

# Learning Decision Trees with Reinforcement Learning: Supplementary Material

Due to the page limit of workshop paper, only a part of the experimental results are illustrated in the main body. This supplementary material demonstrates the complete experimental results.

## A. Results of Experiment 1

Experiment 1 has been introduced in the paper. All the datasets are split into training, validation and test set in proportion of 50%, 25%, 25%. Table 1 shows the final AUC score of RDT and CART on the validation and test set.

Table 1: Final AUC score of RDT and CART.

Datasets	Validation Set			Test Set			Search Space	Iters
	CART	RDT Avg	RDT Top	CART	RDT Avg	RDT Top		
Pima	80.47	<b>80.93 <math>\pm</math> 0.44</b>	81.54	79.09	<b>79.96 <math>\pm</math> 1.70</b>	78.90	$10^{13}$	3000
Heart	91.04	<b>94.73 <math>\pm</math> 1.09</b>	96.98	82.41	<b>85.18 <math>\pm</math> 1.71</b>	87.19	$10^{19}$	2000
Breast	99.18	<b>99.78 <math>\pm</math> 0.27</b>	99.99	95.15	<b>95.22 <math>\pm</math> 2.06</b>	95.93	$10^{45}$	2000
German	68.59	<b>74.84 <math>\pm</math> 0.96</b>	76.19	68.05	<b>73.01 <math>\pm</math> 1.36</b>	74.43	$10^{42}$	2000
HTRU	96.91	<b>96.92 <math>\pm</math> 0.17</b>	97.09	96.12	<b>96.70 <math>\pm</math> 0.23</b>	97.03	$10^{13}$	2000
Credit	75.79	<b>75.88 <math>\pm</math> 0.07</b>	75.98	75.22	<b>75.23 <math>\pm</math> 0.29</b>	75.44	$10^{42}$	3000

Figure 1 shows that RDT’s ROC score on the validation set improves over time and exceeds CART baseline on all the datasets.

Figure 2 demonstrates AUC score on the test set of RDT and CART. The results imply that our method slightly outperforms the baseline on all the datasets.

## B. Results of Experiment 2

In experiment 2, the samples in each dataset are shuffled randomly and split into training, validation and test set in proportion of 60%, 20%, 20%. This experiment is designed to test the reproducibility of the results in experiment 1. The final AUC score of RDT and CART on the validation and test set is shown in Table 2.

Table 2: Final AUC score of RDT and CART.

Datasets	Validation Set			Test Set			Search Space	Iters
	CART	RDT Avg	RDT Top	CART	RDT Avg	RDT Top		
Pima	84.46	<b>87.14 <math>\pm</math> 1.16</b>	89.10	73.61	<b>76.22 <math>\pm</math> 1.96</b>	77.84	$10^{13}$	3000
Heart	87.29	<b>93.46 <math>\pm</math> 0.66</b>	93.89	85.62	<b>87.35 <math>\pm</math> 2.04</b>	89.34	$10^{19}$	2000
Breast	96.72	<b>98.91 <math>\pm</math> 0.58</b>	99.61	94.51	<b>95.61 <math>\pm</math> 1.52</b>	97.52	$10^{45}$	2000
German	73.08	<b>77.97 <math>\pm</math> 0.95</b>	79.31	69.18	<b>71.90 <math>\pm</math> 1.22</b>	73.60	$10^{42}$	2000
HTRU	97.55	<b>97.56 <math>\pm</math> 0.04</b>	97.61	<b>97.36</b>	97.34 $\pm$ 0.00	97.34	$10^{13}$	2000
Credit	75.89	<b>75.93 <math>\pm</math> 0.14</b>	76.14	<b>75.10</b>	74.74 $\pm$ 0.05	74.68	$10^{42}$	3000

