
Margins are Insufficient for Explaining Gradient Boosting

Allan Grønlund ^{*†}

Lior Kamma [†]

Kasper Green Larsen [†]

Abstract

Boosting is one of the most successful ideas in machine learning, achieving great practical performance with little fine-tuning. The success of boosted classifiers is most often attributed to improvements in margins. The focus on margin explanations was pioneered in the seminal work by Schapire et al. (1998) and has culminated in the k 'th margin generalization bound by Gao and Zhou (2013), which was recently proved to be near-tight for some data distributions (Grønlund et al. 2019). In this work, we first demonstrate that the k 'th margin bound is inadequate in explaining the performance of state-of-the-art gradient boosters. We then explain the short comings of the k 'th margin bound and prove a stronger and more refined margin-based generalization bound for boosted classifiers that indeed succeeds in explaining the performance of modern gradient boosters. Finally, we improve upon the recent generalization lower bound by Grønlund et al. (2019).

1 Introduction

Boosting is a powerful technique for producing highly accurate voting classifiers by combining less accurate base learners. Boosting algorithms are typically easy to fine tune and obtain state-of-the-art performance on many learning tasks. Boosting dates back to the seminal work introducing the AdaBoost algorithm [4] and much work has gone into understanding and developing better boosting algorithms. The best performing boosting algorithms are typically variants of gradient boosters [5], such as LightGBM [9] and XGBoost [2], using Regression Trees as base learners.

Classic experiments [14] showed that boosting algorithms tend to improve their test accuracy even when training past the point of perfectly classifying the training data. This may seem counter-intuitive, as adding more base learners, results in a more complex model, that hence might be more prone to overfitting. This phenomenon is often explained by observed improvements in margins. For binary classification with a sample space \mathcal{X} , labels in $\{-1, 1\}$ and a class of base learners $\mathcal{H} \subseteq \mathcal{X} \rightarrow [-1, 1]$, a voting classifier $f : \mathcal{X} \rightarrow \{-1, 1\}$ has the form $f(x) = \text{sign}(\sum_{h \in \mathcal{H}} \alpha_h h(x))$ with all $\alpha_h \geq 0$. A voting classifier thus takes a weighted “vote” among the base learners to obtain its prediction. When speaking of margins, we assume $\sum_h \alpha_h = 1$, which can always be achieved by rescaling the α 's by their sum without changing f . The margin of a training point (x, y) with $x \in \mathcal{X}$ and $y \in \{-1, 1\}$ is then defined as $y \sum_h \alpha_h h(x)$. The margin is thus a value in $[-1, 1]$ which is positive when $f(x) = y$ and negative otherwise. Intuitively, large (positive) margins mean that f is not only correct but very certain in its predictions. Margin theory, starting with the work of Schapire et al. [4], formalized this by proving generalization bounds demonstrating that large margins imply better generalization. It was also shown that the theoretical generalization bounds fit very well with the observed behavior of AdaBoost that tends to keep improving margins even when training past the point of perfectly classifying the training data [14].

^{*}All authors contributed equally, and are presented in alphabetical order.

[†]Department of Computer Science, Aarhus University, {j.allan, lior.kamma, larsen}@cs.au.dk

However, shortly after [4] and [14] was published, Breiman [1] proved a generalization bound based on the minimal margin (the smallest margin achieved by a training point) that was sharper than the generalization bound in Schapire et al. [4]. He then designed a new boosting algorithm, named *Arc-GV*, that provably optimizes the minimal margin, which AdaBoost does not (see [12] for the full story of maximizing the minimal margin). In the same paper, Breiman experimentally showed that Arc-GV produced not just a better minimal margin, but better margins overall, than AdaBoost. However, AdaBoost still obtained a better generalization and test error. This seemed to contradict margin theory, as according to margin theory, all other things being equal, then larger margins should imply better generalization. Later it was shown by Reyzin and Schapire [13] that Breiman’s experiments did not accurately take into account the complexity of the base learner trees created by AdaBoost and Arc-GV, as repeating the experiments showed that Arc-GV produced trees of larger depth than AdaBoost, and deeper trees may be more prone to overfitting. Reyzin and Schapire then considered the same experiments using stumps as base learners, forcing identical depth trees between the algorithms, and in this case, AdaBoost produced better margin distributions than Arc-GV and also generalized better. These findings support the view that better margins provide better generalization as presented in [4, 14].

Later, [16, 10, 6] showed improved generalization bounds that subsumed both the generalization bounds by Schapire et al., and Breiman, providing further theoretical support for margin theory. The current strongest generalization bounds are as follows. Let \mathcal{D} be any distribution over $\mathcal{X} \times \{-1, 1\}$ and define $\mathcal{L}_{\mathcal{D}}(f) = \Pr_{(x,y) \sim \mathcal{D}}[f(x) \neq y]$ as the out-of-sample error of a voting classifier f . Also, for a set $S = \{(x_i, y_i)\}_{i=1}^m$ of m labeled samples drawn i.i.d. from \mathcal{D} , define $\mathcal{L}_S^\theta(f) = \Pr_{(x,y) \sim S}[yf(x) < \theta]$ as the fraction of points in S with margin less than θ (the notation $(x, y) \sim S$ denotes a uniform random point (x, y) in S). With this notation, there are two strongest current generalization bounds. The first [10] uses Rademacher complexity to show that with high probability over the sample set S , it holds for every margin $\theta \in (0, 1]$ and every voting classifier f that:

$$\mathcal{L}_{\mathcal{D}}(f) \leq \mathcal{L}_S^\theta(f) + O\left(\sqrt{\frac{\lg |\mathcal{H}|}{\theta^2 m}}\right). \quad (1)$$

The k ’th margin bound by Gao and Zhou [6] improves this for $\mathcal{L}_S^\theta(f) = o(1/\lg m)$ and is as follows:

$$\mathcal{L}_{\mathcal{D}}(f) \leq \mathcal{L}_S^\theta(f) + O\left(\frac{\lg |\mathcal{H}| \lg m}{\theta^2 m} + \sqrt{\mathcal{L}_S^\theta(f) \cdot \frac{\lg |\mathcal{H}| \lg m}{\theta^2 m}}\right). \quad (2)$$

The k ’th margin bound subsumes both Breiman’s min margin generalization bound and the original generalization bound by Schapire et al. For infinite \mathcal{H} , one may replace $\lg |\mathcal{H}|$ in the above bounds with the VC-dimension of \mathcal{H} times a $\lg m$ factor (as is standard). For simplicity, we focus on the case of finite \mathcal{H} throughout the paper. Moreover, recent work by Grønlund et al. [7] shows that the margin bounds above are near-tight. Formally, they show that for (almost) all margins θ , there exists a data distribution \mathcal{D} and a set of base learners \mathcal{H} , such that with constant probability over the sample set S , there is a voting classifier f such that

$$\mathcal{L}_{\mathcal{D}}(f) \geq \mathcal{L}_S^\theta(f) + \Omega\left(\frac{\lg |\mathcal{H}| \lg m}{\theta^2 m} + \sqrt{\mathcal{L}_S^\theta(f) \cdot \frac{\lg |\mathcal{H}|}{\theta^2 m}}\right). \quad (3)$$

Moreover, the lower bound holds for any value of $\mathcal{L}_S^\theta(f) \leq 49/100$ and any value of $\lg |\mathcal{H}|$ [7].

Remark. Many boosting algorithms produce classifiers $f = \sum_h \alpha_h h$ where $\sum_h \alpha_h \neq 1$ or where base learners output values in \mathbb{R} rather than $[-1, 1]$. To apply margin theory, following [15], such classifiers are rescaled as follows: For each h with output range $[a_h, b_h]$ and coefficient α_h , divide all outputs of h by $\Delta_h = \max\{|a_h|, |b_h|\}$ and multiply α_h by Δ_h . Afterwards, divide all α_h by $\sum_h \alpha_h$.

1.1 Our contribution.

A new margin lower bound: Comparing the current best upper and lower bounds, we see that (2) and (3) match when $\mathcal{L}_S^\theta(f)$ approaches 0. Similarly, we see that (2) and (1) match as $\mathcal{L}_S^\theta(f)$ approaches a constant. But what is the true behavior in-between? The k ’th margin bound (2) gained the factor $\mathcal{L}_S^\theta(f)$ inside the $\sqrt{\cdot}$ but lost a factor $\lg m$ compared to (1). Can the $\lg m$ factor be removed?

What is the correct behavior as $\mathcal{L}_S^\theta(f)$ goes from 0 towards 1? In this work, we show an improved generalization lower bound of:

$$\mathcal{L}_D(f) \geq \mathcal{L}_S^\theta(f) + \Omega \left(\frac{\lg |\mathcal{H}| \lg m}{\theta^2 m} + \sqrt{\mathcal{L}_S^\theta(f) \cdot \frac{\lg |\mathcal{H}| \lg (\mathcal{L}_S^\theta(f)^{-1})}{\theta^2 m}} \right). \quad (4)$$

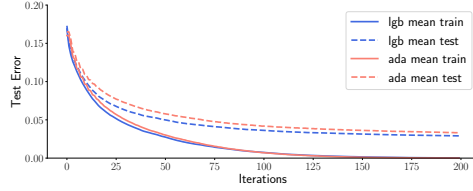
Our lower bound shows that the $\lg m$ factor inside the $\sqrt{\cdot}$ has to show up as $\mathcal{L}_S^\theta(f)$ drops to $m^{-\varepsilon}$ for any constant $\varepsilon > 0$. Moreover, our new lower bound completely settles the generalization performance of boosting in terms of margins whenever $\mathcal{L}_S^\theta(f)$ is outside the range $m^{-o(1)}$ to $o(1)$. It also nicely interpolates between the $\mathcal{L}_S^\theta(f) = 0$ and $\mathcal{L}_S^\theta(f) = 1$ case. We conjecture that the lower bound gives the correct margin-based tradeoff, i.e. that it is possible to improve the upper bounds (1) and (2) to match (4).

