

Introduction to AdaBoost

Balázs Kégl

January 21, 2014

Contents

1	Introduction	9
1.1	Bibliographical notes	9
2	The basic algorithm: binary ADABOOST	11
2.1	Basics	11
2.2	The binary ADABOOST algorithm	11
2.3	Basic analysis	14
2.4	Base classifiers	17
2.4.1	Decision stumps	17
2.5	A simple example	18
2.6	Discussion	19
2.7	Bibliographical notes	19

Notations - Definitions

General

- \mathcal{X} : measurable observation space; usually $\mathcal{X} = \mathbb{R}^d$
- X : random observation vector in \mathcal{X}
- \mathbf{x} : fixed observation vector in \mathcal{X}
- d : dimension of the observation space (if finite)
- K : number of classes in multi-class classification
- \mathcal{Y} : label space; $\mathcal{Y} = \{-1, 1\}$ in binary classification, $\mathcal{Y} \subseteq \mathbb{R}$ in regression, $\mathcal{Y} = \{1, \dots, K\}$ in multi-class classification, and $\mathcal{Y} = \{0, 1\}^K$ in multi-class multi-label classification.
- Y : random label in \mathcal{Y}
- y : fixed label in \mathcal{Y}
- $Z = (X, Y)$
- $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ or $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$: training sample of size n
- n : training sample size
- m : test sample size
- P : distribution of (X, Y)
- $P^{\otimes n}$: distribution of $(X_1, Y_1), \dots, (X_n, Y_n)$
- μ : marginal distribution of X
- η : regression function, $\forall \mathbf{x} \in \mathcal{X}, \eta(\mathbf{x}) = \mathbb{E}\{Y \mid X = \mathbf{x}\}$.
- R : risk
- $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$: loss function
- $f : \mathcal{X} \rightarrow \mathcal{Y}$: predictor
- \mathcal{F} : restricted set of predictors
- Definition: risk of predictor f

$$R(f) = \mathbb{E} \{ \ell(Y, f(X)) \} = \int_{\mathcal{X} \times \mathcal{Y}} \ell(y, f(\mathbf{x})) dP(\mathbf{x}, y)$$

- f^* : Bayes predictor
- $f_{\mathcal{F}}$: best predictor in the set \mathcal{F}
- R^* : Bayes risk
- Definition: empirical risk (training error) of predictor f

$$R_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i))$$

- \hat{f}_n : estimator based on \mathcal{D}_n

B: For the time being I keep \mathcal{F} and f but we may switch to g and \mathcal{G} all above.

N: OK

Classification

- $h : \mathcal{X} \rightarrow \mathcal{Y}$: base classifier
- \mathcal{H} : set of base classifiers
- \mathcal{Z} : discriminant space, $\mathcal{Z} = \mathbb{R}$ in binary classification, $\mathcal{Z} = \mathbb{R}^K$ in multi-class (and multi-class multi-label) classification.
- $f : \mathcal{X} \rightarrow \mathcal{Z}$: discriminant function
- \mathcal{F} : set of discriminant functions
- g_f : classifier obtained from the discriminant function f . In binary classification $g_f = 2\mathbb{I}\{f > 0\} - 1$, in multi-class classification $g_f = \arg \max_{k \in \{1, \dots, K\}} f_k$
- \mathcal{F}_N : set of finite convex combinations of N elements of \mathcal{H}
- \mathcal{F}_λ : set of linear combinations of elements of \mathcal{H} such that the 1-norm of the coefficients is equal to λ .
- $\mathcal{F}_{N,\lambda}$: set of linear combinations of N elements of \mathcal{H} such that the 1-norm of the coefficients is equal to λ .
- α_j : real valued coefficients
- $f = \sum_{j \geq 1} \alpha_j h_j$: element of \mathcal{F}
- $\alpha = (\alpha_1, \dots)$: coefficient vector
- $\|\alpha\|_1 = \sum_{j \geq 1} \alpha_j$: 1-norm of the coefficient vector
- $\tilde{\alpha}_j = \frac{\alpha_j}{\|\alpha\|_1}$: normalized coefficients
- $\tilde{\alpha} = (\tilde{\alpha}_1, \dots)$: normalized coefficient vector (discrete probability distribution over \mathcal{H})
- $\tilde{f} = \sum_{j \geq 1} \tilde{\alpha}_j h_j \in \mathcal{F}_1$: normalized discriminant function
- $g : \mathcal{X} \rightarrow \{-1, 1\}$: final (binary) classifier
- \mathcal{G} : set of classifiers
- V : VC dimension

- $\mathbf{w} = (w_1, \dots, w_n) : \text{distribution over } \mathcal{D}_n$
- $g_f = 2\mathbb{I}\{f > 0\} - 1 : \text{classifier obtained from the real-valued predictor } f$
- Definition: risk of classifier g

$$L(g) = \mathbb{P}\{Y \neq g(X)\} = \int_{\mathcal{X} \times \mathcal{Y}} \mathbb{I}\{y \neq g(\mathbf{x})\} dP(\mathbf{x}, y)$$

- $g^* : \text{Bayes classifier}$
- $g_{\mathcal{G}} : \text{best classifier in the restricted set } \mathcal{G}$
- Definition: empirical risk of classifier g

$$L_n(g) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{Y_i \neq g(X_i)\}$$

- $\hat{g}_n : \text{estimator based on } \mathcal{D}_n$
- $\varphi : (\text{surrogate}) \text{ convex cost function}$
- Definition: convex risk of predictor f

$$A(f) = \mathbb{E}\{\varphi(-Y f(X))\}$$

- Definition: empirical convex risk of predictor f

$$A_n(f) = \frac{1}{n} \sum_{i=1}^n \varphi(-Y_i f(X_i))$$

B: I came back on this issue to an earlier proposition. I know it's not perfect, but I don't know what to do. Clearly, in the algorithmic section the data points will be given, so the margin varies with the discriminant function (so f is the more natural variable), while in the statistical section the function is given, and the natural variable is the data point X .

- $\rho_Z(f) = f(X)Y : \text{margin of a discriminant function } f \text{ on data point } (X, Y) \text{ (multi-class, regression: later)}$
- $\tilde{\rho}_Z(f) = \tilde{f}(X)Y : \text{normalized margin of a discriminant function } f \text{ on data point } (X, Y)$
- $\rho_i(\tilde{\rho}_i) : \text{margin (normalized margin) of a discriminant function (omitted in the notation) on the training data point } (X_i, Y_i) \text{ or } (\mathbf{x}_i, y_i).$
- $\rho_i(\alpha)(\tilde{\rho}_i(\alpha)) : \text{margin (normalized margin) of the discriminant function } \sum_{j \geq 1} \alpha_j h_j \text{ on the training point } (X_i, Y_i) \text{ or } (\mathbf{x}_i, y_i).$
- $\rho(\alpha)(\tilde{\rho}(\alpha)) : \text{margin (normalized margin) vector over the data set } \mathcal{D}_n.$
- Definition: marginal empirical risk of a discriminant function f

$$L_n^{(\rho)}(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{f(X_i)Y_i < \rho\} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\{\rho_i < \rho\}$$