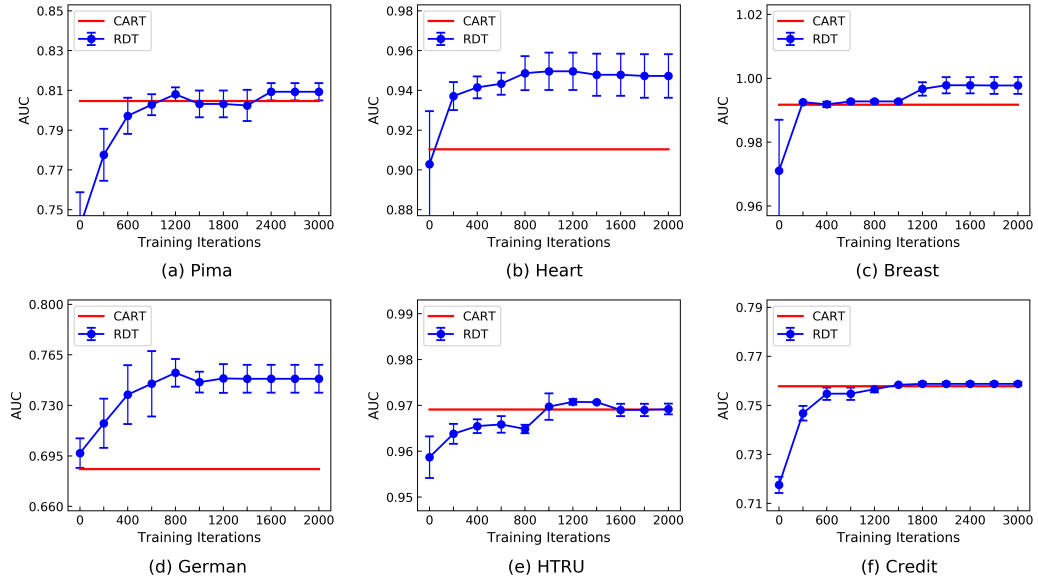


Figure 1: AUC score on the validation set of RDT and CART in experiment 1.

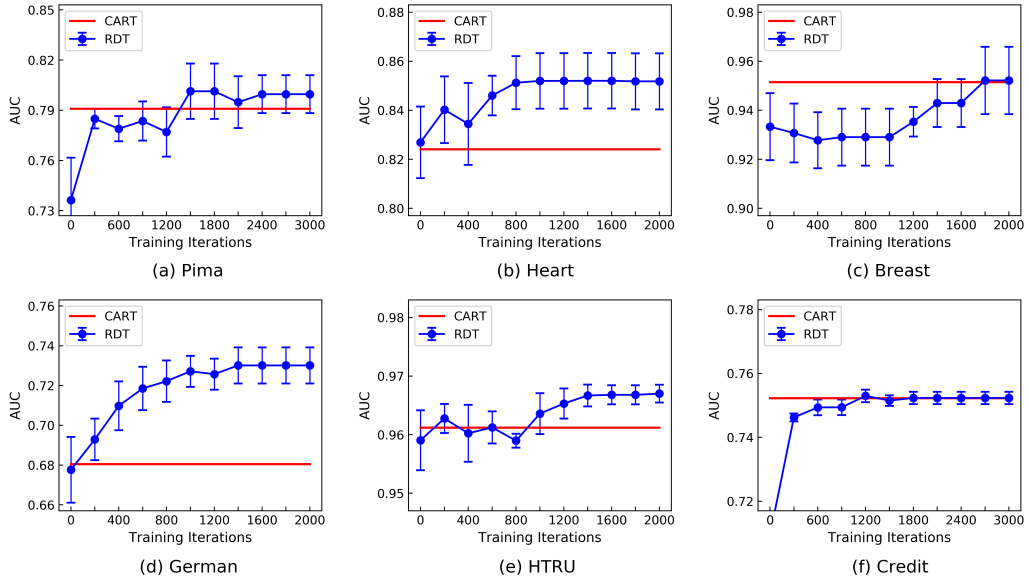


Figure 2: AUC score on the test set of RDT and CART in experiment 1.

Figure 3 and Figure 4 demonstrate AUC score of RDT and CART on the validation and test set respectively. The results demonstrate that RDT acquires competitive performance compared to CART, which is consistent with the conclusion in experiment 1 and validates the effectiveness of our method.