A new refined margin generalization bound: The main part of our paper considers a new refined margin based generalization bound for voting classifiers (boosting algorithms). First, we present experiments showing that the classic margin bounds alone fail to explain the performance of state-of-the-art gradient boosting algorithms. More concretely, we show that gradient boosters actually may produce smaller and smaller margins when run for many iterations, despite the test accuracy staying the same or even improving. We additionally demonstrate that the classic version of AdaBoost may produce significantly better margins than gradient boosters, despite gradient boosters obtaining similar or even better test accuracy and generalization error than AdaBoost. To explain this inconsistency, we observe experimentally that the trees produced by gradient boosters return very small values on all but a few training points, thus making minimal changes to most predictions when added to the voting classifier. We then use this insight to prove a new margin-based generalization bound for boosting algorithms which also take into account the magnitude of predictions by base learners. Finally, we run experiments demonstrating that our refined generalization bounds in fact succeed in explaining and predicting the performance of boosting algorithms. In addition to achieving a better theoretical understanding of boosting algorithms, in particular gradient boosters, these new insights may potentially lead to new algorithms with better accuracy by using regularization inspired by our new generalization bound or more directly optimizing it.

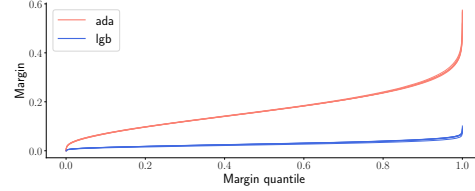
2 Insufficiency of current margin bounds

From the margin-based upper and lower bounds, it may seem that we have all the theory necessary for understanding the generalization performance of boosters. To confirm the theory, we ran experiments with AdaBoost and the state-of-the-art gradient booster LightGBM on standard data sets with the same size trees as base learners. For all experiments we only change the tree size and learning rate of the LightGBM hyperparameters. For AdaBoost we allow the same tree size, unlimited depth, as well as forcing a minimum number of elements in each tree learner to be 20 as is default in LightGBM.

Figure 1b shows a plot of the margin distributions for the two boosters trained on the Forest Cover dataset. From this plot, it is obvious that AdaBoost achieves significantly better margins than LightGBM. Indeed, the k 'th smallest margin of AdaBoost, is much larger than the k 'th smallest margin of LightGBM for all k where at least one of the two margins are non-negative. Thus, from the generalization bounds (1) and (2), AdaBoost should have a much smaller out-of-sample error than LightGBM. However, the corresponding test errors in Figure 1a show a very different story, with LightGBM slightly outperforming AdaBoost. Furthermore, as shown in Section 3, the trees produced by LightGBM are in fact deeper than the trees produced by AdaBoost. This gives rise to some concerns regarding the explanatory power of margins. To further underline the theoretical inconsistency, we examine the two generalization bounds (1) and (2). When applying the generalization bounds to AdaBoost and LightGBM, then for any choice of $p = \mathcal{L}_S^\theta(f) \in [0, 1]$, the only parameter that vary between AdaBoost and LightGBM is θ^{-2} . That is, if we choose θ as the (pm) 'th smallest margin, i.e. fix $\mathcal{L}_S^\theta(f) = p$, then only the value of θ differ between the two boosters and the generalization error grows as θ^{-2} . Figure 2a shows a plot of θ^{-2} as a function of $\mathcal{L}_S^\theta(f)$ for the two boosters. Clearly the penalty in the generalization error is much smaller for AdaBoost, suggesting that AdaBoost should perform much better than LightGBM, despite the test errors in Figure 1a showing that LightGBM outperforms AdaBoost. To investigate this phenomenon further, we have plotted the margin distribution of the two boosters after $t = 10, 20$ and 50 iterations of training, see Figure 2b. It is clear from this plot that the margins of the gradient booster, learned by

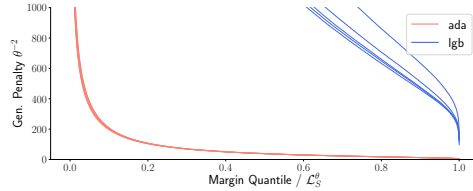


(a) Mean training and test error over five runs. The standard deviation of the final test error is 0.00037 for AdaBoost and smaller for LightGBM.

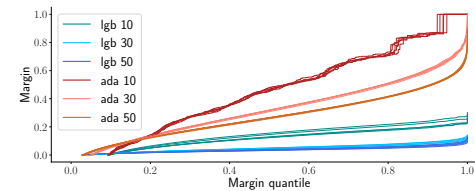


(b) Sorted margin values.

Figure 1: Accuracy and margin plots for AdaBoost and LightGBM on the Forest Cover data set.



(a) Plot of θ^{-2} when choosing θ as the (pm) 'th smallest margin for $p \in [0, 1]$. The margins are those also shown in Figure 1b.



(b) Development in margin distributions for AdaBoost and LightGBM.

Figure 2: Generalization penalties and margin distributions on the Forest Cover data set.

LightGBM, deteriorate quickly with the number of training iterations. To explain why the margins quickly drops towards 0 for the gradient booster, we take a closer look at the trees produced by LightGBM compared to AdaBoost. Figure 3 shows a histogram of the predictions made by the trees produced by LightGBM. It is very striking from this histogram that the trees making up the LightGBM gradient booster makes very small (in absolute value) predictions on most data points, whereas AdaBoost always makes predictions in $\{-1, 1\}$. Note that each tree always has its largest prediction among $\{-1, 1\}$. Thus, LightGBM produces trees that only significantly change the predictions of very few data points, while leaving almost all others unchanged. When training more and more trees, this causes the margins to diminish. To see this, consider as an example a training point $(x, 1)$ and assume the first trained tree h makes a (correct) prediction of $h(x) = 1$ and is assigned a weight of $\alpha_h = 1$. After the first training iteration, the margin of $(x, 1)$ is 1. However, as training progresses, many more trees may be produced that all predict 0 on x while being assigned a weight of 1. Since margins are normalized, $\sum_{h \in \mathcal{H}} \alpha_h = 1$, this means that the margin of x drops to $1/t$ after t rounds of training. The drop in predicted accuracy by the generalization bounds (1) and (2) seem unreasonable if we think about the data point x (the error is expected to grow as t^2 or t). A

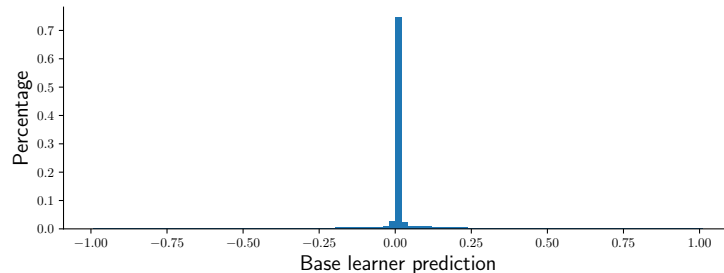


Figure 3: Histogram of base learner predictions for LightGBM on the Forest Cover data set. Only about 1 in 5000 predictions are larger than 0.95 in absolute value.

possible explanation of the shortcomings of current generalization bounds is thus that they simply treat base learners as arbitrary functions in $\mathcal{X} \rightarrow [-1, 1]$. That is, they pay no attention to the fact that base learners trained by gradient boosters make very small predictions on almost all data points.

To further support this claim, we note that the proof of the previous generalization lower bound (3) as well as our improved bound (4) construct a set of base learners \mathcal{H} where all $h \in \mathcal{H}$ make predictions among $\{-1, 1\}$, i.e. they make no predictions of small magnitude. This further supports the belief that an explanation based on the magnitude of predictions may be found, which is the focus of the next section. We have used a tree size of 256 as large tree sizes are used in practice and provide better test errors. Furthermore, the phenomena we are studying is clearer for large tree sizes. In Section 3 we show results for both large trees and stumps. We note that base learners with real valued predictions were first considered by Schapire and Singer [15] that generalized the generalization bound of Schapire et al. [14] to work with real values but without otherwise changing the bound.

3 Refined margin bounds

Motivated by the empirical observations in the previous section, we prove a more refined margin based generalization bound for voting classifiers. Define from a voting classifier f the notation $\Delta(x, h) := |f(x) - h(x)|$. Intuitively, if a voting classifier f has a small margin on a training point x , but this is the result of using mostly base learners h that make small predictions (in absolute value), then $\Delta(x, h)$ will be small for most h in f . Also define from a voting classifier $f = \sum_h \alpha_h h$ the distribution $\mathcal{Q}(f)$ over base learners, which simply returns h with probability α_h . With this notation, our new generalization bound states that for any distribution \mathcal{D} over $\mathcal{X} \times \{-1, 1\}$ and for any margin θ , it holds with high probability over a set $S \sim \mathcal{D}^m$ that all voting classifiers f satisfy:

$$\mathcal{L}_{\mathcal{D}}(f) \leq \mathcal{L}_S^\theta(f) + O\left(\frac{N \lg |\mathcal{H}| \lg m}{m} + \sqrt{\mathcal{L}_S^\theta(f) \cdot \frac{N \lg |\mathcal{H}| \lg m}{m}}\right), \quad (5)$$

where $N = \max\{\theta^{-2} \cdot \left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(16m))/2}\right]\right)^{2/(\lg(16m))}, \theta^{-1}\}$.

Never worse. Comparing our bound to the k 'th margin bound (2), we see that (5) equals the k 'th margin bound when $N = \Theta(\theta^{-2})$. First, we argue that we always have $N = O(\theta^{-2})$, i.e. (5) is never worse than the k 'th margin bound. To see this, observe that $\Delta(x, h) \leq 2$ since all $h \in \mathcal{H}$ produce values in $[-1, 1]$. Thus, $\Delta(x, h)^2 \leq 4$ and $\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2] \leq 4$. This implies $\left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(16m))/2}\right]\right)^{2/\lg(16m)} \leq 4$, hence we always have $N = O(\theta^{-2})$.

