

---

# On Boosting and the Exponential Loss

---

**Abraham J. Wyner**

Statistics Dept.

Wharton School

University of Pennsylvania

Philadelphia, PA 19104

ajw@wharton.upenn.edu

## Abstract

Boosting algorithms in general and AdaBoost in particular, initially baffled the statistical world by posing two questions: (1) Why is it that AdaBoost performs so well? and (2) What makes Boosting methods resistant to overfitting? In response to question (1) Hastie, Tibshirani and Friedman (2000) take a statistical view of Boosting by recasting it as a stagewise approach to the minimization of an exponential loss function by means of an additive model in a process similar to additive logistic regression. This characterization has since been well integrated in the statistics and computer science communities as the best statistical answer to question (1). In this paper, we argue that this well assimilated view is questionable and that perhaps Boosting's success has nothing to do with the minimization of an exponential criterion or indeed any optimization at all. Our argument rests on a constructive theorem that states that for any sequence of classifiers there exists a linear combination for which the exponential criterion equals one. Furthermore, we present a Boosting algorithm which performs empirically like AdaBoost while stabilizing the exponential loss to a constant.

## 1 ON THE EXPONENTIAL LOSS

The performance of Boosting surprised statisticians. AdaBoost (Freund and Schapire, 1996) works well—really well, on a huge variety of problems in various different contexts. Yet it achieves this success in defiance of several important statistical paradigms. Its program of model selection is haphazard, random and unregulated. It shrinks the training set error rate to zero, yet, paradoxically, it continues to improve its out of

sample performance even after this mark is obtained. Finally, It seems remarkably resistant to overfitting, in many, if not all contexts, even though it never validates on hold out data, never estimates an out of sample error rate, or otherwise considers performance on anything other than the training data.

There are primarily three explanations for the success of Boosting. The first explanation is set in the context of PAC learning. While mathematically attractive the theory is practically unappealing since it offers no hint of why AdaBoost should generalize well on test data. Furthermore, from a practical point of view the weak learner hypothesis is unnatural at best and false at worst. A second explanation see Schapire et. al., focuses on AdaBoost's ability to Boost the "margin". This work offers concrete bounds on the generalization error which decrease as the margin increases. Unfortunately, in practice, these bounds are either too loose to be of use, or too restrictive to be applicable. Furthermore, Breiman (2000) has found simple counterexamples for which increased margin leads to higher generalization error. At best, the role of the margins in classification is still unknown and a worthy topic for further research.

Hastie, Tibshirani, and Friedman (Hastie et. al 2000) take great strides to clear up the mystery of boosting to provide statisticians with a statistical view of the subject. The heart of their article is the recasting of boosting as a statistically familiar program for finding an additive model by means of a forward stagewise approximate optimization of an exponential criterion.

This explanation has been widely assimilated and has reappeared in the statistical literature (see for example Friedman et al. 2001) as well as in a plethora of computer science articles (for example Schapire's overview 1999). Further, there has been very interesting work connecting the optimization realized by AdaBoost to other optimizations (see Collins et. al (2000) for ex-

ample).

## 1.1 STATISTICAL VIEWS OF BOOSTING

Let's review this argument. The statistical view in Friedman et al. 2000 begins by asking a statistical question: Given a space of classifiers  $G$ , consider linear combinations of the form

$$f(x) = \sum_{m=1}^M \alpha_m g_m(x) \text{ with } g_m \in G \text{ and } \alpha_m \in R.$$

