

Smooth Boosting Using an Information-Based Criterion

Hatano, Kohei
Department of Informatics, Kyushu University

<https://hdl.handle.net/2324/3057>

出版情報 : DOI Technical Report. 225, 2006-05-19. Department of Informatics, Kyushu University
バージョン :
権利関係 :



Smooth Boosting Using an Information-Based Criterion

Kohei Hatano

Department of Informatics, Kyushu University

hatano@i.kyushu-u.ac.jp

May 19, 2006

Abstract

Smooth boosting algorithms are variants of boosting methods which handle only smooth distributions on the data. They are proved to be noise-tolerant and can be used in the “boosting by filtering” scheme, which is suitable for huge data. However, current smooth boosting algorithms have rooms for improvements: Among non-smooth boosting algorithms, real AdaBoost or InfoBoost, can perform more efficiently than typical boosting algorithms by using an information-based criterion for choosing hypotheses. In this paper, we propose a new smooth boosting algorithm with another information-based criterion based on Gini index. We show that it inherits the advantages of two approaches, smooth boosting and information-based approaches.

1 Introduction

In recent years, huge data have become available due to the development of computers and the Internet. As size of such huge data can reach hundreds of gigabytes in knowledge discovery and machine learning tasks, it is important to make knowledge discovery or machine learning algorithms scalable. Sampling is one of effective techniques to deal with large data. There are many results on sampling techniques [23, 21, 27, 6] and applications to data mining tasks such as decision tree learning [8], support vector machine [2], and boosting [6, 7].

Especially, boosting is simple and efficient learning method among machine learning algorithms. The basic idea of boosting is to learn many slightly accurate hypotheses (or *weak hypotheses*) with respect to different distributions over the data, and to combining them into a highly accurate one. Originally, boosting was invented under the *boosting by filtering* framework (or the filtering framework), where the booster can sample examples randomly from the whole instance space [25, 11]. On the other hand, in the *boosting by subsampling* framework (or, the subsampling framework), the booster is given a bunch of examples in advance. There are two advantages of the filtering framework over the subsampling framework. First, the booster does not have to keep less examples as it “filters” examples and accepts only necessary ones. The second advantage is that the booster can automatically determine the sufficient sample size. Note that it is not trivial to determine the sufficient sample size a priori in the subsampling framework. So the boosting by filtering framework seems to fit learning over huge data. However, early boosting algorithms [25, 11] which work in the filtering framework were not practical,

because they were not “adaptive”, i.e., they need the prior knowledge on the accuracy of weak hypotheses.

Madaboost, a modification of AdaBoost [12], is the first adaptive boosting algorithm which works in the filtering framework [7]. Combining with adaptive sampling methods [6], Madaboost is shown to be more efficient than AdaBoost over huge data, while keeping the prediction accuracy. By its nature of updating scheme, MadaBoost is categorized as one of “smooth” boosting algorithms [13, 29, 15], where the name, smooth boosting, comes from the fact that these boosting algorithms only deal with smooth distributions over data (In contrast, for example, AdaBoost might construct exponentially skew distributions over data). Smoothness of distributions enables boosting algorithms to sample data efficiently. Also, smooth boosting algorithms have theoretical guarantees for noise tolerance in the various noisy learning settings, such as statistical query model [7, 18], malicious noise model [29, 32] and agnostic boosting [15, 20].

However, there seems still room for improvements on smooth boosting. A non-smooth boosting algorithm, InfoBoost [1] (which is a special form of real AdaBoost [26]), performs more efficiently than other boosting algorithms in the boosting by subsampling framework. More precisely, given hypotheses with error $1/2 - \gamma$, typical boosting algorithms take $O(1/\gamma^2)$ iterations to learn sufficiently accurate hypothesis. On the other hand, InfoBoost learns in from $O(1/\gamma)$ to $O(1/\gamma^2)$ iterations by taking advantage of the situation when weak hypotheses have low false positive error [16, 17]. So InfoBoost can be more efficient at most by $O(1/\gamma)$ times.

The main difference between InfoBoost and other boosting algorithms such as AdaBoost or MadaBoost is the criterion for choosing weak hypotheses. Typical boosting algorithms are designed to choose hypotheses whose errors are minimum with respect to given distributions. In contrast, InfoBoost uses an information-based criterion to choose weak hypotheses. The criterion was previously proposed by Kearns and Mansour in the context of decision tree learning [19], and also applied to boosting algorithms using branching programs [22, 30]. But, so far, no smooth algorithm has such the nice property of InfoBoost.

In this paper, we propose a new smooth boosting algorithm using another information-based criterion which is based on *Gini index* [3]. Our boosting algorithm also works in the filtering framework. Preliminary experiments show that our algorithm, which we call GiniBoost, improves MadaBoost in the filtering framework over large data.

2 Preliminaries

2.1 Learning Model

We adapt the PAC learning model [31]. Let \mathcal{X} be an *instance space* and let $\mathcal{Y} = \{-1, +1\}$ be a set of labels. We assume an unknown *target function* $f : \mathcal{X} \rightarrow \mathcal{Y}$. Further we assume that f is contained in a known class \mathcal{F} of functions from \mathcal{X} to \mathcal{Y} . Let D be an unknown distribution over \mathcal{X} . The learner has an access to the *example oracle* $\text{EX}(f, D)$. When given a call from the learner, $\text{EX}(f, D)$ returns an *example* $(\mathbf{x}, f(\mathbf{x}))$ where each instance \mathbf{x} is drawn randomly according to D . Let \mathcal{H} be a hypothesis space, or a set of functions from \mathcal{X} to \mathcal{Y} . We assume that $\mathcal{H} \supset \mathcal{F}$. For any distribution D over \mathcal{X} , *error* of hypothesis $h \in \mathcal{H}$ is defined as $\text{err}_D(h) \stackrel{\text{def}}{=} \Pr_D\{h(\mathbf{x}) \neq f(\mathbf{x})\}$. Let S be a *sample*, a set of examples $((\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_m, f(\mathbf{x}_m)))$. For any sample S , *training error* of hypothesis $h \in \mathcal{H}$ is defined as $\widehat{\text{err}}_S(h) \stackrel{\text{def}}{=} |\{(\mathbf{x}_i, f(\mathbf{x}_i)) \in S \mid h(\mathbf{x}_i) \neq f(\mathbf{x}_i)\}|/|S|$.

We say that learning algorithm A is a *strong learner* for \mathcal{F} if and only if, for any $f \in \mathcal{F}$ and any distribution D , given ε, δ ($0 < \varepsilon, \delta < 1$), a hypothesis space \mathcal{H} , and access to the example oracle $\text{EX}(f, D)$ as inputs, A outputs a hypothesis $h \in \mathcal{H}$ such that $\text{err}_D(h) = \Pr_D\{h(x) \neq f(x)\} \leq \varepsilon$ with probability at least $1 - \delta$. We also consider a weaker learner. Specifically, we say that learning algorithm A is a *weak learner*¹ for \mathcal{F} if and only if, for any $f \in \mathcal{F}$, given a hypothesis space \mathcal{H} , and access to the example oracle $\text{EX}(f, D)$ as inputs, A outputs a hypothesis $h \in \mathcal{H}$ such that $\text{err}_D(h) \leq \frac{1}{2} - \gamma$ for a fixed γ ($0 < \gamma < \frac{1}{2}$).

2.2 Boosting Approach

Schapire proved that the strong and weak learnability are equivalent to each other for the first time [25]. Especially the technique to construct a strong learner by using a weak learner is called “boosting”. Basic idea of boosting is the following: First, the booster trains a weak learner with respect to different distributions D_1, \dots, D_T over the domain \mathcal{X} , and gets different “weak” hypotheses h_1, \dots, h_T such that $\text{err}_{D_t}(h_t) \leq 1/2 - \gamma_t$ for each $t = 1, \dots, T$. Then the booster combines weak hypotheses h_1, \dots, h_T into a final hypotheses h_{final} satisfying $\text{err}_D(h_{final}) \leq \varepsilon$.

In the subsampling framework, the booster calls $\text{EX}(f, D)$ for a number of times and obtains a sample $S = ((\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_m, f(\mathbf{x}_m)))$ in advance. Then the booster constructs the final hypothesis whose training $\widehat{\text{err}}_S(h_{final}) \leq \varepsilon$ by training the weak learner over the given sample S . Then the $\text{err}_D(h_{final})$ can be estimated by using arguments on VC-dimension or margin (E.g., see [12] or [24], respectively). For example, for typical boosting algorithms, $\text{err}_D(h_{final}) \leq \widehat{\text{err}}_S(h_{final}) + \tilde{O}(\sqrt{T \log |\mathcal{W}|/m})$ ² with high probability, where T is the size of the final hypotheses, i.e., the number of weak hypotheses combined in h_{final} .

In the filtering framework, on the other hand, the booster deal with the whole instance space \mathcal{X} through $\text{EX}(f, D)$. By using statistics obtained from calls of $\text{EX}(f, D)$, the booster tries to minimize $\text{err}_D(h_{final})$ directly.

Smooth boosting algorithms only deal such distributions D_1, \dots, D_t that are “smooth” with respect to the original distribution D . We define the following measure of smoothness.

Definition 1. Let D and D' be any distributions over \mathcal{X} . We say that D' is λ -smooth with respect to D if $\sup_{\mathbf{x} \in \mathcal{X}} D'(\mathbf{x})/D(\mathbf{x}) \leq \lambda$.

