

# Orthogonal Decision Trees

Hillol Kargupta<sup>\*‡</sup>, Byung-Hoon Park<sup>†</sup>, Haimonti Dutta<sup>\*</sup>

<sup>\*</sup>Department of Computer Science and Electrical Engineering,

University of Maryland Baltimore County,

1000 Hilltop Circle,

Baltimore, MD 21250

Email: {hillol,hdutta1}@csee.umbc.edu

<sup>†</sup>Computer Science and Mathematics Division,

Oak Ridge National Laboratory,

PO BOX 2008 MS6164,

Oak Ridge, TN 37831-6164.

Email: parkbh@ornl.gov

<sup>‡</sup>The author is also affiliated to Agnik, LLC., Columbia, MD.

## Abstract

This paper introduces orthogonal decision trees that offer an effective way to construct a redundancy-free, accurate, and meaningful representation of large decision-tree-ensembles often created by popular techniques such as Bagging, Boosting, Random Forests and many distributed and data stream mining algorithms. Orthogonal decision trees are functionally orthogonal to each other and they correspond to the principal components of the underlying function space. This paper offers a technique to construct such trees based on Fourier transformation of decision trees and eigen-analysis of the ensemble in the Fourier representation. It offers experimental results to document the performance of orthogonal trees on grounds of accuracy and model complexity.

## Index Terms

Orthogonal Decision Trees, Redundancy Free Trees, Principle Component Analysis, Fourier Transform.

## I. INTRODUCTION

Decision tree [1] ensembles are frequently used in data mining and machine learning applications. Boosting [2], [3], Bagging[4], Stacking [5], and Random Forests [6] are some of the well-known ensemble-learning techniques. Many of these techniques often produce large ensembles that combine the outputs of a large number of trees for producing the overall output. Ensemble-based classification and outlier detection techniques are also frequently used in mining continuous data streams [7], [8]. Large ensembles pose several problems to a data miner. They are difficult to understand and the overall functional structure of the ensemble is not very “actionable” since it is difficult to manually combine the physical meaning of different trees in order to produce a simplified set of rules that can be used in practice. Moreover, in many time-critical applications such as monitoring data streams in resource-constrained environments [9], maintaining a large ensemble and using it for continuous monitoring are computationally challenging. So it will be useful if we can develop a technique to construct a redundancy-free meaningful compact representation of large ensembles. This paper offers a technique to do that and possibly more.

This paper presents a technique to construct redundancy-free decision-tree-ensembles by con-

structuring orthogonal decision trees. The technique first constructs an algebraic representation of trees using multi-variate discrete Fourier bases. The new representation is then used for eigen-analysis of the covariance matrix generated by the decision trees in Fourier representation. The proposed approach then converts the corresponding principal components to decision trees. These trees are defined in the original attributes-space and they are functionally orthogonal to each other. These orthogonal trees are in turn used for accurate (in many cases with improved accuracy) and redundancy-free (in the sense of orthogonal basis set) compact representation of large ensembles.

Section II presents the motivation of this work. Section III presents a brief overview of the Fourier spectrum of decision trees. Section IV describes the algorithms for computing the Fourier transform of a decision tree. Section V offers the algorithm for computing the tree from its Fourier spectrum. Section VI discusses orthogonal decision trees. Section VII presents experimental results using many well-known data sets. Finally, Section VIII concludes this paper.

## II. MOTIVATION

This paper extends our earlier work [10], [9], [11] on Fourier spectrum of decision trees. The main motivation behind this approach is to create an algebraic framework for meta-level analysis of models, produced by many ensemble learning, data stream mining, distributed data mining, and other related techniques. Most of the existing techniques treat the discrete model structures such as decision trees in an ensemble primarily as a black box. Only the output of the models are considered and combined in order to produce the overall output. Fourier bases offer a compact representation of a discrete structure that allows algebraic manipulation of decision trees. For example, we can literally add two different trees, produce weighted average of the trees themselves or perform eigen analysis of an ensemble of trees. Fourier representation of decision trees may offer something that is philosophically similar to what spectral representation of graphs [12] offers—an algebraic representation that allows deep analysis of discrete structures.

Fourier representation allows us to bring in the rich volume of well-understood techniques from Linear Algebra and Linear Systems Theory. This opens up many exciting possibilities for future research, such as quantifying the stability of an ensemble classifier, mining and monitoring mission-critical data streams using properties of the eigenvalues of the ensemble. This paper takes some steps toward achieving these goals.

The main contributions of this paper are listed below:

- 1) It offers several new analytical results regarding the properties of the Fourier spectra of decision trees.
- 2) It presents a detailed discussion on the Tree Construction from Fourier Spectrum (TCFS) algorithm for computing a decision tree from the Fourier coefficients. This includes discussion and experimental evaluation of the TCFS algorithm. New experimental results compare the performance of the trees constructed using the TCFS technique with that of the trees constructed using standard techniques such as C4.5.
- 3) It discusses Orthogonal Decision Trees (ODTs) in details and offers extensive experimental results documenting the performance of ODTs on benchmarked data sets.

The following section reviews the Fourier representation of decision trees.

### III. DECISION TREES AND THE FOURIER REPRESENTATION

This section reviews the Fourier representation of decision tree ensembles, introduced elsewhere [13], [14]. It also presents some new analytical results.

#### A. *Decision Trees as Numeric Functions*

The approach developed in this paper makes use of linear algebraic representation of the trees. In order to do that we first need to convert the tree into a numeric tree just in case the attributes are symbolic. A decision tree defined over a domain of categorical attributes can be treated as a numeric function. First note that a decision tree is a function that maps its domain members to a range of class labels. Sometimes, it is a symbolic function where attributes take symbolic (non-numeric) values. However, a symbolic function can be easily converted to a numeric function by simply replacing the symbols with numeric values in a consistent manner. Since the proposed approach of constructing orthogonal trees uses this representation as an intermediate stage and eventually the physical tree is converted back to the exact scheme for replacing the symbols (if any) does not matter as long as it is consistent.

Once the tree is converted to a discrete numeric function, we can also apply any appropriate analytical transformation as necessary. Fourier transformation is one such interesting possibility. Fourier representation of a function is a linear combination of the Fourier basis functions. The weights, called Fourier coefficients, completely define the representation. Each coefficient is

associated with a Fourier basis function that depends on a certain subset of features defining the domain. This section reviews the Fourier representation of decision tree ensembles, introduced elsewhere [9].

### B. A Brief Review of Multivariate Fourier Basis

Fourier basis set is comprised of orthogonal functions that can be used to represent any discrete function. In other words, it is a functionally complete representation. Consider the set of all  $\ell$ -dimensional feature vectors where the  $i$ -th feature can take  $\lambda_i$  different discrete values. The Fourier basis set that spans this space is comprised of  $\prod_{i=0}^{\ell} \lambda_i$  basis functions. Each Fourier basis function is defined as,

$$\psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x}) = \frac{1}{\sqrt{\prod_{i=1}^{\ell} \lambda_i}} \prod_{m=1}^{\ell} \exp \frac{2\pi i}{\lambda_m} x_m j_m$$

where  $\mathbf{j}$  and  $\mathbf{x}$  are vectors of length  $\ell$ ;  $x_m$  and  $j_m$  are  $m$ -th attribute-value in  $\mathbf{x}$  and  $\mathbf{j}$ , respectively;  $x_m, j_m \in \{0, 1, \dots, \lambda_i\}$  and  $\bar{\lambda}$  represents the feature-cardinality vector,  $\lambda_0, \dots, \lambda_{\ell}$ ;  $\psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$  is called the  $\mathbf{j}$ -th basis function. The vector  $\mathbf{j}$  is called a *partition*, and the *order* of a partition  $\mathbf{j}$  is the number of non-zero feature values it contains. A Fourier basis function depends on some  $x_i$  only when the corresponding  $j_i \neq 0$ . If a partition  $\mathbf{j}$  has exactly  $\alpha$  number of non-zeros values, then we say the partition is of order  $\alpha$  since the corresponding Fourier basis function depends only on those  $\alpha$  number of variables that take non-zero values in the partition  $\mathbf{j}$ .

