

# Boosting

## *Seminar: Foundations of Data Science*

Christian Peters, enrolment no.: 213996  
Faculty of Statistics  
TU Dortmund

January 8, 2021  
winter term 2020/21

### Abstract

A short summary of about two to three sentences that briefly and  
concisely outline the content ...

## 1 Introduction

Training machine learning models in practice is not always as simple as it might seem from a theoretical standpoint. In [1, chapter 2], the empirical risk minimization (ERM) rule was introduced as the learning algorithm of choice. However, this theoretical principle of choosing a hypothesis  $h \in \mathcal{H}$  such that  $L_S(h) = \min_{h \in \mathcal{H}} L_S(h)$ , where  $L_S(h)$  is the error of  $h$  on the training sequence  $S$ , can be impossible to use in practical applications due to the sometimes enormous computational complexity of searching through interesting hypothesis classes  $\mathcal{H}$ .

This problem leads to the question if it is possible to arrive at a strong learning algorithm in a way that doesn't require the computational cost of searching through complex hypothesis classes. Is it perhaps possible to create a strong learner by finding a way to combine "weak" learners that are potentially easier to compute? The algorithmic paradigm of boosting deals with exactly this question, resulting in the widely used AdaBoost algorithm that shows how "weak" learners can be combined in order to obtain a strong learning algorithm.

The goal of this article is to first lay down the foundations of boosting by explaining the concept of weak learnability which will be used to arrive

at the AdaBoost algorithm followed by a discussion of its implications as well as a practical example of how it can be used in the domain of image classification.

## 2 Weak Learnability

Assuming that the realizable assumption [1, chapter 2] holds, the generalization error  $L_{(\mathcal{D},f)}(h)$  of a PAC learner with respect to a distribution  $\mathcal{D}$  and a labeling function  $f$  can by the definition of PAC learning [1, chapter 2] be reduced to an arbitrarily small number (with confidence of  $1 - \delta$ ) by increasing the sample size of the training sequence  $S$ . This, however, can be computationally infeasible in practical applications.

The concept of weak learnability aims to relax the requirement that the generalization error of a hypothesis must become arbitrarily small the more the sample size is increased. For weak learning, it is sufficient that the learning algorithm yields a hypothesis  $h$ , that performs only slightly better than a random guess. One can think of weak learning as applying a simple rule which isn't fully capable of modeling the data generating process, but can still learn a little bit about the underlying problem so that it performs better than random.

Formally, the definition of a weak learner differs only slightly from the definition of PAC learning. In the situation of a two-class classification problem, the definition of a weak learner, can be given as follows:

**Definition 1.** ( *$\gamma$ -weak-learner*) An Algorithm  $A$  is a  $\gamma$ -weak-learner for a hypothesis class  $\mathcal{H}$  if for every  $\delta \in (0, 1)$  there exists a threshold  $m_{\mathcal{H}}(\delta) \in \mathbb{N}$  such that for every distribution  $\mathcal{D}$  over the instance space  $\mathcal{X}$  and for every labeling function  $f : \mathcal{X} \rightarrow \{\pm 1\}$  if running  $A$  on  $m \geq m_{\mathcal{H}}$  training examples, it will yield a hypothesis  $h$  such that with probability of at least  $1 - \delta$  the generalization error  $L_{(\mathcal{D},f)}(h)$  is at most  $\frac{1}{2} - \gamma$ , provided that the realizable assumption holds.

The parameter  $\gamma$  in Definition 1 tells us how well we can expect the weak learner to perform. For example if a learning algorithm is a 1%-weak-learner for a class  $\mathcal{H}$ , then the generalization error can at most be 49% provided that first, the learner didn't fail (which can happen with probability  $\delta$  of drawing a bad sample from  $\mathcal{D}$ ) and second, the learner was run on at least  $m_{\mathcal{H}}(\delta)$  training examples.

It still remains the question of how to obtain a weak learning algorithm for a class  $\mathcal{H}$ . We already saw that applying the ERM rule to complex hypothesis classes can be computationally hard. But what if we choose a simple

hypothesis class  $B$  instead, where the ERM rule can be applied efficiently? It follows directly from Definition 1 that this can only work if the new algorithm  $\text{ERM}_B$  has an error of at most  $\frac{1}{2} - \gamma$  for every sample that was labeled by a hypothesis from  $\mathcal{H}$ . In this case, applying the ERM rule with respect to the simpler class  $B$  would yield a weak learner for  $\mathcal{H}$ .

## 2.1 An Efficient ERM Algorithm for Decision Stumps

One such hypothesis class, where the ERM algorithm can be implemented efficiently, is the class of decision stumps. On the instance space  $\mathcal{X} = \mathbb{R}^d$ , this class is given as follows:

$$\mathcal{H}_{DS} = \{\mathbf{x} \mapsto \text{sign}(\theta - x_i) \cdot b : \theta \in \mathbb{R}, i \in [d], b \in \{\pm 1\}\}$$

The idea behind decision stumps is to divide the instance space  $\mathcal{X}$  along one of its dimensions  $i \in [d]$  at a threshold  $\theta$  into two partitions, such that all the instances in one partition are labeled  $+1$  and all the instances in the other partition are labeled  $-1$ .

Fixing  $b = 1$  without the loss of generality, an ERM algorithm for  $\mathcal{H}_{DS}$  has to find the optimal splitting dimension  $i \in [d]$  as well as the optimal splitting threshold  $\theta$  such that the training error  $L_S(h)$  on a training sequence  $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$  is minimized.

