

Chapter 3

Boosting Algorithms: A Review of Methods, Theory, and Applications

Artur Ferreira and Mário Figueiredo

3.1 Introduction

Boosting is a class of machine learning methods based on the idea that a combination of simple classifiers (obtained by a *weak learner*) can perform better than any of the simple classifiers alone. A *weak learner* (WL) is a learning algorithm capable of producing classifiers with probability of error strictly (but only slightly) less than that of random guessing (0.5, in the binary case). On the other hand, a *strong learner* (SL) is able (given enough training data) to yield classifiers with arbitrarily small error probability.

An ensemble (or committee) of classifiers is a classifier build upon some combination of weak learners. The strategy of boosting, and ensembles of classifiers, is to learn many *weak* classifiers and combine them in some way, instead of trying to learn a single *strong* classifier. This idea of building ensembles of classifiers has gained interest in the last decade [66]; the rationale is that it may be easier to train several simple classifiers and combine them into a more complex classifier than to learn a single complex classifier. For instance, instead of training a large *neural network* (NN), we may train several simpler NNs and combine their individual outputs in order to produce the final output (as illustrated in Fig. 3.1).

Letting $H_m: \mathcal{X} \rightarrow \{-1, +1\}$ be the m -th weak binary classifier (for $m = 1, \dots, M$), and $\mathbf{x} \in \mathcal{X}$ some input pattern to be classified, there are many ways to combine the outputs $H_1(\mathbf{x}), \dots, H_M(\mathbf{x})$ into a single class prediction [66]. For example, assuming that the classifiers err independently of each other, a majority vote combination should yield a lower probability of error than any of the individual classifiers. Con-

Artur Ferreira

Instituto de Telecomunicações, and Instituto Superior de Engenharia de Lisboa, Portugal.

e-mail: arturj@isel.pt

Mário Figueiredo

Instituto de Telecomunicações, and Instituto Superior Técnico, Technical University of Lisbon, Portugal.

e-mail: mario.figueiredo@lx.it.pt

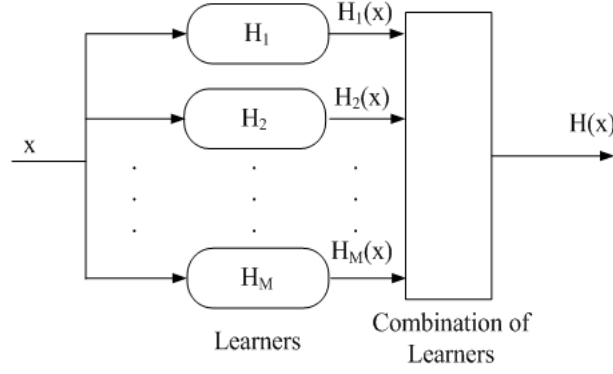


Fig. 3.1 The concept of ensemble of classifiers. The outputs of the weak learners $H_m(\mathbf{x})$ with $m \in \{1, \dots, M\}$ are combined to produce the output of the ensemble of classifiers given by $H(\mathbf{x})$.

sidering a weighted linear combination of the outputs of the weak classifiers, the ensemble prediction function $H : \mathcal{X} \rightarrow \{-1, +1\}$ is given by

$$H(\mathbf{x}) = \text{sign} \left(\sum_{m=1}^M \alpha_m H_m(\mathbf{x}) \right), \quad (3.1)$$

where $\alpha_1, \dots, \alpha_M$ is a set of weights (a simple majority vote results if all the weights are equal).

Among the many different ways in which ensembles of classifiers can be learned and combined [66], boosting techniques exhibit, in addition to good practical performance, several theoretical and algorithmic features that makes them particularly attractive [57], [81], [97]. Essentially, boosting consists in repeatedly using the base weak learning algorithm, on differently weighted versions of the training data, yielding a sequence of weak classifiers that are combined as in (3.1). The weighting of each instance in the training data, at each round of the algorithm, depends on the accuracy of the previous classifiers, thus allowing the algorithm to focus its *attention* on those samples that are still incorrectly classified. The several variants of boosting algorithms differ in their choice of base learners and criterion for updating the weights of the training samples. AdaBoost (which stands for *adaptive boosting*) is arguably the best known boosting algorithm, and was responsible for sparking the explosion of interest in this class of algorithms that happened after the publication of the seminal works of Freund and Schapire [46], [47], [48], [49],

3.1.1 Chapter outline

The remaining sections of this chapter are organized as follows. Section 2 addresses the foundations and origins of boosting algorithms, as a class of methods to im-

prove the accuracy of learning algorithms, by building ensembles of classifiers; the connection of boosting with other machine learning techniques, such as bootstrap and bagging is mentioned. Section 3 describes AdaBoost and discusses some of its theoretical properties, regarding training error, generalization error, and the problem of overfitting. In Section 4, we describe variants of AdaBoost and their properties, including extensions for multi-class problems, while boosting algorithms for semi-supervised learning are discussed in Section 5. Section 6 discusses several successful applications of batch and online boosting algorithms and presents an experimental evaluation of some boosting algorithms, compared to other machine learning techniques on standard benchmark datasets. Section 7 provides a summary and a discussion on boosting algorithms. Finally, Section 8 ends the chapter with some bibliographic and historical remarks.

3.2 The Origins of Boosting and Adaptive Boosting

3.2.1 Bootstrapping and Bagging

Bootstrapping [36, 37] is a general purpose sample-based statistical method in which several (non disjoint) training sets are obtained by drawing randomly, *with* replacement, from a single base dataset. In a dataset with N samples, each instance is selected with probability $1/N$; consequently, after N draws (with large N), the probability that a given instance was not selected is

$$\left(1 - \frac{1}{N}\right)^N \approx \exp(-1) \approx 0.368; \quad (3.2)$$

the validity of this approximation is illustrated in Fig. 3.2, showing that it is quite accurate even with only a moderately large N . This implies that each sample contains roughly 63.2 % of the instances.

Classically, bootstrapping is used to infer some statistic $T(P)$ about a (say infinitely large) population P , from N samples thereof: $Z = \{z_1, \dots, z_N\}$. The idea is to obtain B sets $Z_b^* \subseteq Z$, for $b = 1, \dots, B$, each containing N random samples (with replacement) from Z , from which B estimates of $T(P)$ are obtained. These estimates are then averaged into a final estimate; it is also possible to obtain variance estimates or confidence intervals. The procedure is formally described in Algorithm 1.

Bagging (which stands for *bootstrap aggregation* [11]) is a technique which uses bootstrap sampling to reduce the variance and/or improve the accuracy of some predictor (it may be used in classification and regression). Consider a size- N dataset $Z = \{z_1, z_2, \dots, z_N\}$, where now $z_i = (\mathbf{x}_i, y_i)$, where y_i is a class label, in classification problems, or a real number, in regression problems. The rationale of bagging is to learn a set of B predictors (each from a bootstrap sample $Z_b^* \subseteq Z$, for $b = 1, \dots, B$) and then produce a final predictor by combining (by averaging, in regression, or majority voting, in classification) this set of predictors. The combination of multiple

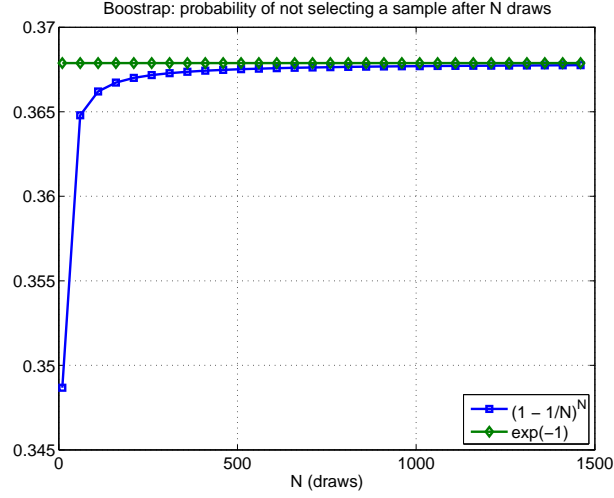


Fig. 3.2 The Bootstrap procedure: probability of not selecting a training sample after N draws and its approximation $\exp(-1)$.

Algorithm 1 Bootstrap Procedure

Input: Size- N sample $Z = \{z_1, z_2, \dots, z_N\}$ of a (potentially infinite) population P .
 B , number of bootstrap samples.

Output: Estimate $\hat{T}(P)$ of the population statistic.

- 1: **for** $b = 1$ to B **do**
 - 2: Draw, with replacement, N samples from Z , obtaining the b -th bootstrap sample Z_b^* .
 - 3: Compute, for each sample Z_b^* , the estimate of the statistic $\hat{T}(Z_b^*)$.
 - 4: **end for**
 - 5: Compute the bootstrap estimate $\hat{T}(P)$ as the average of $\hat{T}(Z_1^*), \dots, \hat{T}(Z_B^*)$.
 - 6: Compute the accuracy of the estimate, using, *e.g.*, the variance of $\hat{T}(Z_1^*), \dots, \hat{T}(Z_B^*)$.
-

predictors decreases the expected error because it reduces the variance component of the bias-variance decomposition [57]. The reduction on this variance component is proportional to the number of classifiers applied in the ensemble. The bagging procedure, for binary classification, is described in Algorithm 2.

As compared to the process of learning a classifier in a conventional way, that is, from the full training set, bagging has two main advantages:

- increases classifier stability and accuracy;
- reduces classifier variance, in terms of the bias-variance decomposition [57].

The use of the bagging technique improves the classification results whenever the base classifiers are unstable, this being the main reasons why the bagging approach works well for classification. Fig. 3.3 depicts the bagging approach for classification.

For further reading on bagging, see [12, 13, 94, 137, 138]. In [137], the authors argue that for very weak learners (*e.g.*, decision stumps, which are tree classifier with

Algorithm 2 Bagging Procedure for Classification

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 B , number of bootstrap samples.

Output: $H : \mathcal{X} \rightarrow \{-1, +1\}$, the final classifier.

- 1: **for** $b = 1$ to B **do**
- 2: Draw, with replacement, N samples from Z , obtaining the b -th bootstrap sample Z_b^* .
- 3: From each bootstrap sample Z_b^* , learn classifier H_b .
- 4: **end for**
- 5: Produce the final classifier by a majority vote of H_1, \dots, H_B , that is, $H(\mathbf{x}) = \text{sign}\left(\sum_{b=1}^B H_b(\mathbf{x})\right)$.

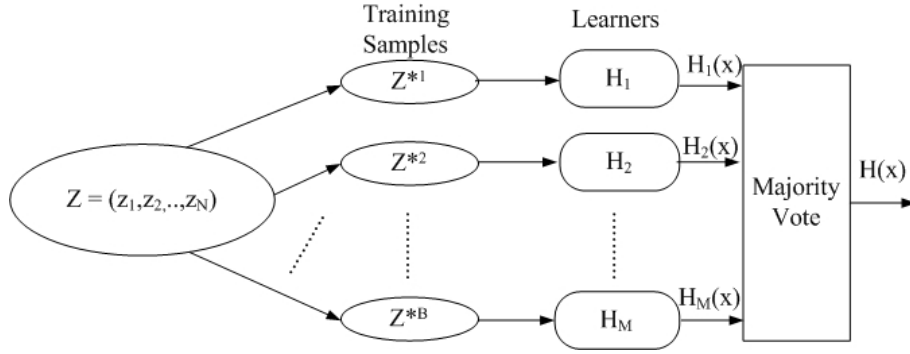


Fig. 3.3 The bagging approach to classification. Using bootstrap, we produce several training samples; each of these samples is fed into a weak learner. The final classification decision is produced by a majority vote on the weak learners output.

only one inner node), the base classifiers built from bootstrap samples are strongly correlated. As a consequence, a simple bagged classifier with these weak learners has very little improvement compared to a single classifier trained from the same data. To overcome this problem, they propose the *local lazy learning based bagging approach* (LLLB), where base learners are trained from a small subset surrounding each test instance. The experimental results on real-world datasets show that the LLLB method significantly outperforms standard bagging.

3.2.2 Strong and Weak Learners

Weak and *strong* learning are fundamental concepts at the heart of boosting algorithms, so we briefly review their formal definitions. These concepts are rooted in the theory of PAC (probably approximately correct) learning [112], where they are defined as follows. Consider an hypothesis, *i.e.*, a classification rule $f : \mathcal{X} \rightarrow \{-1, +1\}$, such that $f \in \mathcal{F}$, where \mathcal{F} is some class of functions from \mathcal{X} to $\{-1, +1\}$. Consider also a set of examples of that hypothesis, *i.e.*, a set of pairs

$\{(\mathbf{x}_i, y_i), i = 1, \dots, N\}$ such that $y_i = f(\mathbf{x}_i)$ and the \mathbf{x}_i are samples of some distribution P . A *strong* learner is capable of, given enough data, producing an arbitrarily good classifier with high probability, that is, for every $P, f \in \mathcal{F}$, $\varepsilon \geq 0$, and $\delta \leq 1/2$, it outputs, with probability no less than $1 - \delta$, a classifier $h : \mathcal{X} \rightarrow \{-1, +1\}$ satisfying $\mathbb{P}_P[h(\mathbf{x}) \neq f(\mathbf{x})] \leq \varepsilon$. Furthermore, the time complexity of the algorithm can be at most polynomial in $1/\varepsilon$, $1/\delta$, N , and the dimension of \mathcal{X} [81].

A weak learner is formally defined in a similar way as a strong one, but with weaker quantification with respect to ε and δ . Given a particular (rather than “for every”) pair $\varepsilon_0 \geq 0$, and $\delta_0 \leq 1/2$, a weak learner outputs, with probability no less than $1 - \delta_0$, a classifier $h : \mathcal{X} \rightarrow \{-1, +1\}$ satisfying $\mathbb{P}_P[h(\mathbf{x}) \neq f(\mathbf{x})] \leq \varepsilon_0$. Underlying the idea of boosting is the fact, proved by Schapire [94], that it is possible to obtain a strong learner by combining weak learners.

3.2.3 Boosting Algorithms

The first boosting procedure was proposed in by Schapire in [94], where the key result is that weak and strong learnability are equivalent, in the sense that strong learning can be performed by combining weak learners. The boosting procedure proposed in [94] is described in detail in Algorithm 3.

Algorithm 3 Boosting Procedure for Classification

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

- 1: Randomly select, without replacement, $L_1 < N$ samples from Z to obtain Z_1^* .
 - 2: Run the weak learner on Z_1^* , yielding classifier H_1 .
 - 3: Select $L_2 < N$ samples from Z , with half of the samples misclassified by H_1 , to obtain Z_2^* .
 - 4: Run the weak learner on Z_2^* , yielding classifier H_2 .
 - 5: Select all samples from Z on which H_1 and H_2 disagree, producing Z_3^* .
 - 6: Run the weak learner on Z_3^* , yielding classifier H_3 .
 - 7: Produce the final classifier as a majority vote: $H(\mathbf{x}) = \text{sign}\left(\sum_{b=1}^3 H_b(\mathbf{x})\right)$.
-

As can be seen in Algorithm 3, the training set is randomly divided without replacement into three partitions, Z_1^* , Z_2^* , and Z_3^* . For a given instance, if the first two classifiers (H_1 and H_2) agree on the class label, this is the final decision for that instance. The set of instances on which they disagree defines the partition Z_3^* , which is used to learn H_3 . Schapire has shown that this learning method is strong, in the sense defined above. Moreover, the error can be further reduced by using this approach recursively, that is, each learner can itself be a boosting procedure. Fig. 3.4 illustrates the boosting approach.

After this proposal by Schapire, Freund [43] proposed a new boosting algorithm based on, and improving, the ideas presented in [94]. That algorithm improves the

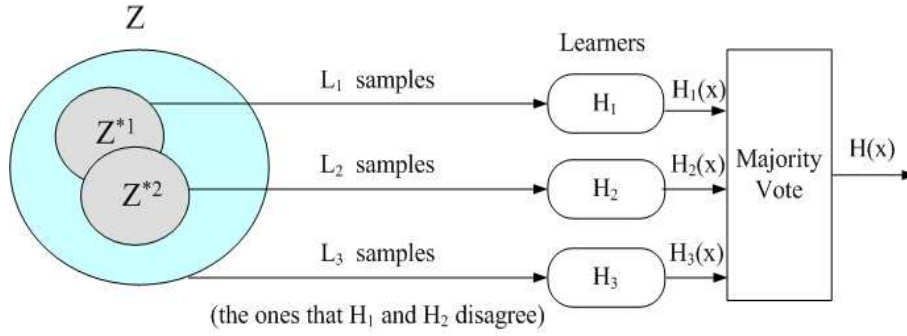


Fig. 3.4 A graphical idea of the first boosting approach proposed in [94]. Notice that each learner can be itself learned by the boosting algorithm in a recursive fashion.

accuracy of algorithms for learning binary classifiers, by combining a large number of classifiers, each of which is obtained by running the given learning algorithm on a different set of examples. As in [94], Freund's new proposals also suffered from several drawbacks, namely the need for a very large training set, due to the fact that this set is divided into subsets.

