

Appendix

A Robust Relabeling Criterion

In Algorithm 1 we give pseudocode for training decision trees using the robust relabeling procedure as a splitting criterion. The algorithm does a depth-first search through all possible decision nodes up to a maximum depth and sets thresholds that locally maximize adversarial accuracy.

The maximum matching step takes worst-case $\mathcal{O}(n^{2.5})$ time. Since this step gets called for every decision node (2^d nodes) and for every possible split value (nm values) the worst-case overall runtime is $\mathcal{O}(2^d mn^{3.5})$. Here n is the number of samples, m the number of features and d the maximum depth of the tree.

Algorithm 1 Robust decision tree learning with robust relabeling criterion

Input: dataset X (n samples, m features), labels y , tree leaves \mathcal{T}_L , maximum depth d

```

1:  $L \leftarrow \{i \mid y_i = 0\}$ 
2:  $R \leftarrow \{i \mid y_i = 1\}$ 
3:  $V_j \leftarrow$  sorted values of feature  $j = 1 \dots m$ 
4:  $\text{loss} \leftarrow n$ 
5: for  $k \in 1 \dots 2^d$  do                                 $\triangleright$  For each decision node up to maximum depth
6:    $\text{best\_loss} \leftarrow \text{loss}$ 
7:    $\mathcal{T}_L^* \leftarrow \mathcal{T}_L$ 
8:   for  $j = 1 \dots m, v \in V_j$  do                         $\triangleright$  Try every split value
9:      $\mathcal{T}_L' \leftarrow \text{ADD\_DECISION\_NODE}(\mathcal{T}_L, j, v, k)$ 
10:     $E' = \{(u, v) \mid u \in L, v \in R, \mathcal{T}_L'^{S(u)} \cap \mathcal{T}_L'^{S(v)} \neq \emptyset\}$ 
11:     $\text{loss}' \leftarrow |\text{MAXIMUM\_MATCHING}(L, R, E')|$          $\triangleright$  Use matching size as loss
12:    if  $\text{loss}' < \text{best\_loss}$  then
13:       $\mathcal{T}_L^* \leftarrow \mathcal{T}_L'$ 
14:       $\text{best\_loss} \leftarrow \text{loss}'$ 
15:    end if
16:  end for
17:   $\mathcal{T}_L \leftarrow \mathcal{T}_L^*$                                  $\triangleright$  Keep the split with lowest loss
18: end for
19:  $\mathcal{T}_L \leftarrow \text{ROBUST\_RELABELING}(X, y, \mathcal{T}_L)$          $\triangleright$  Relabel to set leaf predictions

```

B Cost Complexity Pruning vs. Robust Relabeling

In the paper we compared Cost Complexity Pruning and robust relabeling on the Pima-Indians-diabetes dataset. Below we give the same plots for all 10 datasets that we used in the paper.

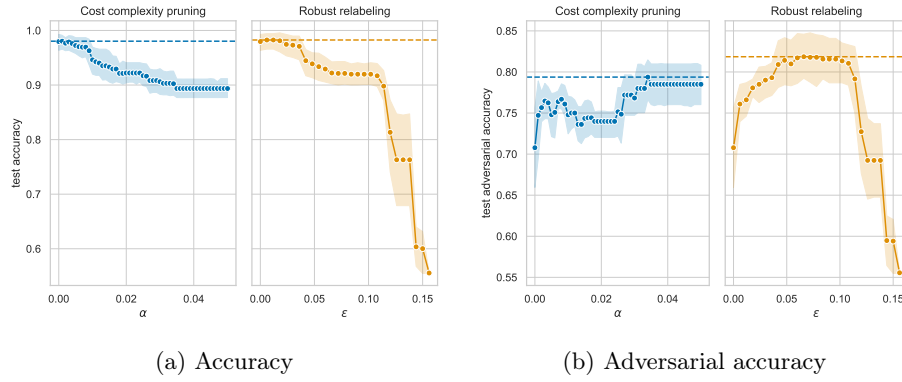


Fig. 1: 5-fold cross validation on the Banknote-authentication dataset.

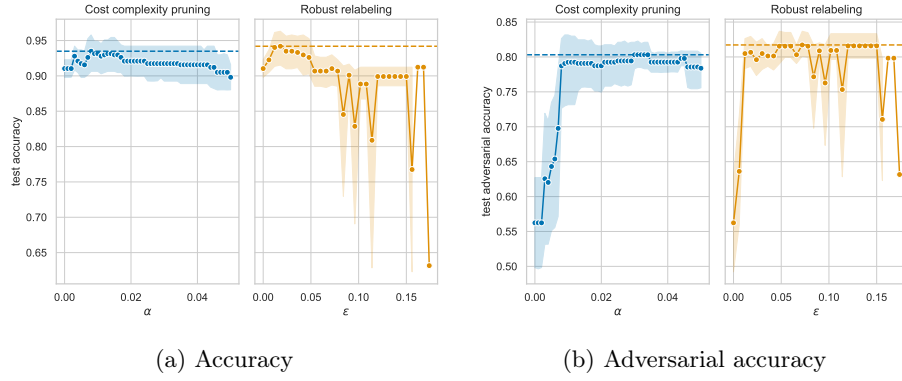


Fig. 2: 5-fold cross validation on the Breast-cancer-diagnostic dataset.