Potentially much better. Next, we demonstrate that our new bound may be significantly better than previous generalization bounds for very natural voting classifiers. For any desired margin $\theta \in (0, 1]$, consider an example of a voting classifier $f(x) = \sum_{i=1}^{1/\theta} \theta h_i(x)$ such that for each training point (x, y) , there is exactly one hypothesis h_i with $h_i(x) = y$ and all others have $h_j(x) = 0$. This example is quite similar to the empirical performance of LightGBM seen in Section 2, where most hypotheses make small predictions on most training points. The voting classifier f has a margin of θ on all training points and thus the k 'th margin bound predicts a generalization error of $O(\lg |\mathcal{H}| \lg m / (m\theta^2))$ (since $\mathcal{L}_S^\theta(f) = 0$ when all points have margin θ). Let us now estimate N in (5). First, fix an $(x, y) \in S$ and consider the expression $\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(16m))/2} = \left(\sum_{i=1}^{1/\theta} \theta \cdot \Delta(x, h_i)^2\right)^{(\lg(16m))/2} = (\theta \cdot (1 - \theta)^2 + (1 - \theta)\theta^2)^{(\lg(16m))/2} < \theta^{(\lg(16m))/2}$. Since this holds for every (x, y) , we have $\left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(16m))/2}\right]\right)^{2/\lg(16m)} < \theta$. Plugging that into the definition of N , we see that $N \leq \max\{\theta^{-2} \cdot \theta, \theta^{-1}\} = \theta^{-1}$. That is, the dependency on the margin has improved by a factor θ and our new generalization bound predicts $\mathcal{L}_{\mathcal{D}}(f) = O(\lg |\mathcal{H}| \lg m / (m\theta))$.

Comparison to earlier work. In recent work, Cortes et al. [3], also proved refined generalization bounds for gradient boosters. Their work shows, that if the q -norm of the vector of leaf predictions for each tree trained by a gradient booster is small, then the trees have smaller VC-dimension and hence the voting classifier has better generalization performance (by using previous generalization bounds). Note that their bound only depends on the leaf predictions and does not take into account the number of training points in each leaf. Our experiment in Figure 3 shows that for each base learner, only a tiny fraction (about 1 in 5000) of training points end in a leaf with large prediction, which our bound takes into account.

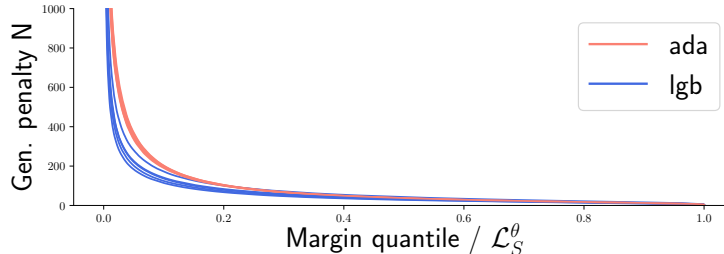


Figure 4: Generalization penalty N on the Forest Cover data set when choosing θ as the (pm) 'th smallest margin for $p \in [0, 1]$.

Table 1: Comparing AdaBoost with LightGBM. In this experiment the trees used as bare learners are of increasing size relative to the data size. Each value shown is the average over several runs and each run use 200 rounds of boosting. Moment is $\left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(m))/2} \right]\right)^{2/\lg(m)}$.

Data Set	Alg.	Train Err	Test Err	Mean Margin	Max Depth	Mean Depth	Moment
Forest	ada	0.0001	0.0331	0.1696	22.0	12.4	0.969
	lgb	0.0002	0.0291	0.0280	23.7	13.9	0.025
Boone	ada	0.00009	0.0589	0.311	17.5	10.2	0.917
	lgb	0.00009	0.0552	0.0818	17.6	10.4	0.0564
Higgs	ada	0.178	0.277	0.0747	24.9	13.5	0.99
	lgb	0.185	0.251	0.018	26	14.7	0.0289
Diabetes	ada	0	0.268	0.148	3.5	2.63	0.973
	lgb	0.0264	0.26	0.142	3.5	2.63	0.214

Empirical evaluation. Our new generalization bound carefully takes the magnitude of predictions made by the base learners into account, thus there is hope that (5) may better explain the experiments in the previous section. To test this, we have run the experiments again, this time plotting the value of N as a function of $p = \mathcal{L}_S^\theta(f)$. That is, we notice that for two voting classifiers produced by AdaBoost and LightGBM, respectively, the only thing that varies in (5) when choosing the (pm) 'th smallest margin, i.e. $p = \mathcal{L}_S^\theta(f)$, is the value of N . Thus smaller values of N imply better generalization according to the theory. Figure 4 shows the result of the experiment. Quite remarkably, the relative ordering of AdaBoost and LightGBM match the observed test errors from Figure 1a much better, i.e. LightGBM slightly outperforms AdaBoost. We have repeated the same experiment on more data sets and summarized the results in Table 1. The parameters for the experiments are shown in Table 2

Table 2: Data sets, all freely available, and parameters considered in the experiments. LR means learning rate as used in LightGBM. For each experiment we randomly split the data set in half to get a training set and a test set of equal size. For the Higgs dataset of size 11 million, we sample a subset of 2 million data points that we randomly split evenly into train and test set. For Forest Cover only the first two classes are used to make it into a binary classification problem.

Data Set	Data Size	Tree Size	LR	Stumps LR	Runs
Diabetes	768	5	0.1	0.1	100
Boone	65032	96	0.2	0.6	5
Forest Cover	495141	256	0.3	0.3	5
Higgs	2000000	512	0.3	0.3	5

In all experiments, the margin distribution, here represented by the mean margin, is much worse for the LightGBM classifier, while the height of the trees used, both the max height and the mean height, is larger. Still the LightGBM classifier generalizes at least as well (in fact, slightly better) than the AdaBoost classifier. Table 1 also shows that the moment value from our generalization bound is significantly better for the LightGBM classifier. When we consider our new generalization bound, the theory nicely matches the observed test errors in the same way as was shown in Figure 4 for all data sets. While not final proof that this is the real or only explanation, it suggests that the success of gradient boosters, despite having poor margins, may be explained by the many small predictions made by the base learner trees. The standard deviations of the test statistics are left out since they are extremely small for the three large data sets (and we have run 100 iterations of the small Diabetes data set). For completeness we have included the same experiment replacing the large trees with stumps and shown the results in Table 3. The results for stumps match those from the larger trees, just with a smaller difference in margins and moment values.

Table 3: Experiments with stumps as base learners. Same setup as in Table 1.

Data Set	Alg.	Train Err	Test Err	Mean Margin	Moment
Forest	ada	0.223	0.224	0.0754	0.987
	lgb	0.217	0.218	0.0225	0.0986
Boone	ada	0.0781	0.0817	0.138	0.975
	lgb	0.0669	0.0744	0.0422	0.239
Higgs	ada	0.309	0.31	0.059	0.986
	lgb	0.301	0.302	0.0309	0.329
Diabetes	ada	0.161	0.246	0.108	0.976
	lgb	0.176	0.238	0.138	0.299

4 Generalization Bound Proof

This section is devoted to the proof of our refined margin based generalization bound for voting classifiers, presented hereafter as Theorem 1. First we recollect some notation. Let \mathcal{X} be some ground set, \mathcal{D} a distribution over $\mathcal{X} \times [-1, 1]$, $\mathcal{H} \subseteq \mathcal{X} \rightarrow [-1, 1]$, and $\mathcal{C} = \mathcal{C}(\mathcal{H})$ be the convex hull of \mathcal{H} . Fix a voting classifier f , then there exists a sequence $(\alpha_h)_{h \in \mathcal{H}} \in \mathbb{R}_+^{\mathcal{H}}$ such that $\sum_{h \in \mathcal{H}} \alpha_h = 1$ and $f = \sum_{h \in \mathcal{H}} \alpha_h \cdot h$. Thus f implicitly defines a distribution $\mathcal{Q} = \mathcal{Q}(f)$ over \mathcal{H} , where $\Pr_{h \sim \mathcal{Q}}[h = h'] = \alpha_{h'}$ for all $h' \in \mathcal{H}$. Finally, let $\Delta : \mathcal{X} \times \mathcal{H} \rightarrow \mathbb{R}$ be defined by $\Delta(x, h) := |f(x) - h(x)|$ for every $x \in \mathcal{X}$, $h \in \mathcal{H}$. We show the following.

Theorem 1. *Let \mathcal{D} be a distribution over $\mathcal{X} \times \{-1, 1\}$ where \mathcal{X} is some ground set, and let $\mathcal{H} \subseteq \mathcal{X} \rightarrow [-1, 1]$. For every $\delta > 0$, it holds with probability at least $1 - \delta$ over a set of m samples $S \sim \mathcal{D}^m$, that for every voting classifier $f \in \mathcal{C}(\mathcal{H})$ and every margin $\theta > 0$, we have*

$$\mathcal{L}_{\mathcal{D}}(f) \leq \mathcal{L}_S^\theta(f) + O\left(\frac{N \lg |\mathcal{H}| + \lg(1/\delta)}{m} + \sqrt{\frac{N \lg |\mathcal{H}| + \lg(1/\delta)}{m} \mathcal{L}_S^\theta(f)}\right), \quad (6)$$

where $N = O\left(\max\{\theta^{-2} \cdot \left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}(f)} [\Delta(x, h)^2]^{(\lg(16m))/2}\right]\right)^{2/(\lg(16m))}, \theta^{-1}\} \lg m\right)$.

Denote by $\mathcal{E} = \mathcal{E}(\delta)$ the event that for every voting classifier f and every margin $\theta > 0$, the bound in (6) holds with N as defined in Theorem 1. In these notations we prove that $\Pr_{S \sim \mathcal{D}^m}[\mathcal{E}] \geq 1 - \delta$.