Algorithm

- T : number of iterations in the algorithm
- t : iteration index
- $\mathbf{w}^{(t)}$: distribution over \mathcal{D}_n in the t th iteration (computed before the t th iteration)
- $h^{(t)}$: base classifier selected in the t th iteration
- $\epsilon^{(t)} = \sum_{i=1}^n w_i^{(t)} \mathbb{I} \left\{ h^{(t)}(\mathbf{x}_i) \neq y_i \right\}$: weighted error of base classifier in the t th iteration
- $\gamma^{(t)} = 1 - 2\epsilon^{(t)}$: edge of $h^{(t)}$
- $\alpha^{(t)}$: coefficient of $h^{(t)}$
- $f^{(t)} = \sum_{t=1}^T \alpha^{(t)} h^{(t)}$: aggregate discriminant function, output of ADABOOST
- $\rho_i^{(t)} = \rho_{f^{(t)}}(Z_i)$: margin of the i th training point after the t th iteration

Chapter 1

Introduction

Definitions: observations, features (aka. attribute, coordinate), labels, classifier, discriminant function, margin, loss, risk, empirical loss, empirical risk, empirical risk minimization, margin risk functions, hyper-parameters, regularization (give just hints, no deep explanation how they are related).

For the time being I'm defining everything where I'm using it. For pedagogical reasons we could keep this policy: give basic definitions in this chapter, so if somebody wants to look them up, everything can be found here, but also define them later when they are first used in the boosting context.

Notations: bold for vectors: \mathbf{x} ; same letter but italic for the elements: x_i ; capital letter for random variables (vector or scalar): X, Y ; capital bold for matrices, random or not: \mathbf{X} . Same letter for indices over the same variable: i over data points, j over base classifiers, t over iterations, ℓ over classes (multiclass).

Due to its simplicity, ADABOOST is also a very good algorithm to introduce machine learning. It can be taught after a first year programming and algorithms class. The key elements of binary AdaBoosting decision stumps is completely elementary. At the same time, it's state-of-the-art on benchmark tests. Implies some choices, e.g., we give the pseudocode of decision stumps in the second chapter so the whole algorithm can be programmed from scratch.

Our main audience is a novice, but we hope to give some novelty (or at least synthesize) for the ML expert (who is not necessarily a boosting expert).

Although several questions remain unanswered (no foolproof explanation why it "works"), it is a very good algo in practice, so deserves a monographie.

1.1 Bibliographical notes

Bishop, Duda-Hart, Devroye-Györfi-Lugosi

Chapter 2

The basic algorithm: binary ADABOOST

In this chapter we describe the basic ADABOOST algorithm used to learn binary (two-class) classifiers. Algorithmic issues, algorithmic convergence results. Section on base classifiers. Simple example.

2.1 Basics

data, error, margin, error minimization

2.2 The binary ADABOOST algorithm

ADABOOST is an *ensemble* (or *meta-learning*) method that constructs a classifier in an iterative fashion. In each iteration, it calls a simple learning algorithm (called the *base learner*) that returns a classifier, and assigns a weight coefficient to it. The final classification will be decided by a weighted “vote” of the base classifiers. The smaller the error of the base classifier, the larger is its weight in the final vote. The base classifiers have to be only slightly better than a random guess (from where their alternative name *weak classifier* derives), which gives great flexibility to the design of the base classifier (or feature) set.

For the formal description of ADABOOST, let the training set be $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$. The algorithm runs for T iterations. T is the only pre-specified hyper-parameter of ADABOOST that can be set by, for example, cross-validation. In each iteration $t = 1, \dots, T$, we choose a base classifier $h^{(t)}$ from a set \mathcal{H} of classifiers and set its coefficient $\alpha^{(t)}$. In the simplest version of the algorithm, \mathcal{H} is a finite set of binary classifiers of the form $h : \mathbb{R}^d \rightarrow \{-1, 1\}$, and the base learner executes an exhaustive search over \mathcal{H} in each iteration. The output of ADABOOST is a *discriminant function* constructed as a weighted vote of the base classifiers

$$f^{(T)}(\cdot) \triangleq \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot). \quad (2.1)$$

The sign of $f^{(T)}(\mathbf{x})$ is then used as the final classification of \mathbf{x} .

The algorithm maintains a weight distribution $\mathbf{w}^{(t)} = (w_1^{(t)}, \dots, w_n^{(t)})$ over the data points¹. The weights are initialized uniformly in line 1, and are updated in each iteration in lines 7-10 (Figure 2.1). The goal of the base learner is to minimize the *weighted (base) error*

$$\epsilon^{(t)} \triangleq \sum_{i=1}^n w_i^{(t)} \mathbb{I} \{h^{(t)}(\mathbf{x}_i) \neq y_i\}.^2 \quad (2.2)$$

Reference to decision stumps.

¹To avoid confusion, from now on we will call the base classifiers' weights $\alpha^{(t)}$ *coefficients*, and keep the term *weight* for the weights w_i of the data points.

²The indicator function $\mathbb{I}\{A\}$ is 1 if its argument A is true and 0 otherwise.

The coefficient $\alpha^{(t)}$ of $h^{(t)}$ is set in line 5 to

$$\alpha^{(t)} \triangleq \frac{1}{2} \ln \left(\frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right). \quad (2.3)$$

```

ADABOOST( $\mathcal{D}_n, \text{BASE}(\cdot, \cdot), T$ )
1  $\mathbf{w}^{(1)} \leftarrow (1/n, \dots, 1/n)$   $\triangleright$  initial weights
2 for  $t \leftarrow 1$  to  $T$ 
3    $h^{(t)} \leftarrow \text{BASE}(\mathcal{D}_n, \mathbf{w}^{(t)})$   $\triangleright$  base classifier
4    $\epsilon^{(t)} \leftarrow \sum_{i=1}^n w_i^{(t)} \mathbb{I} \{ h^{(t)}(\mathbf{x}_i) \neq y_i \}$   $\triangleright$  weighted error of the base classifier
5    $\alpha^{(t)} \leftarrow \frac{1}{2} \ln \left( \frac{1 - \epsilon^{(t)}}{\epsilon^{(t)}} \right)$   $\triangleright$  coefficient of the base classifier
6   for  $i \leftarrow 1$  to  $n$   $\triangleright$  re-weighting the training points
7     if  $h^{(t)}(\mathbf{x}_i) \neq y_i$  then  $\triangleright$  error
8        $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2\epsilon^{(t)}}$   $\triangleright$  weight increases
9     else  $\triangleright$  correct classification
10       $w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{2(1-\epsilon^{(t)})}$   $\triangleright$  weight decreases
11 return  $f^{(T)}(\cdot) = \sum_{t=1}^T \alpha^{(t)} h^{(t)}(\cdot)$   $\triangleright$  weighted “vote” of base classifiers

```

Figure 2.1: The pseudocode of ADABOOST. $\mathcal{D}_n = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ is the training set, $\text{BASE}(\cdot, \cdot)$ is the base learner, and T is the number of iterations.