A function  $f : \mathbf{X}^{\ell} \rightarrow \mathfrak{R}$ , that maps an  $\ell$ -dimensional discrete domain to a real-valued range, can be represented using the Fourier basis functions:  $f(\mathbf{x}) = \sum_{\mathbf{j}} w_{\mathbf{j}} \psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$ . where  $w_{\mathbf{j}}$  is the Fourier Coefficient (FC) corresponding to the partition  $\mathbf{j}$  and  $\bar{\psi}_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$  is the complex conjugate of  $\psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$ ;  $w_{\mathbf{j}} = \sum_{\mathbf{x}} \bar{\psi}_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x}) f(\mathbf{x})$ . The Fourier coefficient  $w_{\mathbf{j}}$  can be viewed as the relative contribution of the partition  $\mathbf{j}$  to the function value of  $f(\mathbf{x})$ . Therefore, the absolute value of  $w_{\mathbf{j}}$  can be used as the ‘‘significance’’ of the corresponding partition  $\mathbf{j}$ . If the magnitude of some  $w_{\mathbf{j}}$  is very small compared to other coefficients, we may consider the  $\mathbf{j}$ -th partition to be insignificant and neglect its contribution. The *order* of a Fourier coefficient is nothing but the order of the corresponding partition. We shall often use terms like *high order* or *low order* coefficients to refer to a set of Fourier coefficients whose orders are relatively large or small respectively. Energy of a spectrum is defined by the summation  $\sum_{\mathbf{j}} w_{\mathbf{j}}^2$ . Let us also define the inner product between two spectra  $\mathbf{w}_{(1)}$  and  $\mathbf{w}_{(2)}$  where  $\mathbf{w}_{(i)} = [w_{(i),1}, w_{(i),2}, \dots, w_{(i),|J|}]^T$  is the column matrix of all Fourier coefficients

in an arbitrary but fixed order. Superscript  $T$  denotes the transpose operation and  $|J|$  denotes the total number of coefficients in the spectrum. The inner product,  $\langle \mathbf{w}_{(1)}, \mathbf{w}_{(2)} \rangle = \sum_{\mathbf{j}} w_{(1),\mathbf{j}} w_{(2),\mathbf{j}}$ . We will also use the definition of the inner product between a pair of real-valued functions defined over some domain  $\Omega$ . This is defined as  $\langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle = \sum_{\mathbf{x} \in \Omega} f_1(\mathbf{x}) f_2(\mathbf{x})$ .

The following section considers the Fourier spectrum of decision trees and discusses some of its useful properties.

### C. Properties of Decision Trees in the Fourier Domain

For almost all practical purposes decision trees have bounded depths. This section will therefore consider decision trees of finite depth bounded by some constant. The underlying functions in such decision trees are computable by a constant depth Boolean AND and OR circuit (or equivalently  $AC^0$  circuit). Linial et al. [15] noted that the Fourier spectrum of  $AC^0$  circuit has very interesting properties and proved the following lemma.

*Lemma 1: (Linial, 1993)* Let  $M$  and  $d$  be the size and depth of an  $AC^0$  circuit. Then

$$\sum_{\{\mathbf{j} \mid o(\mathbf{j}) > t\}} w_{\mathbf{j}}^2 \leq 2M2^{-t^{1/d}/20}$$

where  $o(\mathbf{j})$  denotes the order (the number of non-zero variable) of partition  $\mathbf{j}$  and  $t$  is a non-negative integer. The term on the left hand side of the inequality represents the energy of the spectrum captured by the coefficients with order greater than a given constant  $t$ .

The lemma essentially states the following properties about decision trees:

- 1) High order Fourier coefficients are small in magnitude.
- 2) The energy preserved in all high order Fourier coefficients is also small.

The key aspect of these properties is that the energy of the Fourier coefficients of higher order decays exponentially. This observation suggests that the spectrum of a Boolean decision tree (or equivalently bounded depth function) can be approximated by computing only a small number of low order Fourier coefficients. So Fourier basis offers an efficient numeric representation of a decision tree in terms of an algebraic function that can be easily stored and manipulated.

The exponential decay property of Fourier spectrum also holds for non-Boolean decision trees. The complete proof is given in the appendix which is available as supplementary material from the publisher.

There are two additional important characteristics of the Fourier spectrum of a decision tree that we will use in this paper:

- 1) The Fourier spectrum of a decision tree can be efficiently computed [9].
- 2) The Fourier spectrum can be directly used for constructing the tree.

In other words, we can go back and forth between the tree and its spectrum. This is philosophically similar to the switching between the time and frequency domains in the traditional application of Fourier analysis for signal processing. These two issues will be discussed in details later in this paper. However, before that we would like to make a note of one additional property.

Fourier transformation of decision trees preserves inner product. The functional behavior of a decision tree is defined by the class labels it assigns. Therefore, if  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{|\Omega|}\}$  are the members of the domain  $\Omega$  then the functional behavior of a decision tree  $f(\mathbf{x})$  can be captured by the vector  $[f]_{x \in \Omega} = [f(\mathbf{x}_1) f(\mathbf{x}_2) \dots f(\mathbf{x}_{|\Omega|})]^T$ , where the superscript  $T$  denotes the transpose operation. The following lemma proves that the inner product between two such vectors is identical to the same in between their respective Fourier spectra.

*Lemma 2:* Given two functions  $f_1(\mathbf{x}) = \sum_{\mathbf{j}} w_{(1),\mathbf{j}} \overline{\psi_{\mathbf{j}}^{\lambda}}(\mathbf{x})$  and  $f_2(\mathbf{x}) = \sum_{\mathbf{j}} w_{(2),\mathbf{j}} \overline{\psi_{\mathbf{j}}^{\lambda}}(\mathbf{x})$  in Fourier representation. Then  $\langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle = \langle \mathbf{w}_{(1)}, \mathbf{w}_{(2)} \rangle$ .

**Proof:**

$$\begin{aligned} \langle f_1(\mathbf{x}), f_2(\mathbf{x}) \rangle &= \sum_{\mathbf{x} \in \Omega} f_1(\mathbf{x}) f_2(\mathbf{x}) = \sum_{\mathbf{x} \in \Omega} \sum_{\mathbf{j}, \mathbf{i}} w_{(1),\mathbf{j}} \overline{\psi_{\mathbf{j}}^{\lambda}}(\mathbf{x}) w_{(2),\mathbf{i}} \overline{\psi_{\mathbf{i}}^{\lambda}}(\mathbf{x}) \\ &= \sum_{\mathbf{j}, \mathbf{i}} w_{(1),\mathbf{j}} w_{(2),\mathbf{i}} \sum_{\mathbf{x} \in \Omega} \overline{\psi_{\mathbf{j}}^{\lambda}}(\mathbf{x}) \overline{\psi_{\mathbf{i}}^{\lambda}}(\mathbf{x}) = \sum_{\mathbf{j}} w_{(1),\mathbf{j}} w_{(2),\mathbf{j}} = \langle \mathbf{w}_{(1)}, \mathbf{w}_{(2)} \rangle . \end{aligned}$$

■

The fourth step is true since Fourier basis functions are orthonormal.

#### IV. COMPUTING THE FOURIER TRANSFORM OF A DECISION TREE

The Fourier spectrum of a given tree can be computed efficiently by traversing the tree. This section first reviews an algorithm to do that. It discusses aggregation of the multiple spectra computed from the base classifiers of an ensemble. It also extends the technique for dealing with non-Boolean class labels. Kushilevitz and Mansour [16] considered the issue of learning the low order Fourier spectrum of the target function (represented by a Boolean decision tree)

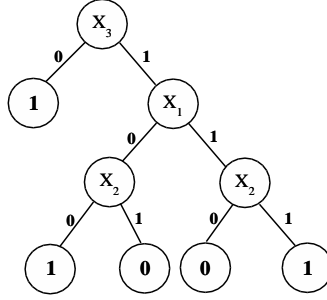


Fig. 1. A Boolean decision tree.

from a data set with uniformly distributed observations. Note that the current contribution is fundamentally different from their goal. This paper does not try to learn the spectrum directly from the data. Rather it considers the problem of computing the spectrum from the decision tree generated from the data.