Proof overview. Inspired by techniques presented by Schapire *et al.* [14] and employed by Gao and Zhou [6], our proof incorporates a discretization of the set of all voting classifiers over \mathcal{H} to a discrete *net* of classifiers, such that, loosely speaking, every voting classifier over \mathcal{H} can be approximated by a classifier that belongs to the net, and in addition, the size of the net is not too big, and thus union bounding over the net yields the desired probability bounds. Thus, intuitively speaking, by randomly rounding every voting classifier f to the net we get an upper bound on the out of sample error for f . More specifically, $N \in \mathbb{N}^+$ be some positive integer. We define a net \mathcal{C}_N of voting classifier

by $\mathcal{C}_N := \left\{ \frac{1}{N} \sum_{j \in [N]} h_j : \langle h_j \rangle_{j \in [N]} \in \mathcal{H}^N \right\}$. For every voting classifier f over \mathcal{H} , we then give a randomized rounding scheme that essentially associates a random net element $g \in \mathcal{C}_N$ with f , and show that with high probability the out of sample error with respect to g well-approximates that of f . By choosing N carefully and union bounding over \mathcal{C}_N we get an upper bound on the out of sample error for all voting classifiers f . The crux of the proof lies in carefully choosing the size of the net, namely N . Loosely speaking, the net size N has to be large enough, so that the net is rich enough to approximate every voting classifier well, but on the other hand small enough, so that union bounding over the net does not incur too large a cost for the probability bound. By subtly choosing N and proving refined bounds on the rounding scheme we get the bound in Theorem 1.

Formally we define for every $N \in \mathbb{N}^+$, the event \mathcal{E}_N to be the set of all samples $S \in (\mathcal{X} \times \{-1, 1\})^m$ satisfying that for all voting classifiers $g \in \mathcal{C}_N$ and integer $\ell \in [0, N]$ it holds that

$$\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \leq \mathcal{L}_S^{\ell/N}(g) + \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} + 4\sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}(g);$$

and

$$\Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] \leq 2 \Pr_{(x,y) \sim S} [|f(x) - g(x)| > \ell/N] + \frac{8 \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m}.$$

Intuitively speaking, for $S \in \mathcal{E}_N$, the first bound ensures a good generalization bound for every voting classifier g in the net, whereas the second bound shows that g approximates f over \mathcal{D} almost as well as it approximates f over S . In turn these two bounds imply that the behavior of f, g over S predicts their behavior over \mathcal{D} . As $\sum_{N=1}^{\infty} \frac{1}{N(N+1)} = 1$, the following lemma implies Theorem 1 by applying a union bound.

Lemma 2. For every $N \in \mathbb{N}^+$ we have $\Pr_{S \sim \mathcal{D}^m} [\mathcal{E}_N] \geq 1 - \frac{\delta}{N(N+1)}$, and moreover, $\bigcap_{N \in \mathbb{N}^+} \mathcal{E}_N \subseteq \mathcal{E}$.

The proof of the lemma is quite involved technically, and most of the proof is thus deferred to the appendix. Our main novelty lies in showing that for our choice of $N = N(f, \theta)$, for every sample set $S \in \text{supp}(\mathcal{D}^m)$, with very high probability over the choice of a point $x \in \mathcal{X}$ and a net-classifier $g \in \mathcal{C}_N$, g approximates f . In turn, this implies that if $S \in \bigcap_{N \in \mathbb{N}^+} \mathcal{E}_N$, then for every voting classifier f and $\theta > 0$, f is well-approximated by a randomized rounding to the net \mathcal{C}_N . Formally we show the following for every f and θ .

Lemma 3. $\Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [\Delta(x, g) > 49\theta/100] \leq \frac{1}{m^2}$, where

$$N = N(f, \theta) := \lg(16m) \cdot \max\{256\theta^{-1} \|\Delta(x, h)\|_{\lg(16m)}, 100/\theta, 128e\theta^{-2} \cdot \left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{(\lg(16m))/2} \right] \right)^{2/(\lg(16m))}\}.$$

Proof. Let $Z = \Delta(x, g)$, then for every integer $r \geq 1$ we conclude from Markov's inequality that

$$\Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [Z > 49\theta/100] = \Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [Z^r > (49\theta/100)^r] \leq \left(\frac{100}{49\theta} \right)^r \|Z\|_r^r. \quad (7)$$

It is therefore enough to show $\|Z\|_r^r \leq \left(\frac{49\theta}{100} \right)^r m^{-2}$ for some positive integer $r \geq 1$. Let $r = 2 \cdot \lceil \lg(4m)/2 \rceil$, then r is an even integer, satisfying $\lg(4m) = 2 \lg(4m)/2 \leq r \leq \lg(4m) + 2 \leq N$. Since r is even, then for $g = \frac{1}{N} \sum_{j \in [N]} h_j$ we get that

$$Z^r = Z(x, g)^r = \left(\frac{1}{N} \sum_{j \in [N]} (f(x) - h_j(x)) \right)^r = \frac{1}{N^r} \sum_{T=(j_i)_{i \in [r]} \in [N]^r} \prod_{i \in [r]} (f(x) - h_{j_i}(x)).$$

For every $T = (j_i)_{i \in [r]} \in [N]^r$ let $D(T) := \{j \in [N] : \exists i \in [r]. j_i = j\}$ be the set of distinct indices occurring in T , and for every $j \in [N]$, let $c_T(j) := |\{i \in [r] : j_i = j\}|$ be the number of times j occurs in T . Then in these notations we have

$$Z^r = \frac{1}{N^r} \sum_{T \in [N]^r} \prod_{j \in D(T)} (f(x) - h_j(x))^{c_T(j)}.$$

As h_1, \dots, h_N are chosen independently, we get that

$$\mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] = \frac{1}{N^r} \sum_{T \in [N]^r} \prod_{j \in D(T)} \mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} \left[(f(x) - h_j(x))^{c_T(j)} \right].$$

Let $T \in [N]^r$, and assume that for some $j \in D(T)$ we have $c_T(j) = 1$, then

$$\mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} \left[(f(x) - h_j(x))^{c_T(j)} \right] = \mathbb{E}_{h \sim \mathcal{Q}} [f(x) - h(x)] = f(x) - \mathbb{E}_{h \sim \mathcal{Q}} [h(x)] = f(x) - \sum_{h \in \mathcal{H}} \alpha_h h(x) = 0,$$

Denote $\mathcal{T} := \{T \in [N]^r : \forall j \in D(T). c_T(j) > 1\}$, then

$$\begin{aligned} \mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] &= \frac{1}{N^r} \sum_{T \in \mathcal{T}} \prod_{j \in D(T)} \mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} \left[(f(x) - h_j(x))^{c_T(j)} \right] \\ &= \frac{1}{N^r} \sum_{T \in \mathcal{T}} \prod_{j \in D(T)} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^{c_T(j)} \right]. \end{aligned} \quad (8)$$

By Lyapunov's Theorem (see, e.g. [11]), $\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^\xi]$ is logarithmic convex for $\xi \in [1, +\infty)$, and as $c_T(j) \geq 2$ for all $j \in D(T)$ we get that

$$\prod_{j \in D(T)} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^{c_T(j)} \right] \leq \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^2 \right]^{|D(T)|-1} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^{r-2|D(T)|+2} \right].$$

Plugging into (8) we get that

$$\mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] \leq \frac{1}{N^r} \sum_{T \in \mathcal{T}} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^2 \right]^{|D(T)|-1} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^{r-2|D(T)|+2} \right]. \quad (9)$$

For every $d \in \mathbb{N}$ denote $\mathcal{T}_d := \{T \in \mathcal{T} : |D(T)| = d\}$. Since for every $T \in \mathcal{T}$ and every $j \in D(T)$, we know that $c_T(j) \geq 2$, then for every $d > r/2$ we get that $\mathcal{T}_d = \emptyset$. Therefore $\mathcal{T} = \bigcup_{d \in [r/2]} |\mathcal{T}_d|$. Moreover, for every $d \in [r/2]$ and every $T \in \mathcal{T}_d$, we have

$$\mathbb{E}_{h \sim \mathcal{Q}} \left[|h(x)|^2 \right]^{|D(T)|-1} \mathbb{E}_{h \sim \mathcal{Q}} \left[|h(x)|^{r-2|D(T)|+2} \right] = \mathbb{E}_{h \sim \mathcal{Q}} \left[|h(x)|^2 \right]^{d-1} \mathbb{E}_{h \sim \mathcal{Q}} \left[|h(x)|^{r-2d+2} \right].$$

We therefore refine (9) to get

$$\mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] \leq \frac{1}{N^r} \sum_{d \in [r/2]} |\mathcal{T}_d| \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^2 \right]^{d-1} \mathbb{E}_{h \sim \mathcal{Q}} \left[\Delta(x, h)^{r-2d+2} \right]. \quad (10)$$

Claim 4. For every $d \in [r/2]$, $|\mathcal{T}_d| \leq r^r \sqrt{2e\pi r} \left(\frac{Ne}{r}\right)^d$.