The smoothness parameter λ has crucial roles in robustness of boosting algorithms [7, 29, 15]. Also, it affects the efficiency of sampling methods. For example, by rejection sampling, we use $1/\lambda$ calls of $\text{EX}(f, D)$ on average to simulate a call of $\text{EX}(f, D')$ for a distribution D' that is λ -smooth w. r. t. D .

2.3 Our Assumption and Technical Goal

In the rest of the paper, we assume that the learner is given a finite set \mathcal{W} of hypotheses such that for any distribution D' over \mathcal{X} , there exists a hypothesis $h \in \mathcal{W}$ satisfying $\text{err}_{D'}(h) \leq 1/2 - \gamma/2$.

¹In the original definition of [25], the weak learning algorithm is allowed to output a hypothesis h with $\text{err}_D(h) > 1/2 - \gamma$ with probability at most δ as well. But in our definition we omit δ to make our discussion simple. Of course, we can use the original definition, while our analysis becomes slightly more complicated.

²In the $\tilde{O}(g(n))$ notation, we neglect $\text{poly}(\log(n))$ terms.

Now our technical goal is to construct an efficient smooth boosting algorithm which works in both the subsampling and the filtering framework.

3 Boosting by Subsampling

In this section, we propose our boosting algorithm in the subsampling framework.

3.1 Derivation

First of all, we derive our algorithm. It is well known that many of boosting algorithms can be explained as greedy minimizers of loss functions [14]. More precisely, it can be viewed that they minimize particular loss functions that bound the training errors. The derivation of our algorithm is also explained simply in terms of its loss function.

Suppose that the learner is given a sample $S = \{(\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_m, f(\mathbf{x}_m))\}$, a set \mathcal{W} of hypotheses, and the current final hypothesis $H_t(\mathbf{x}) = \sum_{j=1}^t \alpha_j h_j(\mathbf{x})$, where each $h_j \in \mathcal{W}$ and $\alpha_j \in \mathbb{R}$ for $j = 1, \dots, t$. The training error of $H_t(\mathbf{x})$ over S is defined by $\widehat{\text{err}}(\text{sign}(H_t)) = \frac{1}{m} \sum_{i=1}^m I(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))$, where $I(a) = 1$ if $a > 1$ and $I(a) = 0$, otherwise. We assume a function $L : \mathbb{R} \rightarrow [0, +\infty)$ such that $I(a) \leq L(a)$ for any $a \in \mathbb{R}$. Then, by definition, $\widehat{\text{err}}(\text{sign}(H_t)) \leq \frac{1}{m} \sum_{i=1}^m L(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))$. If the function L is convex, the upperbound of the training error have a global minimum. Given a new hypothesis $h \in \mathcal{W}$, a typical boosting algorithm assigns α to h that minimizes a particular loss function. For example, AdaBoost solves the following minimization problem:

$$\min_{\alpha \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m L_{\text{exp}}(-f(\mathbf{x}_i)\{H_t(\mathbf{x}_i) + \alpha h(\mathbf{x}_i)\}),$$

where its loss function is given by *exponential loss*, $L_{\text{exp}}(x) = e^x$. The solution is given analytically as $\alpha = \frac{1}{2} \ln \frac{1+\gamma}{1-\gamma}$, where $\gamma = \sum_{i=1}^m f(\mathbf{x}_i)h(\mathbf{x}_i)D_t(\mathbf{x}_i)$, and $D_t(\mathbf{x}_i) = \frac{\exp(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}{\sum_{i=1}^m \exp(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}$. InfoBoost is designed to minimize the same loss function L_{exp} as AdaBoost, but it uses a slightly different form of the final hypothesis $H_t(\mathbf{x}) = \sum_{j=1}^r \alpha [h_j(\mathbf{x})]h_j(\mathbf{x})$. The main difference is that InfoBoost assigns coefficients for each prediction $+1$ and -1 given a hypothesis. Then, the minimization problem of InfoBoost is given as:

$$\min_{\alpha_{[+1]}, \alpha_{[-1]} \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m L_{\text{exp}}(-f(\mathbf{x}_i)\{H_t(\mathbf{x}_i) + \alpha_{[+1]}[h(\mathbf{x}_i)]h(\mathbf{x}_i) + \alpha_{[-1]}[-h(\mathbf{x}_i)]h(\mathbf{x}_i)\}).$$

This problem also has the analytical solution: $\alpha_{[\pm 1]} = \frac{1}{2} \ln \frac{1+\gamma_{[\pm 1]}}{1-\gamma_{[\pm 1]}}$, $\gamma_{[\pm 1]} = \frac{\sum_{i:h(\mathbf{x}_i)=\pm 1} f(\mathbf{x}_i)h(\mathbf{x}_i)D_t(\mathbf{x}_i)}{\sum_{i:h(\mathbf{x}_i)=\pm 1} D_t(\mathbf{x}_i)}$, and $D_t(\mathbf{x}_i) = \frac{\exp(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}{\sum_{i=1}^m \exp(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}$. Curiously, this derivation makes InfoBoost choose a hypothesis that maximizes information gain, where the entropy function is defined not by Shannon's entropy function $E_{\text{Shannon}}(p) = -p \log p - (1-p) \log(1-p)$, but by the entropy function $E_{KM}(p) = 2\sqrt{p(1-p)}$ proposed by Kearns and Mansour [19] (See [30] for details). MadaBoost is formulated as the same minimization problem of AdaBoost, except that its loss function is replaced with $L_{\text{mada}}(x) = e^x$, if $x \leq 0$, $L_{\text{mada}}(x) = x$, otherwise.

Now combining the derivations of InfoBoost and MadaBoost in a straightforward way, our boosting algorithm is given by

$$\min_{\alpha_{[+1]}, \alpha_{[-1]} \in \mathbb{R}} \frac{1}{m} \sum_{i=1}^m L_{mada}(-f(\mathbf{x}_i) \{H_t(\mathbf{x}_i) + \alpha[h(\mathbf{x}_i)]h(\mathbf{x}_i)\}). \quad (1)$$

Since the solution cannot be solved analytically, we minimize an upperbound of the (1). The way of our approximation is a modification of the technique used for AdaFlat [15]. By using Lemma 3 in Appendix we have $L_{mada}(x+a) \leq L_{mada}(a) + L'_{mada}(a)x + \tilde{L}''(a)x^2$, where $\tilde{L}''(x) = \sup_x \max\{\frac{dL_{mada}(x)}{dx dx+}, \frac{dL_{mada}(x)}{dx dx-}\}$. Let

$$\ell(x) = L'_{mada}(x) = \begin{cases} 1, & x \geq 0 \\ e^x, & x < 0. \end{cases}$$

Then we get

$$\begin{aligned} & \frac{1}{m} \sum_{i=1}^m L_{mada}(-f(\mathbf{x}_i)H_t(\mathbf{x}_i)) - \frac{1}{m} \sum_{i=1}^m L_{mada}(-f(\mathbf{x}_i)H_{t+1}(\mathbf{x}_i)) \\ & \geq \frac{1}{m} \sum_{i=1}^m \{f(\mathbf{x}_i)h_t(\mathbf{x}_i)\alpha[h(\mathbf{x}_i)]\ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i)) - \alpha[h(\mathbf{x}_i)]^2\ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))\} \\ & \stackrel{\text{def}}{=} \Delta L_t(h). \end{aligned}$$

By solving the equations $\partial \Delta L_t(h) / \partial \alpha_t[b] = 0$ for $b = \pm 1$, we see that $\Delta L_t(h)$ is maximized if $\alpha_t[b] = \gamma_t[b](h)/2$, where

$$\gamma_t[b](h) = \frac{\sum_{i:h(\mathbf{x}_i)=b} h(\mathbf{x}_i) f(\mathbf{x}_i) D_t(\mathbf{x}_i)}{\sum_{i:h(\mathbf{x}_i)=b} D_t(\mathbf{x}_i)}, \text{ and } D_t(\mathbf{x}_i) = \frac{\ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}{\sum_{i=1}^m \ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}.$$

By substituting $\alpha_t[b] = \gamma_t[b](h)/2$ for $b = \pm 1$, we get

$$\Delta L_t(h) = \frac{\mu_t}{4} \{p_t(h)\gamma_t[+1](h)^2 + (1-p_t(h))\gamma_t[-1](h)^2\} \quad (2)$$

where $\mu_t = \frac{\sum_{i=1}^m \ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))}{m}$, and $p_t(h) = \Pr_{D_t}\{h(\mathbf{x}_i) = +1\}$.

Our derivation suggests that a new criterion to choose a weak hypothesis. That is, we choose $h \in \mathcal{W}$ that maximizes

$$\Delta_t(h) = p_t(h)\gamma_t[+1](h)^2 + (1-p_t(h))\gamma_t[-1](h)^2.$$

We call the quantity *pseudo gain* of hypothesis h with respect to f and D_t . Now we motivate the pseudo gain in the following way. Let $\varepsilon_t[\pm 1](h) = \Pr_{D_t}\{f(\mathbf{x}_i) = \mp 1 | h(\mathbf{x}_i) = \pm 1\}$. Note that $\gamma_t[\pm 1](h) = 1 - 2\varepsilon_t[\pm 1](h)$. Then

$$\begin{aligned} & 1 - \Delta_t(h) \\ & = p_t(h)\{1 - (1 - 2\varepsilon_t[+1](h))^2\} + (1-p_t(h))\{1 - (1 - 2\varepsilon_t[-1](h))^2\} \\ & = p_t(h) \cdot 4\varepsilon_t[+1](h)(1 - \varepsilon_t[+1](h)) + (1-p_t(h)) \cdot 4\varepsilon_t[-1](h)(1 - \varepsilon_t[-1](h)), \end{aligned}$$

which can be interpreted as the conditional entropy of f given h with respect to D_t , where the entropy is defined by *Gini index* $E_{Gini}(p) = 4p(1-p)$ [3] (See other entropy measures in Figure 3.1 for comparison). So, maximizing the pseudo gain is equivalent to maximizing the information gain defined with Gini index.