3.2.4 Relationship Between Boosting, Bagging, and Bootstrapping

Fig. 3.5 shows the connection between bootstrapping, bagging, and boosting, focusing on what they produce and how they handle the training data. The figure emphasizes the fact that these three techniques are all built upon random sampling, being that bootstrapping and bagging perform sampling with replacement while boosting does not. Bagging and boosting have in common the fact that they both use provide a final classifier that is a majority vote of the individual classifiers.

In [28], a comparison of the effectiveness of randomization, bagging, and boosting for improving the performance of the decision-tree algorithm C 4.5 [87] is presented. The experimental results show that for cases with little or no classification noise, randomization is competitive with (and perhaps slightly superior to) bagging but not as accurate as boosting. For situations with substantial classification noise, bagging is much better than boosting, and sometimes better than randomization.

3.3 The AdaBoost Algorithm

After their initial separate work on boosting algorithms, Freund and Schapire proposed the *adaptive boosting* (AdaBoost) algorithm [46], [47], [49]. The key idea behind AdaBoost is to use *weighted* versions of the *same* training data instead of

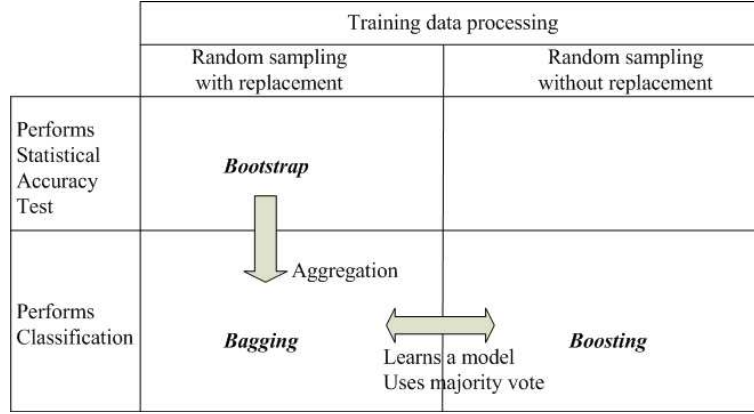


Fig. 3.5 Bootstrapping, bagging, and boosting: what they yield and how they handle the training data.

randomly subsamples thereof. The same training set is repeatedly used and, for this reason, it does need to be very large, as was required by earlier boosting methods.

The AdaBoost algorithm is now a well known and deeply studied method to build ensembles of classifiers with very good performance [57]. The algorithm learns a set of classifiers, using a weak learner, in order to produce the final classifier of the form (3.1). The weak classifiers¹ are obtained sequentially, using re-weighted versions of the training data, with the weights depending on the accuracy of the previous classifiers. The training set is always the same at each iteration, with each training instance weighted according to its (mis)classification by the previous classifiers. This allows the weak learner at each iteration to focus on patterns that were not well classified by the previous weak classifiers. It is important to choose weak learners to obtain the base classifiers, allowing them to learn without decreasing significantly the weight of the previously correctly classified instances. If the base learner is too strong, it may achieve high accuracy, leaving only outliers and noisy instances with significant weight to be learned in the following rounds. Fig. 3.6 depicts the structure of AdaBoost, which is described in detail in Algorithm 4.

The function $h : \mathbb{R} \rightarrow \{0, 1\}$ used in line 4 of the algorithm is the Heaviside function, defined as $h(x) = 1$, if $x \geq 0$, and $h(x) = 0$, if $x < 0$. Consequently, since both y_i and $H_m(\mathbf{x}_i)$ take values in $\{-1, +1\}$, we have that $h(-y_i H_m(\mathbf{x}_i)) = 1$, if $y_i \neq H_m(\mathbf{x}_i)$, and $h(-y_i H_m(\mathbf{x}_i)) = 0$, if $y_i = H_m(\mathbf{x}_i)$, and err_m is the weighted error rate of the m -th classifier.

Line 3 requires some explanation: what does it mean to run a weak learning algorithm on a weighted version of the training set? It means that the goal of the weak learner is to obtain a classifier, say H_m , belonging to a given family of classifiers \mathcal{H} , that satisfies

¹ We refer to a classifier learned by a weak learner as a weak classifier.

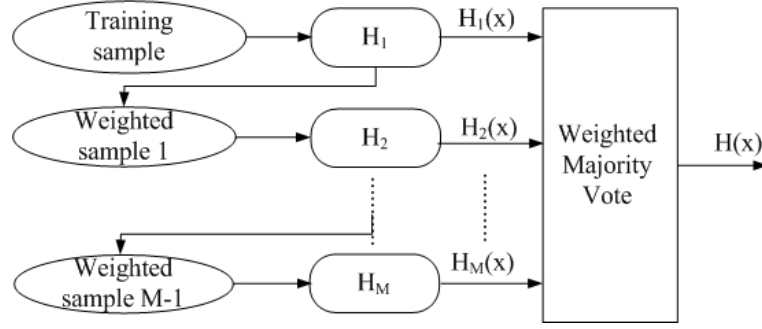


Fig. 3.6 Graphical idea of the adaptive boosting algorithm (adapted from [57]). Each weak learner is trained on a different weighted version of the training data sample. There is no sampling of the training data and the weight of each instance for the following round depends on the performance of the previous learner.

Algorithm 4 (Discrete) AdaBoost Algorithm for Binary Classification

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

- 1: Initialize the weights $w_i^{(1)} = 1/N$, $i \in \{1, \dots, N\}$, and set $m = 1$.
 - 2: **while** $m \leq M$ **do**
 - 3: Run weak learner on Z , using weights $w_i^{(m)}$, yielding classifier $H_m : \mathcal{X} \rightarrow \{-1, +1\}$.
 - 4: Compute $\text{err}_m = \sum_{i=1}^N w_i^{(m)} h(-y_i H_m(\mathbf{x}_i))$, the weighted error of H_m .
 - 5: Compute $\alpha_m = \frac{1}{2} \log \left(\frac{1 - \text{err}_m}{\text{err}_m} \right)$. { /* Weight of the weak learner. */ }
 - 6: For each sample $i = 1, \dots, N$, update the weight $v_i^{(m)} = w_i^{(m)} \exp(-\alpha_m y_i H_m(\mathbf{x}_i))$.
 - 7: Renormalize the weights: compute $S_m = \sum_{j=1}^N v_j$ and, for $i = 1, \dots, N$, $w_i^{(m+1)} = v_i^{(m)} / S_m$.
 - 8: Increment the iteration counter: $m \leftarrow m + 1$
 - 9: **end while**
 - 10: Final classifier: $H(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$.
-

$$\sum_{i=1}^N w_i h(-y_i H_m(\mathbf{x}_i)) \leq \frac{1}{2} - \varepsilon, \quad (3.3)$$

for some small positive ε . Notice that only if $w_i = 1/N$, for $i = 1, \dots, N$ (e.g., at the first iteration of AdaBoost) does the left hand side of (3.3) coincides with the classical error rate on the training set. The existence of such weak classifiers is an important ingredient of boosting, and the interested reader is referred to [81] for more details. Notice that the *weakness* of the classifier is usually controlled by letting \mathcal{H} contain only simple classifiers; for example, when $\mathcal{X} = \mathbb{R}^d$, the family \mathcal{H} may contain only linear rules of the form $H(\mathbf{x}) = \text{sign}(\mathbf{u}^T \mathbf{x} + r)$, where $\mathbf{u} \in \mathbb{R}^d$ and $r \in \mathbb{R}$, which is sometimes known as a perceptron, or rules based on a single

component of the input, *i.e.*, of the form $H(\mathbf{x}) = \text{sign}(u x_j + t)$, where $u \in \{-1, +1\}$ and $t \in \mathbb{R}$, which is called a decision stump.

Notice that the AdaBoost algorithm can actually handle weak classifiers with weighted error rate larger than $1/2$; of course, by simply inverting the output of such a classifier, we obtain a classifier with weighted error rate less than $1/2$. Such an inversion is automatically performed by AdaBoost, because if $\text{err}_m > 1/2$, the corresponding weight α_m is negative, as is clear in Fig. 3.7.

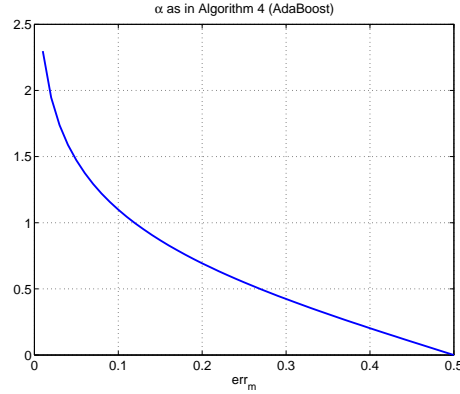


Fig. 3.7 The computation of $\alpha_m (\geq 0)$ as a function of the weighted classification error for each weak learner. As the error tends to 0.5 (random guessing) the contribution (importance) of the weak learner decreases.

Table 3.1 shows the connection between the boosting algorithm (Algorithm 3) and AdaBoost (Algorithm 4). We compare these algorithms in terms of how the training data is processed, the number of classifiers and how the final decision is produced.

Table 3.1 Summary of the main differences between Algorithms 3 (Boosting) and 4 (AdaBoost), regarding how training data is used, the number of classifiers, and the decision mechanism.

	Boosting (Algorithm 3)	AdaBoost (Algorithm 4)
Data usage	Random sampling, no replacement	Weighting (no sampling)
Number of samples	Three	One
Number of classifiers	Three	Up to M
Decision	Majority vote	Weighted majority vote

A key issue when using the weights for the instances, is that the following learner is provided with more information about the importance of each instance and how the previous learners were (or not) able to deal with that instance. This does not happen in bagging nor boosting.

Notice that a straightforward consequence of the instance weighting scheme is that, after the M AdaBoost rounds, the misclassified patterns assigned with higher weights are “hard” patterns to learn; these patterns are probably outliers. This is a kind of side effect of Adaboost, which can be used for outlier detection on given training set.

The AdaBoost algorithm has also been extended for regression tasks. In [3], the prediction error is compared against a threshold to mark it as an error or not and then the AdaBoost version for classification is used. In [30], the probabilities kept by the algorithm are modified based on the magnitude of the error; instances with large error on the previous learners have a higher probability of being chosen to train the following base learner. The median or weighted average is then applied to combine the predictions of the different base learners.

3.3.1 Some Theoretical Properties

We now review several properties of AdaBoost that were shown by Freund and Schapire [49], [99], namely the exponential decay of the training error rate.

The first result shows that the training error of the classifier obtained after M boosting rounds is upper bounded by the product of the normalizing constants of the weights of all the rounds, that is,

$$TE = \frac{1}{N} \sum_{i=1}^N h(-y_i H(\mathbf{x}_i)) \leq \prod_{j=1}^M S_j, \quad (3.4)$$

where S_j is the normalizing constant used in line 7 at iteration j (the proof of this result can be found in Appendix A).

The second result shows how the TE depends on the weighted error rates of the weak classifiers (denoted err_m). Assume that $\text{err}_m = 1/2 - \gamma_m$, with $\gamma_m > \gamma > 0$, for all $m = 1, \dots, M$. Then

$$TE \leq \exp(-2M\gamma^2), \quad (3.5)$$

that is, the training error decreases exponentially with M , and does so at a rate that depends on γ (the proof of this result can be found in Appendix A).

The *expected test error* commonly addressed as the *generalization error* (GE) also has an upper bound, as demonstrated in [49]. The GE of the final classifier, is upper bounded, with high probability by

$$TE + \tilde{O}\left(\sqrt{\frac{Md}{N}}\right), \quad (3.6)$$

where d is the Vapnik-Chervonenkis (VC) [9, 115] dimension of the set of base classifiers. This results shows that there is a tradeoff controlled by the “richness” or “complexity” of the base (weak) classifiers; “stronger” base classifiers allow the TE

to be lower, but correspond to a larger CV dimensions; on the other hand, simpler classifiers have a lower CV dimension, but require more boosting rounds to decrease the TE.

It has been found empirically that the GE usually does not increase as the size of ensemble becomes very large; moreover, it is often observed that the GE continues to decrease even after the training error has reached zero. In [98], it is shown that this behavior is related to the distribution of margins of the training examples with respect to the generated voting classification rule. The margin of an example is defined as the difference between the number of correct votes and the maximum number of votes received by any incorrect label.

3.3.2 Different Views of AdaBoost

It has been argued that one reason for the success of AdaBoost is its ability to increase the *margin* between positive and negative examples [98]. This view provides a connection between margin-based discriminative learning (as in *support vector machines* – SVM) and boosting.

The adaptive boosting techniques can be considered as a greedy optimization method for minimizing the exponential error function

$$\frac{1}{N} \sum_{i=1}^N \exp(-y_i f(\mathbf{x}_i)) = \sum_{i=1}^N \exp\left(-y_i \sum_{m=1}^M \alpha_m H_m(\mathbf{x}_i)\right), \quad (3.7)$$

by learning H_m and choosing the most adequate value of α_m at each round. Detailed analysis of boosting and different views of how this learning procedure behaves can be found in [34, 49, 57, 79, 80].

In [21] a unified view of boosting and logistic regression [57] is described. These learning problems are cast in terms of optimization of Bregman distances, due to their high similarity under this framework. For both problems, new sequential and parallel algorithms are proposed and their potential advantages over existing methods are shown. A general proof of convergence for AdaBoost is also presented. Some connections of AdaBoost with game-theory, linear programming, logistic regression and estimation of probabilities and outliers are discussed in [96, 99].

An evaluation of bagging and boosting using both neural networks and decision trees as learners is carried out in [76]. The experimental results show two important features. The first is that, even though bagging almost always produces a better classifier than any of its individual component classifiers and is relatively impervious to overfitting, it does not generalize any better than a baseline neural-network ensemble method. The second is that, although boosting is a powerful technique that can usually produce better ensembles than bagging, it is more susceptible to noise and overfitting.

In [41], AdaBoost is evaluated on synthetic and real data using two types of weak learners: generative classifiers and radial basis function classifiers. The Ad-

aBoost algorithm with these weak learners shows good convergence properties. On benchmark data, boosting of these weak learners attains results close to the Real AdaBoost algorithm (with decision trees) and SVM, constituting a low computational complexity competitive choice.

3.4 Variants of AdaBoost

In this section, we review several variants of AdaBoost (although we don't claim to have an exhaustive list), both for binary and multi-class supervised learning problems. Many of these variants have proven to be successful in different types of learning scenarios.

The proposal of the AdaBoost algorithm stimulated a significant amount of research on this type learning technique, exploiting its theoretical properties and experimental performance. From this research, several variants of AdaBoost have emerged, some targeted at specific problems, such as, for example, face detection and text categorization. Those variants follow the overall structure of AdaBoost (learn a weak classifier, compute the amount of error, update the weights of the training patterns and repeat the process), but introduce changes on several aspects, such as the weight update expression.

3.4.1 Detailed analysis of some variants

After AdaBoost was introduced, several modified versions (variants) have been proposed, developed, and compared with AdaBoost. This section addresses some of these variants (shown in the timeline of Fig. 3.8) for supervised learning of binary classifiers.

The following subsections describe in detail some of these variants for binary classification. These variants were selected to be presented in more detail, because they are either the first variants to appear after AdaBoost was proposed or they bring quite different new ideas into the adaptive boosting scheme. These variants have in common the fact that all of them were proved to be successful in real-world machine learning problems.

3.4.1.1 Real AdaBoost

The first variant we consider is Real AdaBoost [51, 99], where the term *real* refers to the fact that the algorithm uses real-valued “classifiers” (*i.e.*, before thresholding). This real value can be seen as the probability, or degree of confidence, that a given input pattern belongs to a class, considering the current weight distribution for the training set. The Real AdaBoost algorithm is presented as Algorithm 5.

