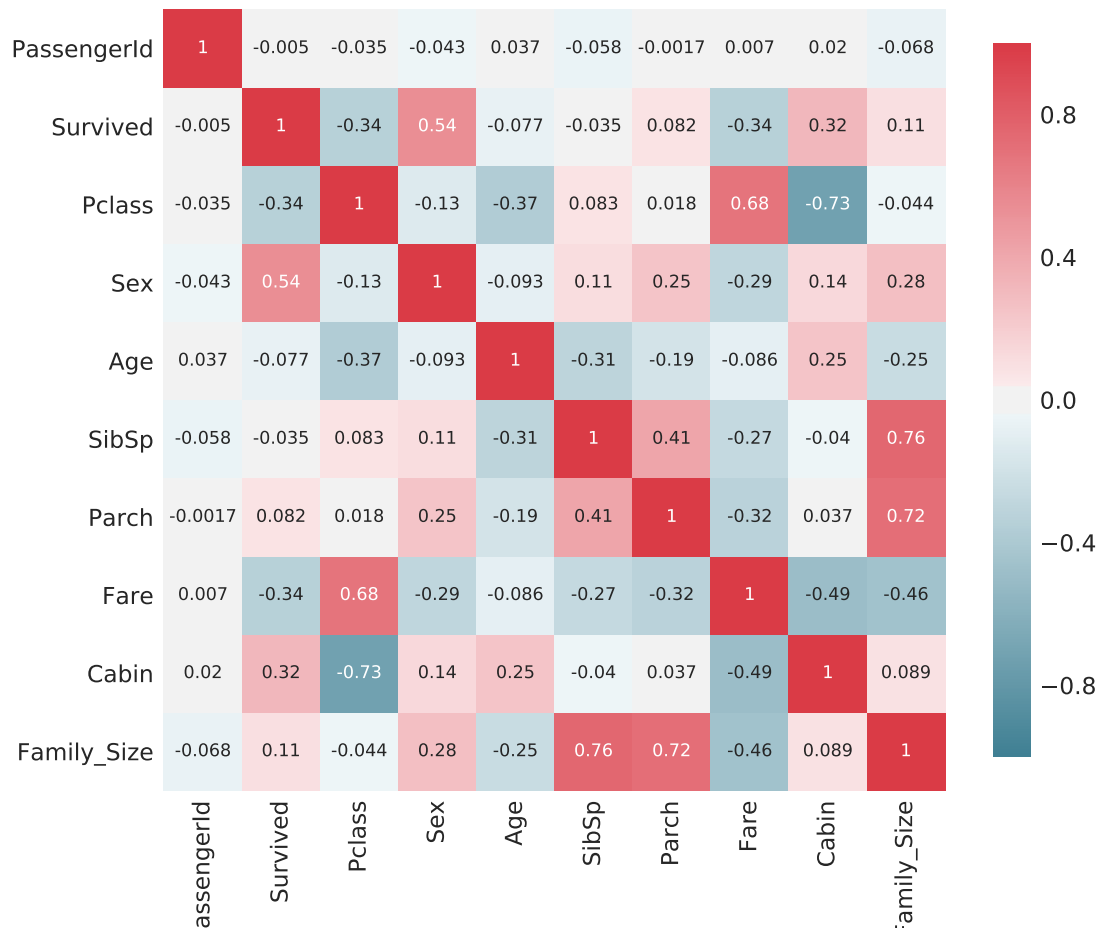
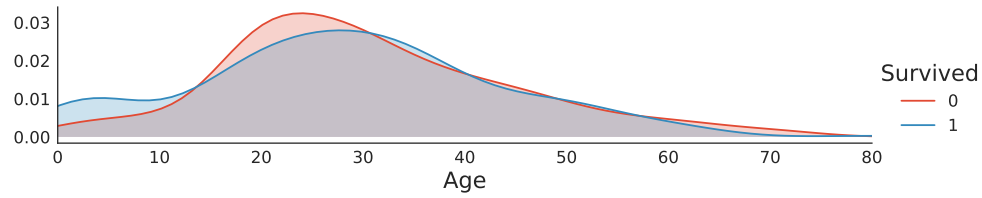
It turned out that the quality of the result heavily depends on the parameter that determines after how many *changeless* runs it terminates. As you can see in figure 8, the algorithm yields a lot support vectors that do not lie within the margin when it stops after one changeless iteration. The quality of the solution increases when the number of passes is larger. The same holds for the runtime unfortunately. This behaviour results from the optimization behaviour of the SVM which is that it only optimizes data points which lay within the margin. Since this "quickly" becomes a small number of data points it is unlikely that the random selection finds a pair to optimize.

4 Summary

We were able to apply our knowledge from class on a simple example problem and thereby extend it with practical skills such as analysing and optimizing features, using the SciKit Framework and dealing with missing data. Additionally we gained deeper insight into one Machine Learning algorithm and were able to point out one critical factor for its performance.

Appendix

A Figures



B Pseudo Code of the SMO algorithm

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**

Figure 11: Pseudo Code of the implemented SMO algorithm. Taken from[3]

References

- [1] Kaggle. *Titanic: Machine Learning from Disaster*. July 13, 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*. July 13, 2017. URL: <http://cs229.stanford.edu/materials/smo.pdf>.