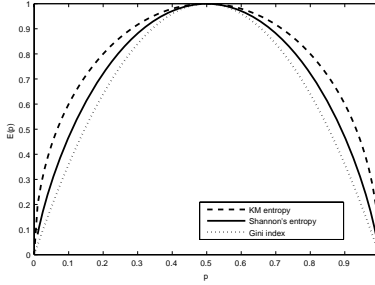


Figure 1: Plots of three entropy functions, KM entropy (upper) $E_{KM}(p) = 2\sqrt{p(1-p)}$, Shannon's entropy (middle) $E_{Shannon}(p) = -p \log p - (1-p) \log(1-p)$, and Gini index (lower) $E_{Gini}(p) = 4p(1-p)$.

3.2 Our Algorithm

Based on our derivation we propose GiniBoost. The description of our modification is given in Figure 3.2. To make our notation simple, we denote $p_t(h_t) = p_t$, $\gamma_t[\pm 1](h_t) = \gamma_t[\pm 1]$, and $\Delta_t(h_t) = \Delta_t$.

First, we show that the smoothness of distributions D_t .

Proposition 1. During the execution of GiniBoost, each distribution D_t ($t \geq 1$) is $1/\varepsilon$ -smooth with respect to D_1 , the uniform distribution over S .

Proof. Note that, during the while-loops, $\mu_t \geq \text{err}_S(h_{final}) > \varepsilon$. Therefore, for any i , $D_t(i)/D_1(i) = \ell(-f(\mathbf{x}_i)H_t(\mathbf{x}_i))/\mu_t < 1/\varepsilon$. \square

It is already shown that smoothness $1/\varepsilon$ is optimal, i.e., there is no boosting algorithm that achieves the smoothness less than $1/\varepsilon$ [29, 15].

Next, we prove the time complexity of GiniBoost.

Theorem 2. Suppose that, during the while-loops, $\text{err}_{D_t}(h_t) \leq 1/2 - \gamma_t/2 \leq 1/2 - \gamma/2$ for some $\gamma > 0$. Then, GiniBoost outputs a final hypothesis h_{final} satisfying $\widehat{\text{err}}_S(h_{final}) \leq \varepsilon$ within $T = O(1/\varepsilon\gamma^2)$ iterations.

Proof. By our derivation of GiniBoost, for any $T \geq 1$, the training error $\widehat{\text{err}}(H_T)$ is less than $1 - \sum_{t=1}^T \Delta L_t(h_t)$. By Jensen's inequality,

$$\Delta_t \geq p_t \gamma_t [+1]^2 + (1-p_t) \gamma_t [-1]^2 \geq \gamma_t^2 \geq \gamma^2.$$

As in the proof of Proposition 1, $\mu_t \geq \varepsilon$. So we have $\Delta L_t(h_t) \geq \varepsilon \gamma^2 / 4$ and thus $\widehat{\text{err}}_S(h_{final}) \leq \varepsilon$ if $T = 4/\varepsilon^2 \gamma^2$. \square

Remark. We discuss the efficiency of other boosting algorithms and GiniBoost. Smooth boosting algorithms MadaBoost [7] and SmoothBoost [29] run in $O(1/\varepsilon\gamma^2)$ iterations as well. However, the former needs a technical assumption that $\gamma_t \geq \gamma_{t+1}$ for each iteration t . Also the latter is not adaptive, i.e., it needs the prior knowledge of $\gamma > 0$. On the other hand, GiniBoost is adaptive and does not need such the technical assumption. AdaFlat [15] is another smooth boosting

GiniBoost

Given: $S = ((\mathbf{x}_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_m, f(\mathbf{x}_m)))$, and ε ($0 < \varepsilon < 1$)

begin

1. $D_1(i) \leftarrow 1/m$; ($i = 1, \dots, m$) $H_0(\mathbf{x}) \leftarrow 0$; $t \leftarrow 1$;

2. **while** $\widehat{\text{err}}_S(h_{\text{final}}) > \varepsilon$ **do**

a) $h_t \leftarrow \arg \max_{h \in \mathcal{W}} \Delta_t(h)$;

b) $\alpha_t[+1] \leftarrow \gamma_t[+1]/2$; $\alpha_t[-1] \leftarrow \gamma_t[-1]/2$;

c) $H_{t+1}(x) \leftarrow H_t(x) + \alpha_t[h_t(x)]h_t(x)$;

d) Define the next distribution D_{t+1} as

$$D_{t+1}(i) = \frac{\ell(-f(\mathbf{x}_i)H_{t+1}(\mathbf{x}_i))}{\sum_{i=1}^m \ell(-f(\mathbf{x}_i)H_{t+1}(\mathbf{x}_i))};$$

e) $t \leftarrow t + 1$;

end-while

3. Output the final hypothesis $h_{\text{final}}(\mathbf{x}) = \text{sign}(H_{t+1}(\mathbf{x}))$.

end.

Figure 2: GiniBoost

algorithm which is adaptive, but it takes $O(1/\varepsilon^2\gamma^2)$ iterations. Finally, AdaBoost [12] achieves $O(\log(1/\varepsilon)/\gamma^2)$ bound and the bound is optimal [11]. But AdaBoost might construct exponentially skew distributions. It is shown that a combination of boosting algorithms (“boosting tandems approach” [10, 15]) can achieve $O(\log(1/\varepsilon)/\gamma^2)$ with smoothness $\tilde{O}(1/\varepsilon)$. Yet, it is still open whether a single adaptive boosting algorithm can learn in $O(\log(1/\varepsilon)/\gamma^2)$ iterations while keeping the optimal smoothness $1/\varepsilon$.

4 Boosting by Filtering

In this section, we propose GiniBoost_{filt} in the filtering framework. Let

$$D_t(\mathbf{x}) = \frac{D(\mathbf{x})\ell(-f(\mathbf{x})H_t(\mathbf{x}))}{\sum_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x})\ell(-f(\mathbf{x})H_t(\mathbf{x}))}.$$

We define $\mu_t = \sum_{\mathbf{x} \in \mathcal{X}} D(\mathbf{x})\ell(-f(\mathbf{x})H_t(\mathbf{x}))$. We denote \hat{a} as the empirical estimate of the parameter a given a sample S_t . The description of GiniBoost_{filt} is given in Figure 3.

The following property of FiltEX can be immediately verified.

Proposition 3. Fix any iteration t , (i) FiltEX outputs $(x, f(x))$, where x is drawn according to D_t , and (ii) the probability that FiltEX outputs an example is at least $\mu_t \geq \text{err}_D(\text{sign}(H_t))$.

Then, we prove a multiplicative tail bound on the estimate $\hat{\Delta}_t(h)$ of the pseudo gain.

Lemma 1. Fix any $t \geq 1$. Let $\hat{\Delta}_t(h) = \hat{p}_t(h)\hat{\gamma}_t[+1](h)^2 + (1 - \hat{p}_t(h))\hat{\gamma}_t[-1](h)^2$ be the empirical estimate of $\Delta_t(h)$ given S_t . Then it holds for any ε ($0 < \varepsilon < 1$) that