How does one select  $g_m, \alpha_m$  and  $M$  to optimize an exponential criterion  $\exp(-yf(x))$ ? Here and throughout this discussion, we will consider  $x$  a vector of predictors,  $y \in \{-1, 1\}$  a binary class label, and  $T$  a training set of  $N$  pairs  $\{x_i, y_i\}$  for  $i = 1, \dots, N$ . Since global minimization is usually not possible, a forward stage-wise optimization is proposed. Given

$$f_{m-1} = \sum_{i=1}^{m-1} \alpha_i g_i(x)$$

the goal is transformed to finding  $g_m$  and  $\alpha_m$  that minimize

$$\sum_{i=1}^N \exp(-y_i(f_{m-1}(x_i) + \alpha_m g_m(x_i))).$$

In light of the defining property of the exponential, namely factorization, the previous problem is reduced to a *weighted* minimization. That is

$$\begin{aligned} \{\alpha_m, g_m\} = \arg \min_{\alpha, g \in G} \exp(-\alpha) \sum_{i=1}^N w_{m-1}(i) \\ + (e^\alpha - e^{-\alpha}) \sum_{i=1}^N w_{m-1}(i) 1\{y_i \neq g(x_i)\} \end{aligned}$$

with

$$w_{m-1}(i) = \exp(-y_i f_{m-1}(x_i)) / \sum_{i=1}^N \exp(-y_i f_{m-1}(x_i))$$

defined to be weights. The AdaBoost algorithm is recovered by noting that this minimum is found by letting

$$g_m = \arg \min_{g \in G} \sum_{i=1}^N w_{m-1}(i) 1\{y_i \neq g(x_i)\}$$

while assuming that every  $g \in G$  maps to  $\{-1, 1\}$ . Taking a single derivative leads to the solution

$$g_m = \arg \min_{g \in G} \sum_{i=1}^N w_{m-1}(i) 1\{y_i \neq g(x_i)\}$$

and

$$\alpha_m = \frac{1}{2} \log \frac{1 - \epsilon_m}{\epsilon_m},$$

where

$$\epsilon_m = \sum_{i=1}^N w_{m-1}(i) 1\{y_i \neq g_m(i)\}$$

is the weighted error rate on the training set by classifier  $g_m$ . At the next stage, the weights are recomputed by an updating procedure where each correctly classified training example gets down weighted by the factor  $\exp(-\alpha_m)$  and each incorrect example gets up weighted by the factor  $\exp(-\alpha_m)$ . This is AdaBoost.

Freund and Schapire (1996) arrive at the same solution and calculation but in a different context. Since one can construct a binary classifier by taking the sign of a linear combination of real values functions, the training set error rate  $\epsilon(M)$  of an  $M$  term linear combination  $f_M$ , can be re-written in the following form:

$$\epsilon(M) = \sum_{i=1}^N 1\{y_i f(x_i) < 0\}.$$

Freund and Schapire observe that

$$\sum_{i=1}^M 1\{y_i f_M(x_i) < 0\} \leq \sum_{i=1}^M \exp(y_i f(x_i)).$$

Thus, a stagewise minimization of the left hand side (*training set error rate*) is bounded by a stagewise minimization of the right hand side (the exponential) and the resulting algorithm is AdaBoost.

The exponential function's identification with Adaboost is uncontested. Nevertheless, minimization of the exponential loss is sensible only in light of the *population* version of the algorithm for which it is easily proven (see Friedman et al. 2001) that a one-step optimization of the exponential criterion results in estimates equal to (half) the log odds ratio. Thus, it would appear that taking the sign of such an estimator (as Adaboost does) would result in the optimal Bayes classifier.

In summary, we have recounted the connection between stagewise optimization of the exponential loss and Adaboost. Unanswered is the question of whether optimization of the exponential on *training data* will result in *population* optimization as well.

## 2 STABILIZING THE LOSS

From the perspective of the exponential loss AdaBoost should be a disastrous overfit, since it never validates

or otherwise considers hold out samples. Furthermore, it computes the loss exclusively on the data while ignoring any consideration of its expectation over the population. In more traditional statistical contexts this ignorance may not be problematic. For example, the construction of a linear regression fits a function to the data (the least squares line) which is proven to be good only in expectation (the conditional expectation is the population minimizer). In small parameter regression contexts, this is rarely a concern since the function space is small and a fit to the data is likely to be a good fit out of sample. Not so in our context, where not only are the base classifiers often strong learners, but the resulting linear combination of classifiers are almost always sufficiently rich enough to separate the training data, even in high dimensions. To wit, Yali Amit has observed (in private communication) that the loss can increase exponentially on a hold out sample while decaying exponentially on the training data. Breiman (2000), in a discussion of the statistical view, has expressed his misgivings with this explanation as well.