#### A. Schema Representation of a Decision Path

For the sake of simplicity, let us consider a Boolean decision tree as shown in Figure 1. The Boolean class labels correspond to positive and negative instances of the concept class. We can express a Boolean decision tree as a function  $f : X^\ell \rightarrow \{0, 1\}$ . The function  $f$  maps positive and negative instances to one and zero respectively. A node in a tree is labeled with a feature  $x_i$ . A downward link from the node  $x_i$  is labeled with an attribute value of the  $i$ -th feature. The path from the root node to a successor node represents the subset of data that satisfies the different feature values labeled along the path. These subsets of the domain are essentially similarity-based equivalence classes and we shall call them *schemata* (schema in singular form). If  $\mathbf{h}$  is a schema, then  $\mathbf{h} \in \{0, 1, *\}^\ell$ , where  $*$  denotes a wildcard that matches any value of the corresponding feature. For example, the path  $\{(x_3 \xrightarrow{1} x_1, x_1 \xrightarrow{0} x_2)\}$  in Figure 1 represents the schema  $0 * 1$ , since all members of the data subset at the final node of this path take feature values 0 and 1 for  $x_1$  and  $x_3$  respectively. We shall use the term *order* to represent the number of non-wildcard values in a schema. The following section describes an algorithm to extract Fourier coefficients from a tree.



## B. Extracting and Calculating Significant Fourier Coefficients from a Tree

Considering a decision tree as a function, the Fourier transform of a decision tree can be defined as:

$$\begin{aligned} w_j &= \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \psi_j(\mathbf{x}) = \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in S_{l_1}} f(\mathbf{x}) \psi_j(\mathbf{x}) + \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in S_{l_2}} f(\mathbf{x}) \psi_j(\mathbf{x}) + \dots + \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in S_{l_n}} f(\mathbf{x}) \psi_j(\mathbf{x}) \\ &= \frac{|S_{l_1}|}{|\Lambda|} f(\mathbf{h}_1) \psi_j(\mathbf{h}_1) + \frac{|S_{l_2}|}{|\Lambda|} f(\mathbf{h}_2) \psi_j(\mathbf{h}_2) + \dots + \frac{|S_{l_n}|}{|\Lambda|} f(\mathbf{h}_n) \psi_j(\mathbf{h}_n) \end{aligned} \quad (1)$$

Where  $\Lambda$  denotes the complete instance space,  $S_{l_i}$  is an instance subspace which  $i^{th}$  leaf node  $l_i$  covers and  $\mathbf{h}_i$  is a schema defined by a path to  $l_i$  respectively (Note that any path to a node in a decision tree is essentially a subspace or hyperplane, thus it is a schema).

*Lemma 3:* For any Fourier basis function  $\psi_j$ ,  $\sum_{\mathbf{x} \in \Lambda} \psi_j(\mathbf{x}) = 0$ .

**Proof:** Since Fourier basis functions form an orthogonal set,

$$\sum_{\mathbf{x} \in \Lambda} \psi_j(\mathbf{x}) = \sum_{\mathbf{x} \in \Lambda} \psi_0(\mathbf{x}) \psi_j(\mathbf{x}) = 0.$$

Here,  $\psi_0$  is the zero-th Fourier basis function, which is constant (one) for all  $\mathbf{x}$ .

■

*Lemma 4:* Let  $\mathbf{h}_i$  be a schema defined by the path to a leaf node  $l_i$ . Then if  $\mathbf{j}$  has a non-zero attribute value at a position where  $\mathbf{h}_i$  has no value (wild-card),

$$\sum_{\mathbf{x} \in S_{l_i}} f(\mathbf{x}) \psi_j(\mathbf{x}) = f(\mathbf{h}_i) \sum_{\mathbf{x} \in S_{l_i}} \psi_j(\mathbf{x}) = 0.$$

Where  $S_{l_i}$  is the subset that  $\mathbf{h}_i$  covers.

**Proof:** Let  $\mathbf{j} = (\mathbf{j}_{in}, \mathbf{j}_{out})$ , where  $\mathbf{j}_{in}$  are features which are included  $h_i$  and  $\mathbf{j}_{out}$  are features not in  $h_i$  respectively. Since all values for  $\mathbf{j}_{in}$  are fixed in  $\mathbf{h}_i$ ,  $\psi_{\mathbf{j}_{in}}(x)$  is constant for all  $x \in S_{l_i}$ . And  $S_{l_i}$  forms redundant (multiples of) complete domain with respect to  $\mathbf{j}_{out}$ . Therefore for a leaf node  $l_i$ ,

$$\sum_{\mathbf{x} \in S_{l_i}} f(\mathbf{x}) \psi_j(\mathbf{x}) = \sum_{\mathbf{x} \in S_{l_i}} f(\mathbf{h}_i) \psi_j(\mathbf{x}) = f(\mathbf{h}_i) \sum_{\mathbf{x} \in S_{l_i}} \psi_{\mathbf{j}_{in}}(\mathbf{x}) \psi_{\mathbf{j}_{out}}(\mathbf{x}) = f(\mathbf{h}_i) \psi_j(\mathbf{h}_i) \sum_{\mathbf{x} \in S_{l_i}} \psi_{\mathbf{j}_{out}}(\mathbf{x}) = 0.$$

■

*Lemma 5:* For any Fourier coefficient  $w_j$  whose order is greater than the depth of a leaf node  $l_i$ ,  $\sum_{\mathbf{x} \in S_{l_i}} \psi_j(\mathbf{x}) = 0$ . If the order of  $w_j$  is greater than the depth of tree, then  $w_j = 0$ .

**Proof:** The proof immediately follows from Lemma 4.

■

Thus, for a FC  $w_j$  to be non-zero, there should exist at least one schema  $\mathbf{h}$  that has non-wild-card attributes for all non-zero attributes of  $\mathbf{j}$ . In other words, there exists a set of non-zero FCs associated with a schema  $\mathbf{h}$ . This observation leads us to a direct way of detecting and calculating all non-zero FCs of a decision tree: For each schema  $\mathbf{h}$  (or path) from the root, we can easily detect all non-zero FCs by enumerating all FCs associated with  $\mathbf{h}$ .

Before describing the algorithm, we need to introduce some notations. Let  $\mathbf{h}_{k=i}$  be a vector that is generated by replacing the  $k$ -th position of  $\mathbf{h}$  with value  $i$ . Note that this notation will be used for both schema and partition. Let us consider a non-leaf node  $n$  that has  $d$  children. In other words, there exist  $d$  disjoint subtrees below  $n$ . If  $x_k$  is the feature appearing in  $n$ , then  $F_{x_k}(i)$  denotes the average function value of domain members covered by a subtree accessible through the  $i$ -th child of  $n$ . For example, in Figure 2,  $F_{x_1}(0)$  is  $\frac{1}{2}$  and  $F_{x_2}(1)$  is one. Note that  $F_{x_k}(i)$  is equivalent to the average of schema  $\mathbf{h}$ , where  $\mathbf{h}$  denotes the path (from the root node) to  $i$ -th subtree of the node where  $x_k$  appears.

The algorithm starts with pre-calculating all  $F_{x_k}(i)$ -s (This is essentially recursive ‘‘Tree-Visit’’ operation). Then it incrementally finds non-zero FCs as it traverses the tree. If we let  $\mathbf{S}$  denote the set of partitions that correspond to non-zero FCs, initially,  $\mathbf{S} = \{000\dots 0\}$  and corresponding  $w_{000\dots 0}$  is calculated with overall average of output. In Figure 2, it is:  $\frac{1}{2} \times \frac{1}{4} + \frac{1}{2} \times 1 = \frac{5}{8}$ . The algorithm continues to extract all remaining non-zero FCs in recursive fashion from the root. New non-zero FCs are identified by inducing their corresponding partitions from the existing  $\mathbf{S}$ . For any  $\mathbf{h} \in \mathbf{S}$ , when a node with the feature  $x_k$  is visited, partitions  $\mathbf{h}_{k=1}, \dots, \mathbf{h}_{k=\lambda_k-1}$  are added into  $\mathbf{S}$ , where  $\lambda_k$  is the cardinality of  $x_k$ . For the tree in Figure 2,  $\mathbf{S}$  is initially  $\{000\}$ . Then 010 is added to  $\mathbf{S}$  when  $x_1$  is visited. Note that 010 is found by replacing the first position (starting from zero) with 1, i.e.,  $\mathbf{h}_{1=1} = 010$  is obtained from  $\mathbf{h} = 000$ .  $w_{010}$  is computed using Equation 1:

$$\begin{aligned} w_{010} &= \frac{1}{2} \times f(*0*)\psi_{010}(*0*) + \frac{1}{2} \times f(*1*)\psi_{010}(*1*) \\ &= \frac{1}{2} \times F_{x_1}(0)\psi_{010}(*0*) + \frac{1}{2} \times F_{x_1}(1)\psi_{010}(*1*) \\ &= \frac{1}{2} \times \frac{1}{2} \times 1 + \frac{1}{2} \times 1 \times (-1) = \frac{1}{4} - \frac{1}{2} = -\frac{1}{4} \end{aligned}$$

For  $x_2$ ,  $\{001, 011\}$  will be added into  $\mathbf{S}$ .  $w_{001}$  and  $w_{011}$  are computed similarly as  $w_{010}$ . The pseudo code of the algorithm is presented in Figure 3.