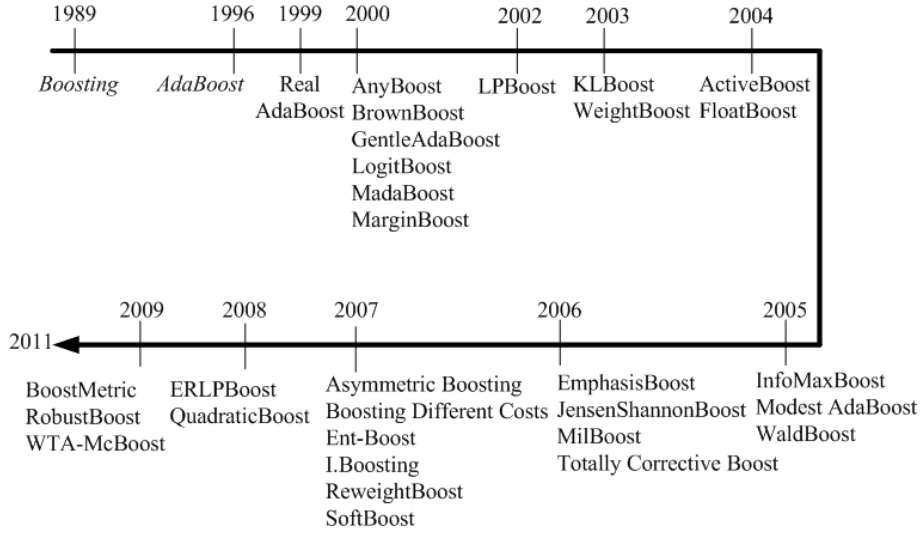


Fig. 3.8 A (possibly incomplete) timeline of AdaBoost variants for supervised learning of binary classifiers, as of 2011.

Algorithm 5 Real AdaBoost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M **do**
 - 3: Fit the class probability estimate $p_m(\mathbf{x}) = \hat{P}_w(y = 1|\mathbf{x})$, using w_i .
 - 4: Set $H_m = \frac{1}{2} \log((1 - p_m(\mathbf{x}))p_m(\mathbf{x})) \in \mathcal{R}$.
 - 5: Update the weights: $w_i \leftarrow w_i \exp(-y_i H_m(\mathbf{x}_i))$
 - 6: Renormalize to weights.
 - 7: **end for**
 - 8: Final classifier: $H(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$.
-

Comparing Real AdaBoost with AdaBoost, we see that the major differences are in lines 3 and 4. In the Real AdaBoost algorithm, these steps consist of computing and using estimates of the probabilities that each training pattern belongs to a class, under the current weight distribution. Standard AdaBoost classifies the input patterns and computes the weighted error rate.

3.4.1.2 Logit Boost

The Logit Boost variant consists of using adaptive Newton steps to fit an additive logistic model [50, 51]. Instead of minimizing the exponential loss, Logit Boost

minimize the logistic loss (negative conditional log-likelihood). Algorithm 6 details the Logit Boost algorithm.

Algorithm 6 Logit Boost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M and while $H_m \neq 0$ **do**
 - 3: Compute the working response $z_i = \frac{y_i^* - p(\mathbf{x}_i)}{p(\mathbf{x}_i)(1 - p(\mathbf{x}_i))}$ and weights $w_i = p(\mathbf{x}_i)(1 - p(\mathbf{x}_i))$.
 - 4: Fit $H_m(\mathbf{x})$ by a weighted least-squares of z_i to \mathbf{x}_i , with weights w_i .
 - 5: Set $H(\mathbf{x}) = H(\mathbf{x}) + \frac{1}{2}H_m(\mathbf{x})$ and $p(\mathbf{x}) = \frac{\exp(H(\mathbf{x}))}{\exp(H(\mathbf{x})) + \exp(-H(\mathbf{x}))}$.
 - 6: **end for**
 - 7: Final classifier: $H(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$.
-

3.4.1.3 Gentle AdaBoost

The Gentle AdaBoost [51] algorithm improves over Real AdaBoost by using Newton steps, providing a more reliable and stable ensemble, since it puts less emphasis on outliers. Instead of fitting a class probability estimate, Gentle AdaBoost (described in Algorithm 7) uses weighted least-squares regression [57] at each iteration. The main difference between Gentle and Real AdaBoost is on the use of the estimates of the weighted class probabilities in order to perform the update. The algorithm is *gentle* because it is considered to be both conservative and more stable as compared to Real AdaBoost. Gentle AdaBoost does not require the computation of log-ratios which can be numerically unstable (since they involve quotients, maybe with the denominator approaching zero). Experimental results on benchmark data show that the conservative Gentle AdaBoost has similar performance to Real AdaBoost and Logit Boost, and in many cases outperforms these other two variants.

3.4.1.4 Modest AdaBoost

The Modest AdaBoost algorithm [117] is known to have less generalization error and higher training error, as compared to Real and Gentle AdaBoost variants. Algorithm 8 shows the details of Modest AdaBoost, which as compared to the previous variants uses a different weighting scheme for the correctly and incorrectly classified patterns, using an “inverted” distribution.

Algorithm 7 Gentle AdaBoost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

-
- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M **do**
 - 3: Train $H_m(\mathbf{x})$ by weighted least-squares of y_i to \mathbf{x}_i , with weights w_i .
 - 4: Update $H(\mathbf{x}) \leftarrow H(\mathbf{x}) + H_m(\mathbf{x})$.
 - 5: Update $w_i \leftarrow w_i \exp(-y_i H_m(\mathbf{x}_i))$ and renormalize to $\sum_i w_i = 1$.
 - 6: **end for**
 - 7: Final classifier: $H(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$.
-

Algorithm 8 Modest AdaBoost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

-
- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M and while $H_m \neq 0$ **do**
 - 3: Train $H_m(\mathbf{x})$ by weighted least-squares of y_i to \mathbf{x}_i , with weights w_i .
 - 4: Compute “inverted” distribution $\bar{w}_i = (1 - w_i)$ and renormalize to $\sum_i \bar{w}_i = 1$.
 - 5: Compute $P_m^{+1} = P_w(y = +1, H_m(\mathbf{x}))$, $\bar{P}_m^{+1} = P_{\bar{w}}(y = +1, H_m(\mathbf{x}))$.
 - 6: Compute $P_m^{-1} = P_w(y = -1, H_m(\mathbf{x}))$, $\bar{P}_m^{-1} = P_{\bar{w}}(y = -1, H_m(\mathbf{x}))$.
 - 7: Set $H_m(\mathbf{x}) = (P_m^{+1}(1 - P_m^{+1}) - P_m^{-1}(1 - P_m^{-1}))$.
 - 8: Update $w_i \leftarrow w_i \exp(-y_i H_m(\mathbf{x}_i))$ and renormalize to $\sum_i w_i = 1$.
 - 9: **end for**
 - 10: Final classifier: $H(\mathbf{x}) = \text{sign} \left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$.
-

The standard distribution w_i assigns high weights to training samples misclassified by earlier steps. On the contrary, \bar{w}_i gives higher weights to samples that are already correctly classified by earlier steps.

Lines 5 and 6 deal with the direct and “inverted” distributions, using the expressions $P_m^{+1} = P_w(y = +1, H_m(\mathbf{x}))$ and $P_m^{-1} = P_w(y = -1, H_m(\mathbf{x}))$; these expressions compute how good is the current weak learner at predicting class labels. On the other hand, the expressions $\bar{P}_m^{+1} = P_{\bar{w}}(y = +1, H_m(\mathbf{x}))$ and $\bar{P}_m^{-1} = P_{\bar{w}}(y = -1, H_m(\mathbf{x}))$ estimate how well our current classifier $H_m(\mathbf{x})$ is working on the data that has been correctly classified by previous steps.

The update $H_m(\mathbf{x}) = (P_m^{+1}(1 - P_m^{+1}) - P_m^{-1}(1 - P_m^{-1}))$ decreases weak classifiers contribution, if it works “too well” on data that has been already correctly classified with high margin. This way, the algorithm is named *Modest* because the classifiers tend to work only in their domain, as defined by w_i .

3.4.1.5 Float Boost

The Float Boost [71, 72] variant is composed of the following stages: 1-initialization; 2-forward inclusion; 3-conditional exclusion; 4-output. All of these stages, with the exception of stage 3, are similar to those of AdaBoost and other variants as discussed so far. The novelty here is the conditional exclusion stage, in which the least significant weak classifier is removed from the set of classifiers, subject to the condition that the removal leads to a error below some threshold. The Float Boost algorithm details are described as Algorithm 9.

3.4.1.6 Emphasis Boost

The Emphasis Boost variant uses a *weighted emphasis* (WE) function [52]. Each input pattern is weighted according to a criterion (parameterized by λ), through the WE function, in such a way that the training process focuses on the “critical” patterns (near the classification boundary) or on the quadratic error of each pattern. Algorithm 10 presents the details of Emphasis Boost.

The WE function is defined by

$$w_i = \exp \left(\lambda \left(\sum_{j=1}^m (\alpha_j H_j(x_i) - y_i)^2 \right) - (1 - \lambda) \left(\sum_{j=1}^m H_j(x_i) \right)^2 \right) \quad (3.8)$$

and controls where the emphasis is placed. This flexible formulation allows choosing how much to consider the *proximity* terms by means of a weighting parameter ($0 \leq \lambda \leq 1$). This way, we have a *boosting by weighting boundary and erroneous samples* technique. Regarding the value of λ , three particular cases are interesting enough to be considered:

- $\lambda = 0$, focus on the “critical” patterns because only the “proximity” to the boundary is taken into account

$$w_i = \exp \left[- \left(\sum_{j=1}^m H_j(\mathbf{x}_i) \right)^2 \right]. \quad (3.9)$$

- $\lambda = 0.5$, we get the classical Real AdaBoost emphasis function

$$w_i = \exp \left[\left(\frac{\sum_{j=1}^m (H_j(\mathbf{x}_i) - y_i)^2}{2} \right) - \frac{(\sum_{j=1}^m H_j(\mathbf{x}_i))^2}{2} \right]. \quad (3.10)$$

- $\lambda = 1$, the emphasis function only pays attention to the quadratic error of each pattern

$$w_i = \exp \left[\sum_{j=1}^m (H_j(\mathbf{x}_i) - y_i)^2 \right]. \quad (3.11)$$

Algorithm 9 Float Boost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.
 N examples $N = a + b$; a examples have $y_i = +1$ and b examples have $y_i = -1$.
 $J(H_M)$, the cost function and the maximum acceptable cost J^* .

Output: A classifier $H : \mathcal{X} \rightarrow \{-1, +1\}$.

```

1:                                     {1 - Initialization stage.}
2: Initialize the weights  $w_i^{(0)} = 1/2a$ , for those examples with  $y_i = +1$ .
3: Initialize the weights  $w_i^{(0)} = 1/2b$ , for those examples with  $y_i = -1$ .
4:  $J_m^{\min} = J^*$   $m = \{1, \dots, M_{\max}\}$ .
5:  $M = 0$ ,  $\mathcal{H}_0 = \{\}$ .
6:                                     {2 - Forward inclusion stage.}
7:  $M \leftarrow M + 1$ .
8: Learn  $H_m(x)$  and  $\alpha_m$ .
9: Update  $w_i^{(M)} \leftarrow w_i^{(M-1)} \exp(-y_i \alpha_m H_m(\mathbf{x}_i))$  and renormalize to  $\sum_i w_i = 1$ .
10:  $\mathcal{H}_M = \mathcal{H}_{M-1} \cup \{H_M\}$ .
11: if  $J_M^{\min} > J(H_M)$  then
12:    $J_M^{\min} = J(H_M)$ .
13: end if
14:                                     {3 - Conditional exclusion stage.}
15:  $h' = \arg \min_{h \in \mathcal{H}_M} J(H_M - h)$ .
16: if  $J(H_M - h') < J_{M-1}^{\min}$  then
17:    $\mathcal{H}_{M-1} = \mathcal{H}_M - h'$ .
18:    $J_{M-1}^{\min} = J(H_M - h')$ 
19:    $M \leftarrow M - 1$ 
20:   if  $h' = h'_m$  then
21:     Recalculate  $w_i^{(j)}$  and  $h_j$  for  $j = \{m', \dots, M\}$ .
22:     Goto line 15.
23:   else
24:     if  $M = M_{\max}$  or  $J(\mathcal{H}_M) < J^*$  then
25:       Goto line 32.
26:     else
27:       Goto line 7.
28:     end if
29:   end if
30: end if
31:                                     {4 - Output stage.}
32: Final classifier:  $H(\mathbf{x}) = \text{sign} \left( \sum_{j=1}^M \alpha_j H_j(\mathbf{x}) \right)$ .
```

The key issue with this algorithm is the choice of λ .

3.4.1.7 Reweight Boost

In the Reweight Boost variant [91], the weak classifiers are stumps (decision trees with a single node). The main idea is to consider as base classifier for boosting,

Algorithm 10 Emphasis Boost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.
 λ , weighting parameter ($0 \leq \lambda \leq 1$).

Output: $H(\mathbf{x})$, a classifier suited for the training set.

-
- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M and while $H_m \neq 0$ **do**
 - 3: Fit a classifier $H_m(\mathbf{x})$ to the training data using weights w_i .
 - 4: Let $\text{err}_m = \sum_{i=1}^N w_i y_i H_m(\mathbf{x}_i) / \sum_{i=1}^N w_i$.
 - 5: Compute $\alpha_m = 0.5 \log((1 + \text{err}_m) / (1 - \text{err}_m))$.
 - 6: Set $w_i = \exp\left(\lambda \left(\sum_{j=1}^m (\alpha_j H_j(\mathbf{x}_i) - y_i)^2\right) - (1 - \lambda) \left(\sum_{j=1}^m H_j(\mathbf{x}_i)\right)^2\right)$.
 - 7: Renormalize to $\sum_i w_i = 1$.
 - 8: **end for**
 - 9: Final classifier: $H(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x})\right)$.
-

not only the last weak classifier, but a classifier formed by the last r selected weak classifiers, using a classifier reuse technique. Algorithm 11 presents the details of the Reweight Boost variant.

Algorithm 11 Reweight Boost

Input: Dataset $Z = \{z_1, z_2, \dots, z_N\}$, with $z_i = (\mathbf{x}_i, y_i)$, where $\mathbf{x}_i \in \mathcal{X}$ and $y_i \in \{-1, +1\}$.
 M , the maximum number of classifiers.
 r , the last r selected weak classifiers.

Output: $H(\mathbf{x})$, a classifier suited for the training set.

-
- 1: Initialize the weights $w_i = 1/N$, $i \in \{1, \dots, N\}$.
 - 2: **for** $m = 1$ to M and while $H_m \neq 0$ **do**
 - 3: Fit a classifier $H_m(\mathbf{x})$ to the training data using weights w_i .
 - 4: Get combined classifier H'_t from $H_t, H_{t-1}, \dots, H_{\max(t-r, 1)}$.
 - 5: Let $\text{err}_m = \sum_{i=1}^N w_i y_i H_m(\mathbf{x}_i) / \sum_{i=1}^N w_i$.
 - 6: Compute $\alpha_m = 0.5 \log((1 - \text{err}_m) / \text{err}_m)$.
 - 7: Set $w_i \leftarrow w_i \exp(-\alpha_m y_i H'_t(\mathbf{x}_i))$.
 - 8: Renormalize to $\sum_i w_i = 1$.
 - 9: **end for**
 - 10: Final classifier: $H(\mathbf{x}) = \text{sign}\left(\sum_{j=1}^M \alpha_j H_j(\mathbf{x})\right)$.
-

3.4.1.8 Other variants

For the sake of both completeness of this chapter and fairness to the many authors of AdaBoost variants, in this subsection we describe further variants for binary classification. We show the name of each variant as well as its main characteristics.

The **KLBoost** [73] variant uses *Kullback-Leibler* (KL) divergence and operates as follows. First, classification is based on the sum of histogram divergences along corresponding global and discriminating linear features. Then, these linear KL features, are iteratively learned by maximizing the projected KL divergence in a boosting manner. Finally, the coefficients to combine the histogram divergences are learned by minimizing the recognition error, once a new feature is added to the classifier. This contrasts with conventional AdaBoost, in which the coefficients are empirically set. Because of these properties, KLBoosting classifier generalizes very well and has been applied to high-dimensional spaces of image data.