Casting Boosting as an optimization has led to interesting new algorithms (see for example, Friedman (2001) and Hastie et. al. (2000)) that are quite similar to AdaBoost since they all involve the averaging of classifiers retrained on perturbed data. In contrast to AdaBoost, these algorithms implement various efforts to protect against overtraining. Broadly, these efforts are of three categories. One method is to limit the complexity of the base learner by using decision stumps instead of trees, another is to shrink the coefficients (see Friedman 2001) and a third is to regularize the algorithm by controlled stopping (Jiang 2001). All three approaches are natural solutions in the context of optimization. AdaBoost works fine without these efforts.

In this paper I will argue that Breiman's intuition was correct: that minimization of an exponential criterion has less to do with the success of Adaboost than it appears. My case rests on one simple theorem and an empirical observation. We start with the theorem.

Consider, as before, a set of binary classifiers  $G$ , such that every  $g \in G$  is a mapping of vectors  $x$  in some  $p$ -dimensional feature space to the binary variable  $y \in \{-1, 1\}$ .

**Theorem A:** Let  $g_m \in G$  from  $m = 1$  to  $M$  be any arbitrary sequence of classifiers. Let  $T = \{x_i, y_i\}_{i=1}^N$  be a training set of  $N$  observations. Starting with  $f_0 = 0$  recursively define  $f_m = \sum_{i=1}^m \alpha_i g_i$  and  $\alpha_m =$

$\beta \log \frac{1-\epsilon_m}{\epsilon_m}$  where

$$\epsilon_m = \sum_{i=1}^N \frac{1}{N} \exp(-y_i f_{m-1}(x_i)) 1\{y_i \neq g_m(x_i)\}.$$

It follows that for all  $M$

$$\sum_{i=1}^N \frac{1}{N} \exp\left(\frac{1}{\beta} y_i f_M(x_i)\right) = 1.$$

*Note:* The value of  $\epsilon_m$  denotes the weighted error rate of  $g_m$  on training set  $T$  using the accumulated exponential loss at each observation  $x_i$  by the linear combination  $f_{m-1}$ .

Theorem A implies that *any* sequence of classifiers can be combined linearly while holding the exponential loss on the data constant. We exploit this fact to produce a *good* linear combination of classifiers that is indifferent to the exponential criterion. To this end we construct the following algorithm:

**Beta-Boosting:** Let  $w_0 = \frac{1}{N}$ . Then for  $m = 1$  to  $M$  do:

1.  $g_m = \operatorname{argmin}_{g \in G} \sum_{i=1}^N w_{m-1} 1\{y_i \neq g(x_i)\}$
2. Let  $\epsilon_m$  be the weighted error rate achieved by  $g_m$
3.  $\alpha_m = \frac{1}{\beta} \log \frac{1-\epsilon_m}{\epsilon_m}$ .
4.  $w_m(i) = w_{m-1}(i) \exp(-y_i g_m(x_i) \alpha_m)$ .
5. Renormalize  $w_m(i)$ .

Output final classifier

$$F(x) = \operatorname{sign}\left[\sum_{m=1}^M \beta \alpha_m g_m(x)\right].$$

Observe that Beta-Boosting is almost Adaboost. The main difference is that we have generalized the factor of  $1/2$  from the definition of  $\alpha_m$  which results in an updating scheme whose weights are the *beta* powers of Adaboost's weights. Theorem A applies in the special case of  $\beta = 1$  (which we call "square-boosting" (since its weights are the squares of Adaboost's), for which it follows that at every stage the weights always sum to 1 (so normalization is not required) and that the exponential criterion remains constant at 1 regardless of the number of iterations. Beta boosting can be viewed in the context of optimization theory as well. The optimal program with respect to the exponential is  $\beta = 1/2$ . Square boosting is a case of Successive Over Relaxation (SOR). In this particular case, however, the doubling of the step size leads not only to over-relaxation but excessive relaxation. Square-Boosting is thus not an optimization.