The weights remain normalized throughout the algorithm, that is, for all t

$$\sum_{i=1}^n w_i^{(t)} = 1. \quad (2.4)$$

To see that, first note that the initialization guarantees that $\sum_{i=1}^n w_i^{(1)} = 1$. Then, assuming that $\sum_{i=1}^n w_i^{(t)} = 1$, we have

$$\sum_{i=1}^n w_i^{(t+1)} = \sum_{i=1}^n \frac{w_i^{(t)}}{2\epsilon^{(t)}} \mathbb{I} \{ h^{(t)}(\mathbf{x}_i) \neq y_i \} + \sum_{i=1}^n \frac{w_i^{(t)}}{2(1-\epsilon^{(t)})} \mathbb{I} \{ h^{(t)}(\mathbf{x}_i) = y_i \} \quad (2.5)$$

$$= \frac{1}{2\epsilon^{(t)}} \sum_{i=1}^n w_i^{(t)} \mathbb{I} \{ h^{(t)}(\mathbf{x}_i) \neq y_i \} + \frac{1}{2(1-\epsilon^{(t)})} \sum_{i=1}^n w_i^{(t)} (1 - \mathbb{I} \{ h^{(t)}(\mathbf{x}_i) \neq y_i \}) \quad (2.6)$$

$$= \frac{1}{2\epsilon^{(t)}} \epsilon^{(t)} + \frac{1}{2(1-\epsilon^{(t)})} (1 - \epsilon^{(t)}) = \frac{1}{2} + \frac{1}{2} = 1. \quad (2.7)$$

In (2.5) we applied the weight update rules in lines 8 and 10, in (2.6) we used the fact that $\mathbb{I} \{ A \} = 1 - \mathbb{I} \{ \neg A \}$, and in (2.7) we applied the definition of $\epsilon^{(t)}$ in (2.2) and the inductive assumption of $\sum_{i=1}^n w_i^{(t)} = 1$.

Since the weights are normalized, $\epsilon^{(t)}$ is a real number between 0 and 1. Moreover, if \mathcal{H} is closed under negation,³ we can guarantee that $\epsilon^{(t)} \leq \frac{1}{2}$, since otherwise we would flip signs and use $-h^{(t)}$ instead of $h^{(t)}$.

The first implication of $\epsilon^{(t)} \leq \frac{1}{2}$ is that $\alpha^{(t)}$ is always non-negative. The coefficient monotonically increases as $\epsilon^{(t)}$ decreases (Figure 2.2), which means that the better the base classifier $h^{(t)}$, the higher is its “vote” in the final classification. As the base error $\epsilon^{(t)}$ goes to 0, the coefficient $\alpha^{(t)}$ tends to infinity. Formally, a “strong” base classifier with $\epsilon^{(t)} = 0$ will automatically receive an infinite coefficient, so it will dominate the linear combination (2.1). In principle, since ADABOOST’s objective is to minimize the training error, $\epsilon^{(t)} = 0$ means that the objective has been achieved, so it is normal to terminate the algorithm. In certain variants of ADABOOST it will be useful to be able to continue even if $\epsilon^{(t)} = 0$; we will discuss how to deal with this problem in Section ??.

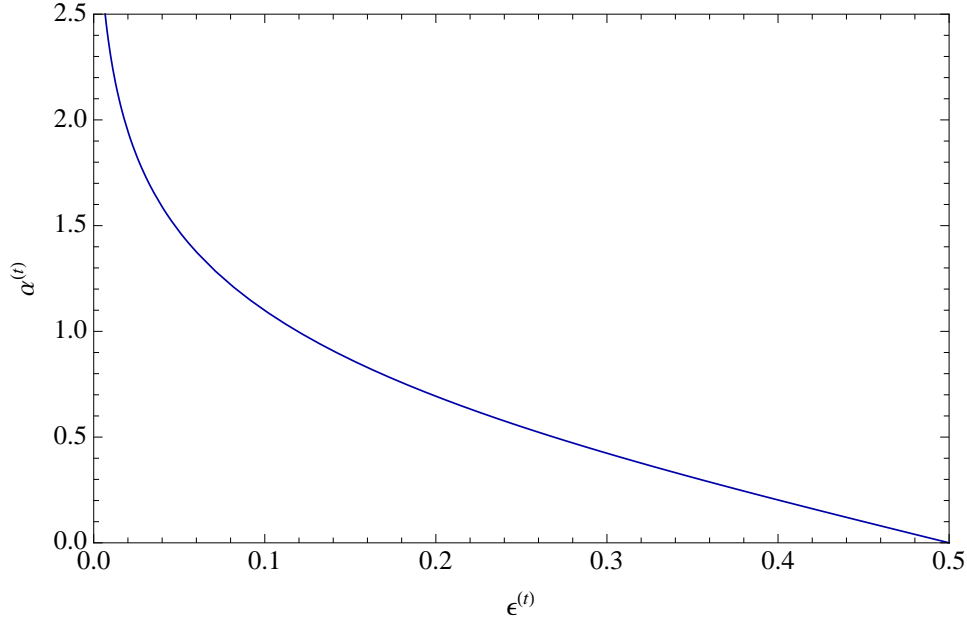


Figure 2.2: The base coefficient $\alpha^{(t)}$ as a function of the base error $\epsilon^{(t)}$. At $\epsilon^{(t)} = \frac{1}{2}$, $\alpha^{(t)} = 0$, and $\alpha^{(t)} \rightarrow \infty$ as $\epsilon^{(t)} \rightarrow 0$.

The second implication of $\epsilon^{(t)} \leq \frac{1}{2}$ is that when the weights w_i of training points are updated in lines 7-10, the weights of misclassified points increase (line 8), and the weights of correctly classified points decrease (line 10). As the algorithm progresses, the weights of frequently misclassified points will tend to be large, so base classifiers will concentrate more and more on these “hard” data points.

Finally, by looking at the evolution of the first summand of (2.5), we can also note that the base error $\epsilon^{(t+1)}$ of the t th base classifier $h^{(t)}$ is $\frac{1}{2}$, so we never select the same base classifier in two consecutive iterations (unless we are in a degenerate case).

³That is, if $h \in \mathcal{H}$ then $-h \in \mathcal{H}$. This property is satisfied by most of the practical classifier classes, and we will assume it from now on. If \mathcal{H} is not closed under negation and we cannot find a base classifier with $\epsilon^{(t)} \leq \frac{1}{2}$, then the algorithm terminates and returns $f^{(t-1)}$. In the degenerate case when $\epsilon^{(t)} = \frac{1}{2}$, we can continue the loop, but since $\alpha^{(\tau)} = 0$ for all $\tau \geq t$, we should also terminate in practice.