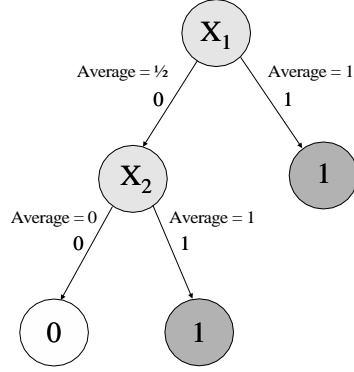


Fig. 2. An instance of Boolean decision tree that shows average output values at each subtree.

### C. Fourier Spectrum of an Ensemble Classifier

The Fourier spectrum of an ensemble classifier that consists of multiple decision trees can be computed by aggregating the spectra of the individual base models. Let  $f(\mathbf{x})$  be the underlying function computed by a tree-ensemble where the output of the ensemble is a weighted linear combination of the outputs of the base tree-classifiers.

$$f(\mathbf{x}) = a_1 f_1(\mathbf{x}) + a_2 f_2(\mathbf{x}) + \dots + a_n f_n(\mathbf{x}) = a_1 \sum_{j \in J_1} w_j^{(1)} \overline{\psi}_j(\mathbf{x}) + \dots + a_n \sum_{j \in J_n} w_j^{(n)} \overline{\psi}_j(\mathbf{x})$$

where  $f_i(\mathbf{x})$  and  $a_i$  are  $i^{\text{th}}$  decision tree and its weight respectively.  $J_i$  is set of non-zero Fourier coefficients that are detected by  $i^{\text{th}}$  decision tree and  $w_j^{(i)}$  is a Fourier coefficient in  $J_i$ . Now equation 2 is written as:  $f(\mathbf{x}) = \sum_{j \in J} w_j \overline{\psi}_j(\mathbf{x})$ , where  $w_j = \sum_{i=1}^n a_i w_j^{(i)}$  and  $J = \cup_{i=1}^n J_i$ . The following section extends the Fourier spectrum-based approach to represent and aggregate decision trees to domains with multiple class labels.

### D. Fourier Spectrum of Multi-Class Decision Trees

A multi-class decision tree has  $k > 2$  different class labels. In general, we can assume that each label is again assigned a unique integer value. Since such decision trees are also functions that map an instance vector to numeric value, the Fourier representation of such tree is essentially not any different. However, the Fourier spectrum cannot be directly applied to represent an

```

1  Function ExtractFS(input: Partition Set  $\mathbf{S}$ , Node  $node$ , Schema  $\mathbf{h}$ )
2     $x_k \leftarrow$  feature appearing in  $node$ 
3     $\mathbf{N} \leftarrow \phi$ 
4    for each  $\mathbf{j} \in \mathbf{S}$ 
5      for each  $i$  from  $(1, \dots, \lambda_k - 1)$ 
6         $\mathbf{N} \leftarrow \mathbf{N} \cup \{\mathbf{j}_{k=i}\}$ 
7      end
8    end
9     $size \leftarrow \frac{|node|}{\lambda_k |\Lambda|}$ 
10   for each  $\mathbf{j} \in \mathbf{N}$ 
11     for each  $i$  from  $(0, \dots, \lambda_k - 1)$ 
12        $w_{\mathbf{j}} \leftarrow w_{\mathbf{j}} + size \times F_{x_k}(i) \psi_{\mathbf{j}}(\mathbf{h}_{k=i})$ 
13     end
14   end
15    $\mathbf{S} \leftarrow \mathbf{S} \cup \mathbf{N}$ 
16   for each  $i$  from  $(0, \dots, \lambda_k - 1)$ 
17     ExtractFS( $\mathbf{S}$ ,  $node_i$ ,  $\mathbf{h}_{k=i}$ )
18   end
19 end

```

Fig. 3. Algorithm for obtaining Fourier spectrum of a decision tree.  $k$  in  $x_k$  implies that  $x_k$  is the  $k$ th feature.  $\lambda_k$  denotes the cardinality of  $x_k$  and  $|node|$  denotes the size of subspace  $node$  covers.  $|\Lambda|$  is the size of the complete instance space.  $node_i$  is the  $i$ th child of  $node$ .

ensemble of decision trees that uses voting as its aggregation scheme. The Fourier spectrum faithfully represents functions in closed forms and ensemble classifiers are not such functions. Therefore, we need a different approach to model a multi-class decision trees with the Fourier basis.

Let us consider a decision tree that has  $k$  classifications. Then let us define  $\mathfrak{S}_i$  to be the Fourier spectrum of a decision tree whose class labels are all set to zero except the  $i$ -th class. In other words, we treat the tree to have a Boolean classification with respect to the  $i$ -th class label. If we define  $f^{(k)}(\mathbf{x})$  to be a partial function that computes the inverse Fourier transform using  $\mathfrak{S}_k$ , classification of an input vector  $\mathbf{x}$  is written as:  $f(\mathbf{x}) = c_1 f^{(1)}(\mathbf{x}) + c_2 f^{(2)}(\mathbf{x}) + \dots + c_l f^{(l)}(\mathbf{x})$ , where each  $c_i$  corresponds to a mapped value for the  $i$ -th classification. Note that if  $\mathbf{x}$  belongs to  $j$ -th class,  $f^{(i)}(\mathbf{x}) = 1$  when  $i = j$ , and 0 otherwise.

Now let us consider an ensemble of decision trees in weighted linear combination form. Then  $f^{(k)}(\mathbf{x})$  can be written as:  $f^{(k)}(\mathbf{x}) = a_1 f_1^{(1)}(\mathbf{x}) + a_2 f_2^{(2)}(\mathbf{x}) + \dots + a_l f_l^{(l)}(\mathbf{x})$ , where  $a_i$  and  $f_i^{(k)}(\mathbf{x})$

represent the weight of  $i$ -th tree in the ensemble and its partial function for the  $k$ -th classification respectively. Finally, the classification of an ensemble of decision tree that adopts voting as its aggregation scheme can be defined as:  $f(\mathbf{x}) = \text{argmax}_k(f^{(k)}(\mathbf{x}))$

In this section, we discussed the Fourier representation of decision trees. We showed that the Fourier spectrum of a decision tree is very compact in size. In particular, we proved that the exponential decay property is also true for a Fourier spectrum of non-Boolean decision trees. In the following section, we will describe how the Fourier spectrum of an ensemble can be used to construct a single tree.

## V. CONSTRUCTION OF A DECISION TREE FROM FOURIER SPECTRUM

This section discusses an algorithm to construct a tree from the Fourier spectrum of an ensemble of decision trees. The following section first shows that the information gain needed to choose an attribute at the decision nodes can be efficiently computed from the Fourier coefficients.