Like Adaboost, the SquareBoost algorithm is a stagewise scheme based on modular reweighting for constructing a linear combination of classifiers fitted on perturbed training sets. Adaboost will drive the exponential criterion to zero exponentially quickly, while Squareboosting keeps the criterion constant at one. More importantly, the criterion is identically one for *all* binary weak learners, not necessarily the one that achieves the lowest weighted training error. Accordingly, we can expect that SquareBoosting should fail to Boost.

**Proof of Theorem A:** Let  $f_0 = 0$  we proceed by induction. For  $m = 0$  the result is trivial. Now, given  $f_{m-1} = \sum_{i=1}^n \alpha_i g_i$ . Let  $w_{m-1}(i) = \exp(-y_i f_{m-1}(x_i)) / \sum_{i=1}^N \exp(-y_i f_{m-1}(x_i))$ . Using the fact that  $y_i g_m(x_i)$  is  $-1$  if  $g_m(x_i) \neq y_i$  and  $1$  otherwise, we have

$$\begin{aligned} & \sum_{i=1}^N \exp(-y_i(f_{m-1} + \alpha_m g_m(x_i))) = \\ & \sum_{i=1}^n w_{m-1}(i) \exp(-y_i \alpha_m g_m(x_i)) \\ & = \sum_{i=1}^N w_{m-1}(i) \exp(\alpha_m) 1\{y_i \neq g_m(x_i)\} + \\ & \sum_{i=1}^N w_{m-1}(i) \exp(-\alpha_m) [1 - 1\{y_i = g_m(x_i)\}] \\ & = \frac{1 - \epsilon_m}{\epsilon_m} \sum_{i=1}^N w_{m-1}(i) 1\{y_i = g_m(x_i)\} + \\ & \frac{\epsilon_m}{1 - \epsilon_m} \sum_{i=1}^N w_{m-1}(i) [1 - 1\{y_i = g_m(x_i)\}] \\ & = \frac{1 - \epsilon_m}{\epsilon_m} \epsilon_m + \frac{\epsilon_m}{1 - \epsilon_m} (1 - \epsilon_m) \\ & = 1 - \epsilon_m + \epsilon_m \\ & = 1 \end{aligned}$$

Thus the exponential criterion equals one for all  $M$  and for *arbitrary* choice of classifiers  $g_m$ , which prove the theorem. Note: If the initial weights are normalized then the sum of the weights will sum to one and renormalization will not be necessary.

### 3 PERFORMANCE

SquareBoost, by construction, will find a linear combination of classifiers by a stagewise process that is indifferent to the exponential criterion. This sets it apart from LogitBoost, GentleBoost (see FHT, 2000)

and Mart (Friedman, 2001). Yet the surprising fact is that SquareBoosting is apparently as successful a classifier as Adaboost. We present SquareBoost tested on examples from the UCI Irvine repository using C4.5 as the base classifier. In all of them, Squareboosting exhibits the properties of AdaBoost: excellent performance (relative to C4.5), training set error rates that reach zero, and test set error rates that continue to drop even after the training error has reached zero. In this draft we include three examples: CAR, SATIMAGE, and SEGMENTATION (details of the simulation appear in Table 1). Other examples (not shown) reveal the same behavior. In all three examples we examine the performance of the algorithm as measured by test set error rate. We then plot the error rate as a function of the number of iterations. Three features clearly stand out:

1. In the examples, the test set error rate of SquareBoost jumps up after the second iteration, but then descends well below the C4.5 error rate. This is likely due to the fact that SquareBoost's reweighting strategy is more aggressive than AdaBoost in the early stages.
2. The test set error rate of SquareBoost is essentially equivalent to that of AdaBoost after 225 iterations.
3. Performance of SquareBoost continues to improve for many iterations after the training set error rate is zero.
4. The performance of AdaBoost and SquareBoosting after 225 iterations is statistically equivalent. AdaBoost slightly outperforms SquareBoost in Satimage (Figure 1), SquareBoost outperforms AdaBoost in Car (Figure 3) and the two are nearly the same in Segmentation (Figure 2)

