

1 **To Reviewer #1: LP-based data selection:** For Fig. 2, we indeed used random sampling to select a subset prior to
2 applying the LP-based data selection on the reduced dataset. This was done because the number of constraints in
3 (7) grows quadratically with the number of samples, N . However, upon further analysis we noticed that the random
4 sampling is unnecessary because the LP in (7) with $f = 0$ and $g = 1$ can in fact be decomposed into N smaller LPs, each
5 with $|\mathcal{N}|$ variables and $d + 2$ constraints. These can be solved in parallel using a LP solver. We implemented this for
6 datasets in Fig. 2 and noticed that the performance of our method is no worse than currently reported, and even better in
7 some cases. We thank the reviewer for pointing out this discrepancy, and will update with new experiments for S10-DS
8 in Fig. 2 with the LP decomposition-based approach (better results without random subsampling preprocessing).

9 **Cutting planes:** Our cutting planes are added once at the beginning before invoking a MILP solver. The cuts can be
10 applied to the other MIP models provided that they use the variables c_i and e_{il} with the same modelling meaning. For
11 example, Fair in [1] can use these cuts, and OCT-H in [4] needs to be modified slightly before they can be applied.

12 **Using data selection for CART, OCT-H (also Reviewer #2):** Thank you for your suggestion. Fig. 1 benchmarks the
13 performance of SVM1-ODT on medium-sized datasets and shows that our ODT formulation outperforms CART,
14 OCT-H, Fair and BinOCT in terms of test accuracy *without using data selection*. Fig. 2 addresses scalability by
15 benchmarking our combined approach (SVM1-ODT and data selection) with CART and other methods such as OCT-H
16 that address scalability. While it is certainly possible to apply our data selection procedure prior to using either CART
17 or OCT-H, our intent in Fig. 2 is largely to benchmark the *scalability* of the combined approach.

18 **Clarification on Eq. (6):** In Sec. 3, we first present the generalized data selection formulation that maximizes the points
19 in the convex hull ($g^T b$) and minimizes the number of selected data points ($f^T a$). In Sec. 3.1, we present a special
20 case of (6) with $f = 0$ and $g = 1$ because choosing these values allows us to decompose it into N smaller LPs while
21 maximizing the points inside the convex hull. The balanced schema is presented in the Supp. without numerical results.

22 **Optimality gap:** For datasets that cannot be solved to optimality within 15 minutes, we observed that the optimality gap
23 could still be large due to the lower bound not improving significantly before the run is terminated. The gap does not
24 indicate time to reach optimality. Thus, we instead use training accuracy to benchmark the performance of our model.

25 **Notation in Section 3.2:** J_N is a subset of extreme points that satisfy a specific condition. K_N contains extreme points,
26 except they are never used with convex combination coefficient larger than $\frac{1}{d+1}$ to express some other points. Since
27 points in K_N carry some information about the distribution of the dataset, we select a subset of K_N , using Alg. 1, for
28 training. “Boundary hyperplanes” refers to the hyperplanes for each leaf node. $\text{dist}(a, h)$ refers to the distance between
29 a to the hyperplane $h^T x = 0$. The term “heuristic” refers to Alg. 1. The LP in (7) is solved in the second step of Alg. 1.

30 **Theorem 2:** Our idea is to disperse the numerical issues between parameters by re-scaling. The numerical instability
31 for a very small ϵ should be easier to handle by an MILP solver than a very big M . *Page 2, line 85:* From (1d, 1e), it
32 follows that w_{il} equals to some u_i which is integral. Then (1c) implies that at the optimal solution \hat{y}_i is also integral. So
33 we do not explicitly enforce \hat{y}_i to be integral. *CART with deeper trees:* We follow the standard benchmarking used in
34 [1, 4] wherein the tree depth for CART and MIPs are the same. Note that a deeper tree for CART is not necessary to
35 outperform the shallow ODT; e.g., for Dermatology the ODT with $D = 2$ has a higher test accuracy (80.7%) compared
36 to the CART tree with $D = 3$ (76.1%) As also observed in [11], for many datasets, even we allow CART to choose the
37 tree depth using its default setting (allowing deep trees), a shallow ODT still outperforms CART (see Tab.13 in [11]).
38 We have incorporated all results in Thm. 1 and 2 into the final model, the final SVM1-ODT model imposes $u_i \in [1, Y]$
39 and $M = 1$. Model parameters ϵ and α_i are tuned via cross-validation for each dataset. We will fix the terminology
40 issues and revise discussions related to “sparsity” for hyperplane, “linear” for the 1-norm, and other suggestions.

41 **To Reviewer #2:** If McCormick linearizations are not used, then there are nonconvex quadratic terms in the constraints.
42 The tractability of the resulting MIQCP is very challenging. CPLEX only supports convex quadratic terms that can be
43 represented as second order cone programs (see shorturl.at/cnDY6). Note that handling logical constraints is not as
44 scalable as the big-M method if we can choose a small value for M . As shown in Theorem 2, we can use $M = 1$.

45 **To Reviewer #3:** The requirements for input features and branching rules from [2] and [11] are different from ours.
46 In particular, [2] and [11] only consider binary features $\mathbf{x}_i \in \{0, 1\}^d$, while our formulation takes numerical features
47 $\mathbf{x}_i \in \mathbb{R}^d$ (and it can be extended to the case with mixed numerical and categorical features, see A.2). Branching rules
48 in these references are binary tests (i.e., univariate splits), while we use a multivariate hyperplane. Hence the authors
49 employ special tools for the specific decision tree, they cannot be generalized to our general tree; for example, [2]
50 proposes a max flow-based MIP formulation. Comparing to [4], we impose additional conditions on hyperplanes so
51 that training samples should be far from the boundary of the cluster at the leaf node by using the multi-hyperplane
52 SVM. Moreover, we use lesser number of binary variables to encode the decision tree ($e_{il}, c_i \in \{0, 1\}$ in our model v.s.
53 $z_{it}, s_{jt}, l_t, d_t \in \{0, 1\}$ in [4]), and we can use a small value for the big-M parameter (i.e., $M = 1$).

54 We did compare with OCT-H from [4] in Figs 1, 2. The model in [2] handles binary input, and 8 data sets were tested.
55 They share 4 datasets with us (Tab. 2 in [2] and Tab. 5 in our Supp.). For same tree depth, we always outperform [2].
56 We have defined \bar{y}_i in lines 118-119. For L124, we will replace it by “... a small constant ϵe_{il} in (2) instead of e_{il} to ...”

57 **To Reviewer #4:** We will give more detailed explanation for cutting planes in Prop. 1 in the revised paper.