

Weighted Oblique Decision Trees*

Bin-Bin Yang, Song-Qing Shen, Wei Gao

National Key Laboratory for Novel Software Technology
Nanjing University, Nanjing, 210023, China
{yangbb, shensq, gaow}@lamda.nju.edu.cn

Abstract

Decision trees have attracted much attention during the past decades. Previous decision trees include axis-parallel and oblique decision trees; both of them try to find the best splits via exhaustive search or heuristic algorithms in each iteration. Oblique decision trees generally simplify tree structure and take better performance, but are always accompanied with higher computation, as well as the initialization with the best axis-parallel splits. This work presents the Weighted Oblique Decision Tree (WODT) based on continuous optimization with random initialization. We consider different weights of each instance for child nodes at all internal nodes, and then obtain a split by optimizing the continuous and differentiable objective function of weighted information entropy. Extensive experiments show the effectiveness of the proposed algorithm.

Introduction

Decision trees have attracted much attention in many real applications such as computer vision (Bosch, Zisserman, and Munoz 2007) and information retrieval (Fuhr and Pfeifer 1994). The classical decision trees include CART (Breiman et al. 1984), ID3 (Quinlan 1986), C4.5 (Quinlan 1993), etc. This motivates a series of studies (Murthy, Kasif, and Salzberg 1994; Loh and Shih 1997; Breiman 2001; Geurts, Ernst, and Wehenkel 2006; Shotton et al. 2013b; Fan 2016; Abuzaid et al. 2016; Zhou and Feng 2017). Recent years have witnessed an increasing popularity on decision trees, for example, Microsoft Kinect makes real time human pose estimation from single depth images by decision trees trained on millions of examples (Shotton et al. 2013a).

The basic idea of decision trees is to separate data with some certain splitting criterion, recursively. This procedure requires an optimization at each internal node of the tree, which partitions the training data in the node into subsets according to some splitting criteria, such as information gain (Quinlan 1993) or Gini impurity index (Breiman et al. 1984). A large number of decision trees have been developed to exploit univariate split functions, according to the feature value below some threshold or not. We call it axis-parallel decision tree since the split at each node can be viewed as

an axis-parallel hyperplane in the feature space, and these trees have made successful applications (Cicalese, Laber, and Saettler 2014; Shotton et al. 2013a).

The axis-parallel decision trees may yield complex tree structure and increase computational cost, when decision boundaries are not parallel to axes. Hence, an oblique split is introduced to make a multivariate linear combination of features followed by binary quantization. Generally, oblique decision trees simplify tree structure and achieve better performance (Breiman et al. 1984; Heath, Kasif, and Salzberg 1993; Brodley and Utgoff 1995; Loh and Shih 1997; Amasyali and Ersoy 2008; Robertson, Price, and Reale 2013; Kotschieder et al. 2015). While training oblique decision trees is always followed with high running-time cost, as well as the initialization with best axis-parallel splits (Murthy, Kasif, and Salzberg 1994; Norouzi et al. 2015).

This work introduces another oblique decision tree based on continuous optimization with random initialization. The main contributions can be summarized as follows:

- Motivated from weighted information entropy, we consider different weights of each instance for child nodes at all internal nodes, and then obtain a split by optimizing a continuous and differentiable objective function. Here we introduce the ‘soft’ splitting instead of ‘hard’ splitting to tackle the intractability for continuous optimization. Our method proceeds with random initialization, whereas previous oblique decision trees require initialization with the best axis-parallel splits.
- Extensive experiments show that our method simplifies tree structure, and achieves significantly better performance than state-of-the-art algorithms of decision trees. Our experiments also show that previous oblique decision trees can not obtain small-size trees without the best axis-parallel splitting initialization. Moreover, our method takes relatively less running-time cost, especially for datasets of dimensionality larger than 200, such as `usps`, `protein` and `mnist`. We finally analyze the tree depths of the proposed WODT method.

The rest of this work is organized as follows: we begin with relevant work and some preliminaries, and then propose our WODT method with empirical supports, and finally conclude with future work.

*Supported by the National Key R&D Program of China (2017YFB1001903), NSFC(61751306, 61876078).
Copyright © 2019, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Relevant Work

Decision trees have a long history from the early work (Messenger and Mandell 1972), which employed a measure of node impurity based on the distribution of class labels for each internal node. Quinlan (1993) and Breiman et al. (1984) introduced the famous C4.5 and CART decision trees based on entropy and Gini index, respectively. Loh and Shih (1997) presented a two-step QUEST tree by splitting each node with significance tests. Due to simplicity and interpretability, a large number of axis-parallel decision trees have been developed in the literature (Domingos and Hulten 2000; Zhou and Chen 2002; Geurts, Ernst, and Wehenkel 2006; Shotton et al. 2013a; Abuzaid et al. 2016), and decision trees have been used as base learners for ensemble algorithms such as boosting and bagging (Breiman 2001; Friedman 2001; Zhou 2012).

Oblique decision trees present another way to construct compact trees and achieve better performance, and the main difference lies in the splits of multivariate linear combinations over features. CART (Breiman et al. 1984) can be applied to oblique decision tree by optimizing the coefficients of oblique splits based on the coordinate descent method. Murthy, Kasif, and Salzberg (1994) made a refinement of CART to find a local optimum by multiple restarts and random perturbations. Some statistical techniques are suggested for oblique decision trees, such as least square method and linear discriminant analysis (Brodley and Utgoff 1995; Loh and Shih 1997; Bennett and Blue 2002; López-Chau et al. 2013).

There are also some heuristic oblique decision trees with good performance under proper assumptions (Amasyali and Ersoy 2008; Manwani and Sastry 2012). In addition, various models have been developed by combining neural networks with decision trees, but with complex structure and high computational cost (Strömberg, Zrida, and Isaksson 1991; Guo and Gelfand 1992; Setiono and Liu 1999; Kotschieder et al. 2015). Norouzi et al. (2015) proposed the oblique decision tree by optimizing a continuous loss, which upper bounds the empirical 0/1 loss with the-best-axis-parallel-split initialization. Our method utilizes the robust sigmoid function to obtain a continuous and differentiable objective function, and proceeds with random initialization.