$$\Pr_{D^m} \{ \hat{\Delta}_t(h) \geq (1 + \varepsilon)\Delta_t(h) \} \leq b_1 e^{-\frac{\varepsilon^2 \Delta_t m}{c_1}}, \quad (3)$$

```

GiniBoostfilt( $\varepsilon, \delta, \mathcal{W}$ )
1. Let  $H_1(\mathbf{x}) = 0$ ;  $t \leftarrow 1$ ;  $\delta_1 \leftarrow \delta/4$ ;
    $S'_1 \leftarrow \frac{18 \log(1/\delta_1)}{\varepsilon}$  random examples drawn by EX( $f, D$ );
2. while  $\widehat{\text{err}}_{S'_t}(\text{sign}(H_t)) \geq \frac{2\varepsilon}{3}$  do
   ( $h_t, S_t$ )  $\leftarrow$  HSelect( $1/2, \delta_t$ );
   ( $\hat{\gamma}_t[+1], \hat{\gamma}_t[-1]$ )  $\leftarrow$  empirical estimates over  $S_t$ ;
    $\alpha_t[+1] \leftarrow \hat{\gamma}_t[+1]/2$ ;  $\alpha_t[-1] \leftarrow \hat{\gamma}_t[-1]/2$ ;
    $H_{t+1}(\mathbf{x}) \leftarrow H_t(\mathbf{x}) + \alpha_t[h_t(\mathbf{x})]h_t(\mathbf{x})$ ;
    $t \leftarrow t + 1$ ;  $\delta_t \leftarrow \delta/(2t(t+1))$ ;
    $S'_t \leftarrow \frac{18 \log(1/\delta_t)}{\varepsilon}$  random examples drawn by EX( $f, D$ );
end-while
3. Output the final hypothesis  $h_{\text{final}}(\mathbf{x}) = \text{sign}(H_t(\mathbf{x}))$ ;

FiltEX()
do
  ( $\mathbf{x}, f(\mathbf{x})$ )  $\leftarrow$  EX( $f, D$ );
   $r \leftarrow$  uniform random number over  $[0, 1]$ ;
  if  $r < \ell(-f(\mathbf{x})H_t(\mathbf{x}))$  then return ( $\mathbf{x}, f(\mathbf{x})$ );
end-do

HSelect( $\varepsilon, \delta$ )
   $m \leftarrow 0$ ;  $S \leftarrow \emptyset$ ;  $i \leftarrow 1$ ;  $\Delta_g \leftarrow 1/2$ ;  $\delta' \leftarrow \delta/(2|\mathcal{W}|)$ ;
  do
    ( $\mathbf{x}, f(\mathbf{x})$ )  $\leftarrow$  FiltEX();
     $S \leftarrow S \cup (\mathbf{x}, f(\mathbf{x}))$ ;  $m \leftarrow m + 1$ ;
    if  $m = \left\lceil \frac{c_1 \ln \frac{b_1}{\delta'}}{\varepsilon^2 \Delta_g} \right\rceil$  then
      Let  $\hat{\Delta}_t(h)$  be the empirical estimate of  $\Delta_t(h)$  over  $S$  for each  $h \in \mathcal{W}$ ;
      if  $\exists h \in \mathcal{W}, \hat{\Delta}_t(h) \geq \Delta_g$  then return  $h$  and  $S$ ;
      else  $\Delta_g \leftarrow \Delta_g/2$ ;  $i \leftarrow i + 1$ ;  $\delta \leftarrow \delta/(i(i+1)|\mathcal{W}|)$ ;
    end-if
  end-do

```

Figure 3: GiniBoost_{filt}

and

$$\Pr_{D^m} \{ \hat{\Delta}_t(h) \leq (1 - \varepsilon) \Delta_t(h) \} \leq b_1 e^{-\frac{\varepsilon^2 \Delta_t m}{c_2}}, \quad (4)$$

where $b_1 \leq 8$, $c_1 \leq 600$, and $c_2 \leq 64$.

The proof of Lemma 1 is given in Appendix. Then, we analyze the adaptive sampling procedure HSelect. Let $\Delta_t^* = \max_{h' \in \mathcal{W}} \Delta_t(h')$. We prove the following lemma.

Lemma 2. Fix any $t \geq 1$. Then, with probability at least $1 - \delta$, (i) HSelect(ε, δ) outputs a hypothesis $h \in \mathcal{W}$ such that $\Delta_t(h) > (1 - \varepsilon)\Delta_t^*$, and (ii) the number of calls of $EX(f, D)$ is

$$O\left(\frac{\log \frac{1}{\delta} + \log |\mathcal{W}| + \log \log \frac{1}{\Delta_t^*}}{\varepsilon^2 \Delta_t^*}\right).$$

Finally we obtain the following theorem.

Theorem 4. With probability at least $1 - \delta$,

- (i) GiniBoost_{filt} outputs the final hypothesis h_{final} such that $\text{err}_D(h_{final}) \leq \varepsilon$,
- (ii) GiniBoost_{filt} terminates in $T = O(1/(\varepsilon\gamma^2))$ iterations, and
- (iii) the number of calls of EX(f, D) is

$$O\left(\frac{\log \frac{1}{\delta} + \log \frac{1}{\varepsilon\gamma} + \log |\mathcal{W}| + \log \log \frac{1}{\gamma}}{\varepsilon^2 \gamma^4} \cdot \left(\log \frac{1}{\delta} + \log \frac{1}{\varepsilon\gamma}\right)\right).$$

Proof. We say that GiniBoost fails at iteration t if one of the following event occurs: (a) HSelect fails, i.e., it does not meet the conditions (i) or (ii) in Lemma 2, (b) FiltEX calls $EX(f, D)$ for more than $(6/\varepsilon)M_t \log(1/\delta_t)$ times at iteration t , where M_t is denoted as the number of calls for FiltEX, (c) $\text{err}_D(\text{sign}(H_t)) > \varepsilon$ and $\widehat{\text{err}}_{S'_t}(\text{sign}(H_t)) < 2\varepsilon/3$, or (d) $\text{err}_D(\text{sign}(H_t)) < \varepsilon/2$ and $\widehat{\text{err}}_{S'_t}(\text{sign}(H_t)) > 2\varepsilon/3$. Note that, by Proposition 3, Lemma 2 and an application of Chernoff bound, the probability of each event (a), ..., (d) is at most δ_t , respectively. So the probability that GiniBoost fails is at most $3\delta_t$ at each iteration t . Then, during T iterations, GiniBoost fails at some iteration is at most $\sum_{t=1}^T 3\delta_t = \delta - \delta/(T+1) < \delta$. Now suppose that GiniBoost does not fail during T iterations. Then, we have $\text{err}_D(h_{final}) \leq 1 - \sum_{i=t}^T (1/8)\Delta_t^*$ by using the similar argument in the proof of Theorem 2, and thus GiniBoost $\text{err}_D(h_{final}) \leq \varepsilon/2$ in $T = 16/(\varepsilon\gamma^2)$ iterations. Then, since GiniBoost does not fail during T iterations, $\widehat{\text{err}}_{S'_t}(\text{sign}(H_t)) < 2\varepsilon/3$ at iteration $T+1$ and GiniBoost outputs h_{final} with $\text{err}_D(h_{final}) \leq \varepsilon/2$ and terminates. The total number of calls of $EX(f, D)$ in $T = O(1/(\varepsilon\gamma^2))$ iterations is $O(T \cdot M_T(1/\varepsilon) \log(1/\delta_T))$ with probability $1 - \delta$ and by combining with Lemma 2, we complete the proof. \square

5 Improvement on Sampling

While Lemma 1 gives a theoretical guarantee without any assumption, the bound has a constant factor $c_1 = 600$, which is too large to apply the lemma in practice. In this section, we derive a practical tail bound on the pseudo gain by using the central limit theorem. We say that a sequence of random variables $\{X_i\}$ is *asymptotically normal* with mean μ_i and variance σ_i^2 (we write X_i is $AN(\mu_i, \sigma_i^2)$ for short) if $(X_i - \mu_i)/\sigma_i$ converges to $N(0, 1)$ in distribution³. The central limit theorem states that, for independent random variables X_1, \dots, X_m from the same distribution with mean μ and variance σ^2 , $\sum_{i=1}^m X_i/m$ is $AN(\mu, \sigma^2/m)$. In particular, we use the multivariate version of the central limit theorem.

³Let $F_1(x), \dots, F_m(x)$, and $F(x)$ be distribution functions. Let X_1, \dots, X_m , and X be corresponding random variables, respectively. X_m converges to X in distribution if $\lim_{m \rightarrow \infty} F_m(x) = F(x)$.

Theorem 5 ([28]). Let $\mathbf{X}_1, \dots, \mathbf{X}_m$ be i.i.d. random vectors with mean $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$. Then, $\sum_{i=1}^m \mathbf{X}_i/m$ is $AN(\boldsymbol{\mu}, \boldsymbol{\Sigma})$.

Fix any hypothesis $h \in \mathcal{W}$, and distribution D_t over \mathcal{X} . Let $X \in \{0, 1\}$ and $Y \in \{-1, +1\}$ be random variables, induced by an independent random draw of $\mathbf{x} \in \mathcal{X}$ under D_t , such that $X = 1$ if $h(\mathbf{x}) = +1$, otherwise $X = 0$ and $Y = f(\mathbf{x})h_t(\mathbf{x})$, respectively. Then the pseudo gain $\Delta_t(h)$ can be written as $E(X) \cdot \{E(XY)/E(X)\}^2 + E(\bar{X}) \cdot \{E(\bar{X}Y)/E(\bar{X})\}^2$, where $\bar{X} = 1 - X$. Our empirical estimate of the pseudo gain is $Z = (\sum_{i=1}^m X_i Y_i/m)^2 / (\sum_{i=1}^m X_i/m) + (\sum_{i=1}^m \bar{X}_i Y_i/m)^2 / (\sum_{i=1}^m \bar{X}_i/m)$. The following theorem guarantees that a combination of sequences of asymptotically normal random variables is also asymptotically normal (Theorem 3.3A in [28]).

Theorem 6 ([28]). Suppose that $\mathbf{X} = (X^{(1)}, \dots, X^{(k)})$ is $AN(\boldsymbol{\mu}, b\boldsymbol{\Sigma})$, with $\boldsymbol{\Sigma}$ a covariance matrix and $b \rightarrow 0$. Let $g(\mathbf{x}) = (g_1(\mathbf{x}), \dots, g_n(\mathbf{x}))$, $\mathbf{x} = (x_1, \dots, x_k)$, be a vector-valued function for which each component function $g_i(\mathbf{x})$ is real-valued and has a nonzero differential at $\mathbf{x} = \boldsymbol{\mu}$. Then, $g(\mathbf{X})$ is $AN(g(\boldsymbol{\mu}), b^2 \mathbf{D}\boldsymbol{\Sigma}\mathbf{D}')$, where

$$\mathbf{D} = \left[\frac{\partial g_i}{\partial x_j} \Big|_{\mathbf{x}=\boldsymbol{\mu}} \right]_{n \times k}.$$

By using Theorem 5 and 6 for $\mathbf{X}_m = (\sum_{i=1}^m X_i/m, \sum_{i=1}^m X_i Y_i/m, \sum_{i=1}^m \bar{X}_i Y_i/m)$ and $g(u, v, w) = v^2/u + w^2/(1-u)$, we get the following result.