2.3 Basic analysis

First we unify the two cases in the weight update formulae (lines 7-10):

$$\begin{aligned}
 w_i^{(t+1)} &= w_i^{(t)} \times \begin{cases} \frac{1}{2\varepsilon^{(t)}} & \text{if } h^{(t)}(\mathbf{x}_i) \neq y_i \\ \frac{1}{2(1-\varepsilon^{(t)})} & \text{if } h^{(t)}(\mathbf{x}_i) = y_i \end{cases} \\
 &= w_i^{(t)} \times \begin{cases} \frac{\sqrt{\frac{1-\varepsilon^{(t)}}{\varepsilon^{(t)}}}}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}} & \text{if } h^{(t)}(\mathbf{x}_i) \neq y_i \\ \frac{\sqrt{\frac{\varepsilon^{(t)}}{1-\varepsilon^{(t)}}}}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}} & \text{if } h^{(t)}(\mathbf{x}_i) = y_i \end{cases} \\
 &= w_i^{(t)} \times \begin{cases} \frac{e^{\alpha^{(t)}}}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}} & \text{if } h^{(t)}(\mathbf{x}_i)y_i = -1 \\ \frac{e^{-\alpha^{(t)}}}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}} & \text{if } h^{(t)}(\mathbf{x}_i)y_i = 1 \end{cases} \quad (2.8)
 \end{aligned}$$

$$= w_i^{(t)} \frac{\exp\left(-\alpha^{(t)}h^{(t)}(\mathbf{x}_i)y_i\right)}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}}, \quad (2.9)$$

where in (2.8) we use the definition (2.3) of the coefficient $\alpha^{(t)}$ and the fact that both y_i and $h^{(t)}(\mathbf{x}_i)$ are $\{-1, 1\}$ -valued. Using the unified equation (2.9), the definition (2.1) of the final discriminant function $f^{(T)}$, and the uniform initialization of the weights $w_i^{(1)}$ in line 1, we can express the weights in a non-recursive way as

$$\begin{aligned}
 w_i^{(t+1)} &= w_i^{(t)} \frac{\exp\left(-\alpha^{(t)}h^{(t)}(\mathbf{x}_i)y_i\right)}{2\sqrt{\varepsilon^{(t)}(1-\varepsilon^{(t)})}} = \dots = \frac{1}{n} \times \frac{\exp\left(-y_i \sum_{\tau=1}^t \alpha^{(\tau)} h^{(\tau)}(\mathbf{x}_i)\right)}{\prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})}} \\
 &= \frac{\exp\left(-f^{(t)}(\mathbf{x}_i)y_i\right)}{n \prod_{\tau=1}^t 2\sqrt{\varepsilon^{(\tau)}(1-\varepsilon^{(\tau)})}}. \quad (2.10)
 \end{aligned}$$

First, this non-recursive formula tells us that the weights are proportional to $\exp\left(-f^{(t)}(\mathbf{x}_i)y_i\right)$. The negative exponent

$$\rho_i^{(t)} \triangleq f^{(t)}(\mathbf{x}_i)y_i$$

is the *unnormalized margin*, defined in general for a discriminant function f and a data point (\mathbf{x}, y) as

$$\rho_{\mathbf{x}, y}(f) \triangleq f(\mathbf{x})y. \quad (2.11)$$

This quantity will play a crucial role in the margin-based analysis of the generalization error (Section ??). For now it is clear, that 1) the sign of the margin $\rho_i^{(t)}$ indicates whether (\mathbf{x}_i, y_i) is correctly classified by $f^{(t)}$ or not, and 2) the larger $\rho_i^{(t)}$, the more “confident” the classification of (\mathbf{x}_i, y_i) is, so intuitively it makes sense that the weight $w_i^{(t+1)}$ decreases monotonically with $\rho_i^{(t)}$.

Second, from (2.7) we know that the weights $w_i^{(t+1)}$ sum to 1, so from the non-recursive formula (2.10) it follows that

$$w_i^{(t+1)} = \frac{\exp\left(-f^{(t)}(\mathbf{x}_i)y_i\right)}{\sum_{l=1}^n \exp\left(-f^{(t)}(\mathbf{x}_l)y_l\right)}. \quad (2.12)$$

By comparing (2.10) and (2.12) we automatically obtain the key equation

$$\frac{1}{n} \sum_{i=1}^n \exp\left(-f^{(t)}(\mathbf{x}_i)y_i\right) = \prod_{\tau=1}^t 2\sqrt{\epsilon^{(\tau)}(1-\epsilon^{(\tau)})} \quad (2.13)$$

of the convergence theorem:

Theorem 1 (Theorem 6 in [FS97]). Assuming that there exists a positive constant δ for which the base error is upper bounded by

$$\epsilon^{(t)} \leq \frac{1}{2} - \delta \quad (2.14)$$

(*** maybe $\frac{1-\delta}{2}$ would be better: lower bound on the edge) for all t , the training error $R_n\left(f^{(T)}\right)$ becomes 0 after at most

$$T = \left\lceil \frac{\ln n}{2\delta^2} \right\rceil + 1$$

iterations.

Proof. First we upper bound the training error by

$$\begin{aligned} R_n\left(f^{(T)}\right) &\triangleq \frac{1}{n} \sum_{i=1}^n \mathbb{I}\left\{\text{sign}\left(f^{(T)}(\mathbf{x}_i)\right) \neq y_i\right\} = \frac{1}{n} \sum_{i=1}^n \mathbb{I}\left\{f^{(T)}(\mathbf{x}_i)y_i < 0\right\} \\ &\leq \frac{1}{n} \sum_{i=1}^n \exp\left(-f^{(T)}(\mathbf{x}_i)y_i\right) \end{aligned} \quad (2.15)$$

$$= \prod_{t=1}^T 2\sqrt{\epsilon^{(t)}(1-\epsilon^{(t)})} = \prod_{t=1}^T \sqrt{1 - (1-2\epsilon^{(t)})^2} \quad (2.16)$$

$$\leq \left(\sqrt{1-4\delta^2}\right)^T \quad (2.17)$$

$$\leq \exp(-2T\delta^2). \quad (2.18)$$

In (2.15) we use $\mathbb{I}\{x < 0\} \leq e^{-x}$, (2.16) was shown in (2.13), in (2.17) we apply the assumption that $\epsilon^{(t)} < \frac{1}{2} - \delta$, and (2.18) follows from $1-x \leq e^{-x}$. The theorem then follows from the discrete nature of the training error: once the upper bound $\exp(-2T\delta^2)$ falls below $\frac{1}{n}$, the error $R_n\left(f^{(T)}\right)$ must be 0. \square

The main significance of Theorem 1 is that it gives a positive answer to PAC learning's (***) explain in the intro) boosting conjecture: it is possible to construct a strong learner (with risk close to 0) by combining a small number of weak learners (with risks only slightly smaller than $\frac{1}{2}$). The exponential convergence speed of the training error, implying that the number of weak learners is $\sim \ln n$, plays a crucial role in proving the conjecture. We will give the formal statement and further discuss its significance in Chapter ??.

The upper bound (2.16) actually explains the base learner's objective of minimizing $\epsilon^{(t)}$. It is also clear that by minimizing $\epsilon^{(t)}$, we do not directly minimize the training error $R_n\left(f^{(T)}\right)$, only through its upper bound

$$R_n^{(e)}\left(f^{(T)}\right) \triangleq \frac{1}{n} \sum_{i=1}^n \exp\left(-f^{(T)}(\mathbf{x}_i)y_i\right) \quad (2.19)$$

called the *exponential empirical risk*. Minimizing such convex upper bounds instead of the training error is a common approach in machine learning for two reasons. First, minimizing $R_n\left(f^{(T)}\right)$ is often computationally hard even over simple function classes. Second, minimizing monotonically decreasing margin losses is often justified by probabilistic arguments: we achieve larger margins and better generalization error (risk) by “smoothing” the empirical risk.