Preliminaries

Let $\mathcal{X} \subset \mathbb{R}^d$ and $\mathcal{Y} = \{1, 2, \dots, C\}$ denote the instance and label space, respectively. Suppose that \mathcal{D} is an (unknown) underlying distribution over the product space $\mathcal{X} \times \mathcal{Y}$. Let $S_m = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)\}$ be a training data, where each example is drawn i.i.d. from the distribution \mathcal{D} .

An oblique decision tree is generally constructed as follows. An instance $\mathbf{x} \in \mathcal{X}$ is directed from the root of the tree down through internal nodes to a leaf node. Each leaf node specifies a distribution over the label space \mathcal{Y} , and each internal node performs a binary test by evaluating a split function $s_\theta(\mathbf{x}): \mathbb{R}^d \rightarrow \mathbb{R}$. If $s_\theta(\mathbf{x}) < 0$, then \mathbf{x} is directed to the left child node; and the right child otherwise. Here, $s_\theta(\mathbf{x}) = \theta^T \mathbf{x}$ is parameterized by $\theta \in \mathbb{R}^d$, and we further incorporate an offset parameter to obtain split functions of the form $\theta^T \mathbf{x} + b$ by appending a constant "1" to the feature vector.

For each internal node, we aim to find a parameter θ for a good oblique split. However, it is difficult to make direct continuous optimization w.r.t. θ , since the indicator function $I[s_\theta(\mathbf{x}) < 0]$ is a discontinuous and piecewise-constant function (Breiman et al. 1984; Murthy, Kasif, and Salzberg 1994; Norouzi et al. 2015). Here, $I[\cdot]$ denotes the indicator function, which returns 1 if the argument is true and 0 otherwise.

Our WODT Method

Motivated from weighted information entropy (Guaşu 1971), this section introduces another oblique decision tree based on instance weights. The basic idea is to consider different weights of each instance for child nodes when searching for split parameters, and obtain a good oblique split by optimizing a continuous and differentiable objective function. We use the 'soft' splitting instead of 'hard' splitting to tackle the intractability for gradient-based optimization.

At an internal node, let S be the set of training examples in this node. Without loss of generality, we assume

$$S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\} \quad \text{for } n \leq m.$$

Given split parameter θ , we compute the weight w.r.t. the left child node by selecting sigmoid function over the negative of the split function $s_\theta(\mathbf{x})$, i.e.,

$$\sigma(-s_\theta(\mathbf{x})) = \sigma(-\theta^T \mathbf{x}) = 1/(1 + e^{\theta^T \mathbf{x}}),$$

and the weight w.r.t. right child node is given by

$$1 - \sigma(-s_\theta(\mathbf{x})) = 1 - \sigma(-\theta^T \mathbf{x}) = \sigma(\theta^T \mathbf{x}).$$

Let S_L and S_R denote the set of training examples and the corresponding weights w.r.t. the left and right child nodes, respectively, i.e.,

$$S_L = \{((\mathbf{x}_i, y_i), w_i^L) | w_i^L = \sigma(-\theta^T \mathbf{x}_i), (\mathbf{x}_i, y_i) \in S\},$$

$$S_R = \{((\mathbf{x}_i, y_i), w_i^R) | w_i^R = \sigma(\theta^T \mathbf{x}_i), (\mathbf{x}_i, y_i) \in S\}.$$

Let W_L and W_R be the sum of weights w.r.t. the left and right child nodes, respectively, that is,

$$W_L(\theta) = \sum_{((\mathbf{x}_i, y_i), w_i^L) \in S_L} w_i^L,$$

$$W_R(\theta) = \sum_{((\mathbf{x}_i, y_i), w_i^R) \in S_R} w_i^R.$$

We further denote by W_L^k and W_R^k the sum of weights in left and right child nodes w.r.t. each class $k \in \mathcal{Y}$, respectively, by

$$W_L^k(\theta) = \sum_{((\mathbf{x}_i, y_i), w_i^L) \in S_L} I[y_i = k] w_i^L,$$

$$W_R^k(\theta) = \sum_{((\mathbf{x}_i, y_i), w_i^R) \in S_R} I[y_i = k] w_i^R.$$

We finally have the objective function $E(\theta)$ as

$$E(\theta) = W_L(\theta) H_L(\theta) + W_R(\theta) H_R(\theta),$$

where H_L and H_R are the left and right weighted information entropies, respectively. More precisely, we have

$$H_L(\theta) = - \sum_{k=1}^C \frac{W_L^k(\theta)}{W_L(\theta)} \log_2 \frac{W_L^k(\theta)}{W_L(\theta)},$$

$$H_R(\theta) = - \sum_{k=1}^C \frac{W_R^k(\theta)}{W_R(\theta)} \log_2 \frac{W_R^k(\theta)}{W_R(\theta)}.$$

Algorithm 1 InduceSubtree(S, D, d) of Weighted Oblique Decision Tree (WODT)

Input: Training data S , maximum tree depth D , depth of the current node d

- 1: Create a node p based on data S
- 2: **if** Examples in data S all belong to class $k \in \mathcal{Y}$ **then**
- 3: Node p is a leaf node labelled with class k
- 4: **end if**
- 5: **if** $d > D$ **then**
- 6: Node p is a leaf node labelled with the majority class $k' \in \mathcal{Y}$
- 7: **end if**
- 8: Calculate the split parameter θ by the L-BFGS algorithm according to Eqns. (1) and (2)
- 9: Obtain the training data for the left-child and right-child nodes according to Eqns.(3) and (4), respectively.
- 10: The left subtree of node p : InduceSubtree($L, D, d + 1$)
- 11: The right subtree of node p : InduceSubtree($R, D, d + 1$)