Proof. Fix some $d \in [r/2]$. There are at most $\binom{N}{d}$ ways to choose a subset $Y \subseteq [N]$ such that $|Y| = d$. Once such a set Y is fixed, there are at most $\binom{d+(r-2d)-1}{r-2d}$ solution to the equation $\sum_{j \in Y} y_j = r$ under the constraint that $y_j \in \mathbb{N} \setminus \{0, 1\}$ for all $j \in Y$. Moreover, once $\{y_j\}_{j \in Y}$ is fixed, there are $r! \cdot \prod_{j \in Y} (y_j!)^{-1}$ ways to form a sequence T satisfying that $D(T) = Y$, $c_T(j) = y_j$ for all $j \in Y$ and $c_T(j) = 0$ otherwise. Note that $\prod_{j \in Y} (y_j!) \geq ((r/d)!)^d$ for every choice of $\{y_j\}_{j \in Y}$, and therefore

$$\begin{aligned} |\mathcal{T}_d| &\leq \binom{N}{d} \cdot \binom{r-d-1}{r-2d} \cdot \frac{r!}{((r/d)!)^d} \leq \left(\frac{Ne}{d}\right)^d \cdot 2^{r-d} \cdot \frac{\sqrt{2e\pi r} (r/e)^r}{(\sqrt{2\pi} (r/d) (r/(ed)))^{r/d}} \\ &\leq \sqrt{2e\pi r} (Ne)^d \cdot r^{r-d} \leq r^r \sqrt{2e\pi r} \left(\frac{Ne}{r}\right)^d \end{aligned}$$

□

Plugging into (10) we conclude that

$$\begin{aligned} \mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] &\leq \frac{1}{N^r} \sum_{d \in [r/2]} r^r \sqrt{2e\pi r} \left(\frac{Ne}{r}\right)^d \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{d-1} \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^{r-2d+2}] \\ &= \sqrt{2e\pi r} \left(\frac{r}{N}\right)^r \sum_{d \in [r/2]} \left(\frac{Ne}{r}\right)^d \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{d-1} \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^{r-2d+2}] \end{aligned}$$

As $\left(\frac{Ne}{r}\right)^\xi$, $\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{\xi-1}$, $\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^{r-2\xi+2}]$ are all logarithmic convex for $\xi \in [1, r/2]$, their product is also logarithmic convex over that range, and thus gets its maximum on either 1 or $r/2$. Concluding we get that

$$\mathbb{E}_{(h_k)_{k \in [N]} \sim \mathcal{Q}^N} [Z^r] \leq \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{r}{N}\right)^r \left(\left(\frac{Ne}{r}\right) \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^r] + \left(\frac{Ne}{r}\right)^{r/2} \mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right).$$

Taking the expectation over $(x, y) \sim \mathcal{D}$ gives

$$\|Z\|_r^r \leq \frac{r}{2} \sqrt{2e\pi r} \left(\frac{r}{N}\right)^r \left(\left(\frac{Ne}{r}\right) \|\Delta(x, h)\|_r^r + \left(\frac{Ne}{r}\right)^{r/2} \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right] \right) \quad (11)$$

To finish the proof of Lemma 3, we show that our bound on N implies that $\|Z\|_r^r \leq \left(\frac{49\theta}{100}\right)^r m^{-2}$. Denote

$$\begin{aligned} \Psi_1 &= \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{r}{N}\right)^r \cdot \left(\frac{Ne}{r}\right) \|\Delta(x, h)\|_r^r = \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{r \|\Delta(x, h)\|_r}{N}\right)^r \cdot \left(\frac{Ne}{r}\right) \\ \Psi_2 &= \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{r}{N}\right)^r \left(\left(\frac{Ne}{r}\right)^{r/2} \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right] \right) \end{aligned}$$

Plugging into (11) we get that $\|Z\|_r^r \leq \Psi_1 + \Psi_2$.

We will show that $\max\{\Psi_1, \Psi_2\} \leq \left(\frac{49\theta}{100}\right)^r \cdot \frac{1}{2m^2}$, which proves the claim. To bound Ψ_1 , note first that Ψ_1 decreases as a function of N (since $r \geq 2$). Since $N \geq 256\theta^{-1} \lg(16m) \cdot \|\Delta(x, h)\|_{\lg(16m)}$ we get that

$$\Psi_1 \leq \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{r \cdot \|\Delta(x, h)\|_r}{256\theta^{-1} \lg(16m) \cdot \|\Delta(x, h)\|_{\lg(16m)}} \right)^r \cdot \left(\frac{256\theta^{-1} \lg(16m) \cdot \|\Delta(x, h)\|_{\lg(16m)} \cdot e}{r} \right)$$

Since $r < \lg(16m)$, and by monotonicity of norms, $\|\Delta(x, h)\|_r \leq \|\Delta(x, h)\|_{\lg(16m)} \leq 2$, where the last inequality is due to the fact that $|f(x) - h(x)| \leq 2$ for all $h \in \mathcal{H}$, $x \in \mathcal{X}$. Moreover, $\lg(4m) \leq r \leq \lg(16m) \leq 2(\lg(4m))$, therefore

$$\begin{aligned} \Psi_1 &\leq \frac{r}{2} \cdot \sqrt{2e\pi r} \left(\frac{\theta}{256}\right)^r \cdot 1024e\theta^{-1} \\ &\leq \left(\frac{49\theta}{100}\right)^r \cdot 3r^{3/2} 125^{-r} \cdot (1024e\theta^{-1}) \leq \left(\frac{49\theta}{100}\right)^r \cdot 3r^{3/2} 64^{-\lg m} 125^{-2} \cdot (1024e\theta^{-1}) \\ &\leq \left(\frac{49\theta}{100}\right)^r \cdot \frac{1}{5} \lg^{3/2}(4m) \cdot m^{-6}\theta^{-1} \leq \left(\frac{49\theta}{100}\right)^r \cdot \frac{1}{2m^2} \cdot \frac{1}{2} (\lg(4m)/m)^{3/2} (m^{5/2}\theta)^{-1} \end{aligned}$$

For large enough m , we have that $\lg(4m)/m \leq 5/8$, and therefore $(\lg(4m)/m)^{3/2} \leq 1/2$. Since $\theta \geq 1/m$ we get that $\Psi_1 \leq \left(\frac{49\theta}{100}\right)^r \cdot \frac{1}{2m^2}$. We now turn to bound Ψ_2 . Recall that $N \geq 128e\theta^{-2} \lg(16m) \cdot$

$\left(\mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{\lg(16m)/2} \right]\right)^{2/\lg(16m)}$, and therefore

$$\begin{aligned} \Psi_2 &\leq 3r^{3/2} \left(\frac{er \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right]^{2/r}}{128e\theta^{-2} \lg(16m) \left(\mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{\lg(16m)/2} \right]\right)^{2/\lg(16m)}} \right)^{r/2} \\ &\leq \left(\frac{49\theta}{100}\right)^r \cdot 3r^{3/2} \left(\frac{r \mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right]^{2/r}}{30 \lg(16m) \left(\mathbb{E}_{(x, y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{\lg(16m)/2} \right]\right)^{2/\lg(16m)}} \right)^{r/2} \end{aligned}$$

Since $r < \log(16m)$, and by monotonicity of norms of random variables, we get that

$$\mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{r/2} \right]^{2/r} \leq \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{\log(16m)/2} \right]^{2/\log(16m)}.$$

Therefore

$$\Psi_2 \leq \left(\frac{49\theta}{100} \right)^r \cdot 3r^{3/2} (30)^{-r/2} \leq \left(\frac{49\theta}{100} \right)^r \cdot 3r^{3/2} (30)^{-(\lg m)/2-1} \leq \left(\frac{49\theta}{100} \right)^r \cdot \frac{1}{2m^2} \cdot \frac{1}{5} r^{3/2} m^{-2/5}$$

Similarly to before, for large enough m , $\lg^{3/2}(4m) \cdot m^{-2/5} \leq 5$, and therefore we conclude that $\Psi_2 \leq \left(\frac{49\theta}{100} \right)^r \cdot \frac{1}{2m^2}$, which completes the proof of the lemma. \square

5 Generalization lower bound

In this section we state and prove our new generalization lower bound, presented as Theorem 5.

Theorem 5. *For every large enough integer N , every $\theta \in (1/N, 1/40)$, $\tau \in [0, 1]$ and every $(\theta^{-2} \ln N)^{1+\Omega(1)} \leq m \leq 2^{N^{O(1)}}$, if $\frac{\ln N \ln m}{m\theta^2} \leq \tau \leq 1$, then there exist a set \mathcal{X} , a hypothesis set \mathcal{H} over \mathcal{X} and a distribution \mathcal{D} over $\mathcal{X} \times \{-1, 1\}$ such that $\ln |\mathcal{H}| = \Theta(\ln N)$ and with probability at least $1/100$ over the choice of samples $S \sim \mathcal{D}^m$ there exists a voting classifier $f_S \in C(\mathcal{H})$ such that*

1. $\mathcal{L}_S^\theta(f_S) \leq \tau$; and
2. $\mathcal{L}_D(f_S) \geq \mathcal{L}_S^\theta(f_S) + \Omega\left(\frac{\ln |\mathcal{H}| \ln m}{m\theta^2} + \sqrt{\tau \ln(\tau^{-1})} \cdot \frac{\ln |\mathcal{H}|}{m\theta^2}\right)$.

Our proof makes use of the lemma, whose proof can be found in [8].

Lemma 6. *For every $\theta \in (0, 1/40)$, $\delta \in (0, 1)$ and integers $d \leq u$, there exists a distribution $\mu = \mu(u, d, \theta, \delta)$ over hypothesis sets $\mathcal{H} \subset \mathcal{X} \rightarrow \{-1, 1\}$, where \mathcal{X} is a set of size u , such that the following holds for $N = \Theta\left(\theta^{-2} \ln d \ln(\theta^{-2} d \delta^{-1}) e^{\Theta(\theta^2 d)}\right)$.*

1. For all $\mathcal{H} \in \text{supp}(\mu)$, we have $|\mathcal{H}| = N$; and
2. For every labeling $\ell \in \{-1, +1\}^u$, if no more than d points $x \in \mathcal{X}$ satisfy $\ell(x) = -1$, then

$$\Pr_{\mathcal{H} \sim \mu} [\exists f \in C(\mathcal{H}) : \forall x \in \mathcal{X}. \ell(x)f(x) \geq \theta] \geq 1 - \delta,$$

We start by describing the outlines of the proofs. To this end fix some integer N , and fix $\theta \in (1/N, 1/40)$. Let u be an integer, and let $\mathcal{X} = \{\xi_1, \dots, \xi_u\}$ be some set with u elements. The distribution \mathcal{D} over $\mathcal{X} \times \{-1, 1\}$, is simply the uniform distribution over $\mathcal{X} \times \{1\}$. That is for every $i \in [u]$ and $y \in \{-1, 1\}$, $\Pr_{\mathcal{D}}[(\xi_i, y)] = \frac{1+y}{2u}$. The following claim is straightforward.

