# 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 $f : \mathcal{X} \to \{\pm 1\}$ 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 $-1$. 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 $S = ((\mathbf{x}_1, y_1), ..., (\mathbf{x}_m, y_m))$ is minimized.

Let $b = 1$ without the loss of generality. Then the $\text{ERM}_{\mathcal{H}_{DS}}$ algorithm has to solve the following optimization problem that minimizes the sum of misclassifications:

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

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

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

# 3 AdaBoost

Here we work with the above definitions and notations and derive important results such as the following Theorem on the least-squares solution

**Theorem 1.** *(useful theorem) Let $X \in \mathbb{R}^{n \times d}, Y \in \mathbb{R}^n$. Further define $\beta^* = \operatorname{argmin}_{\beta \in R^d} \|X\beta - Y\|^2$. Then*

$$\|Y\|^2 = \|X\beta^*\|^2 + \|X\beta^* - Y\|^2.$$

*Proof.* The proof is left as an exercise. $\square$

Sometimes figures help to illustrate a formalism. This is completely unrelated to Theorem 1.

# 4 Conclusion

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

**Note:** BibTeX entries can often be found in the DBLP collection. Google Scholar also offers BibTeX entries, which can be copied into the .bib file and 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.