An extension of the empirical risk  $L_S(h)$  that will be useful when introducing the AdaBoost algorithm later in section 3, is the weighted empirical risk, which is given as

$$L_{\mathbf{D}}(h) = \sum_{i=1}^m D_i \mathbb{1}_{[h(\mathbf{x}_i) \neq y_i]}.$$

The vector  $\mathbf{D} \in \mathbb{R}^m$  is used to give a weight to each training example, where each weight  $D_i$  is nonnegative and  $\sum_{i=1}^m D_i = 1$ . For the special case that  $D_i = \frac{1}{m}$ , the weighted empirical risk is equal to the unweighted empirical risk.

In order to find a decision stump  $h \in \mathcal{H}_{DS}$  that minimizes the weighted empirical risk  $L_{\mathbf{D}}(h)$ , the  $\text{ERM}_{\mathcal{H}_{DS}}$  algorithm has to solve the following optimization problem that minimizes the weighted sum of misclassifications:

$$\min_{j \in [d]} \min_{\theta \in \mathbb{R}} \left( \sum_{i: y_i = 1}^m D_i \mathbb{1}_{[x_{i,j} > \theta]} + \sum_{i: y_i = -1}^m D_i \mathbb{1}_{[x_{i,j} \leq \theta]} \right)$$

Taking a closer look at this problem it becomes clear, that it is not necessary to try every  $\theta \in \mathbb{R}$  during the optimization. In fact, if we first sort the

92 examples for a fixed dimension  $j \in [d]$  so that  $x_{1,j} \leq x_{2,j} \leq \dots \leq x_{m,j}$ , we  
 93 can see that we only have to consider a single splitting threshold in between  
 94 two examples, which leaves us with  $m + 1$  values for  $\theta$  that the  $\text{ERM}_{\mathcal{H}_{DS}}$  has  
 95 to consider for a fixed dimension  $j \in [d]$ . Since the sum of misclassifications  
 96 can be computed in  $\mathcal{O}(m)$  by passing through the data once, the algorithm  
 97 can find optimal values for  $j$  and  $\theta$  in  $\mathcal{O}(dm^2)$  by simple enumeration.

98 This time complexity can be reduced even more by avoiding to recalculate  
 99 the sum of misclassifications for every new value of  $\theta$ . It can be shown [1] that  
 100 it is possible to update the sum for a new value of  $\theta$  in constant time, requiring  
 101 only a single pass through the data for a fixed dimension  $j$ . This leaves us  
 102 with a time complexity of  $\mathcal{O}(dm)$  of solving the optimization problem after  
 103 the sorting step is applied as preprocessing.

### 104 3 AdaBoost

105 In the previous section it was shown how an efficient weak learner can be  
 106 constructed by applying the ERM rule to the class of decision stumps. This  
 107 section presents the AdaBoost (Adaptive Boosting) algorithm, a procedure  
 108 that shows how weak learners such as decision stumps can be used efficiently  
 109 to find a hypothesis with an arbitrarily low empirical error  $L_S(h)$  on a training  
 110 sequence  $S$ .

111 The goal of AdaBoost is to invoke the weak learner multiple times on  
 112 the training data and then to combine the resulting hypotheses similar to a  
 113 weighted majority vote. Let  $T$  be the number of times that AdaBoost invokes  
 114 the weak learner on the training data. Then the resulting output hypothesis  
 115 of AdaBoost has the following form:

$$116 \quad h(x) = \text{sign} \left( \sum_{t=1}^T w_t h_t(x) \right)$$

117 Here,  $h_t(x)$  is the output hypothesis of the weak learner in iteration  $t$  and  $w_t$   
 118 is the corresponding weight that AdaBoost assigns to this hypothesis.

119 In the first iteration  $t = 1$ , the AdaBoost algorithm assigns an equal  
 120 weight  $D_i^{(1)} = \frac{1}{m}$  to each example in the training sequence  $S$  and then invokes  
 121 the weak learner on the weighted training sequence. The error of the resulting  
 122 hypothesis is computed according to

$$123 \quad \epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[h_t(\mathbf{x}_i) \neq y_i]}$$

124 and is at most  $\frac{1}{2} - \gamma$ .

125 The weight  $w_t$  that is assigned to a hypothesis in boosting round  $t$  is  
 126 computed as follows:

$$127 \quad w_t = \frac{1}{2} \log \left( \frac{1}{\epsilon_t} - 1 \right)$$

128 At the end of each iteration, the weights  $D_i^{(t)}$  of each training example  
 129 are updated according to

$$130 \quad D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-w_t y_i h_t(\mathbf{x}_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-w_t y_j h_t(\mathbf{x}_j))}$$

131 This update will assign a higher weight to those examples, that weren't  
 132 correctly classified by  $h_t$ .

## 133 4 Conclusion

134 Even after centuries of research in the field of data science, there is nothing  
 135 more versatile than the useful theorem of chapter 3. It is used everywhere  
 136 and has led to the greatest and most intriguing results, cf. [2]. By the way,  
 137 the book for the seminar [1] is a great reference and should be cited. Further  
 138 literature can be found in the respective *Bibliographic Remarks* sections and  
 139 of course you are welcome to search and add your own references.

140 **Note:** BibTeX entries can often be found in the DBLP collection. Google  
 141 Scholar also offers BibTeX entries, which can be copied into the .bib file and  
 142 may need some minor adjustments.

## References

- [1] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning - From Theory to Algorithms*. Cambridge University Press, 2014.
- [2] J. Someone and J. Someoneelse. Useful theorems, 2003.