

1 **Author Response ("Margins are Insufficient for Explaining Gradient Boosting")**

2 We thank the reviewers for the time and expertise invested in these reviews.

3 **Response to Reviewer 2.** Addressing the reviewer’s concerns, we wish to stress that while for Gradient Boosters
4 (GB) the sigmoid function may be used to define the loss function and to transform the (raw) output prediction of a
5 (learned) ensemble into a probability, the sigmoid function is not used when studying margin theory. We compute the
6 margin of an ensemble of base learners ($\{\alpha_i, h_i \mid h_i(x) \rightarrow [-1, 1]\}$) following classic margin theory (see Section
7 2). Loosely speaking, the margin of a point depends on the output of the voting classifier, and does not involve the
8 sigmoid function. Formally, for data point x_i the margin is $\frac{\sum_{j=1}^n \alpha_j h_j(x_i)}{\sum_{j=1}^n |\alpha_j|}$. Note the final remark of Section 2 of how
9 we transform the output of GB algorithms, e.g. LightGBM, to fit this framework. This is unencumbered by the loss
10 function minimized in training or how sigmoid may be used to transform the output of the ensemble into to a probability
11 distribution by setting $P(y = 1 \mid x) = \sigma(\sum_{j=1}^n \alpha_j h_j(x_i))$. We also note that as opposed to the margins, this probability
12 is not scale-invariant as scaling the weights or adding copies of the base learners provide more extreme probabilities.

13 Albeit unconnected to margin theory, the sigmoid function cannot make large changes to the output if the input is only
14 changed by a small amount. In fact small changes to the input to the sigmoid function makes even smaller changes to
15 the output since the derivative of the sigmoid function is at most 1/4 (and that is only at zero and the norm of derivative
16 goes to zero as $|x|$ goes to infinity very fast).

17 Regarding the reason the margins are reduced with the number of base regressors in GB, in short this is due to the fact
18 that gradient boosters may decide to focus on only a small subset of data in each iteration Consider the margin of a
19 fixed data point x , and assume for simplicity that all weights (α_j) are 1. Then the margin x is $\frac{1}{n} \sum_i h_i(x)$, where n is
20 the number of base learners, which drops towards zero as n increases unless $h_i(x)$ remains large (that is, close to 1)
21 for all $i = 1, \dots, n$. Hence, if each round only focuses on a small subset of data points, meaning that the new base
22 learner only gives non-negligible prediction on a small subset, then the margin of the remaining points decreases. If this
23 happens in each round of boosting, the result is a smaller and smaller margin distribution, and this is what happens in
24 the experiments shown in our paper. We will update our description of this argument in Section 3 of our paper to make
25 it more clear (also see the arguments presented in Section 4 under Potentially Much Better).

26 We thank the reviewer for the specific comments as well. We will certainly take care of them for the final version.

27 **Response to Reviewer 1.** As the reviewer suggests, the scope of our paper and theoretical contributions are more
28 general than gradient boosted trees, as it applies to all voting classifiers irrespective of learning algorithm, including loss
29 function optimized, and base learners. The paper investigates how the actual predictions of voting classifiers on the data
30 may be very different than $\{-1, 1\}$, which is the standard assumption when proving margin bounds. Our theory, and
31 margin theory in general, is orthogonal to specifics of the learning algorithm. We will think about ways of rephrasing
32 our contributions to make it clearer that it is not only specific to gradient boosters with a concrete loss function. In
33 practice, when conducting the numerical experiments, we chose to focus on comparing the classic AdaBoost algorithm
34 where the base learner is restricted to map inputs to $\{-1, 1\}$ with Gradient Boosters that very often returns negligible
35 predictions, as this exactly highlight and fits the phenomena we are investigating. We stress, however, that this is a
36 special case, and the theoretical results we present are significantly broader.

37 For base learners, the same size means the same number of leaves (and no restriction on depth for both algorithms
38 compared). We do not consider decision stumps as these are usually not used by GB in practice and lead to inferior
39 performance (the default number of leaves for LightGBM is size 31, while XGBoost use a max depth default equal to
40 6). Furthermore, with decision stumps, the base classifiers are too weak for the phenomenon with small predictions to
41 even occur. That is, with decision stumps, most predictions are very close to $\{-1, 1\}$. This is most likely due to the
42 fact that decision stumps are incapable of “focusing” on a small subset of training points as discussed above. We will
43 elaborate more on this in the final version and if space allows it, we will also include a histogram of predictions when
44 using decision stumps.

45 In the supplemental material, submitted along with the paper, we included the same experiment on three more data
46 sets to give 4 data sets of increasing size to analyze and demonstrate our new theoretical bound on. The observed
47 behaviour was completely consistent. The mean validation error and standard deviation for the Forest Cover dataset
48 example from the paper are (0.0298, 0.00037) for LightGBM and (0.0327, 0.00053) for AdaBoost. The standard
49 deviation was so small that we chose to only show 3 runs on the plots. We will make sure to comment on this in a
50 final version. For data sets included in the supp. material the mean, and standard deviation are as follows **Boone:**
51 Lgb: mean=0.0574, std=0.00015, Ada: mean=0.0631, std=0.00068 **Higgs:** Lgb: mean=0.2530, std=0.00031, Ada:
52 mean=0.2777, std=0.00009 **Diabetes:** Lgb: mean=0.2532, std=0.0255, Ada: mean=0.2692, std=0.0194 . for Boone
53 and Higgs we used three runs and for Diabetes we used 10 (due to much smaller data size)