Corollary 7. $Z = \frac{(\sum_{i=1}^m X_i Y_i/m)^2}{\sum_{i=1}^m X_i/m} + \frac{(\sum_{i=1}^m \bar{X}_i Y_i/m)^2}{\sum_{i=1}^m \bar{X}_i/m}$ is $AN(\mu_z, \sigma_z^2)$, where $\mu_z = \frac{E(XY)^2}{E(X)} + \frac{E(\bar{X}Y)^2}{E(\bar{X})}$, and $\sigma_z^2 \leq 4\mu_z/m$.

Now we assume that the given sample is large enough to apply the central limit theorem. Then

$$\Pr \left\{ \frac{Z - \mu_z}{\sigma_z} \leq \varepsilon \right\} \approx \Phi(\varepsilon),$$

where $\Phi(x) = \int_{-\infty}^x (1/\sqrt{2\pi})e^{-\frac{1}{2}y^2} dy$. Since $1 - \Phi(x) \leq 1/(x\sqrt{2\pi})e^{-\frac{1}{2}x^2}$ (see, e.g., [9]),

$$\begin{aligned} \Pr \{Z - \mu_z > \varepsilon\mu_z\} &= \Pr \left\{ \frac{Z - \mu_z}{\sigma_z} > \frac{\varepsilon\mu_z}{\sigma_z} \right\} \lesssim \frac{\sigma_z}{\varepsilon\mu_z\sqrt{2\pi}} e^{-\frac{\varepsilon^2\mu_z^2}{2\sigma_z^2}} \\ &< \frac{2}{\sqrt{2\pi}\varepsilon^2\mu_z m} e^{-\frac{\varepsilon^2\mu_z m}{8}}. \end{aligned} \quad (5)$$

Substituting

$$m = \frac{8 \left(\ln \frac{1}{\delta\sqrt{2\pi}} - \frac{1}{2} \ln \ln \frac{1}{\delta\sqrt{2\pi}} \right)}{\varepsilon^2\mu_z}$$

to inequality (5), we obtain

$$\begin{aligned} \Pr \{Z - \mu_z > \varepsilon\mu_z\} &\leq \frac{1}{\sqrt{4\pi \left(\ln \frac{1}{\delta\sqrt{2\pi}} - \frac{1}{2} \ln \ln \frac{1}{\delta\sqrt{2\pi}} \right)}} \cdot \delta\sqrt{2\pi} \cdot \sqrt{\ln \frac{1}{\delta\sqrt{2\pi}}} \\ &< \frac{1}{\sqrt{2\pi \ln \frac{1}{\delta\sqrt{2\pi}}}} \cdot \delta\sqrt{2\pi} \cdot \sqrt{\ln \frac{1}{\delta\sqrt{2\pi}}} = \delta. \end{aligned}$$

Note that the same argument holds for $\Pr\{Z \leq (1-\varepsilon)\mu_z\}$. Therefore, we can replace the estimate of sample size $m = \frac{c_1 \ln(b_1/\delta)}{\varepsilon^2 \Delta_g}$ in HSelect with $m = \frac{8 \left(\ln \frac{1}{\delta \sqrt{2\pi}} - \frac{1}{2} \ln \ln \frac{1}{\delta \sqrt{2\pi}} \right)}{\varepsilon^2 \Delta_g}$ and this modification makes HSelect more practical.

6 Experimental Results

In this section, we show some experimental results on both artificial and real data sets. Our experiments consists of two parts.

6.1 Experiments in the Subsampling Framework

In the first part, we compare GiniBoost, AdaBoost, InfoBoost, and MadaBoost. in the subsampling framework.

For real data, we use some datasets from UCI machine learning repository [5]. Also, we prepare artificial data in order to examine behavior of boosting algorithms in details. The sizes of data we use vary from about 3,000 to 10,000. In particular, for artificial data, we use r -of- k function as the target function. An r -of- k function f over boolean domain $\{-1, +1\}^N$ consists of k relevant variables and $f(\mathbf{x}) = +1$ if at least r of the k relevant variables takes $+1$, otherwise $f(\mathbf{x}) = -1$. Note that 1-of- k function and $k/2$ -of- k function correspond to k -disjunction and k -majority, respectively. In [17], it is shown that, when boolean literals are used as weak hypotheses, InfoBoost can learn r -of- k functions in $O(rk)$ steps by taking advantage of the fact that weak hypotheses have low false positive error, whereas AdaBoost needs $O(k^2)$ steps. Our aim to conduct experiments over these artificial data sets is to see if GiniBoost behaves like InfoBoost. For $r = 1, 3, 5$ and $k = 10$, we fix a r -of- k function over $N = 100$ boolean variables as the target function, and we generate 10,000 random examples labeled by each r -of- k function, where the random examples are drawn so that positive and negative examples are equally likely.

For each dataset, we prepare decision stumps and the constant hypothesis $+1$ (i.e. the hypothesis that always answers $+1$) as weak hypotheses. In each dataset, each record have numeric attributes or binary attributes. For each numeric attribute, we construct a decision stump with a threshold, which predicts $+1$ or -1 depending on whether the value of the attribute is below the threshold or not. The threshold is chosen so that the training error of the decision stump is minimized. For each binary attribute, we prepare the decision stump which answers the value of the attribute.

We consider two versions of GiniBoost in our experiments. The first version is the original one which we described in Section 3. The second version is a slight modification of the original one, in which we use $\alpha_t[\pm 1] = \gamma_t[\pm 1]$. We call this version GiniBoost2.

We evaluate the boosting algorithms by cross validation. We split each data randomly 100 times, where each example is put into a training set with probability 0.7 and a test set with probability 0.3. For each training set, we run the boosting algorithms in 100 steps and evaluate their final hypotheses on the test data.

The results are summarized in Table 6.1. As shown in Table 6.1, performance of GiniBoost and GiniBoost2 appear to be comparable to those of others on real datasets. Also, GiniBoost2 behaves closely to InfoBoost.

dataset	Ada.	InfoB.	Mada.	Gini.	Gini2.
kr-vs-kp	5.2	5.6	5.2	5.6	5.7
hypothyroid	2.5	2.5	2.5	2.5	2.5
sick-euthoroid	5.8	5.6	5.8	5.6	5.5
spambase	23	23	23	23	23
10 of 70	19	6.2	19	10	7.5
20 of 70	8.8	6.1	8.8	10	8.3
30 of 70	7.3	6.7	7.3	10	7.9

Table 1: Test errors (%) of boosting algorithms in the subsampling framework.

6.2 Experiments in the Filtering Framework

In the second part, we compare MadaBoost and GiniBoost in the filtering framework. Basic settings of our experiments in the filtering framework are the same as those in the subsampling framework, except the following: First of all, in order to obtain large datasets, as is done in [6], we inflate the datasets by preparing 100 copies of each record in the data and changing their order randomly. Consequently, the sizes of the inflated data vary from 300,000 to 1,000,000. Second, instead of running each algorithm in 100 steps, we run them until they sample 10,000 examples. More precisely, we run GiniBoost with HSelect(ε, δ), where parameter $\varepsilon = 0.75$ and $\delta = 0.1$ are fixed. Also, we run MadaBoost with geometric AdaSelect [6] whose parameters are $s = 2$, $\varepsilon = 0.5$ and $\delta = 0.1$. Note that, in this setting, we demand both HSelect and AdaSelect to output a weak hypothesis h_t with $\gamma_t^2 \geq (1/4) \max_{h' \in \mathcal{W}} \gamma_t(h')^2$. In the following experiments, we use the approximation based on the central limit theorem, described in Section 5.

The results are summarized in Table ???. Compared to the results in the subsampling framework, both GiniBoost and GiniBoost work slightly better than MadaBoost, especially for artificial data.

Finally, in addition, we apply MadaBoost and GiniBoost for text categorization tasks on a collection of Reuters news (Reuters-21578⁴). We use the modified Apte(“ModApte”) split which contains about 10,000 news documents labeled with topics. We choose five major topics and for each topics, we let boosting algorithms classify whether a news document belongs to the topic or not. As weak hypotheses, we prepare about 30,000 decision stumps corresponding to words. This experiment is done in the same setting of previous ones, except that we do not inflate this dataset. Each algorithm are run until they sample 500,000 examples in total. We conduct 10 fold cross varidation.

The results are shown in Table 3 and Figure 4, As indicated, GiniBoost and GiniBoost2 outperform MadaBoost. We also run AdaBoost (without sampling) for 50 iterations, where Adaboost processes about 500,000 examples. Then, GiniBoost is about three times faster than AdaBoost, while improving the accuracy. The main reason why filtering-based algorithms save time would be that they use rejection sampling. By using rejection sampling, filtering-based algorithms keep only accepted examples in hand. Since the number of accepted example is much smaller than that of the whole given sample, we can find weak hypotheses faster over

⁴<http://www.daviddlewis.com/resources/testcollections/reuters21578>.

dataset	MadaBoost	GiniBoost	GiniBoost2
kr-vs-kp	6.1	6.1	6.1
hypothyroid	2.1	1.9	2.0
sick-euthoroid	5.6	5.6	5.6
spambase	24	22	22
10 of 70	41	29	30
20 of 70	35	30	30
30 of 70	30	30	30

Table 2: Test errors (%) of boosting algorithms in the filtering framework.

accepted examples than over the given sample.