Claim 7. *For every $f : \mathcal{X} \rightarrow \mathbb{R}$ we have $\Pr_{(x,y) \sim \mathcal{D}} [yf(x) < 0] = \frac{1}{u} \sum_{i \in [u]} \mathbb{1}_{f(\xi_i) < 0}$.*

We will show that with some constant probability over a random choice $S \sim \mathcal{D}^m$, an adversarial voting classifier has a high generalization probability. We additionally show existence of a hypothesis set $\hat{\mathcal{H}}$ such that with very high (constant) probability over a random choice of $\ell \in \{-1, 1\}^u$, $C(\hat{\mathcal{H}})$ contains a voting classifier that attains high margins with ℓ over the entire set \mathcal{X} . Finally, we conclude that with positive probability over a random choice of $S \sim \mathcal{D}^m$ both properties are satisfied.

To prove existence of a ‘‘rich’’ yet small enough hypothesis set $\hat{\mathcal{H}}$ we apply Lemma 6 together with Yao’s minimax principle. In order to ensure that the hypothesis sets constructed using Lemma 6 is small enough, and specifically has size $N^{O(1)}$, we need to focus our attention on sparse labelings $\ell \in \{-1, 1\}^u$ only. That is, the labelings cannot contain more than $\frac{\ln N}{\theta^2}$ entries equal to -1 . To this end we will focus on d -sparse vectors. More formally, we define a set of labelings of interest $\mathcal{L}(u, d)$ as follows.

$$\mathcal{L}(u, d) := \{\ell \in \{-1, 1\}^u : |\{i \in [u] : \ell_i = -1\}| \leq d\}. \quad (12)$$

We next show that there exists a small enough (with respect to N) hypothesis set $\hat{\mathcal{H}}$ that is rich enough. That is, with high probability over $\ell \in \mathcal{L}(u, d)$, there exists a voting classifier $f \in C(\hat{\mathcal{H}})$ that attains high minimum margin with ℓ over the entire set \mathcal{X} . Note that the following result, similarly to Lemma 6 does not depend on the size of \mathcal{X} , but only on the sparsity of the labelings in question.

Claim 8. *If $u \leq 2^{N^{O(1)}}$ and $d \leq \frac{\ln N}{\theta^2}$ then there exists a hypothesis set $\hat{\mathcal{H}}$ such that $\ln |\hat{\mathcal{H}}| = \Theta(\ln N)$ and*

$$\Pr_{\ell \in_R \mathcal{L}(u, d)} [\exists f \in C(\hat{\mathcal{H}}) : \forall i \in [u]. \ell_i f(\xi_i) \geq \theta] \geq 1 - 1/N.$$

Proof. Let $\mu = \mu(u, d, \theta, 1/N)$, be the distribution whose existence is guaranteed in Lemma 6. Then for every labeling $\ell \in \mathcal{L}(u, d)$, with probability at least 99/100 over $\mathcal{H} \sim \mu$, there exists a voting classifier $f \in C(\mathcal{H})$ that has minimal margin of θ . That is, for every $i \in [u]$, $\ell_i f(\xi_i) \geq \theta$. By Yao's minimax principle, there exists a hypothesis set $\hat{\mathcal{H}} \in \text{supp}(\mu)$ such that

$$\Pr_{\ell \in_R \mathcal{L}(u, d)} [\exists f \in C(\hat{\mathcal{H}}) : \forall i \in [u]. \ell_i f(x_i) \geq \theta] \geq 1 - 1/N.$$

Moreover, since $\hat{\mathcal{H}} \in \text{supp}(\mu)$, then $|\hat{\mathcal{H}}| = \Theta\left(\theta^{-2} \ln u \cdot \ln(N\theta^{-2} \ln u) \cdot e^{\Theta(\theta^2 d)}\right)$. Since $\theta \geq 1/N$, $\ln u \leq N^{O(1)}$, and $d \leq \frac{\ln N}{\theta^2}$, and thus $e^{\theta^2 d} = N$ we get that there exists some universal constant $C > 0$ such that $|\hat{\mathcal{H}}| = \Theta(N^C)$, and thus $\ln |\hat{\mathcal{H}}| = \Theta(\ln N)$. \square

Let $u = \frac{\ln N}{16\tau\theta^2}$, and let $d = \frac{\ln N}{16e^{28}\theta^2}$. We next introduce some notation. With every set $T \subseteq [u]$ we associate the classifier $h_T : \mathcal{X} \rightarrow \{-1, 1\}$ satisfying that for every $x \in \mathcal{X}$, $h_T(x) = -1$ if and only if $x \in T$. For every m -point sample $S \in (\mathcal{X} \times \{1\})^m$ and every $i \in [u]$, let b_i^S be the number of times ξ_i is sampled into S . If the set S is clear from context, we simply denote b_i . In these notations, $\mathcal{L}_S(h_T) = \frac{1}{m} \sum_{i \in T} b_i^S$ for every $T \subseteq [u]$. Given a sample set S Let $T^* = T^*(S) \subseteq [u]$ be a random set of size d that minimizes $\mathcal{L}_S(h_{T^*(S)}) = \sum_{i \in T^*(S)} b_i^S$. We will show the following.

Lemma 9. *With probability at least 1/100 over the choice of sample $S \sim \mathcal{D}^m$, the following holds.*

1. *There exists a voting classifier $f_S \in C(\hat{\mathcal{H}})$ such that $f_S(\xi_i)h_{T^*(S)}(\xi_i) \geq \theta$ for all $i \in [u]$, and*
2. $\mathcal{L}_S(h_{T^*(S)}) \leq \frac{d}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right)$.

Note that as $\tau \geq \frac{\ln N \ln m}{m\theta^2}$ we know that $u = \frac{\ln N}{16\tau\theta^2} \leq \frac{m}{16 \ln m}$ and therefore $\frac{\ln(u/2d)}{9m/u} \leq \frac{u \ln(e^{28}/\tau)}{9m} \leq \frac{\ln(e^{28}/\tau)}{144 \ln m} \leq \frac{1}{2}$ for large enough N , and therefore the bound in the second part of Lemma 9 is meaningful. We first show that the lemma implies Theorem 5.

Proof of Theorem 5. Fix some $\frac{\ln N \ln m}{m\theta^2} \leq \tau \leq 1$. From Lemma 9 with probability 1/100 over the choice of a sample $S \sim \mathcal{D}^m$ there exists a voting classifier $f_S \in C(\hat{\mathcal{H}})$ such that $f_S(\xi_i)h_{T^*(S)}(\xi_i) \geq \theta$ for all $i \in [u]$ and moreover $\mathcal{L}_S(h_{T^*(S)}) \leq \tau$. Consider f_S , and note first that

$$\mathcal{L}_{\mathcal{D}}(f_S) = \frac{1}{u} \sum_{i \in [u]} \mathbb{1}_{f_S(\xi_i) < 0} = \frac{1}{u} \sum_{i \in [u]} \mathbb{1}_{h_{T^*(S)}(\xi_i) < 0} = \frac{|T^*(S)|}{u} = \frac{d}{u}.$$

Additionally, since for every $i \in [u]$, $f_S(\xi_i) \leq 0$ if and only if $f_S(\xi_i) \leq \theta$, then

$$\mathcal{L}_S^\theta(f_S) = \mathcal{L}_S(f_S) = \mathcal{L}_S(h_{T^*(S)}) \leq \frac{d}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right) \leq \frac{d}{2u} \leq \tau.$$

Summing up we get also that

$$\mathcal{L}_{\mathcal{D}}(f_S) - \mathcal{L}_S^\theta(f_S) \geq \frac{d}{u} \sqrt{\frac{\ln(u/2d)}{9m/u}} = \Omega\left(\tau \sqrt{\frac{u \ln(\tau^{-1})}{m}}\right) = \Omega\left(\sqrt{\frac{\ln N \tau \ln(\tau^{-1})}{m\theta^2}}\right).$$

\square

For the rest of the section we therefore prove Lemma 9. First note that since \mathcal{D} is uniform over $\mathcal{X} \times \{1\}$, and since given $S \sim \mathcal{D}^m$, T^* is sampled uniformly over all subsets $T \in \binom{[u]}{d}$ such that the sum $\sum_{i \in T} b_i^S$ is minimized, we get that for every $T \in \binom{[u]}{d}$, $\Pr_{S \sim \mathcal{D}^m}[T^*(S) = T] = \binom{u}{d}^{-1}$. In other words, for every $h \in \mathcal{L}(u, d)$, $\Pr_{S \sim \mathcal{D}^m}[h_{T^*(S)} = h] = \binom{u}{d}^{-1}$. Therefore $h_{T^*(S)}$ is uniformly distributed over $\mathcal{L}(u, d)$. From claim 8 it follows that for large enough N , the probability over the choice of $S \sim \mathcal{D}^m$ that there exists $f_S \in C(\hat{\mathcal{H}})$ such that $f_S(\xi)h_{T^*(S)}(\xi_i) \geq \theta$ for all $i \in [u]$ is at least 99/100. In order to prove Lemma 9, it is therefore enough to show that with probability at least 1/50 over the choice of $S \sim \mathcal{D}^m$, $\mathcal{L}_S(h_{T^*(S)}) \leq \frac{d}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right)$. We will show that with probability at least 1/50 over the choice of S there exist $i_1, \dots, i_d \in [u]$ such that for every $j \in [d]$, $b_{i_j}^S \leq \frac{m}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right)$. Since $T^*(S)$ minimizes $\sum_{i \in T^*(S)} b_i^S$, it follows that

$$\mathcal{L}_S(h_{T^*(S)}) = \frac{1}{m} \sum_{i \in T^*(S)} b_i^S \leq \frac{1}{m} \sum_{j \in [d]} b_{i_j}^S \leq \frac{d}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right).$$

To this end, fix some $i \in [u]$. For every $j \in [m]$, let I_j^S be an indicator for the event that the j th element selected into S is $(\xi_i, 1)$. Then $b_i^S = \sum_{j \in [m]} I_j^S$, and as \mathcal{D} is uniform, we get that $\mathbb{E}[b_i^S] = \sum_{j \in [m]} \mathbb{E}[I_j^S] = m/u$. We will use the following reverse Chernoff bound and show that with good enough probability, b_i^S is far from its expectation.