SquareBoosting is actually a special case of Beta-Boosting with  $\beta = 1$ , that is weights that are simply the squares of AdaBoost's weights. This choice is interesting since it seems to perform as well as AdaBoost while stabilizing the exponential criterion. We now consider two other special cases: "root" Boosting, that is Beta-Boosting for  $\beta = 1/4$  (weights that are the square roots of AdaBoost) and Quad-Boosting which uses weights that are the fourth powers of AdaBoost. Qualitatively, Root boosting is akin to GentleBoost since the step size in the "optimization" is halved, in many examples a much more gradual fitting to the training set is observed, unfortunately the overall performance is many example (see for example Figure 4) tends to be quite poor relative to AdaBoost. Quad boosting is too aggressive. In all examples considered (See Figure 5, other examples not shown) its

Table 1: Description of three data sets from the UC Irvine Repository

DATASET	DATASET SIZE	TRAINING SET SIZE	ATTRIBUTES CONT/DISCR	CLASSES
CAR	1728	1000	0/6	4
SATIMAGE SEGMENTATION	4435	2000	36/0	6
	2310	1155	19/0	7

performance is weakest relative to the other examples, although it may be possible to improve its performance by means of a regularization algorithm.

## 4 CONCLUSIONS

We have proved that there exists a stagewise forward program for constructing a linear combination of classifiers that is indifferent to the exponential criterion. Nevertheless, this algorithm is nearly AdaBoost in definition and in practice. This forces a reconsideration of the connection between Boosting algorithms and statistical optimization. Another view suggests that despite the difference in their respective treatment of the exponential criterion, both algorithms are actually programs for sequentially perturbing the training data, fitting classifiers and averaging the result. This view is supported by several arguments. One argument is that averaging is in fact a smoothing process that prevents overfitting (see Krieger et al 2001). A closer look at other stagewise programs for finding a linear combination, for example GentleBoost and LogitBoost, reveal the same underlying structure of perturbation and averaging. Given that all the algorithms proposed in this paper (RootBoost, QuadBoost and SquareBoost) are empirically similar to AdaBoost while radically different with respect to optimization, suggests that loss function minimization can only be part of the story. What we have accomplished in this paper is a mathematical and empirical justification for the lingering doubt among some in the Boosting community that optimization is an explanation for AdaBoost.

### Acknowledgements

The author would like to acknowledge Andreas Buja, Yoram Singer and Abba Krieger for many helpful conversations and discussions.

### References

Breiman, L. (2000), Discussion on Statistical View of Boosting, *The Annals of Statistics*, Vol. 28, No. 2.

Collins M., Robert E. Schapire, Yoram Singer, Logistic Regression, AdaBoost and Bregman Distances (2000) *Computational Learning Theory*.

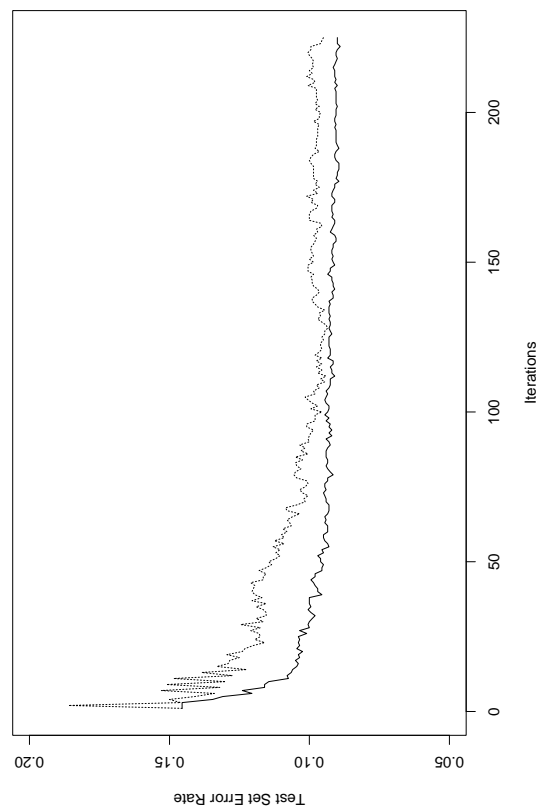


Figure 1: Comparison of SquareBoost (dashed line) and AdaBoost (solid line) applied to Satimage dataset

Freund, Y. and Schapire, R.E. . (1996) Experiments with a new Boosting Algorithm. *Machine Learning: Proceedings of the Thirteenth International Conference* 148-156, Morgan Kaufman, San Francisco.