The upper bound (2.16) also illustrates another important feature of ADABOOST: the greedy nature of its underlying optimization. Indeed, we do not *globally* minimize the exponential risk $R_n^{(e)}(f^{(T)})$ over the linear combinations of \mathcal{H} , rather, in each iteration t we add the base learner $h^{(t)}$ and its coefficient $\alpha^{(t)}$ which is *optimal given the classifier $f^{(t-1)}$* constructed so far. This general interpretation of ADABOOST of greedily optimizing the exponential risk yielded several generalizations and extensions. We will prove the greedy optimization statement formally in Section ??, and show that under certain conditions this local optimization does lead to the global minimum of $R_n^{(e)}(f^{(T)})$.

The condition (2.14) is actually stronger than it seems to be: even for a given data set it requires that the base error is bounded away from $\frac{1}{2}$ for *any* weight distribution $\mathbf{w}^{(t)}$. In fact it can happen easily that (2.14) cannot be satisfied: all we need is a data set \mathcal{D}_n that cannot be separated by any linear combination over \mathcal{H} . In this case the training error can obviously not go to 0, and the base errors $\epsilon^{(t)}$ actually tend to $\frac{1}{2}$. In the limit, the weight distribution is such that for *all* $h \in \mathcal{H}$, the base error is $\frac{1}{2}$. On the other hand, if \mathcal{D}_n is separable by a linear combination over \mathcal{H} , then we can prove that (2.14) can always be satisfied, and the actual value of δ (and so the convergence speed) can be related to the separating *margin* between the two classes. We will formalize this statement in Section ??.

The term $1 - 2\epsilon^{(t)}$ that appears in (2.16) is called the *edge* of the base classifier $h^{(t)}$, denoted by $\gamma^{(t)}$. In a certain sense it is a more natural quantity than the error $\epsilon^{(t)}$, and it will ease the notation in several statements and derivations. It is easy to see that

$$\begin{aligned} \gamma^{(t)} &\triangleq 1 - 2\epsilon^{(t)} = 1 - 2 \sum_{i=1}^n w_i^{(t)} \mathbb{I} \{h^{(t)}(\mathbf{x}_i) \neq y_i\} \\ &= \sum_{i=1}^n w_i^{(t)} \left(1 - 2\mathbb{I} \{h^{(t)}(\mathbf{x}_i) \neq y_i\}\right) \end{aligned} \quad (2.20)$$

$$= \sum_{i=1}^n w_i^{(t)} h^{(t)}(\mathbf{x}_i) y_i, \quad (2.21)$$

where (2.20) follows from $\sum_{i=1}^n w_i^{(t)} = 1$ and (2.21) is true because both y_i and $h^{(t)}(\mathbf{x}_i)$ are $\{-1, 1\}$ -valued. The edge of a random decision is 0 (on average), so $\gamma^{(t)}$ quantifies how much $h^{(t)}$ is better than a random decision (its “edge” over a random decision). Several formulae of the basic algorithm can either be simplified or they become more “symmetric” using the edge instead of the error. For example, the base coefficient (2.3) can be rewritten as

$$\alpha^{(t)} = \frac{1}{2} \ln \left(\frac{1 + \gamma^{(t)}}{1 - \gamma^{(t)}} \right), \quad (2.22)$$

the weight update in lines 7-10 simplifies to

$$w_i^{(t+1)} \leftarrow \frac{w_i^{(t)}}{1 + \gamma^{(t)} h^{(t)}(\mathbf{x}_i) y_i},$$

and the exponential risk (2.13) can be expressed as

$$R_n^{(e)}(f^{(T)}) = \frac{1}{n} \sum_{i=1}^n \exp(-f^{(T)}(\mathbf{x}_i) y_i) = \prod_{t=1}^T \sqrt{(1 + \gamma^{(t)})(1 - \gamma^{(t)})} = \prod_{t=1}^T \sqrt{1 - \gamma^{(t)2}} \quad (2.23)$$

The definition of the margin (2.11) provides another interpretation of (2.21): the edge is the (weighted) *average margin* of the base classifier $h^{(t)}$. In fact the edge is intimately related to margin, and we will elaborate the game-theoretical interplay between the two in Section ??.

As a final remark, note that the upper bound (2.18) can be sharpened by using the a-posteriori edges $\gamma^{(t)}$ instead of their a-priori lower bound 2δ to

$$R_n(f^{(T)}) \leq \exp \left(-\frac{1}{2} \sum_{t=1}^T \gamma^{(t)2} \right), \quad (2.24)$$

so ADABOOST is guaranteed to achieve 0 training error once $\sum_{t=1}^T \gamma^{(t)^2}$ exceeds $2 \ln n$.

2.4 Base classifiers

The algo doesn't have to be perfect: we don't need the minimum. The design of the base classifier set can be determined by a-priori, domain dependent knowledge.

2.4.1 Decision stumps

If there is no particular a-priori knowledge available on the domain of the learning problem, small decision trees or, in the extreme case, *decision stumps* (decision trees with two leaves) are often used. A decision stump can be defined by three parameters, the index j of the feature⁴ that it cuts, the threshold θ of the cut, and the sign of the decision. Formally, a *positive stump* is defined by

$$h_{j,\theta+}(\mathbf{x}) \triangleq 2I_{\{x^{(j)} \geq \theta\}} - 1 = \begin{cases} 1 & \text{if } x^{(j)} \geq \theta, \\ -1 & \text{otherwise,} \end{cases} \quad (2.25)$$

and a *negative stump* is defined by

$$h_{j,\theta-}(\mathbf{x}) \triangleq -h_{j,\theta+}(\mathbf{x}) = 2I_{\{x^{(j)} < \theta\}} - 1 = \begin{cases} 1 & \text{if } x^{(j)} < \theta, \\ -1 & \text{otherwise.} \end{cases} \quad (2.26)$$

Although decision stumps may seem very simple, when boosted, they yield good classifiers in practice. The algorithm we describe here can also be used as a basic building block for more complex base learners, such as trees or products.

The optimal decision stump (that minimizes the the base error $\epsilon^{(t)}$ (2.2)) can be found in $O(nd)$ time using exhaustive search (Figure 2.3). The pseudocode is slightly simpler if we maximize the edge $\gamma^{(t)}$ (2.21) instead of minimizing the error $\epsilon^{(t)}$, so we adopt this notation. First, all feature vectors $\mathbf{x}^{(j)} = (x_1^{(j)}, \dots, x_n^{(j)})$ must be sorted in increasing order. Technically this should be done by creating d indirect index sets and storing the order permutations. This step precedes the boosting loop, so the overall complexity of boosting decision stumps is $O(dn \log n + Tdn) = O(dn(\log n + T))$. To ease the notation, we assume that the feature vectors are sorted, that is, $x_1^{(j)} \leq \dots \leq x_n^{(j)}$ for all $j = 1, \dots, d$.