### A. Schema Average and Information Gain

Consider a classification problem with Boolean class labels— $\{0, 1\}$ . Recall that a schema  $\mathbf{h}$  denotes a path to a node  $n_k$  in a decision tree. In order to compute the information gain introduced by splitting the node using a particular attribute, we first need to compute the entropy of the class distribution at that node. We do that by introducing a quantity called *schema average*. Let us define the *schema average* function value as follows:

$$\phi(\mathbf{h}) = \frac{1}{|\mathbf{h}|} \sum_{\mathbf{x} \in \mathbf{h}} f(\mathbf{x}), \quad (2)$$

where  $f(\mathbf{x})$  is the classification value of  $\mathbf{x}$  and  $|\mathbf{h}|$  denotes the number of members in schema  $\mathbf{h}$ . Note that the schema average  $\phi(\mathbf{h})$  is nothing but the frequency of all instances of the schema  $\mathbf{h}$  with a classification value of 1. Similarly, note that the frequency of the tuples with classification value of 0 is  $(1 - \phi(\mathbf{h}))$ . It can therefore be used to compute the entropy at the node  $n_k$ .

$$\text{confidence}(\mathbf{h}) = \max(\phi(\mathbf{h}), 1 - \phi(\mathbf{h}))$$

$$\text{entropy}(\mathbf{h}) = -\phi(\mathbf{h}) \log \phi(\mathbf{h}) - (1 - \phi(\mathbf{h})) \log(1 - \phi(\mathbf{h}))$$

The computation of  $\phi(\mathbf{h})$  using the above expression for a given ensemble is not practical since we need to evaluate all  $\mathbf{x} \in \mathbf{h}$ . Instead we can use the following expression that computes  $\phi(\mathbf{h})$  directly from the given FS:

$$\phi(\mathbf{h}) = \sum_{p_1} \dots \sum_{p_m} w_{(0, \dots, p_1, \dots, p_m, \dots, 0)} \exp^{2\pi i (\frac{p_1 b_1}{\lambda_{j_1}} + \dots + \frac{p_m b_m}{\lambda_{j_m}})} \quad (3)$$

where  $\mathbf{h} = \dots * b_1 \dots * b_2 \dots * b_m \dots$  that has  $m$  non-wildcard values  $b_i$  at position  $j_i$  and  $p_i \in \{0, 1, \dots, \lambda_{j_i} - 1\}$ . A similar Walsh analysis-based approach for analyzing the behavior of genetic algorithms can be found elsewhere [17]. Note that the summations in Equation 3 are defined only for the fixed (non-wild-card) positions that correspond to the features defining the path to the node  $n_k$ .

Using Equation 3 as a tool to obtain information gain, it is relatively straight-forward to come up with a version of ID3 or C4.5-like algorithms that work using the Fourier spectrum. However, a naive approach may be computationally inefficient. The computation of  $\phi(\mathbf{h})$  requires an exponential number of FCs with respect to the order of  $\mathbf{h}$ . Thus, the cost involved in computing  $\phi(\mathbf{h})$  increases exponentially as the tree becomes deeper. Moreover, since the Fourier spectrum of the ensemble is very compact in size, most Fourier coefficients involved in computing  $\phi(\mathbf{h})$  are zero. Therefore, the evaluation of  $\phi(\mathbf{h})$  using Equation 3 is not only inefficient but also involves unnecessary computations.

Construction of a more efficient algorithm to compute  $\phi(\mathbf{h})$  is possible by taking advantage of the recursive and decomposable nature of Equation 3. When computing the average of an order  $l$  schema  $\mathbf{h}$ , we can reduce some computational steps if any of the order  $l-1$  schemata which subsumes  $\mathbf{h}$  is already evaluated. For a simple example in the Boolean domain, let us consider the evaluation of  $\phi(*1*0**)$ . Let us also assume that  $\phi(*1***)$  is pre-calculated. Then,  $\phi(*1*0**)$  is obtained by simply adding  $w_{000100}$  and  $-w_{010100}$  to  $\phi(*1***)$ . This observation leads us to an efficient algorithm to evaluate schema averages. Recall that the path to a node from the root in a decision tree can be represented as a schema. Then, choosing an attribute for the next node is essentially the same as selecting the best schema among those *candidate* schemata that are subsumed by the current schema whose orders are just one more than that of this schema. In the following section, we describe a tree construction algorithm that is based on these observations.

```

1  Function TCFS(input: Fourier Spectrum FS)
2    Initialize Candidate Feature Set CFSET
3    create root node
4    h ← (***...**)
5    root ← Build(h, FS, SFSET)
6    return root
7  end

```

Fig. 4. Algorithm for constructing a decision tree from Fourier spectrum (TCFS).

### B. Bottom-up Approach to Construct a Tree

Before describing the algorithm, we need to introduce some notations. Let  $\mathbf{h}_{k=i}$  and  $\mathbf{h}$  be two schemata. The order of  $\mathbf{h}_{k=i}$  is one higher than that of  $\mathbf{h}$ . Schema  $\mathbf{h}_{k=i}$  is identical to  $\mathbf{h}$  except at one position—the  $k$ -th feature is set to  $i$  (Note that we use similar notation for ExtractFS). For example, consider schemata  $\mathbf{h} = (*1**2)$  and  $\mathbf{h}_{3=1} = (*1*12)$ . Here we use an integer number-based indexing of the features (zero for the leftmost feature).  $\pi(\mathbf{h})$  denotes a set of partitions that are required to compute  $\phi(\mathbf{h})$  (See Equation 3). A  $k$ -fixed partition is a partition with a non-zero value at the  $k$ -th position. Let  $\xi(k)$  be a set of order one  $k$ -fixed partitions;  $\gamma(\mathbf{h}_{k=i})$  be the partial sum of  $\phi(\mathbf{h}_{k=i})$  which only includes  $k$ -fixed partitions. Now the information gain achieved by choosing the  $k$ -th feature with a given  $\mathbf{h}$  is redefined using these new notations:

$$\begin{aligned}
\text{Gain}(\mathbf{h}, k) &= \text{entropy}(\mathbf{h}) - \frac{1}{\lambda_k} \sum_{i=0}^{\lambda_k-1} \text{entropy}(\mathbf{h}_{k=i}) \\
\text{entropy}(\mathbf{h}_{k=i}) &= -\phi(\mathbf{h}_{k=i}) \log(\phi(\mathbf{h}_{k=i})) - (1 - \phi(\mathbf{h}_{k=i})) \log(1 - \phi(\mathbf{h}_{k=i})) \\
\phi(\mathbf{h}_{k=i}) &= \phi(\mathbf{h}) + \gamma(\mathbf{h}_{k=i}) \\
\gamma(\mathbf{h}_{k=i}) &= \sum_{\mathbf{j} \in \pi(\mathbf{h}) \otimes \xi(k)} \bar{\psi}_{\mathbf{j}}(\mathbf{h}_{k=i}) w_{\mathbf{j}}
\end{aligned}$$

where  $\otimes$  is the Cartesian product and  $\lambda_k$  is the cardinality of the  $k$ -th feature, respectively.

Now we are ready to describe the Tree Construction from Fourier Spectrum (TCFS) algorithm, which essentially notes the decomposable definition of  $\phi(\mathbf{h}_{k=i})$  and focuses on computing  $\gamma(\mathbf{h}_{k=i})$ -s. Note that with a given  $\mathbf{h}$  (the current path), selecting the next feature is essentially identical to choosing the  $k$ -th feature that achieves the maximum  $\text{Gain}(\mathbf{h}, k)$ . Therefore, the basic idea of TCFS is to associate most up-to-date  $\phi(\mathbf{h}_{k=i})$ -s with the  $k$ -th feature. In other words, when TCFS selects the next node (after some  $i$  is chosen for  $\mathbf{h}_{k=i}$ ),  $\mathbf{h}_{k=i}$  becomes the new  $\mathbf{h}$ . Then, it

```

1  Function Build(input: Schema h,
   Fourier Spectrum FS, Candidate Feature Set CFSET)
2  create root node
3  odr  $\leftarrow (1, \text{order}(\mathbf{h}) + 1)$ 
4  Marked  $\leftarrow \phi$ 
5  for each Fourier Coefficient  $w_i$  within odr from FS
6  ft = intersect(h, i, CFSET)
7  if ft is not  $\phi$ 
8  for each value  $j$  of ft
9  update  $\gamma(\mathbf{h}_{ft=j})$  with  $w_i$ 
10 end
11 add  $w_i$  to Marked
12 end
13 end
14 if Marked is  $\phi$ 
15 set label for root using average of h
16 return root
17 end
18 for each feature  $f_i$  in CFSET
19  $gain_i \leftarrow \text{Gain}(\mathbf{h}, f_i)$ 
20 end
21 remove  $k$  with the maximum  $gain_i$  from CFSET
22 root  $\leftarrow k$ 
23 FS  $\leftarrow \mathbf{FS} - \mathbf{Marked}$ 
24 for each possible branch  $br_i$  of  $k$ 
25  $\mathbf{h}_{k=i} \leftarrow$  update h with  $k = i$ 
26  $br_i \leftarrow \text{Build}(\mathbf{h}_{k=i}, \mathbf{FS}, \mathbf{CFSET})$ 
27 end
28 add  $k$  into CFSET
29 add Marked into FS
30 return root
31 end

