

Supervised classification exercise

Introduction to Machine Learning

Marc Oliu & Joan T. Matamalas

Abstract—Nulla vitae elit libero, a pharetra augue. Donec id elit non mi porta gravida at eget metus. Etiam porta sem malesuada magna mollis euismod. Cras justo odio, dapibus ac facilisis in, egestas eget quam. Nulla vitae elit libero, a pharetra augue. Nullam id dolor id nibh ultricies vehicula ut id elit.

I. INTRODUCTION

The problem of classifying a set of images with two possible classes by means of a binary classifier is a long studied problem which has many techniques for its resolution. The goal of this paper is to analyze the performance of three commonly used techniques for two distinct data sets with a varying number of features and different data distribution. More precisely, the algorithms will be run for a linearly separable dataset with 13 features, and for a non-linearly separable dataset with 34 features.

The three algorithms to be analyzed are the linear Support Vector Machine (SVM), a kernelized version of the Support Vector Machine using the Radial Basis Function kernel (RBF-SVM), and the Adaboost meta-classifier.

The SVM algorithm is a linear classifier which defines a delimiting vector to separate the two classes of the data. For this implementation, a soft-margin version of the classifier which allows for data with overlapping classes is used in order to prevent the algorithm from destabilizing in non-linearly separable data. The RBF-SVM algorithm is a kernelized version of the SVM which transforms the data to be classified by projecting the points into a new dimension by means of a radial basis function, which allows for a better separability of the data.

Finally, the Adaboost meta-classifier is used in conjunction with the linear SVM. The Adaboost (or adaptive boosting) is an algorithm consisting on the generation of a set of weighted weak models which then classify the data by means of a weighted voting of the results of

the model. Multiple weak algorithms were considered for the meta-classifier, but finally SVM was chosen for being the one with the lowest out of sample error.

II. SPECIFICATION OF THE PROBLEM

The general algorithm of the problem consists on a 10-fold cross-validation of the test data sets to calculate the out of sample error for the used algorithms, and for each fold a 3-fold cross-validation is performed over the train data to select the parameters of the models. Three folds are used instead of 10 because of the computational cost of the system, but a 10-fold cross-validation would be preferred to increase the accuracy of the estimation of the out of sample error for the given parameters.

In this section the specification of each algorithm will be explained, as well as the general algorithm used for the obtention and analysis of the out of sample error.

A. The Linear Support Vector Machine (SVM) algorithm

The SVM algorithm consists on the optimization of a hyperplane which is then used for the classification of the data. The two parameters which describe the hyperplane are its director vector, represented by W , and the offset of the hyperplane, represented by b .

The hyperplane resulting from the training of the model is then used as a separator between the two classes, where all elements which fall above the hyperplane belong to one of the classes, and the ones falling below belong to the other. These two classes are represented by a +1 and a -1 value correspondingly, which simplifies the training and classification tasks.

The equation 1 represents the maximization function used to select the alpha values for the different support vectors in the dual form of the SVM with soft margin. The alpha values represent the weight of each support vector, and has a value between 0 and C , being C the parameter to optimize. A higher value for alpha means

that a miss-classification of the related support vector is more heavily penalized, so increasing the value of C increases the maximum weight of each support vector.

$$\begin{aligned} \max_{\alpha_i} \quad & L(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{s.t.} \quad & 0 \leq \alpha_i \leq C \wedge \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned} \quad (1)$$

Once the alpha values are optimized, the values for the director vector W of the hyperplane and its offset 'b' are calculated out of the support vectors, considering as support vectors only those elements with a value above 0 for alpha, as it can be seen on the equations at figure 2.

$$\begin{aligned} w &= \sum_i \alpha_i y_i x_i = 0 \\ b &= \frac{1}{\#SV} \sum_{i \in SV} y_i - \left(\sum_{m \in SV} \alpha_m y_m x_m \cdot x_i \right) \end{aligned} \quad (2)$$

Once the training of the model is completed, the classification of one element can be performed by using only the W and 'b' parameters, which represent the model. To do so, the equation of figure 3 is used. What this equation does is to calculate the position of a data point with respect to the hyperplane and get the sign of the result. If the element is above the hyperplane, it belongs to the +1 class, and if below, to the -1 class.