Friedman, J.H. (2001). Greedy function Approximation: A gradient Boosting Machine. *Annals of Statistics*, Vol. 29, No. 5.

Hastie T., Tibshirani R., and Friedman J. (2000). A Statistical View of Boosting, *The Annals of Statistics* 2000, Vol. 28, No. 2.

Hastie T., Tibshirani R., and J. Friedman. (2001) *The Elements of Statistical Learning: data mining, inference and prediction*. Springer, 2001.

Jiang, W., (2001) Is Regularization Unnecessary for boosting. *Proceedings of the 8th annual conference of A.I. and Statistics*, Morgan Kaufmann, January, 2001.

Krieger A., Long C. and A. Wyner (2001) Boosting Noisy Data, *Proceedings of International Conference on Machine Learning*, Morgan Kaufman, July, 2001 .

Schapire R.E. (1999), *Theoretical views of boosting*

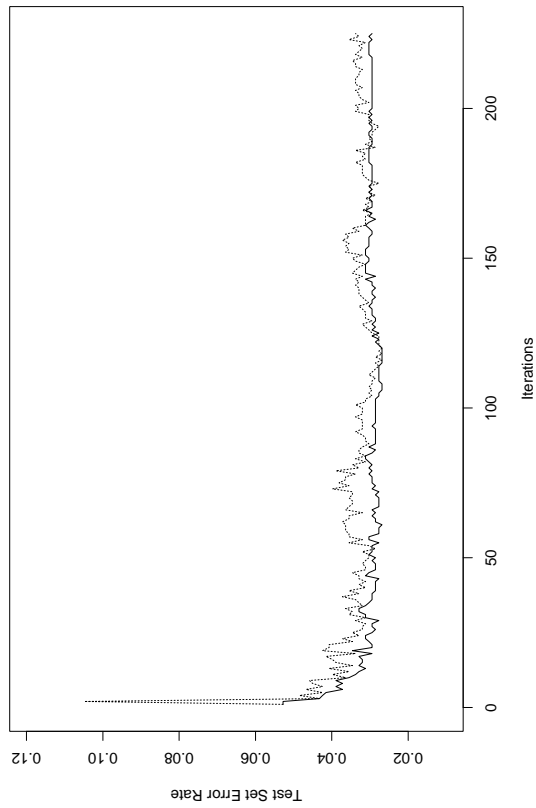


Figure 2: Comparison of SquareBoost (dashed line) and Adaboost(solid line) applied to Segmentation dataset

and applications. In Tenth International Conference on Algorithmic Learning Theory, 1999.

Schapire R.E. , Freund, Y., Bartlett, P. and Lee, W.S. (1998) Boosting the margin: a new explanation for the effectiveness of voting methods. Annals of Statistics 26(5): 1651-1686.