Output: A decision subtree with the root node p

The objective function $E(\theta)$ can be further expressed as

$$E(\theta) = W_L \log_2 W_L + W_R \log_2 W_R - \sum_{k=1}^C W_L^k \log_2 W_L^k - \sum_{k=1}^C W_R^k \log_2 W_R^k. \quad (1)$$

It is easy to observe that the objective function $E(\theta)$ is continuous and differentiable w.r.t. the split parameter θ , and we could make use of standard optimization techniques to find a good split parameter θ , such as gradient-descent and quasi-Newton method. It is important to calculate the gradient function $g(\theta)$ of our objective function $E(\theta)$, that is,

$$g(\theta) = dE(\theta)/d\theta.$$

We have

$$g(\theta) \cdot \ln 2 = (1 + \ln W_L) \frac{dW_L}{d\theta} + (1 + \ln W_R) \frac{dW_R}{d\theta} - \sum_{k=1}^C (1 + \ln W_L^k) \frac{dW_L^k}{d\theta} - \sum_{k=1}^C (1 + \ln W_R^k) \frac{dW_R^k}{d\theta}.$$

From $\sigma(-z) = 1 - \sigma(z)$, $\sigma'(z) = \sigma(z)(1 - \sigma(z))$, we have

$$\begin{aligned} \frac{dW_L}{d\theta} &= \sum_{i=1}^n \sigma(-\theta^T \mathbf{x}_i) [1 - \sigma(-\theta^T \mathbf{x}_i)] (-\mathbf{x}_i) \\ &= - \sum_{i=1}^n [1 - \sigma(\theta^T \mathbf{x}_i)] \sigma(\theta^T \mathbf{x}_i) \mathbf{x}_i = - \frac{dW_R}{d\theta}. \end{aligned}$$

For simplicity, we denote by

$$\beta_i = \sigma(\theta^T \mathbf{x}_i) [1 - \sigma(\theta^T \mathbf{x}_i)] \mathbf{x}_i,$$

and we have $dW_R/d\theta = \sum_{i=1}^n \beta_i$. Similarly, we have $dW_L^k/d\theta = -dW_R^k/d\theta$ and $dW_R^k/d\theta = \sum_{i=1}^n I[y_i = k] \beta_i$.

Table 1: Benchmark datasets

dataset	#instance	#feature	dataset	#instance	#feature
iris	150	4	satimage	6435	36
wine	178	13	usps	9298	256
glass	214	9	pendigits	10992	16
heart	270	13	letter	20000	16
breast	683	10	protein	24387	357
diabetes	768	8	shuttle	58000	9
vehicle	846	18	connect4	67557	126
fourclass	862	2	mnist	70000	780
segment	2310	19	ijcnn1	141691	22
dna	3186	180	cod-rna	331152	8

This follows that

$$\begin{aligned} g(\theta) \cdot \ln 2 &= \ln \frac{W_R}{W_L} \cdot \frac{dW_R}{d\theta} - \sum_{k=1}^C \ln \frac{W_R^k}{W_L^k} \cdot \frac{dW_R^k}{d\theta} \\ &= \ln \frac{W_R}{W_L} \cdot \sum_{i=1}^n \beta_i - \sum_{i=1}^n \beta_i \ln \frac{W_R^{y_i}}{W_L^{y_i}} \\ &= \sum_{i=1}^n \beta_i \ln \frac{W_R W_L^{y_i}}{W_L W_R^{y_i}}, \end{aligned}$$

which yields

$$g(\theta) = \log_2 \frac{W_R}{W_L} \cdot \sum_{i=1}^n \beta_i \log_2 \frac{W_L^{y_i}}{W_R^{y_i}}. \quad (2)$$

In the implementation, we can use vectorization methods to accelerate our WODT method according to Eqns. (1) and (2). We also optimize the objective function $E(\theta)$ based on the L-BFGS algorithm, where we initialize the parameter θ with random vectors. This is quite different from previous oblique decision trees which require the initialization with the best axis-parallel splits. The parameter θ is not initialized to the zero vector so as to avoid the zero gradient.

Given the split parameter θ , we partition data as follows:

$$L = \{(\mathbf{x}, y) \in S | \theta^T \mathbf{x} < 0\}, \quad (3)$$

$$R = \{(\mathbf{x}, y) \in S | \theta^T \mathbf{x} \geq 0\}. \quad (4)$$

As can be seen, we make use of the direction of split parameter θ to partition data, while the objective function $E(\theta)$ is related with norm and direction of parameter θ simultaneously. A natural idea is to make an additional constraint over the norm of θ and then utilize the projection or Lagrange multiplier as in the work (Norouzi et al. 2015). Here, we do not make any additional constrain since $\sigma(\theta^T \mathbf{x})$ will be approximated by floating-point numbers in the implementation of the proposed method.

Algorithm 1 presents the detailed description of the weighted oblique decision subtree. Given an instance $\mathbf{x} \in \mathcal{X}$ and weighted oblique decision tree of Algorithm 1, we predict the label of instance \mathbf{x} according to the leaf node at the end of the path traversed by instance \mathbf{x} .

Experiments

This section empirically evaluates our WODT method on extensive datasets. We begin with the experimental settings,

Table 2: Comparison of test accuracies (mean \pm std.) on benchmark datasets. \bullet/\circ indicates that WODT is significantly better/worse than the corresponding method (pairwise t -tests at 95% significance level). ‘N/A’ means that no results were obtained after running out 250000 seconds (about 3 days).

