

IN4320 Machine Learning Computer Exercise

Version 0.6

March 1, 2017

Computational Learning Theory: boosting

- a. Prove that $e^{-x} \geq (1 - x)$.

Before you answer the following questions, you need to read the paper 'A decision-theoretic generalization of on-line learning and an application to boosting' by Y. Freund and R.E. Schapire, 1995 (also available from Blackboard, under Course Documents, Reading Material, Computation Learning Theory). You do not have to focus too much on sections 1, 2 and 3, but section 4 and 4.1 are important.

- b. Implement a 'weak learner': the decision stump. The decision stump is a very simple classifier that chooses one feature f from a dataset, and checks if the feature value x_f is larger (or smaller) than a certain threshold θ . If $x_f > (<)\theta$ then the object is assigned to class ω_1 and otherwise to ω_2 .

To find the optimal feature f and threshold θ , you have to do an exhaustive search for a training set. So, try all values for f and θ and the sign y of $<$ or $>$, and remember those values f^*, θ^*, y^* for which the classification error on the trainingset is minimum.

Make sure that the function accepts a dataset with labels as input, and that the function outputs the optimal f^*, θ^* and y^* .

Show your (Matlab) code.

- c. Test the implementation on the dataset generated by `gendats` from Prtools. (If you don't want to use Prtools, just generate the two classes from two Gaussian distributions, where the means are $\mu_1 = [0, 0]^T$ and $\mu_2 = [2, 0]^T$, and the covariance matrices are identity matrices.) Make a scatterplot of the data, and give the optimal parameters obtained by your decision stump.

Does the decision stump change if you rescale one of the features (for instance, when you multiply feature 2 by a factor of 10)?

- d. Test the implementation on dataset `optdigitssubset` of last week (available from Blackboard). It consists of the pixel values of small 8×8 images, which are ordered in $N = 1125$ rows of 64 columns wide. The first 554 rows contains the values of 8×8 images of zeros, while the remaining block of 571 rows contains the 64 pixel values of 571 ones. Use the first 50 objects for each class for training, and the rest for testing. What is the classification error on the test objects? How much does this vary when you take other random subsets of 50 for training? Show the mean and standard deviation of the error. Is the performance with training on the first 50 object unexpectedly good or bad?

- e. Extend the implementation of the weak learner such that it accepts a weight per object. Therefore, next to a dataset and labels, the function should accept a weight $w_i > 0$ per object \mathbf{x}_i . The weighted weak learner should now minimize the weighted classification error.

Test the code by using a simple classification problem again, and convince yourself that the code is Bug Free. (That means not only that 'it does not crash', but also show convincingly that objects with heigher weight have a larger influence on the solution.)

- f. Implement the algorithm that is described in Figure 2 of the paper: AdaBoost. Show the code.

Also implement the code to classify new and unseen data. Show the code.

- g. Test your implementation on some dataset. Not only use a simple dataset like `gendats` but also a more complicated dataset like `gendatb`. Find

out which objects obtain a large weight w_i^t .¹ Keep the number of iterations T fixed to, say, $T = 100$.

- h.** Test the implementation on dataset `optdigitssubset` of last week. Use the first 50 objects for each class for training, and the rest for testing. What is the classification error on the test objects? How does this error depend on the number of iterations T ? If you take the classifier with the optimal T , which training objects get a high weight? Show them!

¹And no, it is *not* sufficient to just say 'object 13 and 7 have high weight'. Show me, do the more simple or more difficult objects get high weight?