One of the experimental drawbacks of AdaBoost is that it can not improve the performance of *naïve Bayes* (NB) [33, 128] classifier as expected. **Active-Boost** [122] overcomes this difficulty by using active learning to mitigate the negative effect of noisy data and introduce instability into the boosting procedure. Empirical studies on a set of natural domains show that ActiveBoost has clear advantages with respect to the increasing of the classification accuracy of Naive Bayes when compared against Adaboost.

The **Jensen-Shannon Boosting** [60] incorporates *Jensen-Shannon* (JS) divergence into AdaBoost. JS divergence is advantageous in that it provides a more appropriate measure of dissimilarity between two classes and it is numerically more stable than other measures such as KL divergence.

Infomax Boosting [75] is an efficient feature pursuit scheme for boosting. It is based on the infomax principle, which seeks optimal feature that achieves maximal mutual information with class labels. Direct feature pursuit with infomax is computationally prohibitive, so an efficient gradient ascent algorithm is proposed, based on the quadratic mutual information, nonparametric density estimation and fast Gauss transform. The feature pursuit process is integrated into a boosting framework as infomax boosting. It is similar to Real AdaBoost, but with the following exceptions:

- features are general linear projections;
- generates optimal features;
- uses KL divergence to select features;
- finer tuning on the coefficients.

Ent-Boost [67] uses entropy measures. The class entropy information is used to automatically subspace splitting and optimal weak classifier selection. The number of bins is estimated through a discretization process. KL divergence is applied to probability distribution of positive and negative samples, to select the best weak classifier in the weak classifier set.

The **MadaBoost** [29] algorithm consists on a modification of the weighting scheme of AdaBoost. This variant mitigates the problems that AdaBoost suffers from noisy data, improving its performance.

The **SoftBoost** algorithm [124] is a totally corrective algorithm which optimizes the soft margin and tries to produce a linear combination of hypotheses with the maximum one. The term *soft* means that the algorithm does not concentrate too much on outliers and hard to classify examples. It allows them to lie below the margin (with wrong predictions) but penalizes them linearly via slack variables. Soft-

Boost tries to avoid the problem of overfitting as in AdaBoost when using training data with high degree of noise.

The *linear programming boosting (LPBoost)* [26] algorithm maximizes the margin between training samples of different classes; this way, it belongs to the class of margin-maximizing supervised classification algorithms. The boosting task consists of constructing a learning function in the label space that minimizes misclassification error and maximizes the soft margin, formulated as a linear program which can be efficiently solved using column generation techniques, developed for large-scale optimization problems. Unlike gradient boosting algorithms, which may converge in the limit only, LPBoost converges in a finite number of iterations to a global solution, being computationally competitive with AdaBoost. The optimal solutions of LPBoost are very sparse in contrast with gradient based methods. Empirical findings show that LPBoost converges quickly, often faster than other formulations.

LPBoost performs well on natural data, but there are cases where the number of iterations is linear in the number of training samples instead of logarithmic. By simply adding a relative entropy regularization to the linear objective of LPBoost, we get *entropy regularized LPBoost* ERLPBoost [125], for which there is a logarithmic iteration bound. As compared to a previous algorithm, named SoftBoost, it has the same iteration bound and better generalization error. ERLPBoost does not suffer from this problem and has a simpler motivation. A detailed theoretical and experimental comparison between LPBoost and AdaBoost can be found in [68].

The **MarginBoost** algorithm [80] is a variant of the more general algorithm **AnyBoost** [80]. MarginBoost is also a general algorithm. It chooses a combination of classifiers to optimize the sample average of any cost function of the margin. MarginBoost performs gradient descent in function space, at each iteration choosing a base classifier to include in the combination so as to maximally reduce the cost function. As in AdaBoost, the choice of the base classifier corresponds to a minimization problem involving weighted classification error. That is, for a certain weighting of the training data, the base classifier learning algorithm attempts to return a classifier that minimizes the weight of misclassified training examples.

The general class of algorithms named AnyBoost consists of gradient descent algorithms for choosing linear combinations of elements of an inner product space so as to minimize some functional cost. Each component of the linear combination is chosen to maximize a certain inner product. In MarginBoost, this inner product corresponds to the weighted training error of the base classifier.

Brown Boost [44] uses a non-monotonic weighting function such as examples far from the boundary decrease in weight, trying to achieve a given target error rate. It de-emphasizes outliers when it seems clear that they are too hard to classify correctly, being an adaptive version of Freund's boost-by-majority algorithm [43]. This variant reveals an intriguing connection between boosting and Brownian motion.

The **Weight Boost** algorithm [62] uses input-dependent weighting factors for weak learners. It tries to cope with two possible problems of AdaBoost: suffer from overfitting, especially for noisy data; the assumption that the combination weights are fixed constants and therefore does not take particular input patterns into consideration. A learning procedure which is guaranteed to minimize training errors is

devised. Empirical studies show that Weight Boost almost always achieves a considerably better classification accuracy than AdaBoost. Furthermore, experiments on data with artificially controlled noise indicate that the Weight Boost is more robust to noise than AdaBoost.

Asymmetric Boosting [78] is a cost-sensitive extension of boosting. It is derived from decision-theoretic principles, which exploit the statistical interpretation of boosting to determine a principled extension of the boosting loss. Similarly to AdaBoost, the cost-sensitive extension minimizes this loss by gradient descent on the functional space of convex combinations of weak learners, and produces large margin detectors. Asymmetric boosting is fully compatible with AdaBoost, in the sense that it becomes the latter when errors are weighted equally.

In [58] we have an asymmetric boosting method, ***Boosting with Different Costs***. The motivation is as follows; traditional boosting methods assume the same cost for misclassified instances from different classes, and in this way focus on good performance with respect to overall accuracy. This method is more generic than AdaBoost, and is designed to be more suitable for problems where the major concern is a low false positive (or negative) rate, such as SPAM filtering.

The **Quadratic Boost** [85] algorithm improves AdaBoost with a quadratic combination of base classifiers. It operates by constructing an intermediate learner on the combined linear and quadratic terms. A new method for iterative optimization is proposed; first a classifier is trained by randomizing the labels of the training examples. Subsequently, the input learner is called repeatedly with a systematic update of the labels of the training examples in each round. The quadratic boosting algorithm converges under the condition that the given base learner minimizes the empirical error. The experimental results show that quadratic boosting compares favorably with AdaBoost on large data sets at the cost of the training time.

The **WaldBoost** [104] variant has near optimal time and error rate trade-off. It integrates the AdaBoost algorithm for measurement selection and ordering and the joint probability density estimation, with the optimal sequential probability ratio test decision strategy. It is suited for computer vision classification problems, in which both the error and time characterize the quality of a decision.

In [74] feature reweighting is integrated into the boosting scheme, which not only weights the samples but also weights the features iteratively; it is named **I.Boosting**. To avoid overfitting problems, a relevance feedback mechanism is applied into the boosting framework. I.Boosting is implemented using *adaptive discriminant analysis* (ADA) as base classifiers. The experimental results show the superior performance of I.Boosting over AdaBoost.

In [45] we have a new boosting algorithm, motivated by the large margins theory for boosting. The experimental results point out that the new algorithm is significantly more robust against label noise than existing boosting algorithm.

The algorithm proposed in [84] combines the base learners with symmetric functions. Among its properties of practical relevance, we have significant resistance against noise, and its efficiency even in an agnostic learning setting. Experimental results show the reliability of the classifiers built.

The **MilBoost** [121] variant uses cost functions from the *multiple instance learning* (MIL) literature combined with the AnyBoost framework. The feature selection criterion of MILBoost is modified to optimize the performance of the Viola-Jones cascade method for object detection (see subsection 3.6.1). Experiments with this variant show improvement on the detection rate, as compared to previous approaches. This increased detection rate is a consequence of simultaneously learning the locations and scales of the objects in the training set along with the parameters of the classifier.

The **totally corrective boosting** [126], the weight update of each patterns is analyzed as the minimization of the relative entropy, subject to linear constraints. The algorithm is “totally corrective” in the sense that it takes into account the outputs of all the past weak learners; the “corrective” versions only take into account the last weak learner results. A connection with margin maximization is also shown for totally corrective versions. The experimental results show that the totally corrective versions of AdaBoost attain smaller combinations of weak learners than the corrective ones, being competitive with LPBoost (itself a totally corrective boosting algorithm with no regularization, for which there is no iteration bound known). An asymmetric totally corrective boosting approach for real-time object detection is proposed in [123].

In [103] the **BoostMetric** algorithm is proposed. The goal of this algorithm is to learn a semidefinite metric using boosting techniques. It is a generalization of AdaBoost in the sense that the weak learner is a matrix instead of a classifier, being simple and efficient. It attains better performance than many existing metric learning methods.

A boosting algorithm called *winner-take-all multiple category boosting* (**WTA-McBoost**) was proposed in [133]. On the learning process, the example subcategory labels are modified in order to make better object/non-object decision. Multiple subcategory boosting classifiers are learned simultaneously with the assumption that the final classification of an example will only be determined by the highest score of all the subcategory classifiers (the winner will take all). The subcategory labels of the examples are dynamically assigned in this process, reducing the risk of having outliers in each subcategory. The WTA-McBoost algorithm uses confidence-rated prediction with asymmetric cost and is thus very efficient to train and test. The algorithm is successfully applied by building a multi-view face detector.

The standard boosting procedure is extended to train a two-layer classifier dedicated to handwritten character recognition [42]. This learning scheme relies on a hidden layer and an output layer to obtain a final classification decision. The classical AdaBoost procedure is extended to train a multi-layered structure by propagating the error through the output layer. This extension allows for the selection of optimal weak learners by minimizing a weighted error, in both the output layer and the hidden layer.

3.4.2 Multi-class Variants

Since the first binary classification versions of Adaboost, several generalization of this algorithm to the multi-class case have been proposed. As a result of this direction of research, several multi-class AdaBoost variants have been proposed, as depicted in the timeline shown in Fig. 3.9. Similarly to what we did in subsection 3.4.1.8, for each variant we will point out its main features as well as the connections among them. Many variants address multi-class classification problems as the concatenation of binary problems (a multi-class classification problem can be reformulated as a set of binary problems), while other (more recent) variants apply and combine multi-class classifiers directly.

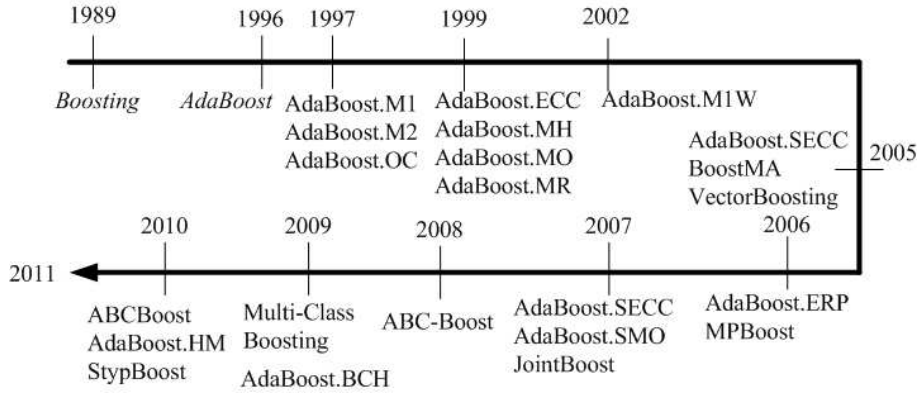


Fig. 3.9 A (possibly incomplete) timeline of AdaBoost variants for supervised learning, on multi-class problems, as of 2011.

AdaBoost.M1 and **AdaBoost.M2** [49] are multi-class extensions of (Discrete) AdaBoost (Algorithm 4). They differ between themselves in the way they treat each class. In the M1 variant, the weight of a base classifier is a function of the error rate. In M2, the sampling weights are increased for instances for which the pseudo-loss exceeds 0.5. The **AdaBoost.M1W** [38] algorithm changes AdaBoost.M1 as follows. In AdaBoost.M1 the weight of a base classifier is a function of the error rate. For AdaBoost.M1W this function is such that it gets positive, if the error rate is less than the error rate of random guessing. **BoostMA** [39] is also a simple modification of AdaBoost.M2 with the advantage that the base classifier minimizes the confidence-rated error, whereas for AdaBoost.M2 the base classifier should minimize the pseudo-loss. This makes BoostMA more easily applicable to already existing base classifiers; it also tends to converge faster than AdaBoost.M2.

AdaBoost.MH [99] is a multi-class and *multi-label*² version of AdaBoost based on Hamming loss [99]. AdaBoost.MH generalizes AdaBoost, being tailored for

² A given instance can be classified into one or more classes.

multi-label *text categorization* (TC) tasks with decision stumps as weak learners. **MPBoost** [40] further improves AdaBoost.MH augmenting its efficiency by performing a multiple pivot selection at each boosting iteration. Both these algorithms use binary features.

The **AdaBoost.MO** algorithm [99] performs a stage-wise functional gradient descent procedure on a given cost function. **AdaBoost.MR** [99] is a multi-class, multi-label version of AdaBoost based on ranking loss. **AdaBoost.OC** [95], where OC stands for output codes and **AdaBoost.ECC** [55], where ECC stands for error-correcting codes are similar algorithms; AdaBoost.OC is a shrinkage version of AdaBoost.ECC, which performs a stage-wise functional gradient descent procedure on an exponential loss cost function.

The **Vector Boosting** [59] algorithm is an extension of the Real AdaBoost in which both its weak learner and its final output are vectors rather than scalars. The idea of Vector Boosting comes from the *multi-class multi-label* (MCML) version of the Real AdaBoost, which assigns a set of labels for each sample and decomposes the original problem into k orthogonal binary ones. The major problem of this algorithm is that for each binary classification problem, a sample is regarded as either positive or negative. However in many complicated cases, it is not tenable since some samples are neither positive nor negative for certain binary classification problems of which they are independent, which makes the MCML version of Real AdaBoost inapplicable.

AdaBoost.ERP [69] is AdaBoost.ECC with repartitioning. This algorithm improves two well-known issues of the quality of the ensemble learned by AdaBoost.ECC: the performance of the base learner; the error-correcting ability of the coding matrix. A coding matrix with strong error-correcting ability may not be overall optimal if the binary problems are too hard for the base learner. A trade-off between error-correcting and base learning is then proposed. The coding matrix is modified according to the learning ability of the base learner.

In [107] shrinkage is applied as regularization in AdaBoost.MO and AdaBoost.ECC and leads two new algorithms named **AdaBoost.SMO** and **AdaBoost.SECC** (the shrinkage versions of MO and ECC, respectively). A similar proposal for AdaBoost.SECC can also be found in [108].

In [136] we have a new algorithm named **Multi-class AdaBoost** that directly extends the AdaBoost algorithm to the multi-class case without reducing it to multiple two-class problems. The algorithm is equivalent to a forward stagewise additive modeling algorithm that minimizes a novel exponential loss for multi-class classification. The algorithm is highly competitive in terms of misclassification error rate.

The **AdaBoost.BCH** algorithm [64] is a multi-class boosting algorithm which solves a C class problem by using $C - 1$ binary classifiers arranged by a hierarchy that is learned on the classes based on their closeness to one another. AdaBoost is then applied to each binary classifier. AdaBoost.BCH requires less computation than AdaBoost.MH, with better or comparable generalization.

In [70] the concept of *adaptive base class boost* (**ABC-Boost**) for multi-class classification is addressed deriving **ABC-MART**, a concrete implementation of ABC-Boost. For binary classification, ABC-MART recovers MART and for multi-

class classification, ABC-MART considerably improves MART, as evaluated on several public datasets.

AdaBoost.HM was proposed in 2010 [63]. It is based on *hypothesis margin* and directly combines multi-class weak classifiers, instead of learning binary weak learners. The hypothesis margin maximizes the output about the positive class and minimizes the maximal outputs about the negative classes. Upper bounds on the training error of AdaBoost.HM are derived and compared against AdaBoost.M1 upper bounds. The weak learners are feedforward neural networks. AdaBoost.HM yields higher classification accuracies than the AdaBoost.M1 and the AdaBoost.MH algorithms, being computationally efficient in training.

Recently, a totally corrective multi-class boosting was proposed [56]. After an analysis of some methods that extend two-class boosting to multi-class, a column-generation based totally-corrective framework for multi-class boosting learning is derived, using the Lagrange dual problems. Experimental results show that the new algorithms have comparable generalization capability but converge much faster than their counterparts.

StypBoost [127] is bilinear boosting algorithm, which extends the multi-class boosting framework of JointBoost to optimize a bilinear objective function. This allows style parameters to be introduced to aid classification, where style is any factor which the classes vary with systematically, modeled by a vector quantity. The algorithm allows learning with different styles. It is applied successfully to two object class segmentation tasks: road surface segmentation and general scene parsing.