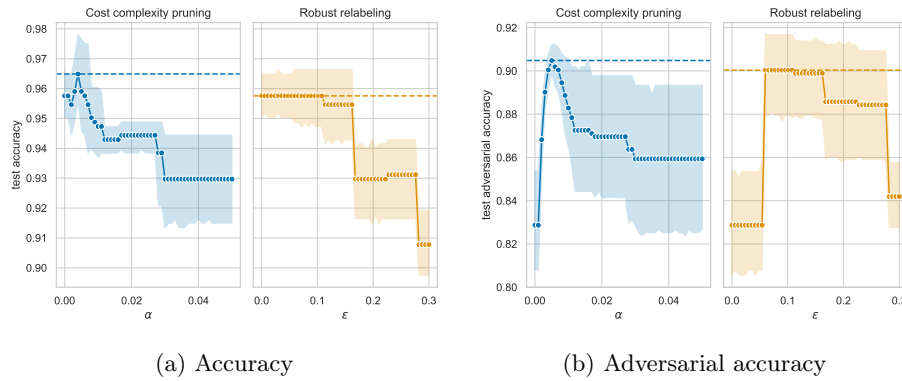


Fig. 3: 5-fold cross validation on the Breast-cancer dataset.

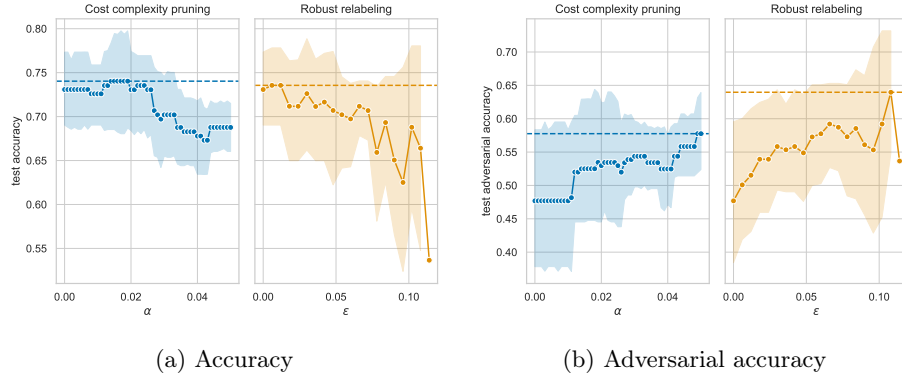


Fig. 4: 5-fold cross validation on the Connectionist-bench-sonar dataset.

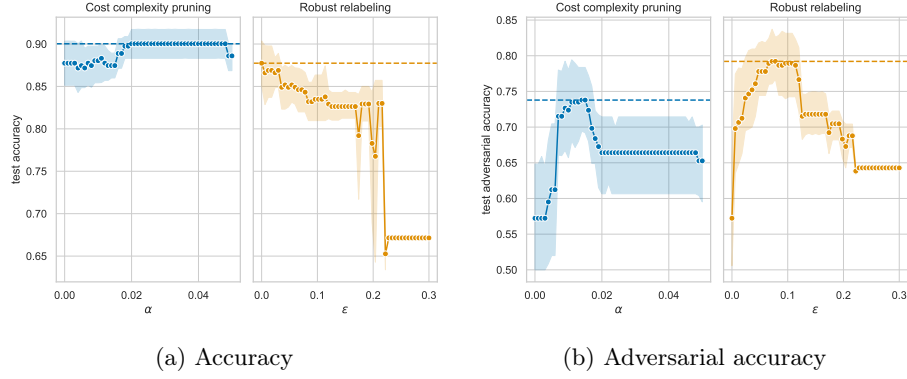


Fig. 5: 5-fold cross validation on the Ionosphere dataset.

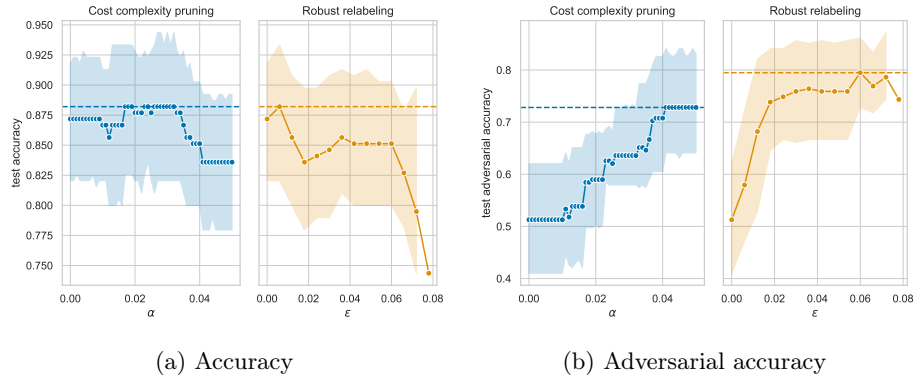


Fig. 6: 5-fold cross validation on the Parkinsons dataset.