The algorithm will examine all possible cutting thresholds θ between non-identical consecutive feature values $x_{i-1}^{(j)}$ and $x_i^{(j)}$. We cut halfway between the feature values, so for all $i = 2, \dots, n$ and $j = 1, \dots, d$,

$$\theta_i^{(j)} = \frac{x_i^{(j)} + x_{i-1}^{(j)}}{2}.$$

Since we have $O(nd)$ thresholds and computing the edge $\gamma^{(t)}$ requires $O(n)$ operations, the naive algorithm would take $O(n^2d)$ time to find the best stump. In Figure 2.3 we use the fact that the edge of “consecutive” stumps can be updated in constant time, so we can find the best stump in linear time using one sweep over data points and features.

For each feature, we start with the constant classifier $h_0(\mathbf{x}) \equiv 1$ (which is equivalent to a positive stump with $\theta < x_1^{(j)}$) and set the edge γ in line 4 to the edge of the constant classifier

$$\gamma_0 = \sum_{i=1}^n w_i y_i.$$

⁴We assume that features are real valued, so observations \mathbf{x} are in \mathbb{R}^d .

```

DECISIONSTUMP( $\mathcal{D}_n, \mathbf{w}$ )
1  $\gamma_0 \leftarrow \sum_{i=1}^n w_i y_i$   $\triangleright$  edge of constant classifier  $h_0(\mathbf{x}) \equiv 1$ 
2  $\gamma^* \leftarrow \gamma_0$   $\triangleright$  best edge
3 for  $j \leftarrow 1$  to  $d$   $\triangleright$  all (numeric) features
4    $\gamma \leftarrow \gamma_0$   $\triangleright$  edge of the constant classifier
5   for  $i \leftarrow 2$  to  $n$   $\triangleright$  all points in order  $x_1^{(j)} \leq \dots \leq x_n^{(j)}$ 
6      $\gamma \leftarrow \gamma - 2w_{i-1}y_{i-1}$   $\triangleright$  update edge of positive stump
7     if  $x_{i-1}^{(j)} \neq x_i^{(j)}$  then  $\triangleright$  no threshold if identical coordinates
8       if  $|\gamma| > |\gamma^*|$  then  $\triangleright$  found better stump
9          $\gamma^* \leftarrow \gamma$   $\triangleright$  update best edge
10         $j^* \leftarrow j$   $\triangleright$  update index of best feature
11         $\theta^* \leftarrow \frac{x_i^{(j)} + x_{i-1}^{(j)}}{2}$   $\triangleright$  update best threshold
12 if  $\gamma^* = \gamma_0$   $\triangleright$  did not beat the constant classifier
13   return  $\text{sign}(\gamma_0) \times h_0$   $\triangleright \pm$  constant classifier
14 else
15   return  $\text{sign}(\gamma^*) \times h_{j^*, \theta^*+}$   $\triangleright$  best stump

```

Figure 2.3: Exhaustive search for the best decision stump. We assume that the coordinates $x_1^{(j)}, \dots, x_n^{(j)}$ are sorted for each feature j . All thresholds θ halfway between non-identical coordinates are examined.

The edge γ is then updated in line 6 by “moving” the next feature value $x_{i-1}^{(j)}$ to the left side of the threshold θ (Figure 2.4). Since we cannot cut between two identical feature values, we only examine the edge if a new stump can be defined (line 7). To use only one sweep, we check for negative and positive stumps at the same time, so we update the edge, the threshold, and the feature index if the absolute value of γ is larger than the absolute value of the currently best γ^* (line 8). At the end of the algorithm we first check if any of the stumps has beaten the constant classifier (line 12), then return either a positive or a negative stump, depending on the sign of γ^* .

*** The same one-sweep approach can be done for other losses such as entropy, Gini-index, quadratic, but, interestingly, not for the absolute error (in case of regression stumps) *exercise*.

2.5 A simple example

In this section we illustrate the theoretical concepts by executing ADABOOST with decision stumps on a simple example. We generate 200 two-dimensional observations $\mathbf{x}_i = (x_i^{(1)}, x_i^{(2)})$ uniformly in a unit square, and, to pose a slight challenge to decision stumps, we define the classes by a diagonal cut (Figure 2.5). Formally, $y_i = 1$ if and only if $x_i^{(1)} + x_i^{(2)} \leq 1$.

The execution of ADABOOST can be followed in Figure 2.6. We plot the data points and the decision regions at $t = 1, 3, 5, 10, 20, 40$ iterations. The areas of the disks are proportional to the weights $w_i^{(t)}$. It can be seen very well that the weights of “hard” data points near the decision boundary will dominate as the ADABOOST progresses. Figure 2.7 shows the same situation after $t = 40$ iterations, but this time we modulate the color intensities of the decision regions by the margin $\rho_{\mathbf{x}, y}(f) = f^{(t)}(\mathbf{x})y$. Naturally, far from the decision boundary the classification is very confident, whereas the light color indicates small margins close to the diagonal decision boundary. The plot

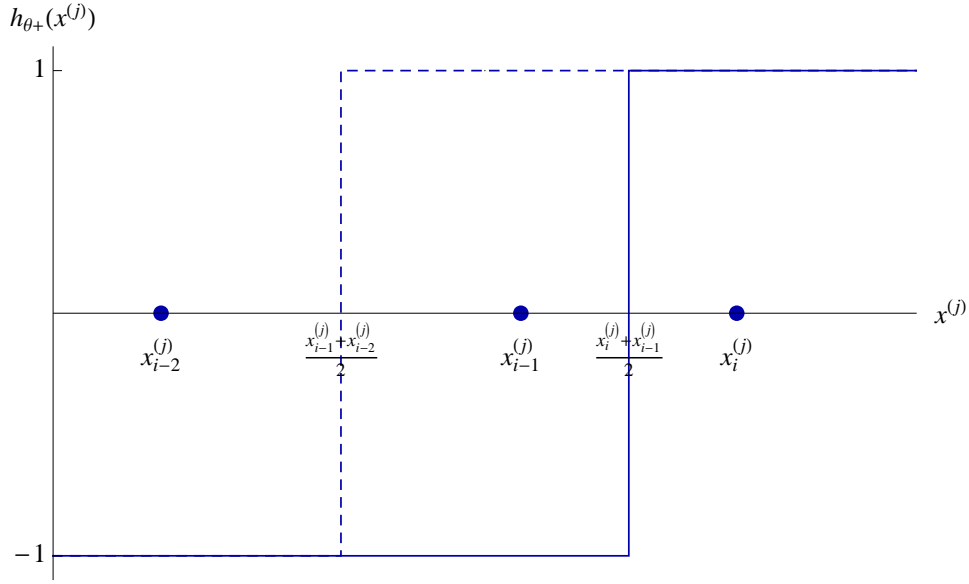


Figure 2.4: Updating the edge γ in line 6 of DECISIONSTUMP($\mathcal{D}_n, \mathbf{w}$). If $y_{i-1}^{(j)} = 1$, then γ decreases by $2w_{i-1}$, and if $y_{i-1}^{(j)} = -1$, then γ increases by $2w_{i-1}$.

also illustrates the formal relationship (2.12) between the weights $w_i^{(t)}$ and the margin $\rho_{\mathbf{x}, \mathbf{y}}(f)$: first, the smaller the margin, the larger are the weights, and second, the weights only depend on the margin, so the disks have an identical size in constant intensity rectangles.