JointBoost [105, 111] is a method where boosted one-versus-all classifiers are trained jointly and are forced to share features. It has been demonstrated to lead both to higher accuracy and smaller classification time, compared to using one-versus-all classifiers that were trained independently and without sharing features.

3.4.2.1 Analysis of multi-class boosting algorithms

In [107] the AdaBoost.MO, AdaBoost.OC, and AdaBoost.ECC algorithms are studied. It is shown that MO and ECC perform stage-wise functional gradient descent on a cost function defined over margin values, and that OC is a shrinkage version of ECC. The AdaBoost.SMO and AdaBoost.SECC are the shrinkage versions of MO and ECC, respectively.

An unifying framework for studying the solution of multi-class categorization problems, by reducing them to multiple binary problems, that are then solved using a margin-based binary learning algorithm is presented [2]. The proposed framework unifies some of the most popular approaches in which each class is compared against all others, or in which all pairs of classes are compared to each other, or in which output codes with error-correcting properties are used. A general method for combining the classifiers generated on the binary problems is proposed. A generic empirical multi-class loss bound given the empirical loss of the individual binary learning algorithms is proven. The scheme and the corresponding bounds apply to

many popular classification learning algorithms including SVM, AdaBoost, regression, logistic regression, and decision-tree algorithms. A multi-class generalization error analysis for general output codes with AdaBoost is provided.

The ability of boosting to achieve drastic improvements compared to the individual weak learners has been noticed by several researchers. For two-class problems it has been observed that AdaBoost, is quite unaffected by overfitting. However, for the case of noisy data, it is also known that AdaBoost can be improved considerably by introducing some regularization technique. In speech-related problems one often considers multi-class problems and boosting formulations have been used successfully to solve them. Under this context, [89] reviews and extends the existing multi-class boosting algorithms to derive new boosting algorithms, which are more robust against outliers and noise in the data; these algorithms are also able to exploit prior knowledge about relationships between the classes.

In [108] a new interpretation of AdaBoost.ECC and AdaBoost.OC is presented. AdaBoost.ECC performs stage-wise functional gradient descent on a cost function, defined in the domain of margin values; AdaBoost.OC is a shrinkage version of AdaBoost.ECC. AdaBoostBCH has slower training and higher generalization ability as compared to AdaBoost.ECC [64].

3.5 Boosting for Semi-supervised Learning

Semi-Supervised learning (SSL) [15] as attracted considerable research effort in the last few years. In many learning problems, we have a large amount of data available, but only a subset of it is labeled. In this section, we review several boosting algorithms for SSL, show in the timeline of Fig. 3.10.

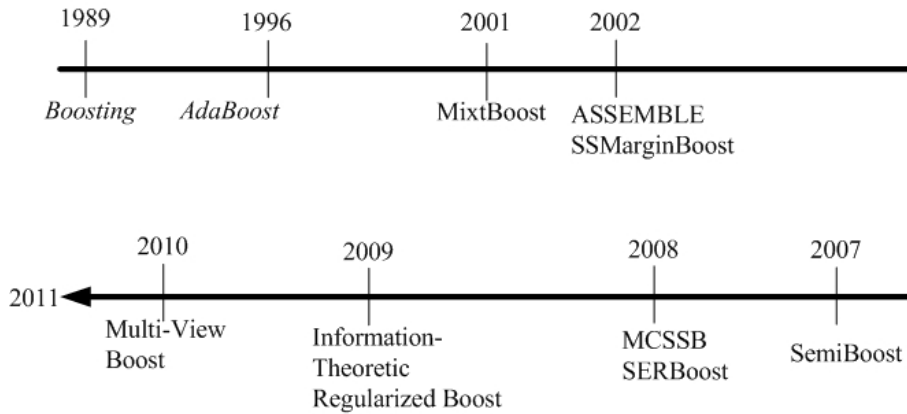


Fig. 3.10 A (possibly incomplete) timeline of AdaBoost variants for semi-supervised learning on binary and multi-class problems, as of 2011.

3.5.1 *MixtBoost*

The MixtBoost algorithm [54] was the first variant of AdaBoost to for SSL; the authors address the question: *can boosting be adapted for SSL learning?* The base classifiers are mixture models, thus MixtBoost can be seen as boosting of mixture models.

The main ingredients of AdaBoost are the *loss* and the *margin*. The simplest way to generalize to SSL is to define these quantities for unlabeled data. This generalization to unlabeled data should not affect the labeled examples and should penalize inconsistencies between the classifier output and the available information. A loss definition is proposed for unlabeled data, from which margins are defined. The missing labels are interpreted as the absence of class information with the key idea that the pattern belongs to a class, but the class is unknown.

3.5.2 *SSMarginBoost*

The SSMarginBoost algorithm is an extension of MarginBoost to SSL [25] that explores the *clustering* assumption of SSL [15] and the large margin criterion. The margin definition is extended to unlabeled data and the gradient descent algorithm for optimizing the resulting margin cost function is derived. SSMarginBoost can be applied with any base classifier able to handle unlabeled data, by means of mixture models trained with an *expectation-maximization* (EM) algorithm [27, 33].

3.5.3 *ASSEMBLE*

In [6], an adaptive semi-supervised ensemble method, named ASSEMBLE, was proposed. The method constructs ensembles based on both labeled and unlabeled data, by alternating between assigning “pseudo-classes” to the unlabeled data using the existing ensemble and constructing the next base classifier using both the labeled and pseudo-labeled data. This algorithm corresponds to maximizing the classification margin in hypothesis space as measured on both the labeled and unlabeled data. Unlike alternative approaches, ASSEMBLE does not require a semi-supervised learning method for the base classifier. It can be used in conjunction with any cost-sensitive classification algorithm for both two-class and multi-class problems. As in SSMarginBoost, ASSEMBLE adopts the MarginBoost notation and strategy adapted to the margin measured on both the labeled and unlabeled data; the key difference is that ASSEMBLE assigns “pseudoclasses” to the unlabeled data.

Moreover, ASSEMBLE using decision trees won the NIPS 2001 Unlabeled Data Competition. It achieves good results on several benchmark datasets using both decision trees and neural networks.

3.5.4 *SemiBoost*

The SemiBoost [77] algorithm was proposed as a boosting framework aiming at improving the classification accuracy of any given supervised learning algorithm by using the available unlabeled examples. The main advantages of SemiBoost over previous approaches are: 1) performance improvement of any supervised learning algorithm, using unlabeled data; 2) efficient computation by the iterative boosting algorithm and 3) exploiting both the SSL *manifold* and *cluster* assumptions [15].

An empirical study on 16 different datasets and text categorization demonstrates that SemiBoost improves the performance of several commonly used supervised learning algorithms, by using a large number of unlabeled examples.

3.5.5 *Multi-Class Semi-Supervised Boosting*

The *multi-class semi-supervised boosting* (MCSSB) algorithm was proposed in [113]. Compared to the existing semi-supervised boosting methods, MCSSB has the advantage to exploit both classification confidence and similarities among examples, when deciding the pseudo-labels for unlabeled examples. This way, it overcomes the shortcoming of the multi-class approach *one-against-the-rest* applied on binary classifiers. Empirical evidence on several datasets shows that the MCSSB algorithm performs better than the previous algorithms for SSL.

3.5.6 *SERBoost*

The problem of bad scaling behavior of many SSL methods on large scale vision problems is addressed in [92]. Based on the *expectation regularization* (ER) principle, the SERBoost SSL boosting algorithm is proposed. It can be applied to large scale vision problems and its complexity is dominated by the base learners. The algorithm provides a margin regularizer for the boosting cost function and shows a principled way of utilizing prior knowledge. As compared to supervised and semi-supervised methods, SERBoost shows improvement both in terms of classification accuracy and computational speed.

3.5.7 *Information Theoretic Regularization Boosting*

An SSL boosting algorithm that incrementally builds linear combinations of weak classifiers through generic functional gradient descent, using both labeled and unlabeled training data, was proposed in [134]. The approach is based on extending the information regularization framework to boosting, bearing loss functions that com-

bine log loss on labeled data with the information-theoretic measures to encode unlabeled data. Even though the information-theoretic regularization terms make the optimization non-convex, a simple sequential gradient descent optimization algorithm is applied. This approach attains good results on synthetic, benchmark and real world tasks as compared to supervised and semi-supervised boosting algorithms.

3.5.8 *Multi-view boosting*

A multi-view boosting algorithm was proposed in [93], that, unlike other approaches, specifically encodes the uncertainties over the unlabeled samples in terms of given priors. Instead of ignoring the unlabeled samples during the training phase of each view, it uses the different views to provide an aggregated prior which is then used as a regularization term inside a semi-supervised boosting method, for multi-class problems. The algorithm uses priors as a regularization component over the unlabeled data. Since the priors may contain a significant amount of noise, a new loss function for the unlabeled regularization is introduced, being robust to noisy priors.

3.5.9 *Extensions on semi-supervised boosting algorithms*

Different strategies have been applied to extend boosting algorithms to SSL problems. Typically, these strategies do not take into account the local smoothness constraints among data into account during ensemble learning. A local smoothness regularizer to semi-supervised boosting algorithms based on the universal optimization framework of margin cost functionals is proposed [17]. This regularizer is applicable to existing SSL boosting algorithms to improve their generalization and speed up their training.

In [18], the problem of using all the three SSL assumptions (*smoothness*, *cluster*, and *manifold*) during boosting is addressed. A novel cost functional consisting of the margin cost on labeled data and the regularization penalty based on unlabeled data is proposed. Thus, minimizing the proposed cost functional with a greedy stagewise optimization procedure leads to a generic boosting framework for SSL.

A local smoothness regularizer for SSL boosting algorithms, based on the universal optimization framework of margin cost functionals, was proposed in [17]. This regularizer is applicable to existing SSL boosting algorithms to improve their generalization and speed up their training.

3.6 Experimental Evaluation

In this section we discuss the application of AdaBoost and its variants on a wide variety of problems. We compare the boosting algorithms with other machine learning techniques, exploiting the theoretical properties of Adaboost.

3.6.1 Successful Applications

Besides its nice theoretical properties, the AdaBoost algorithm and its variants have been found to work very well on problems from different domains. Empirical evidence from many researchers has shown the adequacy of boosting algorithms for real-world problems. This section outlines some successful applications of Adaboost and its variants for binary and multi-class problems. There are many papers which evaluate AdaBoost and its variants for many types of problems; see for instance [5, 28, 32, 61, 76, 88, 102]. It has been shown empirically that AdaBoost with decision trees has excellent performance, being considered the best “off-the-shelf” classification algorithm [5, 57].

The first boosting algorithms were tested on a *optical character recognition* (OCR) problem of optical handwritten digits, with a set of 118 000 instances in boosting multilayer perceptrons [31].

In [120] a **pedestrian detection** system that integrates image intensity information with motion information is proposed. A detection style algorithm scans a detector over two consecutive frames of a video sequence. The detector is trained using AdaBoost to take advantage of both motion and appearance information to detect a walking person. The detector combines two sources of information.

The breast cancer detection problem is addressed in [109]. A data pre-processing, feature selection and Modest AdaBoost algorithm, are applied the breast cancer survival databases in Thailand. For this task, Modest AdaBoost outperforms Real and Gentle AdaBoost variants.

In [100], boosting is applied to **multi-class text categorization** tasks. The approach named BoosTexter has comparable results to other text-categorization algorithms, on a variety of tasks. The BoosTexter system is also applied to **speech categorization** to call-type identification from unconstrained spoken customer responses.

For **face detection**, boosting algorithms have been the most effective of all those developed so far, achieving the best results. They produce classifiers with about the same error rate than neural networks, but they have faster training [71]. Table 3.2 summarizes the use of boosting algorithms for face detection (a *stap* is a decision tree with a single decision node).

In [24] a fast and efficient face detection method has been devised, which relies on the AdaBoost algorithm and a set of Haar Wavelet like features. The face detection problem was been addressed also with **Asymmetric Boosting** [78], where it is shown to outperform a number of previous heuristic proposals for cost-sensitive

Table 3.2 Summary of the use of boosting algorithms for face detection (adapted from [71]).

Face Detector	AdaBoost Variant	Weak Learner
Viola-Jones [118]	Discrete AdaBoost	Stubs
Float Boost	Float Boost [71, 72]	1D Histograms
KLBoost	KLBoost [73]	1D Histograms
Schneiderman	Real AdaBoost [99]	One group of n-D Histograms

boosting. For an updated literature on face detection and the use of boosting and other machine learning techniques, see [131].

A method for selecting edge-type features for **iris recognition** is proposed in [16]. The AdaBoost algorithm is used to select a filter bank from a pile of filter candidates. The decisions of the weak classifiers associated with the filter bank are linearly combined to form a strong classifier. The boosting algorithm can effectively improve the recognition accuracy at the cost of a slight increase on the computation time.

A new approach, proposing two *particle swarm optimisation* (PSO) methods within AdaBoost for **object detection**, for constructing weak classifiers in AdaBoost is proposed in [82]. The experiments show that using PSO for selecting features and evolving associated weak classifiers in AdaBoost is more effective than for selecting features only for this problem.

In [129] a **face recognition** method using AdaBoosted low dimensional and discriminant Gabor features is proposed. AdaBoost is successfully applied to face recognition by introducing the intra-face and extra-face difference space in the Gabor feature space. By using the proposed method, only hundreds of Gabor features are selected. Experiments shown that these hundreds of Gabor features are enough to achieve good performance comparable to that of methods using the complete set of Gabor features.

A **feature selection approach** based on Gabor wavelets and AdaBoost is proposed in [135]. The features are first extracted by a Gabor wavelet transform. For each individual, a small set of significant features are selected by the AdaBoost algorithm from the pool of the Gabor wavelet features. In the feature selection process, each feature is the basis for a weak classifier. In each round of AdaBoost learning, the feature with the lowest error of weak classifiers is selected. The results from the experiment have shown that the approach successfully selects meaningful and explainable features for face verification. The experiments suggest that the feature selection algorithm for face verification selects the features corresponding to the unique characteristics rather than common characteristics, and a large example size statistically benefits AdaBoost feature selection.

The problem of **classifying music by genre** by partitioning songs into smaller pieces and classifying each one separately is addressed in [7]. The choice of features together with an AdaBoost.MH classifier proved to be the most effective method for genre classification at the MIREX 2005 international contest in music information extraction, and the second-best method for recognizing artists.

In [83], 2D cascaded AdaBoost, a novel classifier designing framework, is presented and applied to the **eye localization** problem. There are two cascade classifiers in two directions: the first one is a cascade designed by bootstrapping the positive samples; the second one, as the component classifiers of the first one, is cascaded by bootstrapping the negative samples. The proposed structure is applied to eye localization and evaluated on four public face databases, and extensive experimental results verified the effectiveness, efficiency, and robustness of the proposed method.

AdaBoost can also improve the performance of a strong learning algorithm as proposed in [101]: a neural network based **online character recognition system**. AdaBoost can be used to learn automatically a great variety of writing styles even when the amount of training data for each style varies a lot. The system achieves about 1.4 % error on a handwritten digit database of more than 200 writers.

In [86] the use of boosting and SVM is explored for the **segmentation of white-matter lesions in the MRI scans of human brain**. Simple features are generated from proton density scans. Radial basis function based AdaBoost technique and SVM are employed for this task. The classifiers are trained on severe, moderate and mild cases. The results indicate that the proposed approach can handle MR field inhomogeneities quite well.

A **visual object detection** framework that is capable of processing images extremely rapidly while achieving high detection rates is proposed in [119]. The learning algorithm, based on AdaBoost, selects a small number of critical visual features and yields extremely efficient classifiers. The method combines classifiers in a cascade allowing background regions of the image to be quickly discarded while spending more computation on promising object-like regions. A set of experiments in the domain of **face detection** are presented.

A framework for **classifying face images** using Adaboost and domain-partitioning based classifiers is addressed in [116]. The most interesting aspect of this framework is its ability to build classification systems with high accuracy in dynamical environments, which achieve, at the same time, high processing and training speed. This framework is applied to the specific problem of gender classification using different features, on standard face databases.

In [110] an approach for **image retrieval** using a very large number of highly selective features and efficient online learning is proposed. Our approach is predicated on the assumption that each image is generated by a sparse set of visual “causes” and that images which are visually similar share causes between them. A mechanism for computing a very large number of highly selective features which capture some aspects of this causal structure (with over 45,000 highly selective features). At query time a user selects a few example images, and boosting is used to learn a classification function in this feature space. The boosting procedure learns a simple classifier which only relies on 20 of the features. As a result a very large database of images can be scanned rapidly.

