

Boosting

Seminar: Foundations of Data Science

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

January 14, 2021
winter term 2020/21

Abstract

Boosting is an algorithmic paradigm that aims to create an efficient strong learner by combining multiple weak learners. This paper introduces the concept of weak learnability and explains the AdaBoost algorithm, the first practical implementation of boosting.

1 Introduction

Training machine learning models in practice is not always as simple as it might seem from a theoretical standpoint. In [3, 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, which was first raised by Kearns and Valiant in 1988 [2], resulting in the widely used AdaBoost algorithm that shows how

22 "weak" learners can be combined in order to obtain a strong learning algo-
 23 rithm.

24 The goal of this article is to first lay down the foundations of boosting
 25 by explaining the concept of weak learnability which will be used to arrive
 26 at the AdaBoost algorithm, the first practical implementation of boosting.

27 2 Weak Learnability

28 Assuming that the realizable assumption [3, chapter 2] holds, the generaliza-
 29 tion error $L_{(\mathcal{D},f)}(h)$ of a PAC learner with respect to a distribution \mathcal{D} and
 30 a labeling function f can by the definition of PAC learning [3, chapter 2]
 31 be reduced to an arbitrarily small number (with confidence of $1 - \delta$) by in-
 32 creasing the sample size of the training sequence S . This, however, can be
 33 computationally infeasible in practical applications.

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

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

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

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

58 It still remains the question of how to obtain a weak learning algorithm
59 for a class \mathcal{H} . We already saw that applying the ERM rule to complex hy-
60 pothesis classes can be computationally hard. But what if we choose a simple
61 hypothesis class B instead, where the ERM rule can be applied efficiently?
62 It follows directly from Definition 1 that this can only work if the new algo-
63 rithm ERM_B has an error of at most $\frac{1}{2} - \gamma$ for every sample that was labeled
64 by a hypothesis from \mathcal{H} . In this case, applying the ERM rule with respect
65 to the simpler class B would yield a weak learner for \mathcal{H} .

66 2.1 An Efficient ERM Algorithm for Decision Stumps

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

$$70 \quad \mathcal{H}_{DS} = \{\mathbf{x} \mapsto \text{sign}(\theta - x_i) \cdot b : \theta \in \mathbb{R}, i \in [d], b \in \{\pm 1\}\}$$

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

75 Fixing $b = 1$ without the loss of generality, an ERM algorithm for \mathcal{H}_{DS}
76 has to find the optimal splitting dimension $i \in [d]$ as well as the optimal
77 splitting threshold θ such that the training error $L_S(h)$ on a training sequence
78 $S = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_m, y_m))$ is minimized.

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

$$82 \quad L_{\mathbf{D}}(h) = \sum_{i=1}^m D_i \mathbb{1}_{[h(\mathbf{x}_i) \neq y_i]}.$$

83 The vector $\mathbf{D} \in \mathbb{R}^m$ is used to give a weight to each training example, where
84 each weight D_i is nonnegative and $\sum_{i=1}^m D_i = 1$. For the special case that
85 $D_i = \frac{1}{m}$, the weighted empirical risk is equal to the unweighted empirical
86 risk.

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

$$90 \quad \min_{j \in [d]} \min_{\theta \in \mathbb{R}} \left(\sum_{i: y_i = 1}^m D_i \mathbb{1}_{[x_{i,j} > \theta]} + \sum_{i: y_i = -1}^m D_i \mathbb{1}_{[x_{i,j} \leq \theta]} \right)$$

91 Taking a closer look at this problem it becomes clear, that it is not necessary
 92 to try every $\theta \in \mathbb{R}$ during the optimization. In fact, if we first sort the
 93 examples for a fixed dimension $j \in [d]$ so that $x_{1,j} \leq x_{2,j} \leq \dots \leq x_{m,j}$, we
 94 can see that we only have to consider a single splitting threshold in between
 95 two examples, which leaves us with $m + 1$ values for θ that the $\text{ERM}_{\mathcal{H}_{DS}}$ has
 96 to consider for a fixed dimension $j \in [d]$. Since the sum of misclassifications
 97 can be computed in $\mathcal{O}(m)$ by passing through the data once, the algorithm
 98 can find optimal values for j and θ in $\mathcal{O}(dm^2)$ by simple enumeration.

99 This time complexity can be reduced even more by avoiding to recalculate
 100 the sum of misclassifications for every new value of θ . It can be shown [3] that
 101 it is possible to update the sum for a new value of θ in constant time, requiring
 102 only a single pass through the data for a fixed dimension j . This leaves us
 103 with a time complexity of $\mathcal{O}(dm)$ of solving the optimization problem after
 104 the sorting step is applied as preprocessing.

105 **3 AdaBoost**

106 In the previous section it was shown how an efficient weak learner can be
 107 constructed by applying the ERM rule to the class of decision stumps. This
 108 section presents the AdaBoost (Adaptive Boosting) algorithm, a procedure
 109 that shows how weak learners such as decision stumps can be used efficiently
 110 to find a hypothesis with an arbitrarily low empirical error $L_S(h)$ on a train-
 111 ing sequence S . It was first proposed in 1995 by Yoav Freund and Robert
 112 Schapire [1] and became hugely popular due to its practicability.