```

Fig. 5. Algorithm for constructing a decision tree from Fourier spectrum (TCFS).  $\text{order}(\mathbf{h})$  returns the order of schema **h**.  $\text{intersect}(\mathbf{h}, \mathbf{i})$  returns the feature to be updated using  $w_i$ , if such a feature exists. Otherwise it returns  $\phi$ .

identifies a set of FCs (We call these *appropriate* FCs) that are required to compute all  $\mathbf{h}_{k=i}$ -s for each feature and computes the corresponding entropy. This process can be considered to update each  $\phi(\mathbf{h}_{k=i})$  for the corresponding  $k$ -th feature as if it were selected. The reason is that such computations are needed anyway if a feature is to be selected in the future along the current path. This is essentially updating  $\phi(\mathbf{h}_{k=i})$ -s for a feature  $k$  using bottom-up approach (following the flavor of dynamic programming). Note that  $\phi(\mathbf{h}_{k=i})$  is, in fact, computable by adding  $\gamma(\mathbf{h}_{k=i})$  to



$\phi(\mathbf{h})$ . Here  $\gamma(\mathbf{h}_{k=i})$ -s are partial sums that only current appropriate FCs contribute to. Detection of all appropriate FCs requires a scan over the FS. However, they are removed from the FS once they are used in computation, since they are no longer needed for the calculation of higher order schemata. Thus it takes a lot less time to compute higher order schemata; note that it is just opposite to what we encountered in the naive implementation. The algorithm stops growing a path when either the original FS becomes an empty set or the minimum confidence level is achieved. The depth of the resulting tree can be set to a pre-determined bound. A pictorial description of the algorithm is shown in Figure 6. Pseudo code of the algorithm is presented in Figures 4 and 5.

The TCFS uses the same criteria to construct a tree as that of the C4.5. Both of them require a number of information-gain-tests that grows exponentially with respect to the depth of the tree. In that sense, the asymptotic running time of TCFS is the same as that of the C4.5. However, while the C4.5 uses original data to compute information gains, TCFS uses a Fourier spectrum. Therefore, in practice, a comparison of the running time between the two approaches will depend on the sizes of the original data and that of Fourier spectrum. The following section presents an extension of the TCFS for handling non-Boolean class labels.

### C. Extension of TCFS to Multi-Class Decision Trees

The extension of TCFS algorithm to multi-class problems is immediately possible by redefining the “entropy” function. It should be modified to capture an entropy from the multiple class labels. For this, let us first define  $\phi^{(i)}(\mathbf{h})$  to be a schema average function that uses  $\mathfrak{S}_i$  (See Section IV-D) only. Note that it computes the average occurrence of the  $i$ -th class label in  $\mathbf{h}$ . Then the entropy of a schema is redefined as follows.

$$\text{entropy}(\mathbf{h}) = - \sum_{i=1}^k \phi^{(i)}(\mathbf{h}) \log \phi^{(i)}(\mathbf{h})$$

where  $k$  is the number of class labels.

This expression can be directly used for computing the information gain to choose the decision nodes in a tree for classifying domains with non-Boolean class labels.

In this section, we discussed a way to assign a confidence to a node in a decision tree, and considered a method to estimate information gain using it. Consequently, we showed that

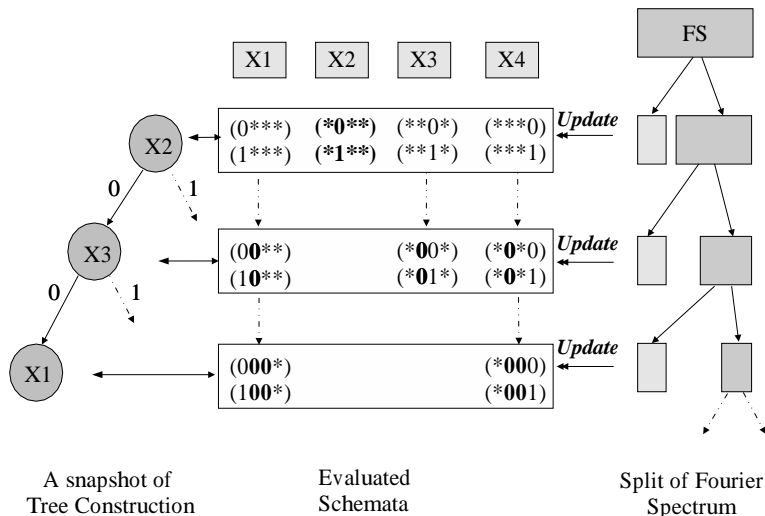


Fig. 6. Illustration of the Tree Construction from Fourier Spectrum (TCFS) algorithm. It shows the constructed tree on the left. The schemata evaluated at different orders are shown in the middle. The rightmost tree shows the splitting of the set of all Fourier coefficients used for making the process of looking up the appropriate coefficients efficient.

a decision tree construction from the Fourier spectrum is possible. In particular, we devised TCFS algorithm that exploits the recursive and decomposable nature of tree building process in spectrum domain, thus constructing a decision tree efficiently. In the following section, we will discuss orthogonal decision trees that can be constructed using the Fourier spectrum of the trees in an ensemble.

## VI. REMOVING REDUNDANCIES FROM ENSEMBLES

Existing ensemble-learning techniques work by combining (usually a linear combination) the output of the base classifiers. They do not structurally combine the classifiers themselves. As a result they often share a lot of redundancies. The Fourier representation offers a unique way to fundamentally aggregate the trees and perform further analysis to construct an efficient representation.

Let  $f_e(\mathbf{x})$  be the underlying function representing the ensemble of  $m$  different decision trees where the output is a weighted linear combination of the outputs of the base classifiers. Then

we can write,

$$f_e(\mathbf{x}) = \alpha_1 \tau_{(1)}(\mathbf{x}) + \alpha_2 \tau_{(2)}(\mathbf{x}) + \cdots + \alpha_m \tau_{(m)}(\mathbf{x}) = \alpha_1 \sum_{j \in \mathcal{J}_1} w_{(1),j} \overline{\psi}_j^\lambda(\mathbf{x}) + \cdots + \alpha_m \sum_{j \in \mathcal{J}_m} w_{(m),j} \overline{\psi}_j^\lambda(\mathbf{x}).$$

Where  $\alpha_i$  is the weight of the  $i^{th}$  decision tree and  $\mathcal{J}_i$  is the set of all partitions with non-zero Fourier coefficients in its spectrum. Therefore,  $f_e(\mathbf{x}) = \sum_{j \in \mathcal{J}} w_{(e),j} \overline{\psi}_j^\lambda(\mathbf{x})$ , where  $w_{(e),j} = \sum_{i=1}^m \alpha_i w_{(i),j}$  and  $\mathcal{J} = \cup_{i=1}^m \mathcal{J}_i$ . Therefore, the Fourier spectrum of  $f_e(\mathbf{x})$  (a linear ensemble classifier) is simply the weighted sum of the spectra of the member trees.

Consider the matrix  $D$  where  $D_{i,j} = \tau_{(j)}(\mathbf{x}_i)$ , where  $\tau_{(j)}(\mathbf{x}_i)$  is the output of the tree  $\tau_{(j)}$  for input  $\mathbf{x}_i \in \Omega$ .  $D$  is an  $|\Omega| \times m$  matrix where  $|\Omega|$  is the size of the input domain and  $m$  is the total number of trees in the ensemble.

An ensemble classifier that combines the outputs of the base classifiers can be viewed as a function defined over the set of all rows in  $D$ . If  $D_{*,j}$  denotes the  $j$ -th column matrix of  $D$  then the ensemble classifier can be viewed as a function of  $D_{*,1}, D_{*,2}, \cdots, D_{*,m}$ . When the ensemble classifier is a linear combination of the outputs of the base classifiers we have  $F = \alpha_1 D_{*,1} + \alpha_2 D_{*,2} + \cdots + \alpha_m D_{*,m}$ , where  $F$  is the column matrix of the overall ensemble-output. Since the base classifiers may have redundancy, we would like to construct a compact low-dimensional representation of the matrix  $D$ . However, explicit construction and manipulation of the matrix  $D$  is difficult, since most practical applications deal with a very large domain. We can try to construct an approximation of  $D$  using only the available training data. One such approximation of  $D$  and its Principal Component Analysis-based projection is reported elsewhere [18]. Their technique performs PCA of the matrix  $D$ , projects the data in the representation defined by the eigenvectors of the covariance matrix of  $D$ , and then performs linear regression for computing the coefficients  $\alpha_1, \alpha_2, \cdots$ , and  $\alpha_m$ .