dataset	our WODT	APDT	CO2	CO2r	OC1	OC1r	CART-LC	CART-LCr
iris	.9733 \pm .0248	.9467 \pm .0499 \bullet	.9467 \pm .0540 \bullet	.9263 \pm .0526 \bullet	.9600 \pm .0442	.9400 \pm .0554 \bullet	.9467 \pm .0499 \bullet	.9000 \pm .1125 \bullet
wine	.9665 \pm .0323	.9271 \pm .0529 \bullet	.9271 \pm .0579 \bullet	.8818 \pm .0565 \bullet	.9213 \pm .0456 \bullet	.8876 \pm .1293 \bullet	.9494 \pm .0450	.9045 \pm .0529 \bullet
glass	.6216 \pm .0256	.6521 \pm .1215	.6521 \pm .0517 \circ	.5837 \pm .0740 \bullet	.6168 \pm .1056	.6075 \pm .0810	.7056 \pm .0763 \circ	.6215 \pm .1142
heart	.7630 \pm .0429	.7370 \pm .0737	.7370 \pm .0300 \bullet	.7167 \pm .0446 \bullet	.7593 \pm .0668	.7815 \pm .0898	.6815 \pm .1285 \bullet	.7556 \pm .0529
breast	.9590 \pm .0099	.9466 \pm .0581	.9346 \pm .0458 \bullet	.9378 \pm .0190 \bullet	.9341 \pm .0608 \bullet	.9356 \pm .0639 \bullet	.9414 \pm .0271 \bullet	.9517 \pm .0173 \bullet
diabetes	.7161 \pm .0191	.7090 \pm .0486	.6908 \pm .0370 \bullet	.6953 \pm .0306 \bullet	.6667 \pm .0387 \bullet	.6914 \pm .0541 \bullet	.7161 \pm .0335	.6784 \pm .0379 \bullet
vehicle	.7069 \pm .0383	.7045 \pm .1603	.7223 \pm .0466	.6501 \pm .0180 \bullet	.6418 \pm .1568 \bullet	.6927 \pm .1050	.6702 \pm .1212	.7069 \pm .0402
fourclass	.9896 \pm .0043	.9843 \pm .0164	.9809 \pm .0191 \bullet	.8778 \pm .0203 \bullet	.9350 \pm .0875 \bullet	.9118 \pm .1139 \bullet	.9548 \pm .1283	.9664 \pm .0252 \bullet
segment	.9623 \pm .0099	.9660 \pm .0869	.9558 \pm .0417	.8653 \pm .0587 \bullet	.8355 \pm .1560 \bullet	.9143 \pm .0809 \bullet	.9580 \pm .0167	.9242 \pm .0214 \bullet
dna	.9250 \pm .0167	.9039 \pm .0241 \bullet	.8820 \pm .0162 \bullet	.7578 \pm .0178 \bullet	.8331 \pm .0114 \bullet	.8516 \pm .1031 \bullet	.8980 \pm .0176 \bullet	.8711 \pm .0102 \bullet
satimage	.8760 \pm .0097	.8485 \pm .0217 \bullet	.8480 \pm .0126 \bullet	.8139 \pm .0103 \bullet	.8485 \pm .0115 \bullet	.8159 \pm .0248 \bullet	.8510 \pm .1034	.8360 \pm .0177 \bullet
usps	.9058 \pm .0050	.8620 \pm .0341 \bullet	N/A	.5879 \pm .0300 \bullet	.8729 \pm .0180 \bullet	.7200 \pm .0320 \bullet	.6542 \pm .0772 \bullet	.5999 \pm .0258 \bullet
pendigits	.9660 \pm .0025	.9177 \pm .0849 \bullet	.9104 \pm .0987 \bullet	.7503 \pm .0343 \bullet	.9374 \pm .0042 \bullet	.9094 \pm .0608 \bullet	.9342 \pm .0062 \bullet	.9248 \pm .0028 \bullet
letter	.8786 \pm .0030	.8672 \pm .0777	.8041 \pm .0578 \bullet	.7215 \pm .0011 \bullet	.8142 \pm .0162 \bullet	.7764 \pm .0254 \bullet	.8530 \pm .0026 \bullet	.7688 \pm .0122 \bullet
protein	.5957 \pm .0052	.4887 \pm .0191 \bullet	.5218 \pm .0084 \bullet	.4285 \pm .0064 \bullet	.5406 \pm .0244 \bullet	.5606 \pm .0162 \bullet	.5248 \pm .0087 \bullet	.4770 \pm .0047 \bullet
shuttle	.9990 \pm .0002	.9999 \pm .0014	.9914 \pm .0135 \bullet	.9177 \pm .0115 \bullet	.9995 \pm .0011	.9982 \pm .0023	.9998 \pm .0005 \circ	.9986 \pm .0008 \bullet
connect4	.7415 \pm .0041	.7527 \pm .0261	.7131 \pm .0201 \bullet	.6457 \pm .0000 \bullet	.7312 \pm .0342	.7060 \pm .0189 \bullet	.7400 \pm .0108	.7249 \pm .0044 \bullet
mnist	.9434 \pm .0013	.8806 \pm .0101 \bullet	N/A	N/A	.7557 \pm .0240 \bullet	.7941 \pm .0203 \bullet	.8890 \pm .0051 \bullet	.8393 \pm .0147 \bullet
ijcnn1	.9703 \pm .0014	.9670 \pm .0020 \bullet	.9634 \pm .0011 \bullet	.9056 \pm .0002 \bullet	.9047 \pm .0045 \bullet	.9005 \pm .0082 \bullet	.9519 \pm .0033 \bullet	.9603 \pm .0023 \bullet
cod-rna	.9543 \pm .0016	.9433 \pm .0969	.8767 \pm .0663 \bullet	.6667 \pm .0016 \bullet	.7838 \pm .0363 \bullet	.9230 \pm .0606 \bullet	.8001 \pm .0017 \bullet	.9365 \pm .0014 \bullet
win/tie/loss		9/11/0	17/2/1	20/0/0	15/5/0	16/4/0	11/7/2	17/3/0

and then make empirical comparisons of our WODT method with state-of-the-art algorithms of decision trees. We further investigate tree sizes based on the cardinality of leaf node, and show the comparisons of running time. We finally analyze the training and generalization performance with respect to different tree depths.

Experimental Setting

We conduct our experiments on twenty benchmark datasets, as summarized in Table 1¹. Most datasets have been well-studied in previous studies on decision trees. The features have been scaled to $[-1, 1]$ for all datasets.

