

1 Thank you reviewers for your positive and helpful suggestions!

2 **R1: Figure 1 no OSDT (warm-start only):** They probably have the same solution; the two points would be identical.

3 **Runtime benchmark against BinOCT and CART:** Good idea. We can add BinOCT and CART execution profiles to
4 figures like that shown in the paper’s Figure 2. (We include examples in Figure 1 below.) In general, OSDT converges
5 and *completes more quickly* than BinOCT; BinOCT *sometimes includes redundant leaves* leading to non-optimal
6 solutions to our problem. CART frequently completes very quickly, but does not always yield results near optimal. For
7 example, on monk-1, monk-2 and tic-tac-toe, the results of CART are far away from the optimal ones.

8 **Show BinOCT and CART trees as well as OSDT trees (Fig 4):** Sure! Figures 5 and 6 in the supplementary materials
9 have example binOCT trees; we can add CART trees as well.

10 **Revised title:** We’ve made it efficient for real-valued features since the submission, owing to a good implementation
11 of Theorem F.1. (see Figure 2 below). **Introduce notation gradually:** easy, will do. **Including more background on**
12 **CORELS:** easy, will do. Thank you for all your suggestions! We appreciate it!

13 **R2: Paper contributions:** Sorry! (1) first practical optimal decision tree algorithm to achieve provably optimal solutions
14 for nontrivial problems. (2) a series of new analytical bounds to reduce the space (Sec 3.2). (3) the first practical
15 algorithmic use of a tree representation using only its leaves (Sec 3). (4) Implementation speedups saving 97% of
16 run time (Sec 4 and appendix). (5) For the important COMPAS and FICO datasets (Sec 5), optimal trees have never
17 previously been published. We present the first ones.

18 **Utility for real and categorical features too:** The COMPAS and FICO datasets used in the paper both have
19 real-valued features, so it is applicable. Since the submission, our implementation of Thm F.1 is more efficient,
20 allowing better scaling with real-valued features (e.g. on the Iris dataset in Figure 2, we reduced the number of tree
21 evaluations by 45%). All possible split points are considered for each real-valued variable in our current implementation.

24 **Scalability:** We have been working on a parallel implementation and expect to have
25 scalability results from it in time for the camera ready submission. Parallelization is
26 *much* less difficult than what we already did.

27 **R3: The authors would need to demonstrate convincingly that improved training**
28 **accuracy also translates into improved test set generalization:** The basis for all of
29 statistical learning theory is that training accuracy *and simplicity* provably lead to
30 better test accuracy, and OSDT’s objective incorporates both accuracy and sparsity.
31 The 10-fold cross validation results in the supplementary materials (Section J) illustrate
32 this, showing in-sample and out-of-sample accuracy.

33 **The results are less conclusive than for training:** Since OSDT found solutions that
34 were as accurate as other methods in testing, but were more sparse, we interpret that
35 as conclusive evidence. We will discuss this in the main body.

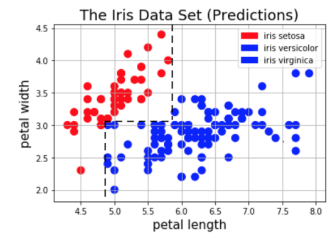


Figure 2: OSDT’s result on Iris dataset. Covariates are real-valued. OSDT considered all possible splits to produce the optimal decision boundary (black).

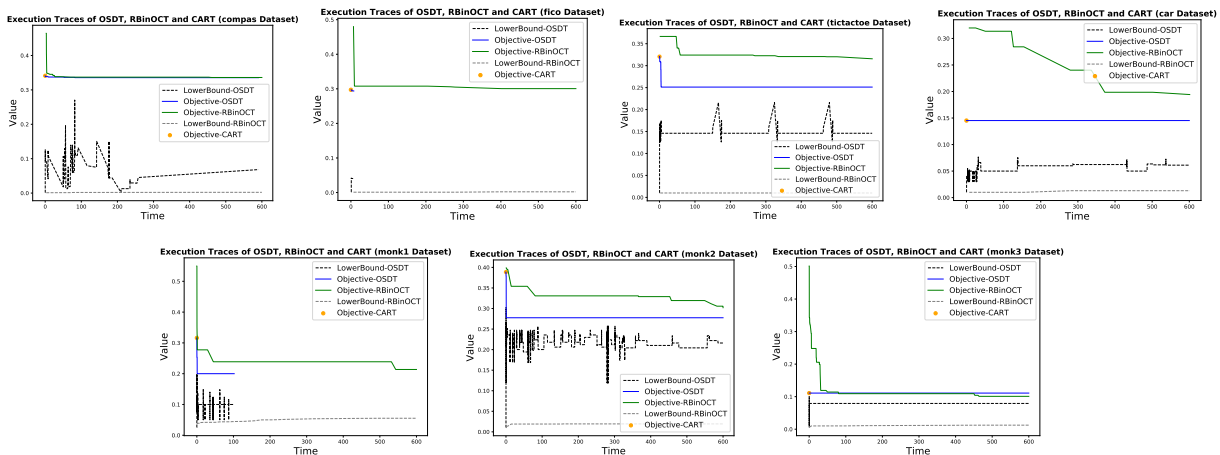


Figure 1: Execution traces of OSDT, CART and regularized BinOCT. OSDT converges much faster than RBinOCT generally. Traces end when execution completes for each algorithm separately. Runs were stopped after 10 min.