In particular, GiniBoost uses fewer accepted examples than MadaBoost. mainly because they use different criteria. Roughly speaking, MadaBoost takes $\tilde{O}(1/\gamma_t^2)$ accepted examples in order to estimate γ_t . On the other hand, in order to estimate Δ_t , GiniBoost takes $\tilde{O}(1/\Delta_t)$ accepted examples, which are smaller than $\tilde{O}(1/\gamma_t^2)$ when we neglect constant factors. This consideration would explain why GiniBoost is faster than MadaBoost.

7 Summary and Future Work

In this paper, we propose a smooth boosting algorithm that uses an information-based criterion based on Gini index for choosing hypotheses. Our preliminary experiments shows that our algorithms perform well in the filtering framework. As future work, we further investigate the connections between boosting and information-based criteria. Also, we will conduct experiments over much huge data in the filtering framework.

Acknowledgments

I would like to thank Prof. Masayuki Takeda in Kyushu University for his various supports. This work is also supported in part by the 21st century COE program at Graduate School of Information Science and Electrical Engineering in Kyushu University.

References

- [1] J. A. Aslam. Improving algorithms for boosting. In *Proc. 13th Annu. Conference on Comput. Learning Theory*, pages 200–207, 2000.
- [2] Jose L. Balcazar, Yang Dai, and Osamu Watanabe. Provably fast training algorithms for support vector machines. In *Proceedings of IEEE International Conference on Data Mining (ICDM'01)*, pages 43–50, 2001.
- [3] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth International Group, 1984.

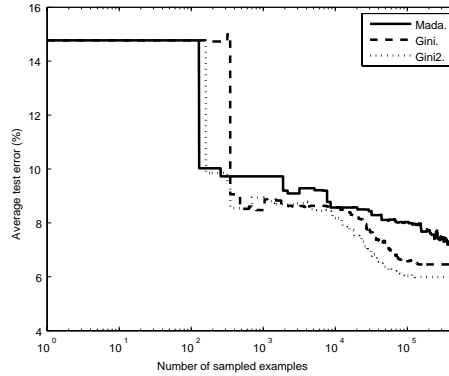


Figure 4: Test errors (%) of boosting algorithms for Reuters-21578 data. The test errors are averaged over topics.

- [4] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Annals of Mathematical Statistics*, 23:493–509, 1958.
- [5] C.L. Blake D.J. Newman, S. Hettich and C.J. Merz. UCI repository of machine learning databases, 1998.
- [6] C. Domingo, R. Gavaldà, and O. Watanabe. Adaptive sampling methods for scaling up knowledge discovery algorithms. *Data Mining and Knowledge Discovery*, 6(2):131–152, 2002.
- [7] C. Domingo and O. Watanabe. MadaBoost: A modification of AdaBoost. In *Proceedings of 13th Annual Conference on Computational Learning Theory*, pages 180–189, 2000.
- [8] P. Domingos and G. Hulten. Mining high-speed data streams. In *Proceedings of the Sixth ACM International Conference on Knowledge Discovery and Data Mining*, pages 71–80, 2000.
- [9] W. Feller. *An introduction to probability theory and its applications*. Wiley, 1950.
- [10] Y. Freund. An improved boosting algorithm and its implications on learning complexity. In *Proc. 5th Annual ACM Workshop on Computational Learning Theory*, pages 391–398. ACM Press, New York, NY, 1992.
- [11] Y. Freund. Boosting a weak learning algorithm by majority. *Information and Computation*, 121(2):256–285, 1995.
- [12] Y. Freund and R. E. Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55(1):119–139, 1997.
- [13] Yoav Freund. An adaptive version of the boost by majority algorithm. In *COLT '99: Proceedings of the twelfth annual conference on Computational learning theory*, pages 102–113, 1999.

	# of sampled examples	# of accepted examples	time (sec.)	test error (%)
Ada.	N/A	N/A	742	6.2
Mada.	527022	83607	314	7.0
Gini.	520408	85548	264	6.4
Gini2.	526657	77177	250	6.0

Table 3: Summary of experiments over Reuters-2158.

- [14] J. Friedman, T. Hastie, and R. Tibshirani. Additive logistic regression: a statistical view of boosting. *Annals of Statistics*, 2:337–374, 2000.
- [15] D. Gavinsky. Optimally-smooth adaptive boosting and application to agnostic learning. *Journal of Machine Learning Research*, 2003.
- [16] K. Hatano and M. K. Warmuth. Boosting versus covering. In *Advances in Neural Information Processing Systems 16*, 2003.
- [17] K. Hatano and O. Watanabe. Learning r-of-k functions by boosting. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory*, pages 114–126, 2004.
- [18] M. Kearns. Efficient noise-tolerant learning from statistical queries. In *Proceedings of 25th Annual ACM Symposium on Theory of Computing*, pages 392–401, 1993.
- [19] M. Kearns and Y. Mansour. On the boosting ability of top-down decision tree learning algorithms. *Journal of Computer and System Sciences*, 58(1):109–128, 1999.
- [20] Michael J. Kearns, Robert E. Schapire, and Linda Sellie. Toward efficient agnostic learning. In *COLT*, pages 341–352, 1992.
- [21] Richard J. Lipton and Jeffrey F. Naughton. Query size estimation by adaptive sampling. *Journal of Computer and System Sciences*, 51(1):18–25, 1995.
- [22] Yishay Mansour and David A. McAllester. Boosting using branching programs. *Journal of Computer and System Sciences*, 64(1):103–112, 2002.
- [23] Oden Maron and Andrew W. Moore. The racing algorithm: Model selection for lazy learners. *Artificial Intelligence Review*, 11(1-5):193–225, 1997.
- [24] R. E. Schapire, Y. Freund, P. Bartlett, and W. S. Lee. Boosting the margin: a new explanation for the effectiveness of voting methods. *The Annals of Statistics*, 26(5):1651–1686, 1998.
- [25] Robert E. Schapire. The strength of weak learnability. *Machine Learning*, 5(2):197–227, 1990.
- [26] Robert E. Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning*, 37(3):297–336, 1999.

- [27] Tobias Scheffer and Stefan Wrobel. Finding the most interesting patterns in a database quickly by using sequential sampling. *Journal of Machine Learning Research*, 3:833–862, 2003.
- [28] R. J. Serfling. *Approximation theorems of mathematical statistics*. Wiley, 1980.
- [29] R. A. Servedio. Smooth boosting and learning with malicious noise. In *14th Annual Conference on Computational Learning Theory*, pages 473–489, 2001.
- [30] Eiji Takimoto, Syuhei Koya, and Akira Maruoka. Boosting based on divide and merge. In *Proceedings of the 15th International Conference on Algorithmic Learning Theory*, pages 127–141, 2004.
- [31] L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- [32] L. G. Valiant. Learning disjunction of conjunctions. In *IJCAI*, pages 560–566, 1985.

Appendix

Lemma 3. Let $L(x) = x + 1$, if $x > 0$ and e^x , otherwise. Then it holds for any $a \in \mathbb{R}$ and any $x \in [-1, +1]$ that

$$L(x + a) \leq L(a) + L'(a)x + L'(a)x^2.$$

Proof. For any $x \in [-1, 1]$, let $g_x(a) = L(a) + L'(a)(x + x^2) - L(x + a)$. We consider the following cases. (Case 1: $x + a, a \leq 0$) We have $g_x(a) = e^a(1 + x + x^2 - e^x) \geq 0$, as $e^x \leq 1 + x + x^2$ for $x \in [-1, 1]$. (Case 2: $x + a, a \geq 0$) It is immediate to see that $g_x(a) = x^2 \geq 0$. (Case 3: $x + a < 0$, and $a > 0$) It holds that $g_x(a) = 1 + a + x + x^2 - e^{x+a} \geq 0$ since $g'_x(a) = 1 - e^{x+a} > 0$ and $g_x(0) = 1 + x + x^2 - e^x \geq 0$. (Case 4: $x + a > 0$, and $a < 0$) By using the fact that $1 + x + x^2 \geq e^x$ for $x \in [-1, 1]$, we have $g_x(a) = e^a(1 + x + x^2) - (x + a + 1) \geq e^{x+a} - (1 + x + a) \geq 0$. \square

Proof of Lemma 1.

Proof. We prove inequality (3) only. The other inequality can be proved in a symmetric way. To make our notation simple, let $\Delta = \Delta_t(h)$, $p = p_t(h)$, $q = \gamma_t[+1](h)$, and $r = \gamma_t[-1](h)$. Similarly, let $\hat{\Delta}$, \hat{p}, \hat{q} , and \hat{r} be the corresponding empirical estimates. Note that, by definition, $q, r, \hat{q}, \hat{r} \in [-1, +1]$.

First we bound the probability that $\hat{p}\hat{q}^2 - pq^2 > \varepsilon\Delta/2$. In the following, we consider three cases.