While the approach is interesting, it has a serious limitation. First of all, the construction of an approximation of  $D$  even for the training data is computationally prohibiting for most large scale data mining applications. Moreover, this is an approximation since the matrix is computed only over the observed data set of the entire domain. In the following we demonstrate a novel way to perform a PCA of the matrix containing the Fourier spectra of trees. The approach works without explicitly generating the matrix  $D$ . It is important to note that the PCA-based regression scheme [18] offers a way to find the weightage for the members of the ensemble. It does not offer any way to aggregate the tree structures and construct a new representation of the ensemble

which the current approach does.

The following analysis will assume that the columns of the matrix  $D$  are mean-zero. This restriction can be easily removed with a simple extension of the analysis. Note that the covariance of the matrix  $D$  is  $D^T D$ . Let us denote this covariance matrix by  $C$ . The  $(i, j)$ -th entry of the matrix,

$$C_{i,j} = \langle D(*, i), D(*, j) \rangle = \langle \tau_{(i)}(\mathbf{x}), \tau_{(j)}(\mathbf{x}) \rangle = \sum_{\mathbf{p}} w_{(i),\mathbf{p}} w_{(j),\mathbf{p}} = \langle \mathbf{w}_{(i)}, \mathbf{w}_{(j)} \rangle \quad (4)$$

The fourth step is true by Lemma 2. Now let us consider the matrix  $W$  where  $W_{i,j}$  is the coefficient corresponding to the  $i$ -th member of the partition set  $\mathcal{J}$  from the spectrum of the tree  $\tau_{(j)}$ . Equation 4 implies that the covariance matrices of  $D$  and  $W$  are identical. Note that  $W$  is an  $|\mathcal{J}| \times m$  dimensional matrix. For most practical applications  $|\mathcal{J}| \ll |\Omega|$ . Therefore analyzing  $W$  using techniques like PCA is significantly easier. The following discourse outlines a PCA-based approach.

PCA of the covariance matrix of  $W$  produces a set of eigenvectors  $V_1, V_2, \dots, V_k$ . The eigenvalue decomposition constructs a new representation of the underlying domain. Note that since the eigenvectors are nothing but a linear combination of the original column vectors of  $W$ , each of them also form a Fourier spectrum and we can reconstruct a decision tree from this spectrum. Moreover, since they are orthogonal to each other, the tree constructed from them also maintain the orthogonality condition and therefore they are redundancy-free. They define a basis set and can be used to represent any given decision tree in the ensemble in the form of a linear combination. Orthogonal decision trees can be defined as an immediate extension of this framework.

A pair of decision trees  $f_1(\mathbf{x})$  and  $f_2(\mathbf{x})$  are orthogonal to each other if and only if  $\langle f_a(\mathbf{x}), f_b(\mathbf{x}) \rangle = 0$  when  $a \neq b$  and  $\langle f_a(\mathbf{x}), f_a(\mathbf{x}) \rangle = 1$  otherwise. The second condition is actually a slightly special case of orthogonal functions—orthonormal condition. A set of trees are pairwise orthogonal if every possible pair of members of this set satisfy the orthogonality condition.

The orthogonality condition guarantees that the representation is not redundant. These orthogonal trees form a basis set that spans the entire function space of the ensemble. The overall output of the ensemble is computed from the output of these orthogonal trees. Specific details

of the ensemble output computation depends on the adopted technique to compute the overall output of the original ensemble. However, for most popular cases considered here boils down to computing the average output. If we choose to go for weighted averages, we may also compute the coefficients corresponding to each  $V_q$  by simply performing linear regression.

## VII. EXPERIMENTAL RESULTS

This section reports the experimental performance of orthogonal decision trees on the following data sets - SPECT, NASDAQ, DNA, House of Votes and Contraceptive Method Usage Data. For each data set, the following three experiments are performed using known classification techniques:

- 1) **C4.5:** The C4.5 classifier is built on training data and validated over test data.
- 2) **Bagging:** A popular ensemble classification technique, bagging, is used to test the classification accuracy of the data set.
- 3) **Random Forest:** Random forests are built on the training data, using approximately half the number of features in the original data set. The number of trees in the forest is identical to that used in the bagging experiment<sup>1</sup>.

We then perform another set of experiments for comparing the techniques described in the previous sections in terms of error in classification and tree complexity.

- 1) **Reconstructed Fourier Tree (RFT):** The training set is uniformly sampled, with replacement and C4.5 trees are built on each sample. The Fourier representation of each individual tree is obtained preserving a certain percentage (e.g. 90%) of the energy. This representation of a tree is used to reconstruct a decision tree using the TCFS algorithm described in Section V. The performance of a reconstructed Fourier tree is compared with the original C4.5 tree. The error in classification and tree complexity of each of the reconstructed trees is reported. The purpose of this experiment is to study the effect of "noise removal" from the ensemble on its classification-accuracy by going to the Fourier domain and then coming back to the tree domain using the TCFS algorithm.
- 2) **Aggregated Fourier Tree(AFT):** The training set is uniformly sampled, with replacement and C4.5 decision trees are built on each sample (This is identical to bagging). A Fourier

<sup>1</sup>We used the WEKA implementation(<http://www.cs.waikato.ac.nz/ml/weka/>) of Bagging and Random Forests

representation of each tree is obtained (preserving a certain percentage of the total energy), and these are aggregated with uniform weighting to obtain the spectrum of an Aggregated Fourier Tree (AFT). The AFT is reconstructed using the TCFS algorithm described before and the classification accuracy and the tree complexity of this aggregated Fourier tree is reported.

- 3) **Orthogonal Decision Trees:** The matrix containing the Fourier coefficients of the decision trees is subjected to principal component analysis. Orthogonal trees are built using the corresponding eigenvectors. In most cases it is found that the first principal eigenvector captures most of the variance, and thus the orthogonal decision tree constructed from this eigenvector is of particular interest. We report the error in classification and tree complexity of the orthogonal decision tree obtained from the most dominant eigenvector. We also perform experiments where we keep  $k^2$  significant eigenvectors. The trees are combined by weighting them according to the coefficients obtained from a Least Square Regression. Each orthogonal decision tree is weighted using coefficients calculated from Least Square Regression. For this, we allow all the orthogonal decision trees to individually produce their classification on the test set. Thus each ODT produces a column vector of its classification estimate. Since the class-labels in the test set are already known, we use the least square regression to obtain the weights to assign to each ODT. The accuracy of the orthogonal decision trees is reported as ODT-LR (ODTs combined using Least Square Regression).

In addition to reporting the error in classification, we also report the tree complexity, the total number of nodes in the tree. Similarly, the term ensemble complexity reflects the total number of nodes in all the trees in the ensemble. A smaller ensemble tree complexity implies a compact representation of an ensemble and therefore it is desirable. Our experiments show that ODTs usually offer significantly reduced ensemble tree complexity without any reduction in the accuracy. The following section presents the results for the SPECT data set.

<sup>2</sup>We select the value of  $k$  in such a manner that the total variance captured is more than 90%. One could potentially do cross-validation to obtain a suitable value of  $k$  as pointed out in [19] but this is beyond the current scope of the work and will be explored in future.

Method of classification	Error Percentage
C4.5	24.5989 (%)
Bagging	20.85 (%)
Random Forest	22.99466 (%)
Aggregated Fourier Tree (AFT)	19.78(%)
ODT from 1st PC	8.02(%)
ODT-LR	8.02(%)

Method of classification	Tree Complexity
C4.5	13
Bagging (average of 40 trees)	5.06
Aggregated Fourier Tree(AFT)(40 trees)	3
Orthogonal Decision Tree from 1st PC	17
Orthogonal Decision Trees (average of 15 trees)	4.3

TABLE I

CLASSIFICATION ERROR(LEFT) AND TREE COMPLEXITY(RIGHT) FOR SPECT DATA.