Lemma 10. *Let $m \in \mathbb{N}^+$ and let I_1, \dots, I_m be independent indicator random variables with success probability $1/u$. Then for every $\sqrt{3/(m/u)} \leq \delta \leq 1/2$ we have*

$$\Pr \left[\sum_{j \in [m]} I_j \leq (1 - \delta)mp \right] \geq e^{-9m\delta^2/u}.$$

Denote $\delta := \sqrt{\frac{\ln(u/2d)}{9m/u}}$. As we have shown earlier, $\delta \leq 1/2$. Moreover, since $\frac{u}{2d} \geq e^{27}\tau^{-1} \geq e^{27}$, we get that $\delta \geq \sqrt{\frac{27}{9m/u}} = \sqrt{\frac{3}{m/u}}$. We can therefore conclude from Lemma 10 that

$$\Pr[b_i^S \geq (1 - \delta)m/u] \geq e^{-9m\delta^2/u} = e^{-\ln(u/2d)} = \frac{2d}{u}.$$

Let B_i^S be the indicator for the event $b_i^S \geq (1 - \delta)m/u$, then $\mathbb{E}[B_i^S] \geq \frac{2d}{u}$. Finally, let $B^S = \sum_{i \in [u]} B_i^S$, then $\mathbb{E}[B^S] \geq 2d$. We will show that with probability at least 1/8 $\geq 1/50$ we have $B^S \geq d$. This implies that there exist i_1, \dots, i_d such that for every $j \in [d]$, $b_{i_j}^S \leq \frac{m}{u} (1 - \delta) = \frac{m}{u} \left(1 - \sqrt{\frac{\ln(u/2d)}{9m/u}}\right)$. To show $B^S \geq d$ with reasonable probability, we use the Paley-Zigmond inequality.

$$\Pr[B^S \geq d] = \Pr \left[B^S \geq \frac{1}{2} \mathbb{E}[B^S] \right] \geq \frac{\mathbb{E}[B^S]^2}{4\mathbb{E}[(B^S)^2]}.$$

Since B_1^S, \dots, B_u^S are negatively correlated, we have that $\mathbb{E}[B_i^S B_j^S] \leq \mathbb{E}[B_i^S] \mathbb{E}[B_j^S] = \mathbb{E}[B_1^S]^2$ for every $i, j \in [u]$. Moreover, as B_1^S, \dots, B_u^S are indicators, $\mathbb{E}[(B_i^S)^2] = \mathbb{E}[B_i^S]$ for all $i \in [u]$. Therefore

$$\begin{aligned} \mathbb{E}[(B^S)^2] &= \sum_{i, j \in [u]} \mathbb{E}[B_i^S B_j^S] \leq (u^2 - u)\mathbb{E}[B_1^S]^2 + u\mathbb{E}[B_1^S] \\ &\leq u^2 \mathbb{E}[B_1^S]^2 + \mathbb{E}[B^S] = \mathbb{E}[B^S]^2 + \mathbb{E}[B^S] \leq 2\mathbb{E}[B^S]^2, \end{aligned}$$

where the last inequality is due to the fact that $\mathbb{E}[B^S] \geq 2d \geq 1$. We conclude that

$$\Pr[B^S \geq d] \geq \frac{\mathbb{E}[B^S]^2}{4\mathbb{E}[(B^S)^2]} \geq \frac{1}{8}.$$

The proof of the lemma, and therefore of Theorem 5 is now complete.

Statement of potential broader impact

In this work, we have empirically shown that gradient boosters produce voting classifiers where many base learners make predictions of small magnitude. We then used this observation to prove stronger generalization bounds that better explain the practical performance of gradient boosters. We hope and believe that our findings may not only advance our theoretical understanding of boosting algorithms, but potentially also lead to algorithms with better accuracy by using regularization inspired by our new generalization bound or more directly optimizing it.

Acknowledgments and Disclosure of Funding

Kasper Green Larsen is supported by DFF Sapere Aude Grant 9064-00068B, a Villum Young Investigator Grant and an AUFF Starting Grant.

References

- [1] Leo Breiman. Prediction games and arcing algorithms. *Neural Computation*, 11(7):1493–1517, 1999.
- [2] Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '16*, page 785–794. Association for Computing Machinery, 2016.
- [3] Corinna Cortes, Mehryar Mohri, and Dmitry Storcheus. Regularized gradient boosting. In *Advances in Neural Information Processing Systems 32*, pages 5449–5458. 2019.
- [4] Yoav Freund and Robert E Schapire. A decision-theoretic generalization of on-line learning and an application to boosting. *J. Comput. Syst. Sci.*, 55(1):119–139, August 1997.
- [5] Jerome H. Friedman. Greedy function approximation: A gradient boosting machine. *Annals of Statistics*, 29:1189–1232, 2000.
- [6] Wei Gao and Zhi-Hua Zhou. On the doubt about margin explanation of boosting. *Artificial Intelligence*, 203:1–18, 2013.
- [7] Allan Grønlund, Lior Kamma, Kasper Green Larsen, Alexander Mathiasen, and Jelani Nelson. Margin-based generalization lower bounds for boosted classifiers. In *Advances in Neural Information Processing Systems 32*, pages 11963–11972. 2019.
- [8] Allan Grønlund, Lior Kamma, Kasper Green Larsen, Alexander Mathiasen, and Jelani Nelson. Margin-based generalization lower bounds for boosted classifiers. In *Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019*, 2019.
- [9] Guolin Ke, Qi Meng, Thomas Finley, Taifeng Wang, Wei Chen, Weidong Ma, Qiwei Ye, and Tie-Yan Liu. Lightgbm: A highly efficient gradient boosting decision tree. In *Advances in Neural Information Processing Systems 30*, pages 3146–3154. 2017.
- [10] V. Koltchinskii and D. Panchenko. Empirical margin distributions and bounding the generalization error of combined classifiers. *Ann. Statist.*, 30(1):1–50, 02 2002.
- [11] Albert W. Marshall, Ingram Olkin, and Barry C. Arnold. *Inequalities: Theory of Majorization and Its Applications*. Springer, 2nd edition, 2011.
- [12] Alexander Mathiasen, Kasper Green Larsen, and Allan Grønlund. Optimal minimal margin maximization with boosting. In *36th International Conference on Machine Learning*, volume 97 of *ICML '19*, pages 4392–4401, 09–15 Jun 2019.
- [13] Lev Reyzin and Robert E. Schapire. How boosting the margin can also boost classifier complexity. In *23rd International Conference on Machine Learning, ICML '06*, page 753–760. Association for Computing Machinery, 2006.
- [14] Robert E Schapire, Yoav Freund, Peter Bartlett, and Wee Sun Lee. Boosting the margin: A new explanation for the effectiveness of voting methods. *The annals of statistics*, 26(5):1651–1686, 1998.
- [15] Robert E. Schapire and Yoram Singer. Improved boosting algorithms using confidence-rated predictions. *Machine Learning Journal*, 37(3):297–336, December 1999.

- [16] Liwei Wang, Masashi Sugiyama, Zhaoxiang Jing, Cheng Yang, Zhi-Hua Zhou, and Jufu Feng. A refined margin analysis for boosting algorithms via equilibrium margin. *Journal of Machine Learning Research*, 12(51):1835–1863, 2011.

A Proof of Lemma 2

We start by handling the first part of the lemma, namely that for every $N \in \mathbb{N}^+$, with high probability over $S \sim \mathcal{D}^m$, $S \in \mathcal{E}_N$.

Claim 11. For every $N \in \mathbb{N}^+$, $g \in \mathcal{C}_N$ and $\ell \in [0, N]$, with probability at least $1 - \frac{\delta}{N(N+1)^2|\mathcal{H}|^N}$ over $S \sim \mathcal{D}^m$ we have

$$\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \leq \mathcal{L}_S^{\ell/N}(g) + \frac{8 \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} + 4\sqrt{\frac{\ln(4N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}; \quad (13)$$

and

$$\Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] \leq 2 \Pr_{(x,y) \sim S} [|f(x) - g(x)| > \ell/N] + \frac{8 \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m}. \quad (14)$$

We draw the reader's attention to the fact that by union bounding over all $g \in \mathcal{C}_N$ and $\ell \in [0, N]$ we get that $\Pr_{S \sim \mathcal{D}^m} [\mathcal{E}_N] \geq 1 - \frac{\delta}{N(N+1)}$ for every $N \in \mathbb{N}^+$, which proves the first part of Lemma 2. The proof of Claim 11 is quite involved technically, and is therefore deferred to the appendix .