The *boosting-based multimodal speaker detection* (BMSD) algorithm is proposed in [130]. It performs **speaker detection, identifying the active speaker** in a video, which can be very helpful for remote participants to understand the dynamics of the meeting. This algorithm fuses audio and visual information at feature level

by using boosting to select features from a combined pool of both audio and visual features simultaneously. It achieves a very accurate speaker detector with extremely high efficiency.

3.6.1.1 Online Boosting

In the recent years, some attention has been given to *online boosting*, in which the training examples become available one at a time [4, 14, 22, 23].

A new family of topic-ranking algorithms for multi-labeled documents is proposed in [22, 23]. The algorithms are simple to implement being both time and memory efficient. Experiments with the proposed family of topic-ranking algorithms on standard corpora, show that these algorithms attain adequate results, outperforms other topic-ranking adaptations of well-known classifiers.

The problem of **online adaptation of binary classifiers for tracking** is addressed in [53]. Online learning allows for simple classifiers since only the current view of the object from its surrounding background needs to be discriminated. However, online adaptation has one key problem: each update of the tracker may introduce an error which, finally, can lead to tracking failure (drifting). A novel online semi-supervised boosting method which significantly alleviates the drifting problem in tracking applications was proposed [53]. This allows to limit the drifting problem while still staying adaptive to appearance changes. The main idea is to formulate the update process in a semi-supervised fashion as combined decision of a given prior and an online classifier without any parameter tuning.

A boosting framework that can be used to derive online boosting algorithms for various cost functions was proposed in [4]. Within this framework, online boosting algorithms for logistic regression, least squares regression, and multiple instance learning are derived.

In [14] a **real-time vision-based vehicle detection** system employing an online boosting algorithm is proposed. It is an online AdaBoost approach for a cascade of strong classifiers instead of a single strong classifier. The idea is to develop a cascade of strong classifiers for vehicle detection that is capable of being online trained in response to changing traffic environments. The proposed online boosting method can improve system adaptability and accuracy to deal with novel types of vehicles and unfamiliar environments.

TransientBoost [106] is an online learning algorithm, which is highly adaptive but still robust. It used an internal multi-class representation and modeling reliable and unreliable data in separate classes. Unreliable data is considered transient, and thus highly adaptive learning parameters are applied to adapt to fast changes in the scene while errors fade out fast. In contrast, the reliable data is preserved completely and not harmed by wrong updates. The algorithm is applied successfully on the tasks of object detection and object tracking.

3.6.2 Comparison with other machine learning techniques

In this section we show some detailed experimental results comparing AdaBoost variants, on well-known public domain datasets. We also describe some public domain software packages with code for AdaBoost and its variants.

Table 3.3 briefly describes the datasets used in the experiments, shown by increasing dimensionality. These datasets have several types of data and represent many different learning problems and are available from the UCI Repository [8]³. We also have some datasets from bioinformatics (micro-array and gene expression data)⁴ as well as the 5 datasets of the NIPS2003 FS Challenge⁵ namely, Arcene, Madelon, Gisette, Dexter, and Dorothea.

Table 3.3 Datasets with binary problems used in the experiments: P and N are the number of features and patterns respectively. The datasets are shown by increasing dimensionality.

Dataset	P	N	Type of data / Classification problem
Crabs	5	200	Classify crabs by gender
Phoneme	5	5404	Speech phoneme classification
Abalone	8	4177	Predict the age of abalone from physical measurements
Pima	8	768	The Pima Indians diabetes detection
Contraceptive	9	1473	Predict the current contraceptive method choice
Hepatitis	19	155	Detect if patients lived or died from hepatitis
WBCD	30	569	Wisconsin breast cancer diagnostic database
Ionosphere	34	351	Radar data - signals returned from the ionosphere
SpamBase	54	4601	Sparse BoW data / classify email as SPAM or not
Madelon	500	4400	Float data / artificial dataset, highly nonlinear and difficult
Colon	2000	62	Colon cancer detection
Gisette	5000	13500	Dense integer / distinguish handwritten digits '4' and '9'
DLBCL	5470	77	Dense integer / Lymphoma detection from medial analysis
Leukemia	7129	72	Cancer detection from medical analysis
Example1	9947	2600	Sparse BoW (subset of Reuters) / text classification
Arcene	10000	900	Dense integer / detect cancer versus normal patterns
Prostate Tumor	10509	102	Cancer detection from medical analysis
Dexter	20000	2600	Same data as Example1 with 10053 distractor features

The Leptograpsus Crabs dataset is from Ripley's book [90], and is publicly available at <http://www.stats.ox.ac.uk/pub/PRNN/>. The Crabs dataset is considered as a two-class problem for male/female detection. The Phoneme⁶ dataset holds log-periodograms to represent speech phonemes as used in [57].

³ <http://archive.ics.uci.edu/ml/datasets.html>

⁴ <http://www.gems-system.org/>

⁵ <http://www.nipsfsc.ecs.soton.ac.uk>

⁶ <http://orange.biolab.si/datasets/phoneme.htm>

The Abalone and Pima Indians are well-known datasets from the UCI Repository; their tasks is to predict the age of abalones from their shell measurements and to predict the presence of Diabetes in the Pima Indians population, respectively.

On the Contraceptive dataset has the task to predict the current contraceptive method choice, for married women who were either not pregnant or do not know if they were at the time of interview. It is a subset of the 1987 National Indonesia Contraceptive Prevalence Survey.

The task of Hepatitis dataset is to classify if patients lived or died from hepatitis, given a set of medical analysis. The WBCD dataset is the well-known Wisconsin breast cancer database. The Ionosphere dataset is a binary classification problem on radar data; “good” radar returns are those showing evidence of some type of structure in the ionosphere and “bad” returns are those that do not; their signals pass through the ionosphere. The SpamBase dataset has sparse bag-of-words floating-point data; the task is to classify email messages as SPAM or non-SPAM. We have considered only the first 54 features which constitute a *bag-of-words* representation. The Madelon dataset is an artificial problem of the NIPS2003 FS challenge⁷; it is a difficult problem, because it is multivariate and highly nonlinear. The Colon, DLBCL, Leukemia, and Prostate Tumor datasets deal with the problem of cancer detection from micro-array data.

In the case of Example1⁸, each pattern is a 9947-dimensional BoW vector. The Dexter dataset has the same data as Example1 (with different train, test, and validation partitions) with 10053 additional distractor features, at random locations; it was created for the NIPS 2003 FS challenge. For both datasets, the task is learn to classify Reuters articles as being about “corporate acquisitions” or not.

The Arcene and Gisette datasets also belong to the NIPS2003 FS challenge, an their tasks are: to distinguish cancer versus normal patterns from mass-spectrometric data; to separate the highly confusable handwritten digits ‘4’ and ‘9’.

3.6.2.1 Software packages for boosting algorithms

There are several software packages, freely available online that include implementations of boosting algorithms:

- The GML AdaBoost Matlab Toolbox⁹ provides implementations of Real, Gentle, and Modest Adaboost.
- The ENT TOOL Matlab Toolbox <http://www.j-wichard.de/entool/> which has many machine learning techniques and includes Real, Gentle, and Modest Adaboost from the GML Toolbox.
- A Java implementation is available at <http://jboost.sourceforge.net/>, including AdaBoost, LogitBoost, RobustBoost, and Boostexter.

⁷ <http://clopinet.com/isabelle/Projects/NIPS2003/#challenge>

⁸ <http://svmlight.joachims.org/>

⁹ <http://graphics.cs.msu.ru/science/research/machinelearning/adaboosttoolbox>

- The well-known WEKA machine learning package includes AdaBoost.M1 and MultiBoost classifiers, and it is available at <http://www.cs.waikato.ac.nz/ml/weka/>.
- A C++ implementation of the MPBoost algorithm, is available at this internet address <http://www.esuli.it/mpboost>.
- An efficient C++ implementation of various boosting algorithms can be found in <http://www.stat.purdue.edu/~vishy/>.
- An open-source implementation of BoosTexter¹⁰ (see subsection 3.6.1) can be found at <http://code.google.com/p/icsiboost/>.
- A BoostMetric implementation as well as other boosting algorithms are available at <http://code.google.com/p/boosting/>.
- In <http://cseweb.ucsd.edu/~yfreund/adaboost/index.html> we have a Java Applet that shows how AdaBoost behaves during training.
- The *generalized boosted regression models* (GBM) implements extensions to AdaBoost and gradient boosting machine. It includes regression methods for least squares, absolute loss, quantile regression, logistic, Poisson, Cox proportional hazards partial likelihood, and AdaBoost exponential loss. It is available at <http://cran.r-project.org/web/packages/gbm/index.html>.

The PRTools toolbox [35] available at <http://www.prtools.org/prtools.html> has many machine learning techniques, but it lacks implementations on boosting algorithms¹¹.

3.6.2.2 Analysis of training and test error

Fig. 3.11 shows the weights of each pattern, after training Real Adaboost with $M \in \{1, 5\}$ rounds, on the Ionosphere dataset using 100 training patterns. Notice that at the beginning of the first round, the weights have an uniform distribution with $1/N$. On the first few iterations, the weight of many patterns is changed in such a way that we get a distribution which is quite different from the uniform.

On Fig. 3.12 we have the training error and the test error of Gentle AdaBoost on the WBCD dataset, as a function of the number of weak learners. We see that the training error drops fast on the first few iterations. Even after the training error reaches zero, the test error continues to drop. Fig. 3.13 shows the test error rates for Real, Gentle, and Modest AdaBoost classifiers, on the WBCD and Pima datasets, as a function of the number of weak learners. For these three classifiers, we have an adequate test set error rate on both datasets. On the WBCD dataset, the best performance is achieved by Gentle AdaBoost and for the Pima dataset Real AdaBoost attains the best results.

¹⁰ <http://www.cs.princeton.edu/~schapire/boostexter.html>

¹¹ As of version PRTools 4.0, available at the time of this writing (July, 2011).

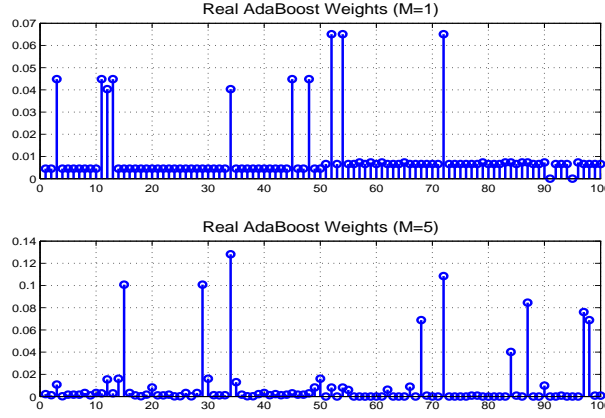


Fig. 3.11 The weights assigned to each pattern, after training Real Adaboost with $M \in \{1, 5\}$ rounds (learners), on the Ionosphere dataset.

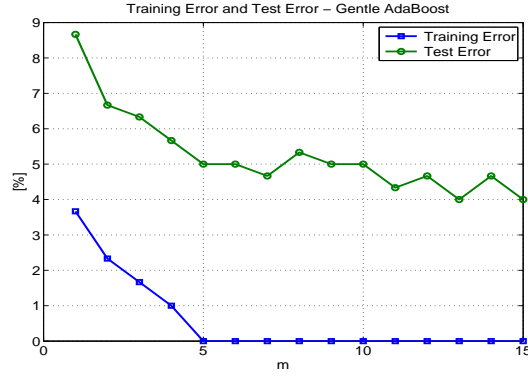


Fig. 3.12 The training error and the test error for Gentle AdaBoost on the WBCD dataset, with $M=15$ learners.

3.6.2.3 Comparison with other classifiers

The reported results in Table 3.4 are averages over ten different random replications of different training/testing partitions, for the standard datasets described in Table 3.3. We compare Real, Gentle, and Modest AdaBoost with linear SVM [10, 20, 114] and KNN [1] classifiers from the PRTools toolbox. The linear kernel SVM classifiers are trained up to 20000 iterations and the KNN classifier uses $K=3$ neighbors. Regarding AdaBoost variants we use the ENTOOL toolbox with $M=15$ weak learners (tree nodes).

In many low and medium dimensional datasets, one of the AdaBoost variants attains adequate results being better than SVM and KNN. However, for the higher-dimensional datasets SVM and KNN tend to perform slightly better than these AdaBoost variants.

3.7 Summary and discussion

The AdaBoost algorithm (and its variants) has many practical advantages, which, combined with theoretical guarantees, makes it a very attractive general purpose learning method. On the practical side, boosting algorithms are simple to implement and debug. The base learner and the number of iterations (learners) are the only two important choices to be made.

AdaBoost had been shown to be resistant to overfitting, despite the fact that it can produce combinations involving very large numbers of base classifiers. However, recent studies have shown that this is not the case, even for base classifiers as simple as decision stumps. The success of the boosting algorithms depends on the amount of data available for training as well as on the type of weak learner.

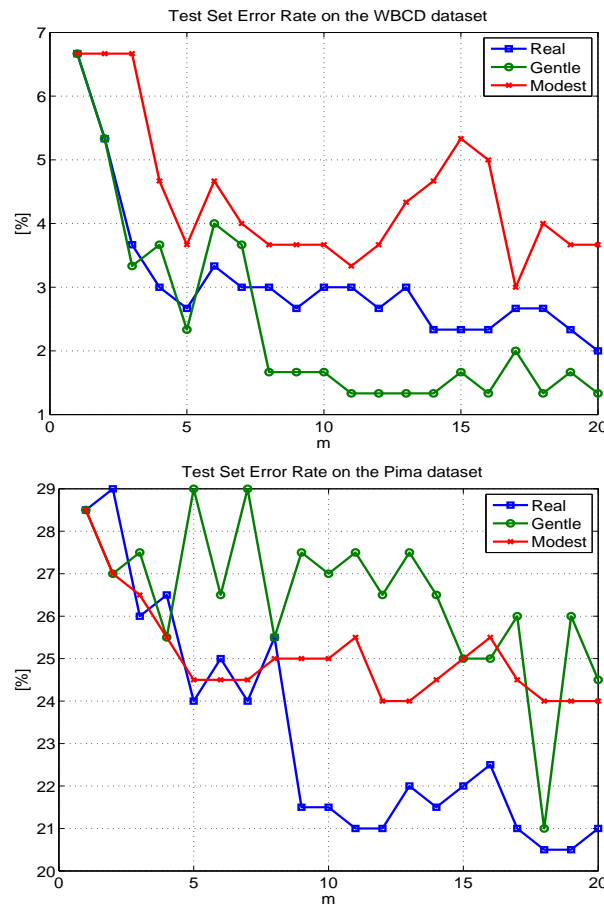


Fig. 3.13 Test error rates for Real, Gentle, and Modest AdaBoost classifiers, on the WBCD and Pima datasets.

Table 3.4 Experimental comparison of three AdaBoost variants with linear SVM and 3-NN classifiers. We show the average \pm standard deviation of the test set error rate, for ten runs with different random training/test partitions on standard datasets of Table 3.3. The best results are in bold face.