(Case 1: $p < \varepsilon\Delta/4$) Observe that $\hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2$ implies $\hat{p} > p + \varepsilon\Delta/4$. So, by applying Chernoff bound [4], we get

$$\Pr_{D^m} \{ \hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2 \} \leq e^{-\frac{(\frac{\varepsilon\Delta}{4p})^2 pm}{3}} \leq e^{-\frac{\varepsilon\Delta m}{12}}.$$

(Case 2: $p \geq \frac{\varepsilon\Delta}{4}$ and $pq^2 < \frac{\varepsilon\Delta}{8}$) Assume that $\hat{p} < 2p$. This implies $q^2 < \frac{\varepsilon\Delta}{4\hat{p}}$. So we get

$$\Pr_{D^m} \{ \hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2 \} \leq \Pr_{D^m} \{ \hat{q}^2 \geq \varepsilon\Delta/(2\hat{p}) \} \leq \Pr_{D^m} \left\{ |\hat{q} - q| \geq \frac{\sqrt{2}-1}{2} \sqrt{\varepsilon\Delta/\hat{p}} \right\}.$$

Notice, for any event A and B , that $\Pr\{A\} = \Pr\{A \wedge B\} + \Pr\{A \wedge \bar{B}\} \leq \Pr\{A|B\} + \Pr\{\bar{B}\}$. By using Hoeffding and Chernoff bounds,

$$\begin{aligned} \Pr_{D^m} \{ \hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2 \} &\leq \Pr_{D^m} \left\{ |\hat{q} - q| \geq \frac{\sqrt{2}-1}{2} \sqrt{\varepsilon\Delta/\hat{p}} \mid \hat{p} < 2p \right\} + \Pr_{D^m} \{ \hat{p} > 2p \} \\ &\leq 2e^{-\frac{(\sqrt{2}-1)^2 \varepsilon\Delta m \hat{p}}{4}} + e^{-\frac{pm}{3}} \leq 3e^{-\frac{(\sqrt{2}-1)^2 \varepsilon\Delta}{4}} \leq 3e^{-\frac{\varepsilon\Delta}{25}}. \end{aligned}$$

(Case 3: $p \geq \frac{\varepsilon\Delta}{4}$ and $pq^2 > \frac{\varepsilon\Delta}{8}$) Let α be the real number such that $pq^2 = \alpha\Delta$. By our assumption on Case 3, it holds that $\alpha > \frac{\varepsilon}{8}$. Then we define the following events:

$$E_1 : \hat{p} \leq (1 + \varepsilon_1)p, \text{ and } E_2 : \hat{q}^2 \leq (1 + \varepsilon_1)q^2,$$

where we let $\varepsilon_1 = \frac{\varepsilon}{\sqrt{24\alpha}}$, so that $2\varepsilon_1 + \varepsilon_1^2 < \frac{\varepsilon}{2\alpha}$. (Note that this inequality might not hold for $\varepsilon > 1$). By definition of α , $\varepsilon_1 < \sqrt{\frac{\varepsilon}{3}} < 1$.

The events E_1 and E_2 imply that

$$\hat{p}\hat{q}^2 - pq^2 \leq (2\varepsilon_1 + \varepsilon_1^2)pq^2 < \frac{\varepsilon}{2\alpha}pq^2 = \frac{\varepsilon}{2}\Delta.$$

Therefore we get

$$\Pr_{D^m} \{ \hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2 \} \leq \Pr_{D^m} \{ \bar{E}_1 \vee \bar{E}_2 \} \leq \Pr_{D^m} \{ \bar{E}_1 \} + \Pr_{D^m} \{ \bar{E}_2 \}. \quad (6)$$

Then, we give an upperbound of the probability that each event E_i does not happen. By Chernoff bound [4], we have

$$\Pr_{D^m} \{ \bar{E}_1 \} \leq e^{-\frac{\varepsilon_1^2 pm}{3}} \leq e^{-\frac{\varepsilon^2 \Delta m}{72}}. \quad (7)$$

Since it holds that $\sqrt{1+x} \geq (\sqrt{2}-1)x + 1$ for $0 \leq x \leq 1$, we have

$$\hat{q}^2 > (1 + \varepsilon_1)q^2 \Leftrightarrow |\hat{q}| > \sqrt{1 + \varepsilon_1}|q| \Rightarrow |\hat{q} - q| > (\sqrt{2}-1)\varepsilon_1|q|.$$

Then, given that $\hat{p} > p/2$, we obtain, by applying Hoeffding bounds that

$$\begin{aligned} \Pr_{D^m} \{ \bar{E}_2 \mid \hat{p} > p/2 \} &\leq \Pr_{D^m} \left\{ |\hat{q} - q| > (\sqrt{2}-1)\varepsilon_1|q| \mid \hat{p} > p/2 \right\} \\ &\leq 2e^{-\frac{(\sqrt{2}-1)^2 \varepsilon_1^2 p q^2 m}{2}} \leq 2e^{-\frac{(\sqrt{2}-1)^2 \varepsilon_1^2 p q^2 m}{4}}. \end{aligned}$$

Again, by using the inequality $\Pr\{A\} \leq \Pr\{A|B\} + \Pr\{\bar{B}\}$ and Chernoff bound, we have

$$\begin{aligned} \Pr_{D^m} \{ \bar{E}_2 \} &= \Pr_{D^m} \{ \bar{E}_2 \mid \hat{p} > p/2 \} + \Pr_{D^m} \{ \hat{p} < p/2 \} \\ &\leq 2e^{-\frac{(\sqrt{2}-1)^2 \varepsilon_1^2 p q^2 m}{4}} + e^{-\frac{pm}{12}} \\ &\leq 2e^{-\frac{(\sqrt{2}-1)^2 \varepsilon^2 \Delta m}{96}} + e^{-\frac{\varepsilon \Delta m}{48}} \leq 3e^{-\frac{\varepsilon^2 \Delta m}{600}}. \end{aligned} \quad (8)$$

Using (6), (7) and (8), we have

$$\Pr_{D^m} \{ \hat{p}\hat{q}^2 - pq^2 \geq \varepsilon\Delta/2 \} \leq \Pr_D \{ \bar{E}_1 \} + \Pr_D \{ \bar{E}_2 \} \leq 4e^{-\frac{\varepsilon^2\Delta m}{600}}. \quad (9)$$

In any case, we have the inequality (9). Applying the similar analysis for the probability that $(1 - \hat{p})\hat{r}^2 - (1 - p)q^2 \geq \varepsilon\Delta/2$, we get

$$\Pr_{D^m} \{ \hat{\Delta} - \Delta > \varepsilon\Delta \} \leq 4e^{-\frac{\varepsilon^2\Delta m}{600}} + 4e^{-\frac{\varepsilon^2\Delta m}{600}} \leq 8e^{-\frac{\varepsilon^2\Delta m}{600}}.$$

□

Proof of Lemma 2

Proof. First, we prove the first statement (i) of Lemma 2. Let i^* be the minimum integer such that $1/2^{i^*} < (1 - \varepsilon/2)\Delta_t^*$. In the following, we show that, with high probability, HSelect outputs h (and S) such that $\Delta_t(h) \geq (1 - \varepsilon)\Delta_t^*$ during i^* trials of the do-loops. For each trial $i = 1, \dots, i^*$, we denote $\Delta_g(i) = 1/2^i$, $\delta(i) = \delta/(|\mathcal{W}|i(i+1))$, $\hat{\Delta}_t^* = \hat{\Delta}_t(h^*)$, where $h^* = \arg \max_{h' \in \mathcal{W}} \Delta_t(h')$, and $m(i) = \left\lceil (c_1 \ln \frac{b_1}{\delta(i)}) / (\varepsilon^2 \Delta_g(i)) \right\rceil$. Let $\mathcal{W}_{bad} = \{h \in \mathcal{W} | \Delta_t(h) < (1 - \varepsilon)\Delta_t^*\}$. We say that Hselect fails if it outputs a hypothesis $h \in \mathcal{W}_{bad}$ during i^* trials. Note that if Hselect fails, either one of the following events happen: (1) At some trial i ($1 \leq i \leq i^*$), some $h \in \mathcal{W}_{bad}$ satisfies $\hat{\Delta}_t(h) \geq \Delta_g$, or (2) at the i^* th trial, some $h \in \mathcal{W}_{bad}$ satisfies $\hat{\Delta}_t(h) > \hat{\Delta}_t^*$ or $\hat{\Delta}_t^* < \Delta_g$. So we can bound the probability that Hselect fails by bounding the probability that either the event (1) or (2) occurs.