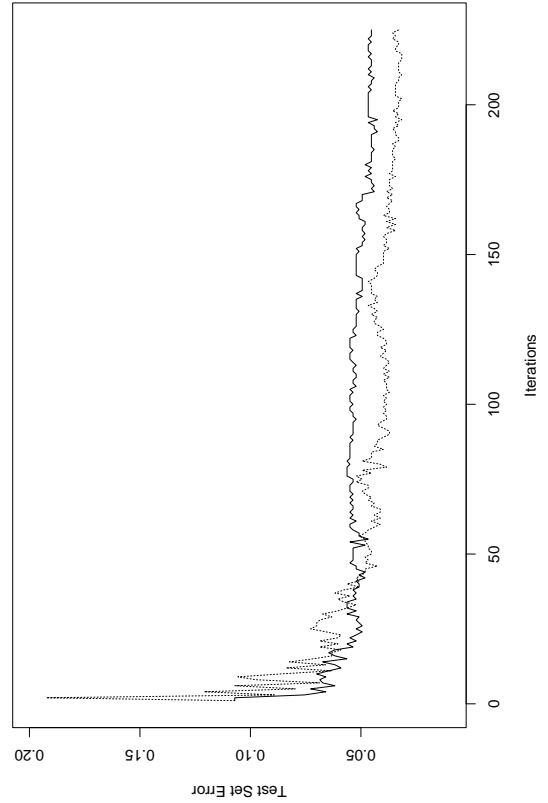


Figure 3: Comparison of SquareBoost (dashed line) and Adaboost (solid line) applied to Car dataset

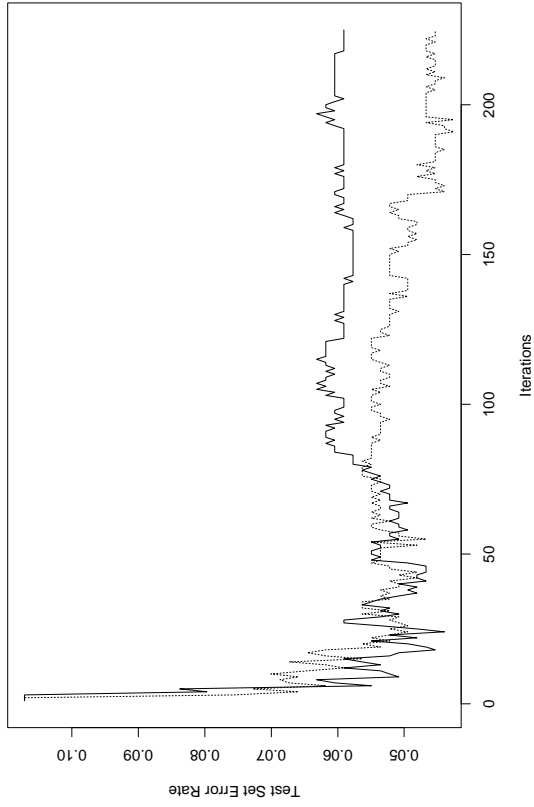


Figure 4: Comparison of RootBoost (Solid line) and Adaboost (Dashed line) applied to Car dataset

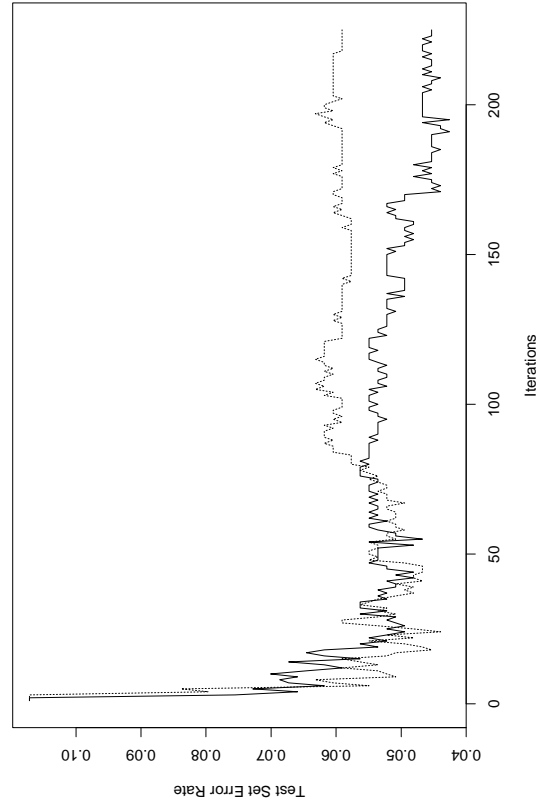


Figure 5: Comparison of QuadBoost (Dashed line) and Adaboost (solid line) applied to Car dataset