We compare our proposed WODT method with state-of-the-art algorithms of decisions tree as follows:

- CART-LC: CART for oblique decision trees with best axis-parallel splitting initialization (Breiman et al. 1984)
- CART-LCr: CART for oblique decision trees with random initialization (Breiman et al. 1984)
- OC1: Oblique decision trees induced by coordinate descent method, multiple restarts and random perturbations with best axis-parallel splitting initialization (Murthy, Kasif, and Salzberg 1994)
- OC1r: Oblique decision trees induced by coordinate descent method, multiple restarts and random perturbations with random initialization (Murthy, Kasif, and Salzberg 1994)
- CO2: Oblique decision trees induced by optimizing a continuous upper bound on the empirical loss with best axis-parallel splitting initialization (Norouzi et al. 2015)

- CO2r: Oblique decision trees induced by optimizing a continuous upper bound on the empirical loss with random initialization (Norouzi et al. 2015)

- APDT: axis-parallel decision trees (Quinlan 1993)

We implement the OC1 method as in the work of (Murthy, Kasif, and Salzberg 1994)², but slightly modify it so that the tree grows up to the fullest extent unless reaching the maximal tree depth. For the CO2 method, we excute 2 trials of 5-cv to select regularization parameter $\nu \in \{0.1, 1, 4, 10, 43, 100\}$ and the learning rate $\eta \in \{0.003, 0.01, 0.03\}$ as in the work (Norouzi et al. 2015). We take the default parameters as in the work of (Murthy, Kasif, and Salzberg 1994) for OC1, OC1r, CART-LC and CART-LCr. We also select information gain as the criterion for APDT, OC1, OC1r, CART-LC and CART-LCr. Our WODT method does not require additional hyper-parameter.

The performance of all methods are evaluated by 10 trials of 5-fold cross validation with different random seeds, and the performance is obtained by averaging over 50 runs. All experiments are performed on a node of computational cluster with 16 CPUs (Intel Xeon Core 3.0GHz) running RedHat Linux Enterprise 5 with 48GB main memory.

Experimental Results

Table 2 shows the test accuracy comparisons of our method with other methods. As can be seen, our WODT method significantly outperforms those oblique decision trees with random initialization, such as CO2r, OC1r and CART-LCr, since the win/tie/loss counts show that our WODT wins for most times and never loses. It is also observable that our WODT achieves better or comparable performance with CO2,

¹<http://www.ics.uci.edu/~mllearn/MLRepository.html>

²The codes of OC1 and CART-LC are downloaded from <http://ccb.jhu.edu/software/oc1/oc1.tar.gz>

Table 3: Comparison of leaves cardinality (mean \pm std.) on benchmark datasets. \bullet/\circ indicates that our WODT generates fewer/more leaf nodes than the corresponding method (pairwise t -tests at 95% significance level). ‘N/A’ means that no results were obtained after running out 250000 seconds (about 3 days), and we adopt the scientific notation $a \pm b(Ec)=a \times 10^c \pm b \times 10^c$.