Dataset	Real	Gentle	Modest	SVM	3-NN
Crabs	23.50 \pm 11.56	23.50 \pm 12.26	19.50 \pm 12.12	5.50 \pm 4.97	33.50 \pm 15.28
Phoneme	24.45 \pm 4.50	23.70 \pm 3.25	23.90 \pm 2.63	26.70 \pm 2.49	22.90 \pm 1.73
Abalone	26.25 \pm 3.26	26.15 \pm 3.12	25.65 \pm 3.98	23.05 \pm 4.56	28.55 \pm 3.24
Pima	23.05 \pm 3.35	22.20 \pm 3.28	23.60 \pm 2.17	25.85 \pm 2.52	28.25 \pm 3.50
Contraceptive	34.35 \pm 4.45	34.50 \pm 4.29	35.50 \pm 3.06	37.95 \pm 3.61	40.05 \pm 3.20
Hepatitis	19.00 \pm 5.80	22.25 \pm 7.77	20.50 \pm 9.92	25.25 \pm 7.77	41.25 \pm 7.38
WBCD	3.57 \pm 1.81	2.40 \pm 1.64	3.83 \pm 2.24	4.80 \pm 1.10	7.10 \pm 1.28
Ionosphere	9.60 \pm 5.02	7.60 \pm 2.99	7.30 \pm 2.36	12.70 \pm 4.42	18.20 \pm 2.86
SpamBase	14.60 \pm 2.68	14.60 \pm 3.11	14.10 \pm 3.04	13.57 \pm 2.39	18.33 \pm 1.52
Madelon	50.12 \pm 2.29	49.78 \pm 1.98	49.45 \pm 1.85	50.97 \pm 1.85	50.58 \pm 1.97
Colon	1.67 \pm 5.27	1.67 \pm 5.27	1.67 \pm 5.27	13.33 \pm 4.10	15.56 \pm 6.70
Gisette	13.67 \pm 1.59	11.06 \pm 0.90	10.62 \pm 1.51	7.79 \pm 0.67	11.39 \pm 1.81
DLBCL	14.00 \pm 10.16	19.00 \pm 9.94	17.00 \pm 10.59	4.67 \pm 4.22	18.67 \pm 7.73
Leukemia	37.00 \pm 17.51	37.00 \pm 17.51	37.00 \pm 17.51	11.00 \pm 7.75	13.50 \pm 9.44
Example1	12.07 \pm 2.62	9.62 \pm 1.57	11.63 \pm 1.42	4.30 \pm 0.61	10.98 \pm 1.13
Arcene	31.10 \pm 6.26	32.00 \pm 5.60	29.70 \pm 6.36	31.00 \pm 0.94	20.90 \pm 6.06
Prost. Tumor	5.75 \pm 4.87	7.25 \pm 3.11	4.75 \pm 4.03	3.50 \pm 2.49	15.63 \pm 4.87
Dexter	18.13 \pm 3.07	15.30 \pm 1.06	14.57 \pm 1.66	10.10 \pm 0.80	25.97 \pm 4.63

There are dozens of variants for binary and multi-class problems that have been proven successful on many problems. In recent years, the research on boosting algorithms has been focused mainly on multi-class and semi-supervised problems. The online boosting algorithms, in which the training examples arrive one at a time, as contrary to the batch mode, is also a fruitful field of research with many successful algorithms. Many of these algorithms are applied to real-time computer vision problems, such as detection or tracking. This is a focus of intensive current research.

AdaBoost and its variants have also been combined with different machine learning techniques, such as random subspaces [65], genetic algorithms [19], and rotation forest [132].

3.8 Bibliographical and historical remarks

There are many papers and tutorials addressing the many boosting algorithms, proposing variants for the multi-class case and/or semi-supervised learning problems. In this section, we point out some of these elements that can be found in the literature.

For the origins of boosting algorithms, one can be interested in reading about the bootstrapping [36, 37] and bagging (bootstrap aggregation) [11] techniques. Bootstrap was initially proposed in 1982, but it regained interest on the decade of

1990-1999, in which bagging was proposed. The use of bagging for classification problems is addressed in many papers, see [12, 13, 94, 137, 138] for many applications. The seminal paper of Schapire [94], proposing the first provable polynomial time boosting procedure is a must read. For a comparison of the effectiveness of randomization, bagging, and boosting see [28].

The seminal papers on the middle of the decade of 1990-1999 introducing *adaptive boosting* (AdaBoost) algorithm [46], [47], [49], with the idea that we can *weight* the data instead of resampling it, are also a must read. On the second half of the decade 1990-1999, the AdaBoost algorithm was also extended for regression tasks, as addressed in [3, 30], for instance.

From 1999 until this date, there are dozens of extensions and variants of AdaBoost for supervised and semi-supervised binary and multi-class problems, covering a wide range of successful applications. In this chapter, besides the theoretical aspects of boosting algorithms, we tried to cover as many variants and successful applications as possible. Section 3.4 covers many variants for supervised learning, whereas Section 3.5 addresses the semi-supervised variants on binary and multi-class problems. On Section 3.6 we have described a wide range of applications. The vast majority of these variants and successful applications were published in the decade of 2000-2009. Online boosting, in contrast with batch boosting, has received a great deal of attention in recent research; we covered many approaches for online boosting on subsection 3.6.1.1.

For further reading on adaptive boosting algorithms, please see [81], which complements well this chapter. Whereas we have aimed at covering a wide range of variants, [81] focuses more on theoretical and practical aspects of boosting and ensemble learning. The webpage http://cbio.mskcc.org/~arvey/boosting_papers.html has many papers, tutorials, and links to software on boosting algorithms. There are also some tutorials about boosting available on-line^{12 13}.

The web pages of boosting and adaptive boosting pioneers R. Schapire¹⁴. and Y. Freund¹⁵ have many useful information about boosting algorithms. This information is useful for both experienced researcher as well as to the new researchers entering this exciting field of research.

¹² <http://www.site.uottawa.ca/~stan/csi5387/boost-tut-ppr.pdf>

¹³ <http://www.stat.purdue.edu/~vishy/>

¹⁴ <http://www.cs.princeton.edu/~schapire/boost.html>

¹⁵ <http://cseweb.ucsd.edu/~yfreund/papers/index.html>

Appendix A: Proofs

A.1: Proof of (3.4)

Consider the final weights, $w_i^{(M+1)}$, and explicitly write the recursion that starts at $w_i^{(1)} = 1/N$ (recall that S_j is the normalizing constant used in line 7 at iteration j)

$$w_i^{(M+1)} = w_i^{(M)} \frac{\exp(-\alpha_M y_i H_M(\mathbf{x}_i))}{S_M} = \frac{\prod_{j=1}^M \exp(-\alpha_j y_i H_j(\mathbf{x}_i))}{N \prod_{j=1}^M S_j}. \quad (3.12)$$

This can be re-written as

$$w_i^{(M+1)} = \frac{\exp\left(-y_i \sum_{j=1}^M \alpha_j H_j(\mathbf{x}_i)\right)}{N \prod_{j=1}^M S_j} = \frac{\exp(-y_i f(\mathbf{x}_i))}{N \prod_{j=1}^M S_j}, \quad (3.13)$$

where $f(\mathbf{x}) = \sum_{j=1}^M \alpha_j H_j(\mathbf{x})$, from which we can conclude that

$$\exp(-y_i f(\mathbf{x}_i)) = w_i^{(M+1)} N \prod_{j=1}^M S_j. \quad (3.14)$$

Now, noticing that $H(\mathbf{x}) = \text{sign}(f(\mathbf{x}))$, and recalling that h denotes the Heaviside function (defined above), we have

$$h(-y_i H(\mathbf{x}_i)) = h(-y_i f(\mathbf{x}_i)) \leq \exp(-y_i f(\mathbf{x}_i)) = w_i^{(M+1)} N \prod_{j=1}^M S_j. \quad (3.15)$$

We can now write and bound the training error (TE) rate of H ,

$$TE = \frac{1}{N} \sum_{i=1}^N h(-y_i H(\mathbf{x}_i)) \leq \frac{1}{N} \sum_{i=1}^N w_i^{(M+1)} N \prod_{j=1}^M S_j = \prod_{j=1}^M S_j, \quad (3.16)$$

because $\sum_{i=1}^N w_i^{(M+1)} = 1$, thus concluding the proof of (3.4).

A.2: Proof of (3.5)

Let us plug the expression for α_m (line 5 of AdaBoost) into the expression of the normalizing factor S_m , and use the fact that $y_i H_m(\mathbf{x}_i) = 1$, if and only if $H_m(\mathbf{x}_i) = y_i$, while $y_i H_m(\mathbf{x}_i) = -1$, if and only if $H_m(\mathbf{x}_i) \neq y_i$,

$$S_m = \sum_{i=1}^N w_i^{(m)} \exp \left(y_i H_m(\mathbf{x}_i) \log \sqrt{\frac{\text{err}_m}{1 - \text{err}_m}} \right) \quad (3.17)$$

$$= \sqrt{\frac{\text{err}_m}{1 - \text{err}_m}} \sum_{i: y_i = H(\mathbf{x}_i)} w_i^{(m)} + \sqrt{\frac{1 - \text{err}_m}{\text{err}_m}} \sum_{i: y_i \neq H(\mathbf{x}_i)} w_i^{(m)} \quad (3.18)$$

$$= 2\sqrt{\text{err}_m(1 - \text{err}_m)}. \quad (3.19)$$

Recalling that $\text{err}_m = 1/2 - \gamma_m$, we have

$$S_m = 2\sqrt{\left(\frac{1}{2} - \gamma_m\right) \left(\frac{1}{2} + \gamma_m\right)} \quad (3.20)$$

$$= \sqrt{1 - 4\gamma_m^2} \quad (3.21)$$

$$= \exp \left(\frac{1}{2} \log (1 - 4\gamma_m^2) \right) \quad (3.22)$$

$$\leq \exp (-2\gamma_m^2), \quad (3.23)$$

where in (3.23) we have used the inequality $\log u \leq u - 1$ (often referred to as the Gibbs inequality). Plugging this inequality into (3.4), and invoking the assumption $\gamma_m > \gamma$, yields

$$TE \leq \exp \left(-2 \sum_{m=1}^M \gamma_m^2 \right) \leq \exp (-2M\gamma^2) \quad (3.24)$$

thus proving inequality (3.5).

References

1. D. Aha, D. Kibler, and M. Albert. Instance-based learning algorithms. In *Machine Learning*, pages 37–66, 1991.
2. E. Allwein, R. Schapire, and Y. Singer. Reducing multiclass to binary: A unifying approach for margin classifiers. *Journal of Machine Learning Research*, 1:113–141, 2000.
3. R. Avnimelech and N. Intrator. Boosting regression estimators. *Neural Computation*, 11:491–513, 1999.
4. B. Babenko, M. Yang, and S. Belongie. A family of online boosting algorithms. In *Learning09*, pages 1346–1353, 2009.
5. E. Bauer and R. Kohavi. An empirical comparison of voting classification algorithms: Bagging, boosting, and variants. *Machine Learning*, 36:105–139, 1999.
6. K. Bennett, A. Demiriz, and R. Maclin. Exploiting unlabeled data in ensemble methods. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '02, pages 289–296, New York, NY, USA, 2002. ACM.
7. J. Bergstra and B. Kégl. Meta-features and adaboost for music classification. In *Machine Learning Journal : Special Issue on Machine Learning in Music*, 2006.
8. C. Blake and C. Merz. Uci repository of machine learning databases. Technical report, University of California, Irvine, Department of Informatics and Computer Science, 1999.
9. A. Blumer, A. Ehrenfeucht, D. Haussler, and M. Warmuth. Learnability and the vapnik-chervonenkis dimension. *Journal of the ACM*, 36:929–965, 1989.

10. B. Boser, I. Guyon, and V. Vapnik. A training algorithm for optimal margin classifiers. In *Proc. of the 5th Annual ACM Workshop on Computational Learning Theory*, pages 144–152, New York, NY, USA, 1992. ACM Press.
11. L. Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996.
12. L. Breiman. Bias, variance, and arcing classifiers. Technical report, UC Berkely, CA, 1996.
13. P. Bühlmann and B. Yu. Analyzing bagging. *Annals of Statistics*, 30:927–961, 2002.
14. W.-C. Chang and C.-W. Cho. Online Boosting for Vehicle Detection. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 40(3):892–902, June 2010.
15. O. Chapelle, B. Schalkopf, and A. Zien. *Semi-Supervised Learning*. The MIT Press, 2006.
16. K. Chen, C. Chou, S. Shih, W. Chen, and D. Chen. Feature selection for iris recognition with adaboost. *Intelligent Information Hiding and Multimedia Signal Processing, 2007. IHHMSP 2007. Third International Conference on*, 2:411–414, 2007.
17. K. Chen and S. Wang. Regularized boost for semi-supervised learning. In *Neural Information Processing Systems*, 2007.
18. K. Chen and S. Wang. Semi-supervised Learning via Regularized Boosting Working on Multiple Semi-supervised Assumptions. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 99(1), 2010.
19. H. Chouaib, O. Terrades, S. Tabbone, F. Cloppet, and N. Vincent. Feature selection combining genetic algorithm and adaboost classifiers. In *ICPR08*, pages 1–4, 2008.
20. N. Christiani and J. Shawe-Taylor. *An Introduction to Support Vector Machines and other kernel based learning methods*. Cambridge University Press, 2000.
21. M. Collins, R. E. Schapire, and Y. Singer. Logistic regression, adaboost and bregman distances. In *Machine Learning*, volume 48, pages 158–169, 2000.
22. K. Crammer and Y. Singer. A new family of online algorithms for category ranking. In *Proceedings of the 25th annual international ACM SIGIR conference on Research and development in information retrieval, SIGIR '02*, pages 151–158, New York, NY, USA, 2002. ACM.
23. K. Crammer, Y. Singer, J. K. T. Hofmann, T. Poggio, and J. Shawe-taylor. A family of additive online algorithms for category ranking. *Journal of Machine Learning Research*, 3:2003, 2003.
24. D. Cristinacce and T. Cootes. Facial feature detection using adaboost with shape constraints. In *In Proc. BMVC03*, pages 231–240, 2003.
25. F. d’Alché Buc, Y. Grandvalet, and C. Ambroise. Semi-supervised marginboost, 2002.
26. A. Demiriz, K. P. Bennett, and J. S. Taylor. Linear Programming Boosting via Column Generation. *Machine Learning*, 46(1-3):225–254, 2002.
27. A. Dempster, N. Laird, and D. Rubin. Maximum likelihood from incomplete data via the EM algorithm (with discussion). *Journal of the Royal Statistical Society, B*, 39:1–38, 1977.
28. T. Dietterich. An experimental comparison of three methods for constructing ensembles of decision trees : Bagging , boosting , and randomization. *Machine Learning*, 40(2):139–157, 2000.
29. C. Domingo and O. Watanabe. Madaboost: A modification of adaboost, 2000.
30. H. Drucker. Improving regressors using boosting techniques. In *Proceedings of the Fourteenth International Conference on Machine Learning, ICML 97*, pages 107–115, San Francisco, CA, USA, 1997. Morgan Kaufmann Publishers Inc.
31. H. Drucker, R. Schapire, and P. Simard. Improving performance in neural networks using a boosting algorithm. In *Advances in Neural Information Processing Systems 5, [NIPS Conference]*, pages 42–49, San Francisco, CA, USA, 1993. Morgan Kaufmann Publishers Inc.
32. H. Drucker and C. Tortes. Boosting decision trees. In *Advances in Neural Information Processing Systems*, volume 8, pages 479–485. MIT Press, 1996.
33. R. Duda, P. Hart, and D. Stork. *Pattern Classification*. John Wiley & Sons, 2nd edition, 2001.
34. N. Duffy and D. Helmbold. Potential boosters? In *Advances in Neural Information Processing Systems 12*, pages 258–264. MIT Press, 2000.