Proof. First note that if $\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \leq 8m^{-1} \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$ then (13) holds for all S , and thus with probability 1 over $S \sim \mathcal{D}^m$. Assume therefore that $\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) > 8m^{-1} \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$. Denote $S = \{(x_j, y_j)\}_{j \in [m]}$, then

$$\mathcal{L}_S^{\ell/N}(g) = \Pr_{(x,y) \sim S} [yg(x) \leq \ell/N] = \frac{1}{m} \sum_{j \in [m]} \mathbb{1}_{yg(x_j) \leq \ell/N}.$$

Moreover $\mathbb{E}[\mathbb{1}_{yg(x_j) \leq \ell/N}] = \mathcal{L}_{\mathcal{D}}^{\ell/N}(g)$ for all $j \in [m]$, and therefore $\mathbb{E}[\mathcal{L}_S^{\ell/N}(g)] = \mathcal{L}_{\mathcal{D}}^{\ell/N}(g)$. Let $\gamma := \sqrt{\frac{2 \ln(4N(N+1)^2|\mathcal{H}|^N/\delta)}{m \mathcal{L}_{\mathcal{D}}^{\ell/N}}}$. Then $\gamma \in (0, 1/2)$, and therefore a Chernoff bound gives the following two inequalities.

$$\Pr_{S \sim \mathcal{D}^m} \left[\mathcal{L}_S^{\ell/N}(g) < (1-\gamma) \mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \right] \leq e^{-\gamma^2 m \mathcal{L}_{\mathcal{D}}^{\ell/N}(g)/2} \leq \frac{\delta}{4N(N+1)^2|\mathcal{H}|^N}$$

$$\Pr_{S \sim \mathcal{D}^m} \left[\mathcal{L}_S^{\ell/N}(g) > 2\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \right] \leq e^{-m \mathcal{L}_{\mathcal{D}}^{\ell/N}(g)/3} \leq \frac{\delta}{4N(N+1)^2|\mathcal{H}|^N},$$

where the last inequality follows from the fact that $\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \geq 8m^{-1} \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$. Therefore with probability at least $1 - \delta/(2N(N+1)^2|\mathcal{H}|^N)$ we get that

$$\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \leq (1-\gamma)^{-1} \mathcal{L}_S^{\ell/N}(g) \leq (1+2\gamma) \mathcal{L}_S^{\ell/N}(g) \leq (1+2\gamma) \mathcal{L}_S^{\ell/N}(g) + \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m}, \quad (15)$$

and moreover

$$\gamma = \sqrt{\frac{2 \ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m \mathcal{L}_{\mathcal{D}}^{\ell/N}(g)}} \leq \sqrt{\frac{4 \ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m \mathcal{L}_S^{\ell/N}(g)}} \quad (16)$$

Plugging (16) into (15) and summing up we get

$$\mathcal{L}_{\mathcal{D}}^{\ell/N}(g) \leq \mathcal{L}_S^{\ell/N}(g) + \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} + 4\sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}(g).$$

Next note once again that if $\Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] \leq 8m^{-1} \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$ then (14) holds for all S , and thus with probability 1 over $S \sim \mathcal{D}^m$. Assume therefore that

$\Pr_{(x,y) \sim \mathcal{D}}[|f(x) - g(x)| > \ell/N] > 8m^{-1} \ln(4\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$. Similarly to the first part of the proof a Chernoff bound gives the following inequality.

$$\begin{aligned} \Pr_{S \sim \mathcal{D}^m} \left[\Pr_{(x,y) \sim S} [|f(x) - g(x)| > \ell/N] > 2 \Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] \right] \\ \leq e^{-m \Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] / 3} \leq \frac{\delta}{4N(N+1)^2|\mathcal{H}|^N}, \end{aligned}$$

where the last inequality follows from the fact that $\Pr_{(x,y) \sim \mathcal{D}} [|f(x) - g(x)| > \ell/N] \geq 8m^{-1} \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)$. Therefore with probability at least $1 - \delta/(2N(N+1)^2|\mathcal{H}|^N)$ we get (14). Union bounding we get that with probability with probability at least $1 - \delta/(N(N+1)^2|\mathcal{H}|^N)$ over the choice of $S \sim \mathcal{D}^m$ we have both (13) and (14). \square

We turn now to prove the second part of Lemma 2, namely that $\bigcap_{N \in \mathbb{N}^+} \mathcal{E}_N \subseteq \mathcal{E}$. To this end, let $S \in \bigcap_{N \in \mathbb{N}^+} \mathcal{E}_N$. Let f be some voting classifier and let $\theta > 0$. As f is a voting classifier, then there exists a sequence $\langle \alpha_h \rangle_{h \in \mathcal{H}} \in \mathbb{R}_+^{\mathcal{H}}$ such that $\sum_{h \in \mathcal{H}} \alpha_h = 1$ and $f = \sum_{h \in \mathcal{H}} \alpha_h \cdot h$. Thus f implicitly defines a distribution $\mathcal{Q} = \mathcal{Q}(f)$ over \mathcal{H} , where $\Pr_{h \sim \mathcal{Q}}[h = h'] = \alpha_{h'}$ for all $h' \in \mathcal{H}$. Recall that $\Delta : \mathcal{X} \times \mathcal{H} \rightarrow \mathbb{R}$ is defined by $\Delta(x, h) := |f(x) - h(x)|$ for every $x \in \mathcal{X}$, $h \in \mathcal{H}$.

Definition 1. Let X be a random variable, and let $r \in \mathbb{N}$, then the r th moment of X is defined by $\|X\|_r^r := \mathbb{E}[X^r]$. The r th norm of X is defined by $\|X\|_r := \sqrt[r]{\mathbb{E}[X^r]}$.

Set hereafter

$$N := \lg(16m) \cdot \max\{256\theta^{-1}\|\Delta(x, h)\|_{\lg(16m)}, 100/\theta\},$$

$$128e\theta^{-2} \cdot \left(\mathbb{E}_{(x,y) \sim S} \left[\mathbb{E}_{h \sim \mathcal{Q}} [\Delta(x, h)^2]^{(\lg(16m))/2} \right]^{2/(\lg(16m))} \right)$$

The product distribution \mathcal{Q}^N defines a distribution over \mathcal{H}^N . By identifying an N -tuple $h_1, \dots, h_N \in \mathcal{H}$ with the corresponding classifier $\frac{1}{N} \sum_{j \in [N]} h_j$ we can think of \mathcal{Q}^N also as a distribution over \mathcal{C}_N . We first observe that

$$\begin{aligned} \mathcal{L}_{\mathcal{D}}(f) &\leq \Pr_{(x,y) \sim \mathcal{D}, g \sim \mathcal{Q}^N} [yf(x) \leq 0 \wedge yg(x) \leq \theta/2] + \Pr_{(x,y) \sim \mathcal{D}, g \sim \mathcal{Q}^N} [yf(x) \leq 0 \wedge yg(x) > \theta/2] \\ &\leq \Pr_{(x,y) \sim \mathcal{D}, g \sim \mathcal{Q}^N} [yg(x) \leq \theta/2] + \Pr_{(x,y) \sim \mathcal{D}, g \sim \mathcal{Q}^N} [|f(x) - g(x)| > \theta/2] \end{aligned} \tag{17}$$

To bound the first summand, let $\ell \in [0, N]$ be the smallest integer such that $\theta/2 \leq \ell/N$. Such ℓ clearly exists as $\theta \in [0, 1]$. Moreover we know that $\theta/2 \leq \ell/N \leq \theta/2 + 1/N \leq 51\theta/100$. Since $S \in \mathcal{E}_N$ we get that

$$\begin{aligned} \Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [yg(x) \leq \theta/2] &\leq \Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [yg(x) \leq \ell/N] = \mathbb{E}_{g \sim \mathcal{Q}^N} \left[\Pr_{(x,y) \sim \mathcal{D}} [yg(x) \leq \ell/N] \right] \\ &\leq \mathbb{E}_{g \sim \mathcal{Q}^N} \left[\Pr_{(x,y) \sim S} [yg(x) \leq \ell/N] + \varepsilon_N(g) \right] \leq \Pr_{(x,y) \sim S} [yg(x) \leq 51\theta/100] + \mathbb{E}_{g \sim \mathcal{Q}^N} [\varepsilon_N(g)], \end{aligned}$$

where $\varepsilon_N(g) = \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} + 4\sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}(g)$. Similarly to (17) we get that

$$\Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [yg(x) \leq 51\theta/100] \leq \Pr_{(x,y) \sim S} [yf(x) \leq \theta] + \Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > 49\theta/100],$$

and therefore

$$\Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [yg(x) \leq \theta/2] \leq \Pr_{(x,y) \sim S} [yf(x) \leq \theta] + \Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > 49\theta/100] + \mathbb{E}_{g \sim \mathcal{Q}^N} [\varepsilon_N(g)]. \tag{18}$$

Moreover, since $S \in \mathcal{E}_N$ we get the following bound over the second summand in (17).

$$\begin{aligned}
\Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > \theta/2] &\leq \Pr_{\substack{(x,y) \sim \mathcal{D} \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > (\ell - 1)/N] \\
&\leq 2 \Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > (\ell - 1)/N] + \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} \\
&\leq 2 \Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > 49\theta/100] + \frac{8 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m}
\end{aligned} \tag{19}$$

Plugging (18) and (19) into (17) we get that

$$\begin{aligned}
\mathcal{L}_{\mathcal{D}}(f) &\leq \Pr_{(x,y) \sim S} [yf(x) \leq \theta] + 3 \Pr_{\substack{(x,y) \sim S \\ g \sim \mathcal{Q}^N}} [|f(x) - g(x)| > 49\theta/100] \\
&\quad + \frac{16 \ln(2\delta^{-1}N(N+1)^2|\mathcal{H}|^N)}{m} + \mathbb{E}_{g \sim \mathcal{Q}^N} \left[\sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}(g) \right]
\end{aligned} \tag{20}$$

From Lemma 3 we get that by Jensen's inequality and sub-additivity of square root

$$\begin{aligned}
\mathbb{E}_{g \sim \mathcal{Q}^N} \left[\sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\ell/N}(g) \right] &\leq \sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathbb{E}_{g \sim \mathcal{Q}^N} \left[\mathcal{L}_S^{51\theta/100}(g) \right] \\
&\leq \sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \left(\mathcal{L}_S^{\theta}(f) + \frac{1}{m^2} \right) \\
&\leq \frac{1}{m} + \sqrt{\frac{\ln(N(N+1)^2|\mathcal{H}|^N/\delta)}{m}} \mathcal{L}_S^{\theta}(f),
\end{aligned} \tag{21}$$

and therefore

$$\mathcal{L}_{\mathcal{D}}(f) \leq \mathcal{L}_S^{\theta}(f) + O \left(\frac{N \lg |H| + \lg(1/\delta)}{m} + \sqrt{\frac{N \lg |H| + \lg(1/\delta)}{m}} \mathcal{L}_S^{\theta}(f) \right),$$

which concludes the proof of Theorem 1.