dataset	our WODT	APDT	CO2	CO2r	OC1	OC1r	CART-LC	CART-LCr
iris	.5200 \pm .1170(E1)	.7800 \pm .1200(E1) \bullet	.8000 \pm .0200(E1) \bullet	.1500 \pm .0000(E2) \bullet	.5400 \pm .0699(E1)	.5589 \pm .1440(E2) \bullet	.6500 \pm .0900(E1)	.1080 \pm .0343(E3) \bullet
wine	.3400 \pm .0490(E1)	.7500 \pm .1100(E1) \bullet	.7500 \pm .5200(E1) \bullet	.2710 \pm .1095(E2) \bullet	.5100 \pm .1200(E1) \bullet	.3730 \pm .3080(E2) \bullet	.5100 \pm .1600(E1) \bullet	.8400 \pm .6500(E1) \bullet
glass	.4620 \pm .0240(E2)	.3620 \pm .0220(E2) \circ	.3620 \pm .0250(E2) \circ	.2976 \pm .0842(E3) \bullet	.8120 \pm .4630(E2) \bullet	.3566 \pm .0462(E3) \bullet	.5280 \pm .2380(E2)	.5901 \pm .0662(E3) \bullet
heart	.2320 \pm .0412(E2)	.3429 \pm .0279(E2) \circ	.3429 \pm .0520(E2) \circ	.9797 \pm .0862(E3) \bullet	.7120 \pm .4700(E2) \bullet	.1949 \pm .0467(E3) \bullet	.9090 \pm .3660(E2) \bullet	.3639 \pm .0598(E3) \bullet
breast	.1400 \pm .0210(E2)	.2320 \pm .0200(E2) \circ	.2390 \pm .0960(E2) \circ	.5391 \pm .1154(E3) \bullet	.3279 \pm .2910(E2) \bullet	.9020 \pm .4100(E2) \bullet	.4610 \pm .1889(E2) \bullet	.1943 \pm .0458(E3) \bullet
diabetes	.1082 \pm .0044(E3)	.1041 \pm .0047(E3) \circ	.1049 \pm .0226(E3)	.1841 \pm .0256(E4) \bullet	.1409 \pm .0940(E3)	.4521 \pm .1076(E3) \bullet	.1565 \pm .0367(E3) \bullet	.9519 \pm .0910(E3) \bullet
vehicle	.1140 \pm .0087(E3)	.1055 \pm .0031(E3) \circ	.2236 \pm .0105(E3) \circ	.1508 \pm .0172(E4) \bullet	.1840 \pm .0895(E3)	.6143 \pm .1993(E3) \bullet	.2148 \pm .0539(E3) \bullet	.1133 \pm .0098(E4) \bullet
fourclass	.1860 \pm .0242(E2)	.2170 \pm .0240(E2) \circ	.2320 \pm .0540(E2) \circ	.3289 \pm .1830(E2) \bullet	.2950 \pm .1920(E2)	.4360 \pm .2680(E2) \bullet	.1989 \pm .1250(E2)	.4520 \pm .2420(E2) \bullet
segment	.6419 \pm .0248(E2)	.5570 \pm .0279(E2) \circ	.2232 \pm .1011(E3) \circ	.3252 \pm .1256(E3) \bullet	.1108 \pm .0537(E3) \bullet	.3972 \pm .1203(E3) \bullet	.8800 \pm .2849(E2) \bullet	.7347 \pm .1109(E3) \bullet
dna	.1601 \pm .0170(E2)	.1194 \pm .0094(E3) \circ	.4130 \pm .1079(E3) \circ	.5216 \pm .1290(E3) \bullet	.1723 \pm .0147(E3) \bullet	.2904 \pm .0222(E3) \bullet	.1242 \pm .0211(E3) \bullet	.3371 \pm .0214(E3) \bullet
satimage	.2103 \pm .0121(E3)	.3298 \pm .0278(E3) \circ	.6102 \pm .0736(E3) \circ	.5791 \pm .0354(E4) \bullet	.5258 \pm .1346(E3) \bullet	.1913 \pm .0144(E4) \bullet	.5137 \pm .0719(E3) \bullet	.2818 \pm .0198(E4) \bullet
usps	.1500 \pm .0096(E3)	.3389 \pm .0517(E3) \bullet	N/A	.5295 \pm .0964(E4) \bullet	.5100 \pm .0073(E3) \bullet	.9744 \pm .0367(E3) \bullet	.6143 \pm .0974(E3) \bullet	.1227 \pm .0022(E4) \bullet
pendigits	.1494 \pm .0407(E3)	.2094 \pm .0597(E3) \bullet	.1508 \pm .0502(E3)	.3315 \pm .0760(E3) \bullet	.3970 \pm .0050(E3) \bullet	.1107 \pm .0150(E4) \bullet	.1941 \pm .0424(E3) \bullet	.1890 \pm .0121(E4) \bullet
letter	.1213 \pm .0023(E4)	.1752 \pm .0063(E4) \bullet	.1198 \pm .0167(E4)	.1787 \pm .0171(E4) \bullet	.2083 \pm .0100(E4) \bullet	.1056 \pm .0020(E5) \bullet	.1804 \pm .0122(E4) \bullet	.1610 \pm .0025(E5) \bullet
protein	.4572 \pm .0173(E3)	.2849 \pm .0050(E4) \bullet	.3044 \pm .0031(E4) \bullet	.4398 \pm .0003(E5) \bullet	.3329 \pm .1120(E4) \bullet	.1502 \pm .0054(E4) \bullet	.3165 \pm .0017(E4) \bullet	.1586 \pm .0058(E4) \bullet
shuttle	.8115 \pm .0363(E2)	.2739 \pm .0409(E2) \circ	.1459 \pm .0589(E3) \circ	.3091 \pm .0855(E3) \bullet	.3979 \pm .2100(E2) \circ	.4869 \pm .1146(E3) \bullet	.3410 \pm .0310(E2) \circ	.8264 \pm .0474(E3) \bullet
connect4	.2976 \pm .0069(E4)	.9700 \pm .0657(E4) \bullet	.1388 \pm .0054(E5) \bullet	.5242 \pm .0000(E5) \bullet	.1041 \pm .0296(E5) \bullet	.6458 \pm .2000(E4) \bullet	.1065 \pm .0074(E5) \bullet	.7490 \pm .0254(E4) \bullet
mnist	.7371 \pm .0270(E3)	.3329 \pm .0167(E4) \bullet	N/A	N/A	.2500 \pm .0119(E4) \bullet	.2032 \pm .0002(E4) \bullet	.3279 \pm .0118(E4) \bullet	.1910 \pm .0020(E4) \bullet
ijcnn1	.7180 \pm .0078(E3)	.8742 \pm .2904(E3) \bullet	.1566 \pm .0218(E4) \bullet	.2327 \pm .0016(E5) \bullet	.1724 \pm .1793(E3) \circ	.1006 \pm .0937(E4) \bullet	.1611 \pm .0181(E4) \bullet	.9714 \pm .0205(E4) \bullet
cod-rna	.7743 \pm .0031(E4)	.2751 \pm .0174(E4) \circ	.9501 \pm .5928(E3) \circ	.1000 \pm .0000(E3) \circ	.7283 \pm .0290(E4) \circ	.1240 \pm .0280(E5) \bullet	.9594 \pm .1264(E3) \circ	.1535 \pm .0121(E5) \bullet
win/tie/loss		14/0/6	15/3/2	19/0/1	13/4/3	20/0/0	15/3/2	20/0/0

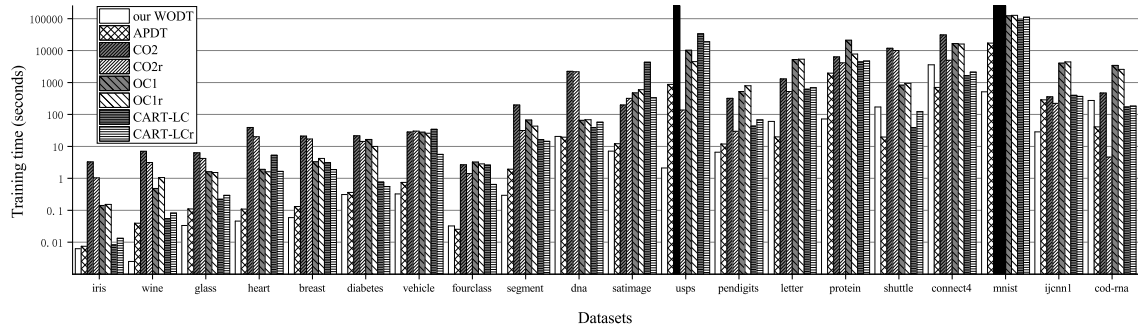


Figure 1: Comparison of the running time (in seconds) of WODT and other decision trees on benchmark datasets. Notice that the y -axis is in log-scale. Full black columns imply that no results were returned after running out 250000 seconds (about 3 days).