35. R. Duin, P. Juszczak, P. Paclik, E. Pekalska, D. Ridder, D. Tax, and S. Verzakov. PRTTools4.1, a Matlab Toolbox for Pattern Recognition. Technical report, Delft University of Technology, 2007.
36. B. Efron. The jackknife, the bootstrap and other resampling plans. *Society for Industrial and Applied Mathematics (SIAM)*, 1982.
37. B. Efron and R. Tibshirani. *An Introduction to the Bootstrap*. Chapman & Hall, New York, 1993.
38. G. Eibl and K. Pfeiffer. How to make adaboost.M1 work for weak classifiers by changing only one line of the code. In *Machine Learning: Thirteenth European Conference*, volume 1, pages 109–120, 2002.
39. G. Eibl and K. Pfeiffer. Multiclass boosting for weak classifiers. *Journal of Machine Learning Research*, 6:189–210, 2005.
40. A. Esuli, T. Fagni, and F. Sebastiani. Mp-boost: A multiple-pivot boosting algorithm and its application to text categorization. In *In Proceedings of the 13th International Symposium on String Processing and Information Retrieval (SPIRE'06)*, 2006.
41. A. Ferreira and M. Figueiredo. Boosting of (very) weak classifiers. In *6th Portuguese Conference on Telecommunications, Confele'07*, Peniche, Portugal, 2007.
42. F. Fleuret. Multi-layer boosting for pattern recognition. *Pattern Recognition Letters*, 30:237–241, February 2009.
43. Y. Freund. Boosting a Weak Learning Algorithm by Majority. *Information and Computation*, 121(2):256–285, 1995.
44. Y. Freund. An adaptive version of the boost by majority algorithm. In *In Proceedings of the Twelfth Annual Conference on Computational Learning Theory*, pages 102–113, 2000.
45. Y. Freund. A more robust boosting algorithm. <http://arxiv.org/abs/0905.2138>, 2009.
46. Y. Freund and R. Schapire. A decision-theoretic generalization of on-line learning and application to boosting. In *European Conference on Computational Learning Theory – EuroCOLT*. Springer, 1994.
47. Y. Freund and R. Schapire. Experiments with a new boosting algorithm. In *Thirteenth International Conference on Machine Learning*, pages 148–156, Bari, Italy, 1996.
48. Y. Freund and R. Schapire. Game theory, on-line prediction and boosting. In *Proceedings of the Ninth Annual Conference on Computational Learning Theory*, pages 325–332. ACM Press, 1996.
49. Y. Freund and R. 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.
50. J. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29:1189–1232, 2000.
51. J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *The Annals of Statistics*, 28(2):337–374, 2000.
52. V. Gómez-Verdejo, M. Ortega-Moral, J. Arenas-García, and A. Figueiras-Vidal. Boosting of weighting critical and erroneous samples. *Neurocomputing*, 69(7-9):679–685, 2006.
53. H. Grabner, C. Leistner, and H. Bischof. Semi-supervised On-Line Boosting for Robust Tracking. In D. Forsyth, P. Torr, and A. Zisserman, editors, *Computer Vision ECCV 2008*, volume 5302 of *Lecture Notes in Computer Science*, chapter 19, pages 234–247. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008.
54. Y. Grandvalet, F. Buc, and C. Ambroise. Boosting mixture models for semi-supervised learning. In *ICANN International Conference on Artificial Neural Networks*, volume 1, pages 41–48, Vienna, Austria, 2001.
55. V. Guruswami and A. Sahai. Multiclass learning, boosting, and error-correcting codes. In *12th Annual Conference on Computational Learning Theory (COLT-99)*, Santa Cruz, USA, 1999.
56. Z. Hao, C. Shen, N. Barnes, and B. Wang. Totally-corrective multi-class boosting. In *ACCV10*, pages IV: 269–280, 2010.
57. T. Hastie, R. Tibshirani, and J. Friedman. *The Elements of Statistical Learning*. Springer, 2nd edition, 2001.

58. J. He and B. Thiesson. Asymmetric gradient boosting with application to SPAM filtering. In *CEAS*, 2007.
59. C. Huang, H. Ai, Y. Li, and S. Lao. Vector boosting for rotation invariant multi-view face detection. In *International Conference on Computer Vision (ICCV)*, volume 1, pages 446–453, 2005.
60. X. Huang, S. Li, and Y. Wang. Jensen-shannon boosting learning for object recognition. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2, pages 144–149, 2005.
61. J. Jackson and M. Craven. Learning sparse perceptrons. In *Advances in Neural Information Processing Systems*, volume 8, pages 654–660. MIT Press, 1996.
62. R. Jin, Y. Liu, L. Si, J. Carbonell, and A. G. Hauptmann. A new boosting algorithm using input-dependent regularizer. In *Proceedings of Twentieth International Conference on Machine Learning (ICML 03)*. AAAI Press, 2003.
63. X. Jin, X. Hou, and C.-L. Liu. Multi-class adaboost with hypothesis margin. In *Proceedings of the 2010 20th International Conference on Pattern Recognition, ICPR '10*, pages 65–68, Washington, DC, USA, 2010. IEEE Computer Society.
64. G. Jun and J. Ghosh. Multi-class boosting with class hierarchies. *Multiple Classifier Systems*, 5519:32–41, 2009.
65. H. Kong and E. Teoh. Coupling adaboost and random subspace for diversified fisher linear discriminant. In *ICARCV06*, pages 1–5, 2006.
66. L. Kuncheva. *Combining Pattern Classifiers: Methods and Algorithms*. Wiley, 2004.
67. D. Le and S. Satoh. Ent-boost: Boosting using entropy measures for robust object detection. *Pattern Recognition Letters*, 2007.
68. H. Li and C. Shen. Boosting the minimum margin: Lpboost vs. adaboost. *Digital Image Computing: Techniques and Applications*, 0:533–539, 2008.
69. L. Li. Multiclass boosting with repartitioning. In *23rd International Conference on Machine Learning (ICML 07)*, Pennsylvania, USA, 2006.
70. P. Li. Abc-boost: adaptive base class boost for multi-class classification. In *International Conference on Machine Learning*, pages 79–632, 2009.
71. S. Li and A. Jain. *Handbook of Face Recognition*. Springer, 2005.
72. S. Li and Z. Zhang. Floatboost learning and statistical face detection. *Transactions on Pattern Analysis and Machine Intelligence*, 26(9):23–38, 2004.
73. C. Liu and H. Shum. Kullback-Leibler boosting. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 587–594, Madison, Wisconsin, USA, 2003.
74. Y. Lu, Q. Tian, and T. S. Huang. Interactive boosting for image classification. In *Proceedings of the 7th international conference on Multiple classifier systems, MCS'07*, pages 180–189, Berlin, Heidelberg, 2007. Springer-Verlag.
75. S. Lyu. Infomax boosting. In *International Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 533–538, 2005.
76. R. Maclin. An empirical evaluation of bagging and boosting. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence*, pages 546–551. AAAI Press, 1997.
77. P. Mallapragada, R. R. Jin, A. Jain, and Y. Liu. SemiBoost: Boosting for Semi-Supervised Learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(11):2000–2014, 2009.
78. H. Masnadi-Shirazi and N. Vasconcelos. Asymmetric boosting. In *Proceedings of the 24th international conference on Machine learning, ICML '07*, pages 609–619, New York, NY, USA, 2007. ACM.
79. L. Mason, J. Baxter, P. Bartlett, and M. Frean. Boosting algorithms as gradient descent in function space, 1999.
80. L. Mason, J. Baxter, P. Bartlett, and M. Frean. Functional gradient techniques for combining hypotheses. *Advances in Large Margin Classifiers*, 1:109–120, 2000.
81. R. Meir and G. Rätsch. An introduction to boosting and leveraging. In S. Mendelson and A. Smola, editors, *Advanced Lectures on Machine Learning*. Springer Verlag, 2006.

82. A. Mohemmed, M. Zhang, and M. Johnston. A PSO Based Adaboost Approach to Object Detection. *Simulated Evolution and Learning*, pages 81–90, 2008.
83. Z. Niu, S. Shan, S. Yan, X. Chen, and W. Gao. 2d cascaded adaboost for eye localization. In *Proc. of the 18th International Conference on Pattern Recognition*, 2006.
84. R. Nock and P. Lefaucheur. A Robust Boosting Algorithm. In T. Elomaa, H. Mannila, and H. Toivonen, editors, *Machine Learning: ECML 2002*, volume 2430 of *Lecture Notes in Computer Science*, pages 319–331, Berlin, Heidelberg, 2002. Springer Berlin / Heidelberg.
85. T. V. Pham and A. W. M. Smeulders. Quadratic boosting. *Pattern Recogn.*, 41:331–341, January 2008.
86. A. Qudus, P. Fieguth, and O. Basir. Adaboost and Support Vector Machines for White Matter Lesion Segmentation in MR Images. In *2005 IEEE Engineering in Medicine and Biology 27th Annual Conference*, pages 463–466. IEEE, 2005.
87. J. Quinlan. *C4.5: Programs for Machine Learning*. Morgan Kaufmann, 1993.
88. J. Quinlan. Bagging, Boosting, and C4.5. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 725–730, 1996.
89. G. Rätsch. Robust multi-class boosting. In *Eurospeech 2003, 8th European Conference on Speech Communication and Technology*, page 9971000, Geneva, Switzerland, 2003.
90. B. Ripley. *Pattern Recognition and Neural Networks*. Cambridge University Press, 1996.
91. J. Rodriguez and J. Maudes. Boosting recombined weak classifiers. *Pattern Recognition Letters*, 29(8):1049–1059, 2007.
92. A. Saffari, H. Grabner, and H. Bischof. Serboost: Semi-supervised boosting with expectation regularization. In D. Forsyth, P. Torr, and A. Zisserman, editors, *Computer Vision ECCV 2008*, volume 5304 of *Lecture Notes in Computer Science*, pages 588–601. Springer Berlin / Heidelberg, 2008.
93. A. Saffari, C. Leistner, M. Godec, and H. Bischof. Robust multi-view boosting with priors. In *ECCV, ECCV'10*, pages 776–789, Berlin, Heidelberg, 2010. Springer-Verlag.
94. R. Schapire. The strength of weak learnability. In *Machine Learning*, volume 5, pages 197–227, 1990.
95. R. Schapire. Using output codes to boost multiclass learning problems. In *14th International Conference on Machine Learning (ICML 97)*, pages 313–321, Tennessee, USA, 1997.
96. R. Schapire. Theoretical views of boosting. In *Proceedings of the 4th European Conference on Computational Learning Theory*, EuroCOLT '99, pages 1–10, London, UK, 1999. Springer-Verlag.
97. R. Schapire. The boosting approach to machine learning: An overview. In *Nonlinear Estimation and Classification*, Berkeley, 2002. Springer.
98. R. Schapire, Y. Freund, P. Bartlett, and W. Lee. Boosting the margin: A new explanation for the effectiveness of voting methods, 1997.
99. R. Schapire and Y. Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297–336, 1999.
100. R. Schapire and Y. Singer. BoosTexter: A Boosting-based System for Text Categorization. *Machine Learning*, 39(2/3):135–168, 2000.
101. H. Schwenk and Y. Bengio. Adaboosting neural networks: Application to on-line character recognition. In *In ICANN'97, LNCS, 1327*, 967–972, pages 967–972. Springer, 1997.
102. H. Schwenk and Y. Bengio. Boosting Neural Networks. *Neural Comp.*, 12(8):1869–1887, 2000.
103. C. Shen, J. Kim, L. Wang, and A. van den Hengel. Positive semidefinite metric learning with boosting. In Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, editors, *Advances in Neural Information Processing Systems (NIPS'09)*, pages 1651–1659, Vancouver, B.C., Canada, December 2009. MIT Press.
104. J. Sochman and J. Matas. Waldboost ” learning for time constrained sequential detection. In *Proceedings of the 2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05) - Volume 2 - Volume 02*, CVPR '05, pages 150–156, Washington, DC, USA, 2005. IEEE Computer Society.
105. A. Stefan, V. Athitsos, Q. Yuan, and S. Sclaroff. Reducing jointboost-based multiclass classification to proximity search. In *CVPR*, pages 589–596. IEEE, 2009.

106. S. Sternig, M. Godec, P. Roth, and H. Bischof. Transientboost: On-line boosting with transient data. In *OLCV10*, pages 22–27, 2010.
107. Y. Sun, S. Todorovic, and J. Li. Unifying multi-class adaboost algorithms with binary base learners under the margin framework. *Pattern Recognition Letters*, 28:631–643, 2007.
108. Y. Sun, S. Todorovic, J. Li, and D. Wu. Unifying the error-correcting and output-code adaboost within the margin framework. In *Proceedings of the 22nd international conference on Machine learning*, ICML 05, pages 872–879, New York, NY, USA, 2005. ACM.
109. J. Thongkam, O. Xu, Y. Zhang, F. Huang, and G. Adaboosts. Breast cancer survivability via adaboost algorithms. In *Proceedings of the second Australasian workshop on Health data and knowledge management - Volume 80*, HDKM '08, pages 55–64, Darlinghurst, Australia, Australia, 2008. Australian Computer Society, Inc.
110. K. Tieu and P. Viola. Boosting image retrieval. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition – CVPR*, volume 1, pages 228–235, 2000.
111. A. Torralba, K. Murphy, and W. Freeman. Sharing visual features for multiclass and multiview object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(5):854 – 869, March 2007.
112. L. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
113. H. Valizadegan, R. Jin, and A. K. Jain. Semi-Supervised Boosting for Multi-Class Classification. In *ECML PKDD '08: Proceedings of the European conference on Machine Learning and Knowledge Discovery in Databases - Part II*, pages 522–537, Berlin, Heidelberg, 2008. Springer-Verlag.
114. V. Vapnik. *The Nature of Statistical Learning Theory*. Springer-Verlag, 1999.
115. V. Vapnik and A. Chervonenkis. On the uniform convergence of relative frequencies of events to their probabilities. *Theory of Probability and its Applications*, 16(2):264–280, 1971.
116. R. Vershae, J. Ruiz-del-solar, and M. Correa. Gender classification of faces using adaboost. In *Lecture Notes in Computer Science (CIARP 2006)* 4225, page 78. Springer, 2006.
117. A. Vezhnevets and V. Vezhnevets. Modest AdaBoost - teaching AdaBoost to generalize better. *Graphicon*, 12(5):987–997, 2005.
118. P. Viola and M. Jones. Rapid object detection using a boosted cascade of simple features. In *Proceeding International Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 511–518, Hawaii, 2001.
119. P. Viola and M. Jones. Robust real-time face detection. *International Journal of Computer Vision*, 57:137–154, 2004.
120. P. Viola, M. Jones, and D. Snow. Detecting pedestrians using patterns of motion and appearance. In *International Conference on Computer Vision – ICCV*, pages 734–741, 2003.
121. P. Viola, J. Platt, and C. Zhang. Multiple instance boosting for object detection. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 1417–1424, Cambridge, MA, 2006. MIT Press.
122. L. Wang, S. Yuan, L. Li, and H. Li. Boosting naïve Bayes by active learning. In *Third International Conference on Machine Learning and Cybernetics*, volume 1, pages 41–48, Shanghai, China, 2004.
123. P. Wang, C. Shen, N. Barnes, H. Zheng, and Z. Ren. Asymmetric totally-corrective boosting for real-time object detection. In *ACCV10*, pages I: 176–188, 2010.
124. M. Warmuth, K. Gloer, and G. Rätsch. Boosting Algorithms for Maximizing the Soft Margin. In *NIPS*, 2007.
125. M. Warmuth, K. Gloer, and S. Vishwanathan. Entropy regularized lpboost. In *Proceedings of the 19th international conference on Algorithmic Learning Theory*, ALT '08, pages 256–271, Berlin, Heidelberg, 2008. Springer-Verlag.
126. M. Warmuth, J. Liao, and G. Rätsch. Totally corrective boosting algorithms that maximize the margin. In *Proceedings of the 23rd international conference on Machine learning*, ICML '06, pages 1001–1008, New York, NY, USA, 2006. ACM.
127. J. Warrell, P. Torr, and S. Prince. Styp-boost: A bilinear boosting algorithm for learning style-parameterized classifiers. In *BMVC10*, 2010.
128. J. Webb, J. Boughton, and Z. Wang. Not so naïve Bayes: Aggregating one-dependence estimators. *Machine Learning*, 58(1):5–24, 2005.

129. P. Yang, S. Shan, W. Gao, S. Z. Li, and D. Zhang. Face recognition using ada-boosted gabor features. In *in: Proceedings of the 16th International Conference on Face and Gesture Recognition*, pages 356–361, 2004.
130. C. Zhang, P. Yin, Y. Rui, R. Cutler, P. Viola, X. Sun, N. Pinto, and Z. Zhang. Boosting-based multimodal speaker detection for distributed meeting videos. *IEEE Trans. on Multimedia*, 10(8):1541–1552, December 2008.
131. C. Zhang and Z. Z. *Boosting-Based Face Detection and Adaptation*. Morgan and Claypool Publishers, 2010.
132. C. Zhang and J. Zhang. Rotboost: A technique for combining rotation forest and adaboost. *Pattern Recognition Letters*, 29(10):1524–1536, July 2008.
133. C. Zhang and Z. Zhang. Winner-take-all multiple category boosting for multi-view face detection. Technical report, One Microsoft Way, Redmond, WA 98052 USA, 2010.
134. L. Zheng, S. Wang, Y. Liu, and C.-H. Lee. Information theoretic regularization for semi-supervised boosting. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, KDD '09, pages 1017–1026, New York, NY, USA, 2009. ACM.
135. M. Zhou, H. Wei, and S. Maybank. Gabor wavelets and adaboost in feature selection for face verification.
136. J. Zhu, H. Zou, S. Rosset, and T. Hastie. Multi-class adaboost. *Statistics and Its Interface*, 2:349–360, 2009.
137. X. Zhu, C. Bao, and W. Qiu. Bagging very weak learners with lazy local learning. In *International Conference on Pattern Recognition (ICPR)*, pages 1–4, 2008.
138. X. Zhu and Y. Yang. A lazy bagging approach to classification. *Pattern Recognition*, 41:2980–2992, 2008.