

1 We would like to thank all reviewers for their invaluable feedback. The next revision of the paper will include fixes for
2 all typos that were mentioned. Responses for questions raised by each reviewer are below.

3 **Reviewer #1**

4 We would like to clarify questions about the theory you've raised. PolyTree framework represent tree and ensembles not
5 as a sum of weights over **leaves**, but over **right split indicator functions** ($c(x) = I\{x > b\}$, see line 57). This fact is a
6 core idea and should be kept in mind to understand most of the propositions: A) l74: *I do not understand condition 2. A*
7 *leaf can be bounded from the left and the right, so can't a feature be present twice?* The feature will still be represented
8 only once. Let look at leaf $x \in (0, 1]$. This leaf is product of two indicator function $I\{x > 0\}I\{x \leq 1\}$, this is equal
9 to $I\{x > 0\}(1 - I\{x > 1\}) = 1 - I\{x > 0\}I\{x > 1\} = 1 - I\{x > 1\}$. $I\{x > 1\}$ is a stronger condition, thus
10 $I\{x > 0\}$ could be removed. We will add this example to the next revision of the paper to clarify this property. *Theorem*
11 *1: Wouldn't it be possible to construct an ensemble H' from an ensemble H by splitting one leaf l and assigning the*
12 *weigh of the old leaf in the new leaves in H' ?* When the ensemble is represented as a sum over leaves, there indeed will
13 be two different ensembles. But, both ensembles will generate the same PolyTree representation, because we perform a
14 reduction of weights for monomials with the same set of conditions (such leaf will generate 2 monomials with the same
15 absolute weights, but different signs, after reduction terms will disappear). This comes from condition 2 (line 74). *l105:*
16 *In this section, we (...) set up a task of tree shape change in an ensemble. We represent a tree ensemble as a sum of trees*
17 *of fixed shape. —> unclear to me whether the shape is changed or not.* By fixed shape we mean a restriction on a set of
18 possible tree structures. For example, tree of depth 6; balanced trees; symmetric oblivious trees and so on. *equation*
19 *2 seems wrong...* Thank you, this is a typo and will be fixed. *Comments about notations* By 2^C we denote set of all
20 possible monomial structures, and $M \in 2^C$ represent some monomial from this set. We use letter d to denote the depth
21 of decision tree, and in the context of line 85 this is a maximum depth of tree in the ensemble. We'll add all this to
22 paragraph with notation. *Provide code* We have Java-based implementation on GitHub and will include the link with
23 non-anonymous version. We are also working on an implementation of proposed methods as a part of one of the major
24 GBDT libraries and it'll be finished before NeurIPS2019.

25 **Reviewer #2**

26 - *In the last paragraph of section 3, the execution time of an original ensemble with the transformed one using the*
27 *proposed algorithm is compared; nevertheless, the overhead time added by this transformation is not mentioned.* Our
28 greedy algorithm is quite slow and could take several seconds to complete. Also, this experiment is a proof-of-concept
29 and we did not optimize it to run as fast as possible. Mostly, such transformations are interesting for production tasks
30 where the same model could be used for days. Transformation time could be big, but it is done once, thus overhead is
31 negligible. *...Symmetric trees can highly benefit from model reduction because of having many zero leaves...* For more
32 clarity, there were balanced trees, like on in XGBoost, but symmetric oblivious trees used in CatBoost. Such trees are
33 often not the best choice to learn on one-hot-encoded datasets, so exploring pruning strategy on this type of trees also
34 requires an in-depth study of what type of trees and when we need to use for such data. It is interesting, but there is
35 enough space and we will divert from the main topic — PolyTree framework introduction.

36 **Reviewer #3**

37 *Why the level of dependency grows with the number of conditions in M .* PolyTree decomposition is similar to n -way
38 ANOVA decomposition, where dependence for factors x, y, z is decomposed to main effects (x, y, z) and their iterations
39 (xy, yz, xyz) . In ANOVA x, y, z are categorical factors, while in PolyTree it is right split indicator functions. - *Authors*
40 *do not mentioned if the cross validation is applied or not. This latter should be applied.* The cross-validation is used
41 during parameter tuning, but the final metric is computed on the fixed independent test set (experiment simulates
42 practical usage of GBDT — choose hyper-parameters / model via cross-validation and then use it in production). *For*
43 *section 3 relative to theoretical analysis, experiments remain poor; testing on only one set (Higg dataset) is insufficient to*
44 *make reliable conclusions. More data sets should be tested.* This example shows, that there are certain type of problems,
45 where to get the better model we need to use one type of tree shapes and are able to benefit from the other shape during
46 exploitation. We don't insist that this situation is universal, but just show its existence. Systematic study what tree shape
47 needed for what problem, and what benefits could be achieved for the specific problem at hand, requires space; we
48 are working on this as follow-up work. *Why what is tested in Table 1 is only sets of binary classes. What about sets*
49 *with multiple classes.* Its a topic for an individual study. There are several ways of how to perform multiclassification
50 in boosting. For some types (like one-vs-all) the problem reduces to binary one, while for others not. There could
51 be several ways to generalized PolyTree for multiclass, and each one should be explored. But this article introduces
52 PolyTree framework and we don't have enough space to explore this problem in details. *In Tab 1, How to explain that*
53 *there is no significant improvement In AUC criterion between trained ensembles and pruned ensembles. In fact, you get*
54 *almost the same AUC for both methods* For the trained ensemble, we have used the early-stopping strategy. GBDT with
55 early stopping usually provided the state-of-the-art model. Thus, this model provides a very strong baseline that is hard
56 to beat with the smaller model.