OC1 and CART-LC, though these oblique decision trees are accompanied with the best axis-parallel splitting initialization in the implementation. In comparison with the axis-parallel decision tree (APDT) method, our proposed WODT method achieves comparable performance for small-size datasets, and shows its superior for datasets of size larger than 5000, such as *satimage* and *pendigits*.

Table 3 shows the comparisons of leaves cardinality of our WODT with other methods. As can be seen, our WODT has fewer leaves than the oblique decision trees with random initialization such as CO2r, OC1r and CART-LCr, except for dataset *cod-rna*. We think that CO2r gets invalid partitions repeatedly until this method runs up to the maximal depth and returns with a deep and thin tree, which results in such few leaves. This analysis is also supported in Table 2 that our WODT method and CO2r achieve accuracy rates of 0.9543 and 0.6667, respectively.

It is also observable, from Table 3, that our WODT method has fewer, or as many leaves as those oblique decision trees with the best axis-parallel splitting initialization. Moreover, we notice that the previous oblique decision trees are not easy to obtain small-size trees without the best axis-parallel

splitting initialization, and the oblique decision trees with random initialization could produce more leaves, which may cause overfitting and take relatively poor performance as shown in Table 2.

We also compare the running time of WODT and the other decision tree methods, and the average CPU time (in seconds) is shown in Figure 1. As can be seen, our WODT takes comparable running time with the traditional axis-parallel decision tree (APDT) method. It is also observable that our WODT takes much less running time than the state-of-the-art oblique decision trees with the best axis-parallel splitting initialization or random initialization in most cases. In particular, our WODT is about 100 times faster than the other decision tree methods for the datasets of feature dimensionality larger than 200, such as *usps*, *protein* and *mnist*.

Depth Analysis

We finally analyze the influence of tree depths. Due to page limitation, we only present empirical results on four datasets *usps*, *protein*, *connect4* and *mnist*, while the trends are similar for the other datasets.

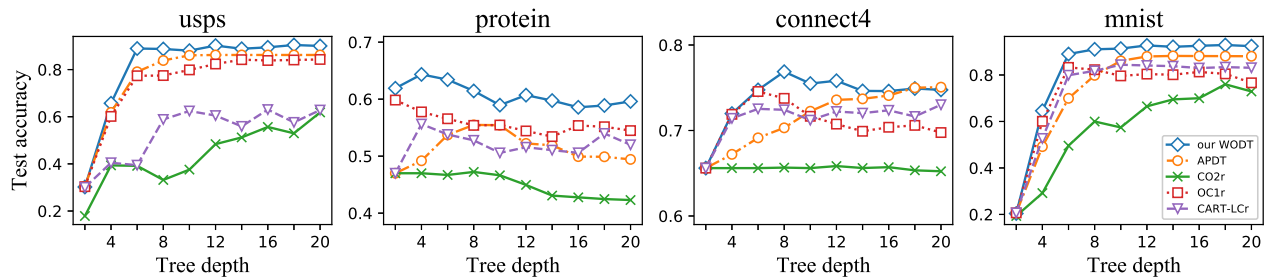


Figure 2: The influence on test accuracy with respect to tree depths for different methods.

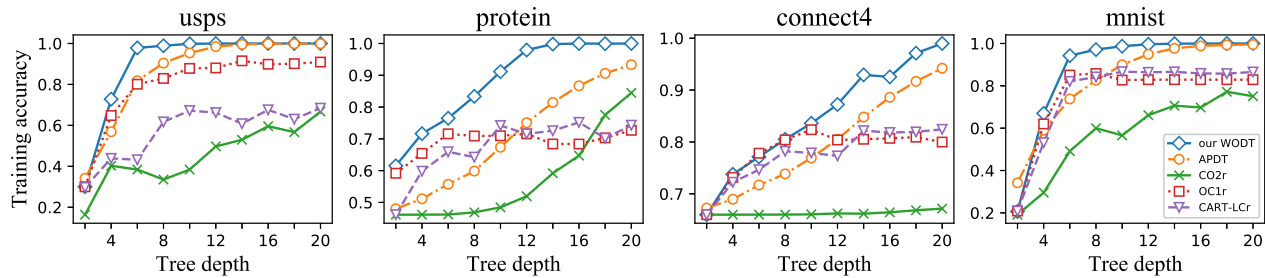


Figure 3: The influence on training accuracy with respect to tree depths for different methods.

Figure 2 shows the relationships between test accuracies and tree depths. Here the tree depth ranges from 2 to 20 with interval 2, and we compare our WODT with four decision trees: APDT, CO2r, OC1 and CART-LCr for simplicity, and the trends are similar to other methods such as CO2, OC1 and CART-LC. As can be seen, our WODT achieves the best performance at different depths in comparison with the other four methods, and it is also observable that our WODT method tends to achieve better generalization performance only with a shallow decision tree.

Figure 3 shows the relationships between the training accuracies and tree depths, where the range and compared methods are similar to those of Figure 2. As can be seen, our WODT method achieves better training accuracies than the others, which shows its stronger ability to fit data of our method. It is noteworthy that our WODT method and axis-parallel decision tree (APDT) method well fit dataset *usps* and *mnist* at a depth of 20 while WODT achieves better test accuracies as shown in Figure 2. This illustrates that optimizing the objective function in this work could yield a robust split by considering weighted entropy.

Conclusions

Oblique decision trees have attracted much attention during the past decades, and previous decision trees rely on the best axis-parallel splitting initialization with high computational cost. This work presents new Weighted Oblique Decision Tree (WODT). The basic idea is motivated from the weighted entropy, and we optimize the continuous and differentiable objective function with random initialization to find the split for each internal node. Extensive experiments show the ef-

fectiveness and robustness of our proposed method. In the future, an interesting work is to construct oblique decision trees based on ‘sparse’ splits by optimizing our objective function under the ℓ_1 penalty, which may present better interpretability and efficiency for prediction. Another interesting future work is to construct ensemble methods, such as random forests and boosting, based on our WODT method.