$$w \cdot x + b \quad (3)$$

B. The Radial Basis Function kernelized SVM (RBF-SVM) algorithm

The RBF-SVM algorithm is a variation of the SVM algorithm which uses a non-linear kernel different from the dot product of data points. In particular, the kernel used in this algorithm for the kernelized SVM is the radial basis function, which is calculated by using the function represented at figure 4. As it can be seen in the function, there is an sigma parameter, which defines the width of the gaussian distribution over the new dimension. As the sigma value increases, the gaussian form is wider, while it gets sharper as sigma decreases.

$$k_{rbf}(x_1, x_2) = e^{-\frac{\|x_1 - x_2\|^2}{2\sigma^2}} \quad (4)$$

The kernel functions allow for the calculation of the dot-product between two points in a higher dimensional space without the need of first increasing the dimensionality of the data points, simplifying the process of optimizing the function, since the same optimization function seen in the previous section

(figure 1) is valid, needing only to replace the dot product by the kernel function.

Since the dot product is performed on the lower-dimensional data, and the dataset isn't transformed to the higher-dimensional space, it's not possible to calculate the hyperplane director vector from the support vector, and thus the support vectors, their alphas and labels must be added to the model in order to classify other data points. The 'b' value, which represents the offset of the hyperplane, is calculated in the same way as with the linear SV, as seen in the second function of figure 2.

To classify the new instances, the equation seen at figure 5 is used. This equation first calculates the position of the data point with respect to the hyperplane without taking into account the offset by using directly the support vector, and then adds the offset to the result. This is the same procedure used in the normal SVM, but since now the equation of the hyperplane isn't defined, the support vectors that define this hyperplane are used directly to check the position of the new data point with respect to the hyperplane.

$$\sum_{i=1} \alpha_i y_i k(x_i, x) \quad (5)$$

As with the SVM models, the sign of the result defines the class of the data point.

C. The Adaboost Algorithm

The adaboost algorithm is a meta-classifier that uses a linear combination of weak classifiers of the same type to classify a data point, as it can be seen in figure 6. As it can be seen, each weak classifier has a weight assigned, which is inversely proportional to the classification error of the weak model.

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right) \quad (6)$$

An weak classifier can be used with the adaboost algorithm, and some have been considered for this work. More specifically, the threshold, perceptron and linear SVM classifiers have been considered for being used with the adaboost algorithm. Finally the selected one has been the SVM algorithm, since it is the one with the highest predictive power, and because it's a good option to check the improvement of the Adaboost

algorithm over the normal SVM classifier.

$$\begin{aligned} D_{t+1}(i) &= \frac{D_t(i)e^{-\alpha_t y_t h_t(x_i)}}{Z_t} \\ Z_t &= \sum_i D_t(i)e^{\alpha_t y_t h_t(x_i)} \end{aligned} \quad (7)$$

The training of the adaboost algorithm is performed in a cyclic way, where at the beginning the weights of each data point to classify are all set to the same value $1/N$, and on each iteration a new model is generated and weighted by following the formula seen at figure 8, and the weights are updated for the misclassified data points on the new model by following the function 7, normalizing the weights afterwards.

$$\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon}{\epsilon} \quad (8)$$

D. Main algorithm

The main algorithm used to run the experiment consists on a 10-fold cross-validation over the data set, performing an optimization of the parameters of each model at each one of the folds, by performing a 3-fold of the training data.

First the data is divided into 10 folds, and an iteration over the folds selects one fold at each iteration to estimate the out of sample error of the models to generate. Then, the remaining ten folds are used to train the three models. These remaining 9 folds are merged into three folds, and a new iteration over the three folds is performed, selecting one fold at each iteration, which will be used to estimate the error of the tested parameters.

The two training folds are used to train the three kinds of models for different values of the parameters, and the out of sample error for the selected values is estimated from the other fold. Once the out of sample error has been calculated from the three folds, the errors for each model type and set of parameters are averaged, and the set of parameters for each model that gives the lowest mean of out of sample errors are selected. Then the models are trained with the data of the three folds (which is the data of the 9 training folds of the main 10-fold cross-validation) and the selected parameters.

Finally, the out of sample errors tested for the optimal parameters at each fold of the 10-fold cross validation are averaged for each model, giving the mean

out of sample error. Also, the standard deviation and confidence intervals are calculated for the ten folds.

As for the error surface, it's obtained by averaging the error surfaces of each fold of the 10-fold cross-validation, which in turn are obtained as the average of the error surfaces for the 3-fold cross validation performed at each fold.

III. ORGANIZATION OF THE SOFTWARE

IV. EXPERIMENTS

A. Experimental methodology

B. Ionosphere dataset

C. Heart-statlog dataset

V. CONCLUSIONS

VI. FUTURE WORK