In Figure 2.8 we plot the training error and its two bounds as a function of the number of iterations. Figure 2.8(a) indicates that the training error becomes 0 after 26 iterations. In Section 2.3 we showed that ADABOOST attempts to minimize the exponential risk $R_n^{(e)}(f^{(t)})$ (2.23) in a greedy fashion; the red curve in Figure 2.8(b) indicates that $R_n^{(e)}(f^{(t)})$ is a rather tight upper bound of the training error. On the other hand, the a-priori bound⁵ (2.17), used to establish the convergence speed independently of the algorithm's dynamics, may be quite loose, as indicated by the green curve in Figure 2.8(b).

2.6 Discussion

For the neurally inclined: it's like constructing a one-hidden-layer neural net incrementally.

The entropy – info-theoretical view.

The game-theoretical view.

Boosting as feature selector.

2.7 Bibliographical notes

⁵In the case of a finite base classifier class \mathcal{H} , the constant δ that guarantees a universal upper bound $\frac{1}{2} - \delta$ on the base error $\epsilon^{(t)}$ can be determined using linear programming; see Section ??.

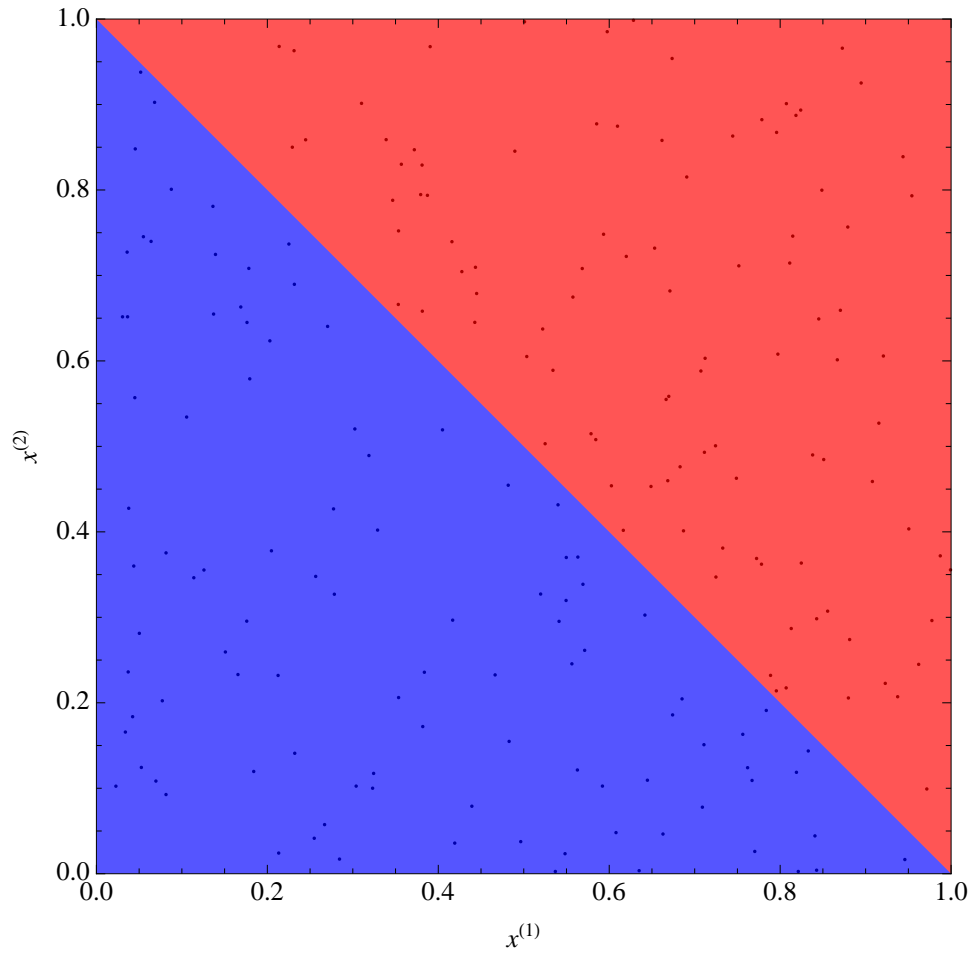
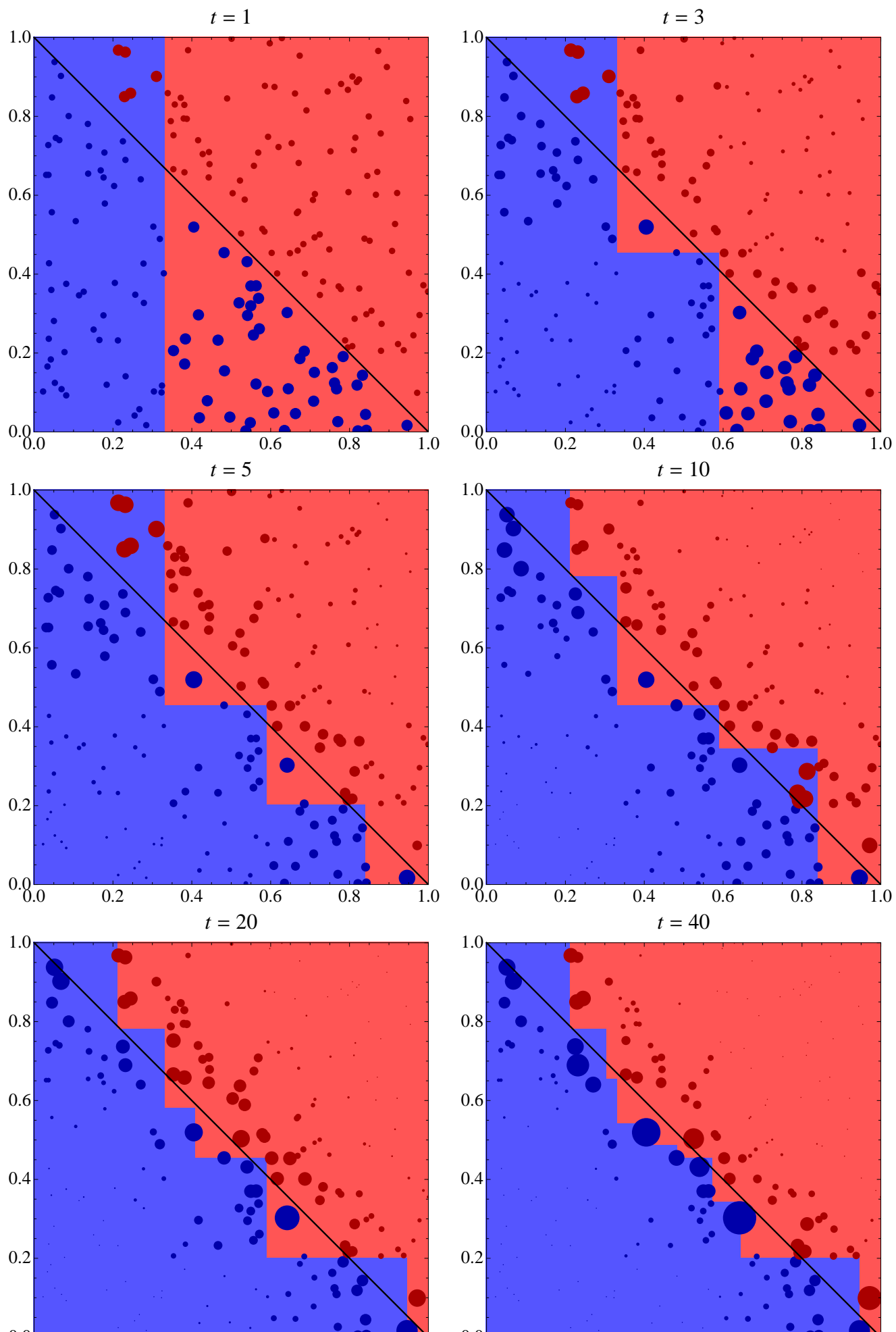