References

- Abuzaid, F.; Bradley, J. K.; Liang, F. T.; Feng, A.; Yang, L.; Zaharia, M.; and Talwalkar, A. S. 2016. Yggdrasil: An optimized system for training deep decision trees at scale. In *Advances in Neural Information Processing Systems 31*, 3817–3825.
- Amasyali, M. F., and Ersoy, O. 2008. Cline: A new decision-tree family. *IEEE Transactions on Neural Networks* 19(2):356–363.
- Bennett, K. P., and Blue, J. A. 2002. A support vector machine approach to decision trees. In *Proceedings of the IEEE International Joint Conference on Neural Networks*, 2396–2401.
- Bosch, A.; Zisserman, A.; and Munoz, X. 2007. Image classification using random forests and ferns. In *Proceedings of the IEEE International Conference on Computer Vision*, 1–8.
- Breiman, L. I.; Friedman, J. H.; Olshen, R. A.; and Stone, C. J. 1984. Classification and Regression Trees (CART). *Encyclopedia of Ecology* 40(3):582–588.
- Breiman, L. 2001. Random Forests. *Machine Learning* 45(1):5–32.

- Brodley, and Utgoff, P. 1995. Multivariate decision trees. *Machine Learning* 19(1):45–77.
- Cicalese, F.; Laber, E.; and Saettler, A. M. 2014. Diagnosis determination: Decision trees optimizing simultaneously worst and expected testing cost. In *Proceedings of the 31st International Conference on Machine Learning*, 414–422.
- Domingos, P., and Hulten, G. 2000. Mining high-speed data streams. In *Proceedings of the 6th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 71–80.
- Fan, L. 2016. Accurate robust and efficient error estimation for decision trees. In *Proceedings of the 33rd International Conference on Machine Learning*, 239–247.
- Friedman, J. H. 2001. Greedy function approximation: A gradient boosting machine. *Annals of Statistics* 29(5):1189–1232.
- Fuhr, N., and Pfeifer, U. 1994. Probabilistic information retrieval as a combination of abstraction, inductive learning, and probabilistic assumptions. *ACM Transactions on Information Systems* 12(1):92–115.
- Geurts, P.; Ernst, D.; and Wehenkel, L. 2006. Extremely randomized trees. *Machine Learning* 63(1):3–42.
- Guiasu, S. 1971. Weighted entropy. *Reports on Mathematical Physics* 2(3):165–179.
- Guo, H., and Gelfand, S. B. 1992. Classification trees with neural network feature extraction. *IEEE Transactions on Neural Networks* 3(6):923–933.
- Heath, D.; Kasif, S.; and Salzberg, S. 1993. Induction of oblique decision trees. In *Proceedings of the 13th International Joint Conference on Artificial Intelligence*, 1002–1007.
- Kontschieder, P.; Fiterau, M.; Criminisi, A.; and Rota Bulò, S. 2015. Deep neural decision forests. In *Proceedings of the IEEE International Conference on Computer Vision*, 1467–1475.
- Loh, W.-Y., and Shih, Y.-S. 1997. Split selection methods for classification trees. *Statistica Sinica* 7(4):815–840.
- López-Chau, A.; Cervantes, J.; López-García, L.; and Lamont, F. G. 2013. Fisher’s decision tree. *Expert Systems with Applications* 40(16):6283–6291.
- Manwani, N., and Sastry, P. S. 2012. Geometric decision tree. *IEEE Transactions on Systems, Man, and Cybernetics* 42(1):181–92.
- Messenger, R., and Mandell, L. 1972. A modal search technique for predictive nominal scale multivariate analysis. *Journal of the American Statistical Association* 67(340):768–772.
- Murthy, S. K.; Kasif, S.; and Salzberg, S. 1994. A system for induction of oblique decision trees. *Journal of Artificial Intelligence Research* 2(1):1–32.
- Norouzi, M.; Collins, M. D.; Fleet, D. J.; and Kohli, P. 2015. CO2 forest: Improved random forest by continuous optimization of oblique splits. *arXiv preprint arXiv:1506.06155*, 2015.
- Quinlan, J. R. 1986. Induction of decision trees. *Machine Learning* 1(1):81–106.
- Quinlan, J. R. 1993. *C4.5: Programs for Machine Learning*. Morgan Kaufmann Publishers.
- Robertson, B.; Price, C.; and Reale, M. 2013. CARTopt: A random search method for nonsmooth unconstrained optimization. *Computational Optimization and Applications* 56(2):291–315.
- Setiono, R., and Liu, H. 1999. A connectionist approach to generating oblique decision trees. *IEEE Transactions on Systems, Man, and Cybernetics* 29(3):440–444.
- Shotton, J.; Girshick, R.; Fitzgibbon, A.; Sharp, T.; Cook, M.; Moore, R.; Moore, R.; Kohli, P.; Criminisi, A.; and Kipman, A. 2013a. Efficient human pose estimation from single depth images. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35(12):2821–2840.
- Shotton, J.; Sharp, T.; Kohli, P.; Nowozin, S.; Winn, J.; and Criminisi, A. 2013b. Decision jungles: Compact and rich models for classification. In *Advances in Neural Information Processing Systems* 28, 234–242.
- Strömberg, J.-E.; Zrida, J.; and Isaksson, A. 1991. Neural trees-using neural nets in a tree classifier structure. In *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing*, 137–140.
- Zhou, Z.-H., and Chen, Z.-Q. 2002. Hybrid decision tree. *Knowledge-Based Systems* 15(8):515–528.
- Zhou, Z.-H., and Feng, J. 2017. Deep forest: Towards an alternative to deep neural networks. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, 3553–3559.
- Zhou, Z.-H. 2012. *Ensemble Methods: Foundations and Algorithms*. Chapman and Hall/CRC.