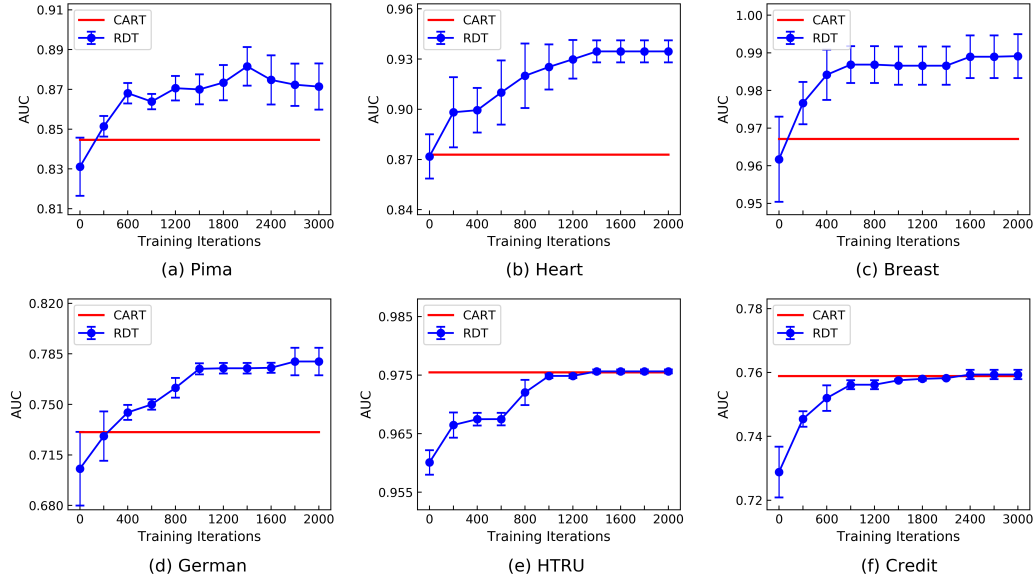


Figure 3: AUC score on the validation set of RDT and CART in experiment 2.

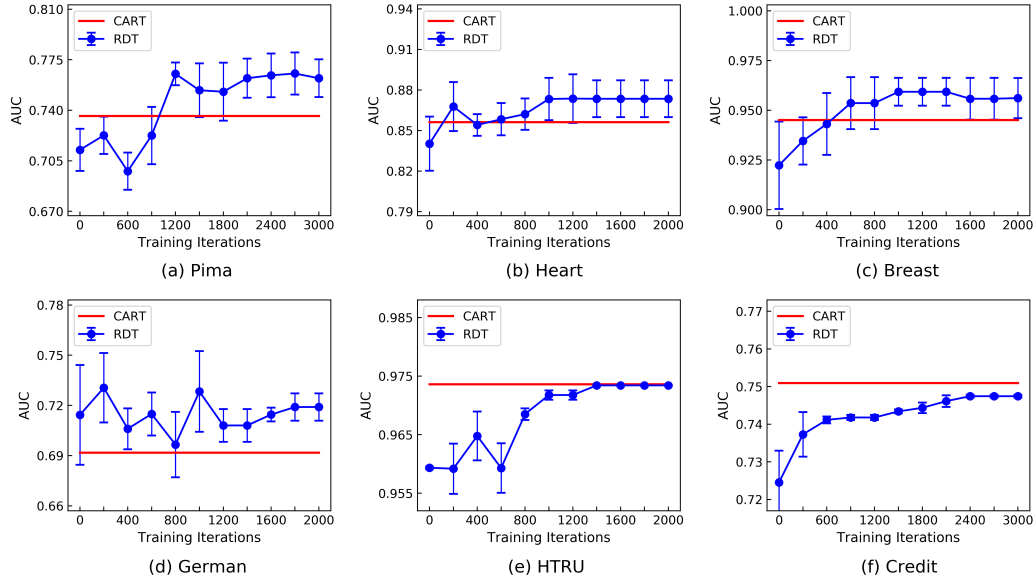


Figure 4: AUC score on the test set of RDT and CART in experiment 2.

### C. Design of Reward Signal

In the experiment, we find that if we directly use the performance score on the validation set as reward signal, the controller may find decision trees which perform extremely well on the validation set but poorly on the test set. This problem is more significant on datasets with relatively small sample number than on large datasets. A possible explanation to this phenomenon is that as the sample number decreases, the variance of the performance score on the validation set will increase, which makes the reward signal more noisy and leads the controller to overfit on the validation set. To solve this problem, we design two methods to improve the quality of the reward signal.

Firstly, we randomly choose 80% of the samples in the validation set and calculate the performance score on this subset. We repeat this process for 5 times and use their mean value as the final score, which can help reduce the variance of the performance score on the validation set.

Secondly, as a model is more likely to perform well on unseen test data if it has achieved good performance on both training and validation set, we design reward signal  $R$  as the F1 score of the performance score on both training and validation set:

$$R = \frac{2 \times PV \times PT}{PV + PT} \quad (1)$$

where PT is the performance score on the training set, and PV is the performance score on the validation set calculated in the way presented in the last paragraph.