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

$$118 \quad h(x) = \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right)$$

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

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

124 hypothesis is computed according to

$$125 \quad \epsilon_t = \sum_{i=1}^m D_i^{(t)} \mathbb{1}_{[h_t(\mathbf{x}_i) \neq y_i]}$$

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

127 The weight w_t that is assigned to a hypothesis in boosting round t is
128 computed as follows:

$$129 \quad w_t = \frac{1}{2} \log \left(\frac{1}{\epsilon_t} - 1 \right)$$

130 As we can see, the smaller the error ϵ_t of the hypothesis h_t is, the bigger the
131 weight w_t will be.

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

$$134 \quad D_i^{(t+1)} = \frac{D_i^{(t)} \exp(-w_t y_i h_t(\mathbf{x}_i))}{\sum_{j=1}^m D_j^{(t)} \exp(-w_t y_j h_t(\mathbf{x}_j))}$$

135 This update assigns a higher weight to those examples, that weren't correctly
136 classified by h_t .

137 In short, the Adaboost algorithm performs the following steps in each of
138 the T iterations:

- 139 1. Invoke the weak learner on the training sequence S weighted by $\mathbf{D}^{(t)}$
- 140 2. Assign a weight w_t to the output hypothesis h_t of the weak learner.
141 Hypotheses with a smaller training error ϵ_t will get a higher weight.
- 142 3. Compute a new weight vector $\mathbf{D}^{(t+1)}$ that gives a higher weight to
143 incorrectly classified examples.

144 The computational complexity of AdaBoost essentially consists of invoking
145 the weak learning algorithm T times on the training data. Thus, if the
146 weak learner can be implemented efficiently (as it is the case for decision
147 stumps), AdaBoost is also efficient.

148 On top of that, it can be shown [3], that the training error $L_S(h)$ of the
149 AdaBoost hypothesis h decreases exponentially in the number of boosting
150 rounds T :

$$151 \quad L_S(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{[h(\mathbf{x}_i) \neq y_i]} \leq e^{-2\gamma^2 T}$$

152 Here, γ describes the γ -weak-learner as defined in Definition 1. This means,
153 that AdaBoost will achieve an arbitrarily low training error and is still an
154 efficient algorithm.

155 In practical applications however, it is more important to also achieve
 156 a good out of sample error. The next section will show that the true risk
 157 $L_{\mathcal{D}}(h)$ of the AdaBoost hypothesis will also be small by taking a look at
 158 the hypothesis class that resembles all the possible output hypotheses of
 159 AdaBoost.

160 3.1 AdaBoost Out-of-Sample Performance

161 The hypothesis class of AdaBoost is parameterized by the hypothesis class
 162 B of the weak learner it uses as well as by the number of boosting rounds T .
 163 Formally, it is given as follows:

$$164 \quad L(B, T) = \left\{ x \mapsto \text{sign} \left(\sum_{t=1}^T w_t h_t(x) \right) : w \in \mathbb{R}^T, h_t \in B \right\}$$

165 The fundamental theorem of statistical learning [3, chapter 6] states
 166 that a hypothesis class is PAC-learnable if its VC dimension is finite. This
 167 means that if the VC dimension of $L(B, T)$ is finite, there exists a threshold
 168 $m_{\mathcal{H}}(\epsilon, \delta) \in \mathbb{N}$ for every $\epsilon, \delta \in (0, 1)$, such that the true error $L_{\mathcal{D}}(h)$ of the
 169 AdaBoost hypothesis $h \in L(B, T)$ is at most ϵ (with confidence $1 - \delta$) for
 170 every distribution \mathcal{D} when training AdaBoost on at least $m_{\mathcal{H}}(\epsilon, \delta)$ training
 171 examples. This tells us that even the Out-of-Sample error of AdaBoost can
 172 be reduced arbitrarily (if the realizable assumption holds), by increasing the
 173 amount of training examples.

174 It can be shown that an upper bound of the VC dimension of $L(B, T)$ is
 175 linear in the VC dimension of the hypothesis class B of the weak learner and
 176 also linear in T , i.e. $\text{VCdim}(L(B, T))$ is in $\mathcal{O}(\text{VCdim}(B)T)$. This means that
 177 if the VC dimension of B is finite, so is the VC dimension of $L(B, T)$.

178 In the case of B being the hypothesis class of all decision stumps, the
 179 VC dimension of B is 2, so it is clearly finite. It follows from the funda-
 180 mental theorem of statistical learning, that when using decision stumps as
 181 the hypothesis class for the weak learner, AdaBoost is a PAC learner for the
 182 class $L(B, T)$. Furthermore, if the weak learner for B can be implemented
 183 efficiently (as we have shown for decision stumps), AdaBoost is also efficient.

184 4 Conclusion

185 In this article it was shown how a weak learning algorithm can be boosted
 186 into a strong PAC learner using the AdaBoost algorithm. The hypothesis
 187 class of decision stumps was identified to be efficiently learnable using an

ERM algorithm, so it followed that the AdaBoost algorithm is also efficient using decision stumps as the base learner.

Further, it was also shown that AdaBoost is a PAC learner for the hypothesis class $L(B, T)$, which tells us that the true error of AdaBoost can be arbitrarily decreased by increasing the size of the training set. All of these properties make AdaBoost a useful algorithm for practical applications of machine learning.

References

- [1] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119 – 139, 1997.
- [2] M. Kearns and L. G. Valiant. Learning boolean formulae or finite automata is as hard as factoring. Technical Report TR 14-88, Harvard University Aiken Computation Laboratory, 1988.
- [3] S. Shalev-Shwartz and S. Ben-David. *Understanding Machine Learning - From Theory to Algorithms*. Cambridge University Press, 2014.