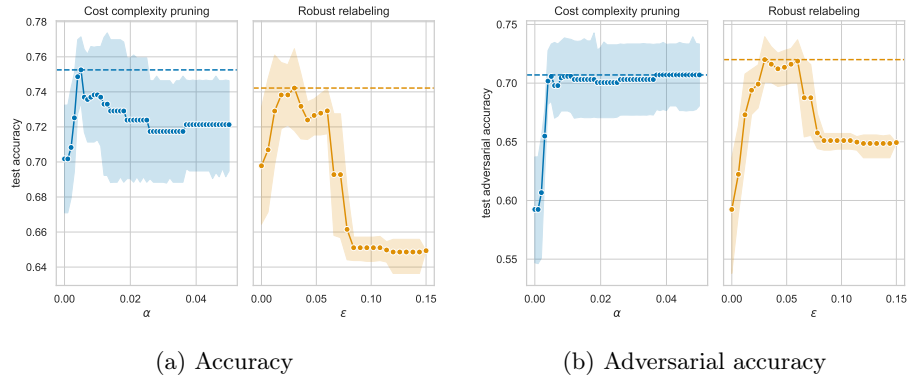


Fig. 7: 5-fold cross validation on the Pima-Indians-diabetes dataset.

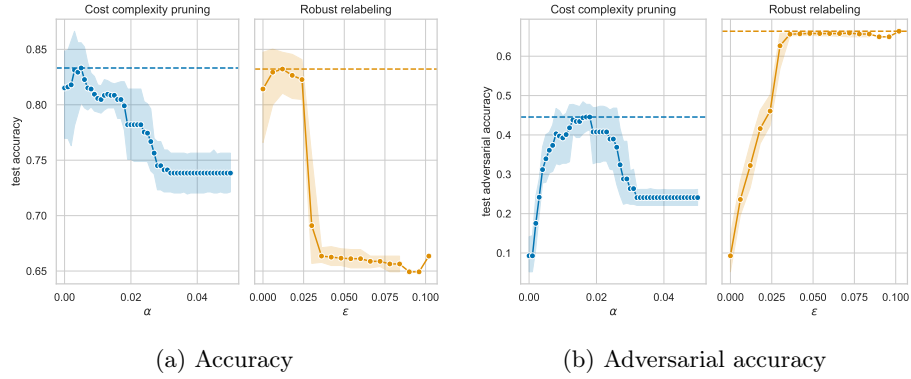


Fig. 8: 5-fold cross validation on the Qsar-biodegradation dataset.

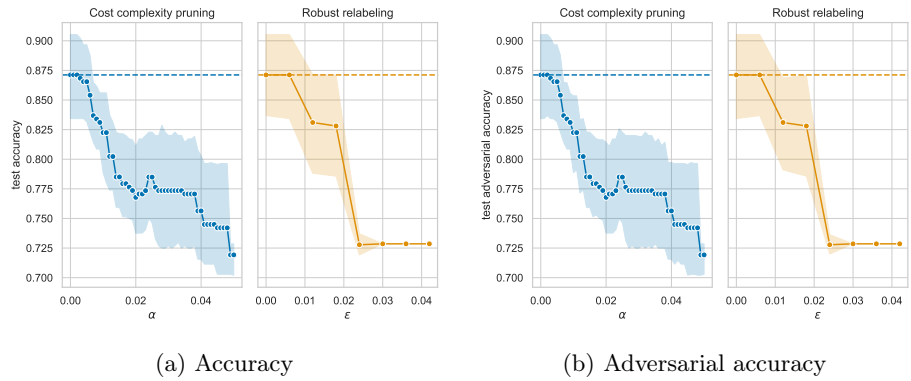
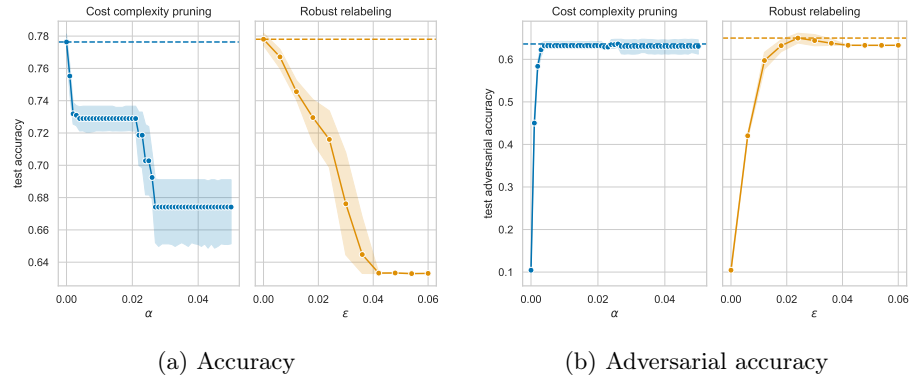


Fig. 9: 5-fold cross validation on the Spectf-heart dataset.



(a) Accuracy

(b) Adversarial accuracy

Fig. 10: 5-fold cross validation on the Wine dataset.