First, consider the event (1). For any i such that $1 \leq i < i^*$ and for any $h \in \mathcal{W}_{bad}$, it holds, by the definition of i^* , that $\Delta_g(i) \geq (1 - \varepsilon/2)\Delta_t^* \geq (1 - \varepsilon/2)\Delta_t(h)/(1 - \varepsilon)$. Note that $\Delta_g(i) \geq (1 - \varepsilon/2)\Delta_t(h)/(1 - \varepsilon)$ if and only if $\Delta_g(i) \geq \Delta_t(h) + \varepsilon/(2 - \varepsilon) \cdot \Delta_g(i)$. Therefore, by Lemma 1, for any i ($1 \leq i < i^*$) and any $h \in \mathcal{W}_{bad}$,

$$\begin{aligned} \Pr_{D^{m(i)}} \{ \hat{\Delta}_t(h) > \Delta_g(i) \} &\leq \Pr_{D^{m(i)}} \left\{ \hat{\Delta}_t(h) > \Delta_t(h) + \frac{\varepsilon}{2 - \varepsilon} \Delta_g(i) \right\} \\ &\leq b_1 e^{-\frac{\left(\frac{\varepsilon}{2 - \varepsilon} \frac{\Delta_g(i)}{\Delta_t(h)}\right)^2 m(i) \Delta_t(h)}{c_1}} \leq \delta(i). \end{aligned}$$

Then the probability of the event (1) is bounded by

$$\begin{aligned} \sum_{1 \leq i < i^*} \sum_{h \in \mathcal{W}_{bad}} \Pr_{D^{m(i)}} \{ \hat{\Delta}_t(h, S_i) > \Delta_g(i) \} &\leq \sum_{1 \leq i < i^*} \sum_{h \in \mathcal{W}_{bad}} \delta(i) \\ &= \sum_{1 \leq i < i^*} \frac{\delta}{i(i+1)} = \delta \left(1 - \frac{1}{i^*} \right). \end{aligned}$$

Next, we bound the probability of the event (2). Observe that, if the event (2) holds, then (a) for some $h \in \mathcal{W}_{bad}$, $\hat{\Delta}_t(h) > (1 - \varepsilon/2)\Delta_t^*$ or (b) $\hat{\Delta}_t^* < (1 - \varepsilon/2)\Delta_t^*$. Further, the event (a) implies $\hat{\Delta}_t(h) > \Delta_t(h) + (\varepsilon/2)\Delta_t^*$. By Lemma 1, we see that the probability that the event (a) occurs

is at most

$$\begin{aligned}
(n-1) \Pr_{D^{m(i^*)}} \{ \hat{\Delta}_t(h) > \Delta_t(h) + (\varepsilon/2)\Delta_t^* \} &\leq (n-1)b_1 e^{-\frac{\varepsilon^2 m(i^*) \Delta_t^{*2}}{4c_1 \Delta_t(h)}} \\
&\leq (n-1)b_1 e^{-\frac{\varepsilon^2 m(i^*) \Delta_t^*}{4c_1(1-\varepsilon)}} \\
&\leq (n-1)b_1 e^{-\frac{\varepsilon^2 m(i^*)(1-\varepsilon/2)\Delta_g(i^*)}{4c_1(1-\varepsilon)}} \\
&\leq (n-1)\delta(i).
\end{aligned}$$

Also, the probability of the event (b) occurs is at most $\Pr_{D^{m(i^*)}} \{ \hat{\Delta}_t^* < (1-\varepsilon/2)\Delta_t^* \} \leq b_1 \exp\{-\frac{\varepsilon^2 \Delta_{*} m(i^*)}{4c_1}\} \leq \delta(i^*)$. So we obtain,

$$\Pr_{D^{m(i^*)}} \{ \text{the event (2) occurs} \} \leq n\delta(i^*) \leq \frac{\delta}{i^*(i^*+1)}.$$

Finally, we obtain that

$$\begin{aligned}
\Pr_{D^{m(i^*)}} \{ \text{HSelect fails} \} &\leq \Pr_{D^{m(i^*)}} \{ \text{event (1) or (2) happens} \} \\
&\leq \delta \left(1 - \frac{1}{i^*} \right) + \frac{\delta}{i^*(i^*+1)} < \delta.
\end{aligned}$$

Next, we consider the sample complexity of HSelect. Note that, by the definition of i^* , $i^* < \log(1/\Delta_t^*) + 2$. Then, it immediately follows that with probability at least $1 - \delta$, HSelect calls FiltEX at most $m(i^*) = O((\log(1/\delta) + \log |\mathcal{W}| + \log \log(1/\Delta_t^*))/\Delta_t^*)$ times. \square

8 Proof of Corollary 7

Proof. Let $g(u, v, w) = v^2/u + w^2/(1-u)$. Then its partial derivatives are

$$\frac{\partial g}{\partial u} = -\frac{v^2}{u^2} + \frac{w^2}{(1-u)^2}, \quad \frac{\partial g}{\partial v} = \frac{2v}{u}, \quad \text{and} \quad \frac{\partial g}{\partial w} = \frac{2w}{1-u}.$$

We denote $U = \sum_{i=1}^m X_i/m$, $\bar{U} = \sum_{i=1}^m \bar{X}_i/m$, $V = \sum_{i=1}^m X_i Y_i/m$, and $W = \sum_{i=1}^m \bar{X}_i Y_i/m$. For any random variables A and B , let $\mu_a = E(A)$ and $\sigma_a^2 = \text{Var}(A)$ and let C_{ab} be the covariance between A and B . By applying Theorem 6, Z is $AN(\mu_z, \sigma_z^2)$, where

$$\begin{aligned}
\sigma_z^2 &= \left(\begin{pmatrix} -\frac{\mu_v^2}{\mu_u^2} + \frac{\mu_w^2}{\mu_{\bar{u}}^2} & \frac{2\mu_v}{\mu_u} & \frac{2\mu_w}{\mu_{\bar{u}}} \end{pmatrix} \begin{pmatrix} \sigma_u^2 & C_{uv} & C_{uw} \\ C_{uv} & \sigma_v^2 & C_{vw} \\ C_{uw} & C_{vw} & \sigma_w^2 \end{pmatrix} \begin{pmatrix} \left(-\frac{\mu_v}{\mu_u} + \frac{\mu_w}{\mu_{\bar{u}}} \right) \\ \frac{2\mu_v}{\mu_u} \\ \frac{2\mu_w}{\mu_{\bar{u}}} \end{pmatrix} \right) \\
&= \left(-\frac{\mu_v^2}{\mu_u^2} + \frac{\mu_w^2}{\mu_{\bar{u}}^2} \right)^2 \sigma_u^2 + \frac{4\mu_v^2}{\mu_u^2} \sigma_v^2 + \frac{4\mu_w^2}{\mu_{\bar{u}}^2} \sigma_w^2 + \frac{8\mu_v \mu_w}{\mu_u \mu_{\bar{u}}} C_{vw} \\
&\quad + \frac{4\mu_v}{\mu_u} \left(-\frac{\mu_v^2}{\mu_u^2} + \frac{\mu_w^2}{\mu_{\bar{u}}^2} \right) C_{uv} + \frac{4\mu_w}{\mu_{\bar{u}}} \left(-\frac{\mu_v^2}{\mu_u^2} + \frac{\mu_w^2}{\mu_{\bar{u}}^2} \right) C_{uw}. \tag{10}
\end{aligned}$$

One can verify that

$$\begin{aligned}
C_{uv} &= \frac{\mu_{\bar{x}} \mu_{xy}}{m}, \quad C_{uw} = -\frac{\mu_x \mu_{\bar{x}y}}{m}, \quad C_{vw} = -\frac{\mu_x \mu_{\bar{x}y}}{m} \\
\sigma_u^2 &= \frac{\mu_x \mu_{\bar{x}}}{m}, \quad \sigma_v^2 = \frac{\sigma_{xy}}{m}, \quad \text{and} \quad \sigma_w^2 = \frac{\sigma_{\bar{x}y}}{m}. \tag{11}
\end{aligned}$$

By plugging (11) into (10), we obtain

$$\begin{aligned}
\sigma_z^2 &= \left(\frac{\mu_{xy}^4}{\mu_x^4} - \frac{2\mu_{xy}^2\mu_{\bar{x}y}^2}{\mu_x^2\mu_{\bar{x}}^2} + \frac{\mu_{\bar{x}y}^4}{\mu_{\bar{x}}^4} \right) \frac{\mu_x\mu_{\bar{x}}}{m} + \frac{4\mu_{xy}^2\sigma_{xy}^2}{\mu_x^2 m} + \frac{4\mu_{\bar{x}y}^2\sigma_{\bar{x}y}^2}{\mu_{\bar{x}}^2 m} - \frac{8\mu_{xy}^2\mu_{\bar{x}y}^2}{\mu_x\mu_{\bar{x}}m} \\
&\quad + \frac{4\mu_{xy}}{\mu_x} \left(-\frac{\mu_{xy}^2}{\mu_x^2} + \frac{\mu_{\bar{x}y}^2}{\mu_{\bar{x}}^2} \right) \frac{\mu_{\bar{x}}\mu_{xy}}{m} - \frac{4\mu_{\bar{x}y}}{\mu_{\bar{x}}} \left(-\frac{\mu_{xy}^2}{\mu_x^2} + \frac{\mu_{\bar{x}y}^2}{\mu_{\bar{x}}^2} \right) \frac{\mu_x\mu_{\bar{x}y}}{m} \\
&\leq \frac{\mu_{xy}^2}{\mu_x m} \left\{ \frac{\mu_{xy}^2\mu_{\bar{x}}}{\mu_x^2} + \frac{4\sigma_{xy}^2}{\mu_x} - \frac{4\mu_{xy}^2\mu_{\bar{x}}}{\mu_x^2} \right\} + \frac{\mu_{\bar{x}y}^2}{\mu_{\bar{x}} m} \left\{ \frac{\mu_{\bar{x}y}^2\mu_x}{\mu_{\bar{x}}^2} + \frac{4\sigma_{\bar{x}y}^2}{\mu_{\bar{x}}} - \frac{4\mu_{\bar{x}y}^2\mu_x}{\mu_{\bar{x}}^2} \right\} \\
&\leq \frac{4\mu_z}{m},
\end{aligned}$$

as desired. □