Figure 2.5: Observations $\mathbf{x}_i = (x_i^{(1)}, x_i^{(2)})$ are generated uniformly in a unit square; class labels are defined by $y_i = 1$ iff $x_i^{(1)} + x_i^{(2)} \leq 1$. Positive points (and region) are blue, and negative points are red.



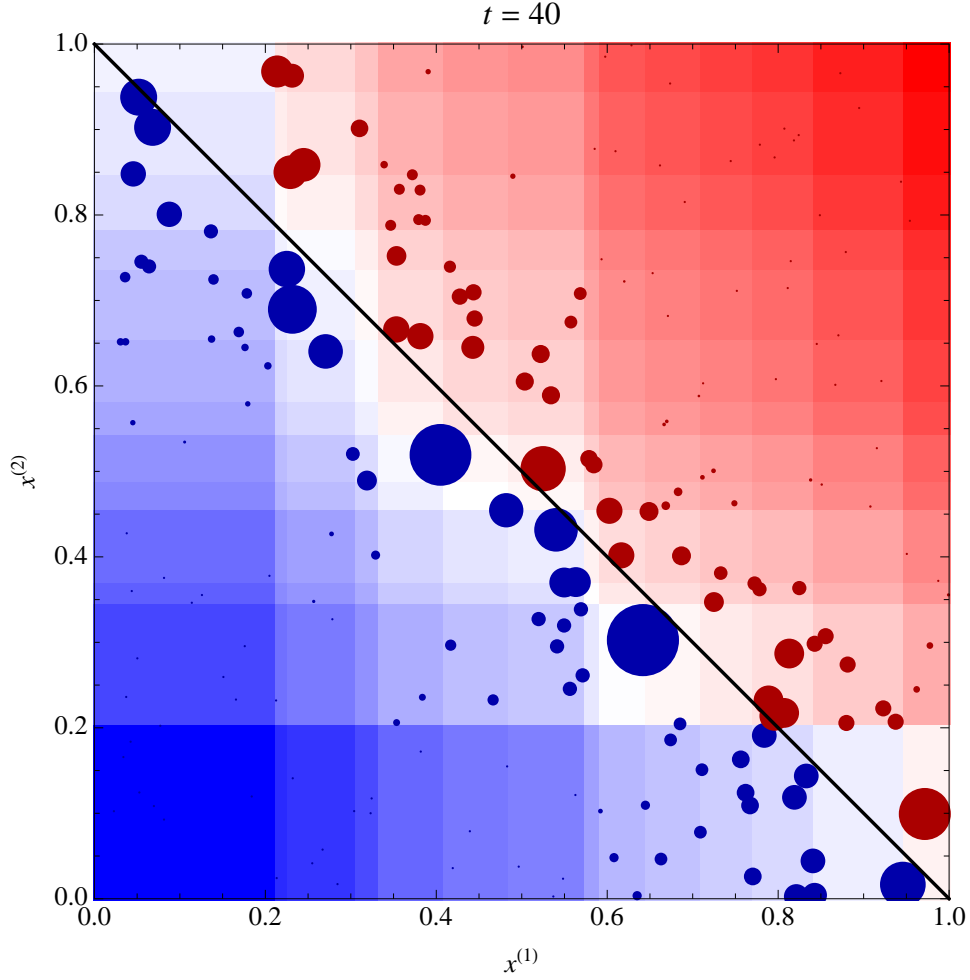


Figure 2.7: The weights and the margins after $t = 40$ iterations. The color intensity is proportional to the margin $\rho_{\mathbf{x},y}(f) = f^{(t)}(\mathbf{x})y$.

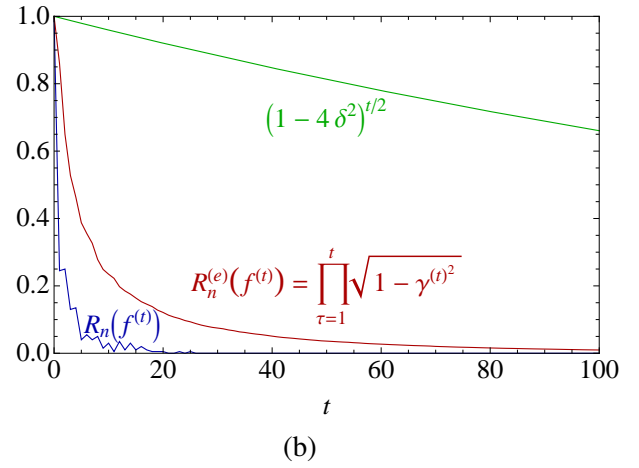
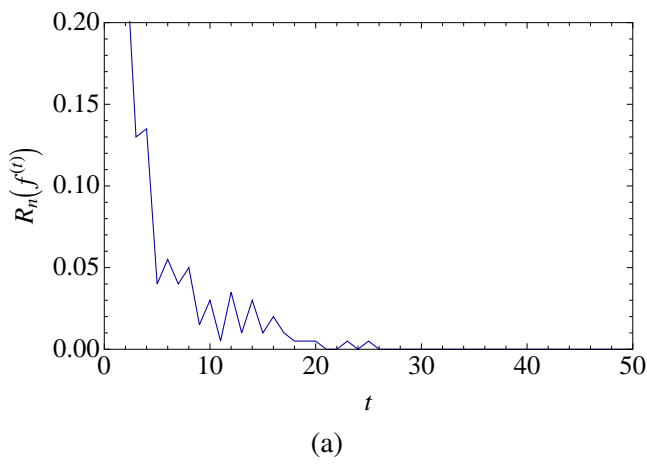


Figure 2.8: (a) The training error $R_n(f^{(t)})$ and (b) its bounds. The green curve is the a-priori bound (2.17) and the red curve is the exponential risk $R_n^{(e)}(f^{(t)})$ (2.23).

Bibliography

- [FS97] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55:119–139, 1997.