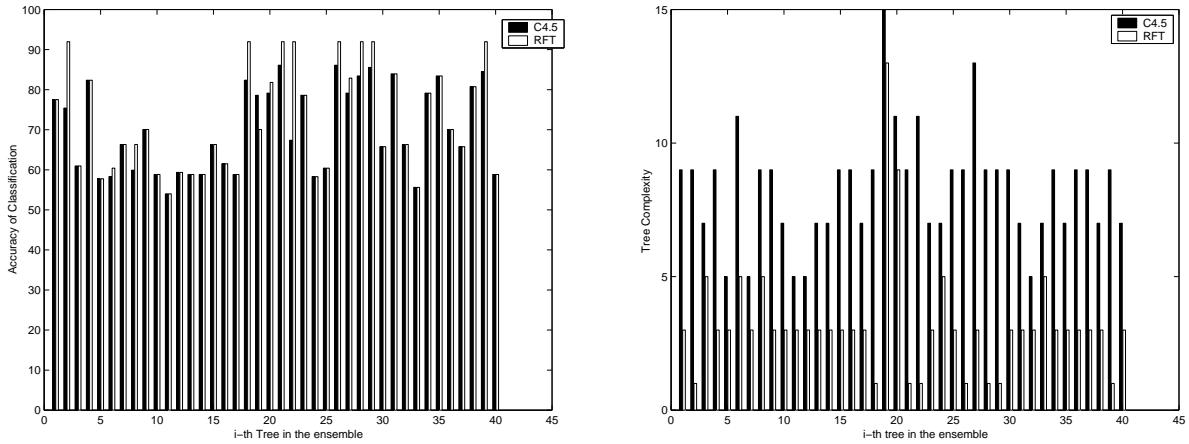


Fig. 7. The accuracy and tree complexity of C4.5 and RFT for SPECT data

### A. SPECT Data set

This section illustrates the idea of orthogonal decision trees using a well known binary data set. The dataset, available from the University of California Irvine, Machine Learning Repository, describes diagnosing of cardiac Single Proton Emission Computed Tomography (SPECT) images into two categories, normal or abnormal. The database of 267 SPECT image sets (patients) is processed to extract features that summarize the original SPECT images. As a result, 44 continuous feature patterns are obtained for each patient, which are further processed to obtain 22 binary feature patterns. The training data set consists of 80 instances and 22 attributes. All the features are binary, and the class label is also binary (depending on whether a patient is deemed normal or abnormal). The test data set consists of 187 instances and 22 attributes.

Table I(Left) shows the error percentage obtained in each of the different classification schemes. The root mean squared error for the 10 fold cross validation in the C4.5 experiment is found to be 0.4803 and the standard deviation is 2.3862. For Bagging, the number of trees in the ensemble is chosen to be forty. Our experiments reveal that further increase in number of trees in the ensemble causes a decrease in accuracy of classification of the ensemble possibly due to over-fitting of the data.

For experiments with Random Forests, forest of 40 trees, each constructed while considering 12 random features is built. The average Out of bag error is reported to be 0.3245.

Figure 7(Left) compares the accuracy of the original C4.5 ensemble with that of the Reconstructed Fourier Tree(RFT) ensemble preserving 90% of the energy of the spectrum. The results reveal that if all of the spectrum is preserved, the accuracy of the original C4.5 tree and RFT are identical. When the higher order Fourier coefficients are removed, this becomes equivalent to pruning a decision tree. This explains the higher accuracy of the reconstructed Fourier tree preserving 90% of the energy of the spectrum. Figure 7(Right) compares the tree complexity of the original C4.5 ensemble with that of the RFT ensemble.

In order to construct the orthogonal decision trees, the coefficient matrix is projected onto the first fifteen most significant principal components. The most significant principal component captures 85.1048% of the variance and the tree complexity of the ODT constructed from this component is 17 with an accuracy of 91.97%. Figure 8 shows the variance captured by all the fifteen principal components.

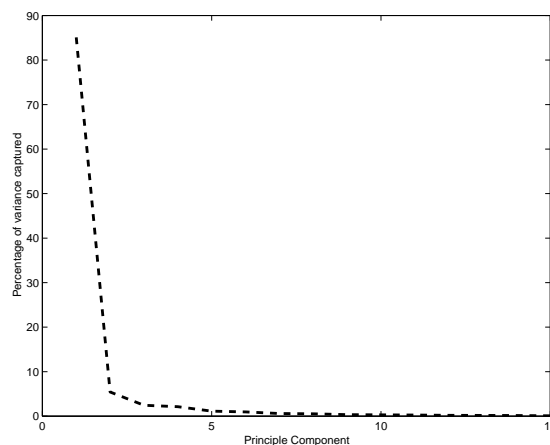


Fig. 8. Percentage of variance captured by principal components for SPECT Data.



Table I(Right) illustrates the tree complexity for this data set. The orthogonal trees are found to be smaller in complexity, thus reducing the complexity of the ensemble.

### B. NASDAQ Data set

The NASDAQ data set is a semi-synthetic data set with 1000 instances and 100 discrete attributes. The original data set has three years of NASDAQ stock quote data. It is preprocessed and transformed to discrete data by encoding percentages of changes in stock quotes between consecutive days. For these experiments we assign, 4 discrete values, that denote levels of changes. The class labels, predict whether the Yahoo stock is likely to increase or decrease based on attribute values of the 99 stocks. We randomly select 200 instances for training and the remaining 800 instances forms the test data set.

Method of classification	Error Percentage
C4.5	24.63 (%)
Bagging	32.75 (%)
Random Forest	25.75 (%)
Aggregated Fourier Tree(AFT)	34.51 (%)
ODT from 1st PC	31.12(%)
ODT-LR	31.12(%)

Method of classification	Tree Complexity
C4.5	29
Bagging(average of 60 trees)	17
Aggregated Fourier Tree(AFT) (60 trees)	15.2
Orthogonal Decision Tree from 1st PC	3
Orthogonal Decision Trees (average of 10 trees)	6.2

TABLE II

CLASSIFICATION ERROR(LEFT) AND TREE COMPLEXITY(RIGHT) FOR NASDAQ DATA.

Table II(Left) illustrates the classification accuracies of different experiments performed on this data set. The root mean squared error for the 10 fold cross validation in the C4.5 experiment is found to be 0.4818 and the standard deviation is 2.2247. C4.5 has the best classification accuracy, though the tree built has the highest tree complexity also. For the bagging experiment, C4.5 trees are built on the dataset, such that the size of each bag (used to build the tree) as a percentage of the data set is 40%. Also, Random forest of 60 trees, each constructed while considering 50 random features is built on the training data and tested with the test data set. The average out of bag error is reported to be 0.3165.

Figure 9(Left) compares the accuracy of the original C4.5 ensemble with that of the Reconstructed Fourier Tree(RFT) ensemble preserving 90% of the energy of the spectrum. Fig-

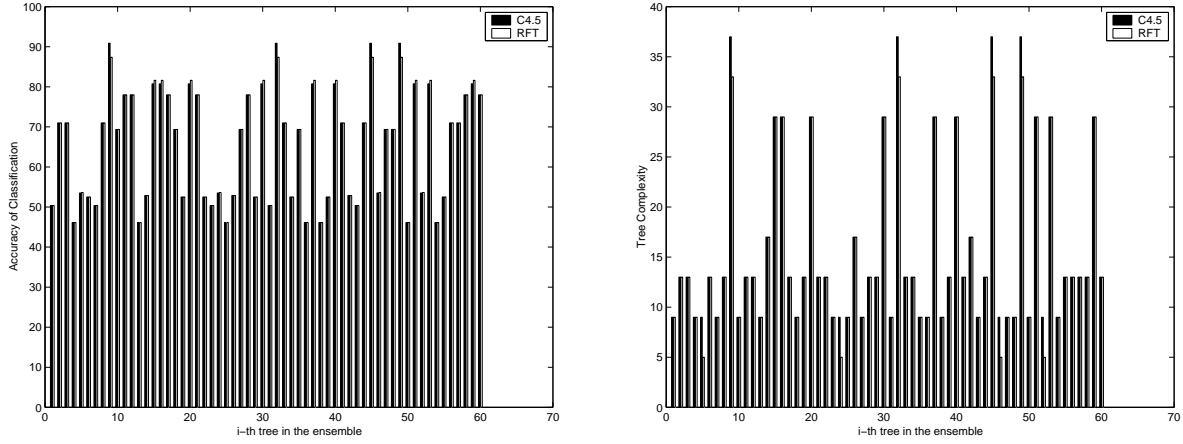


Fig. 9. The accuracy and tree complexity of C4.5 and RFT for Nasdaq data

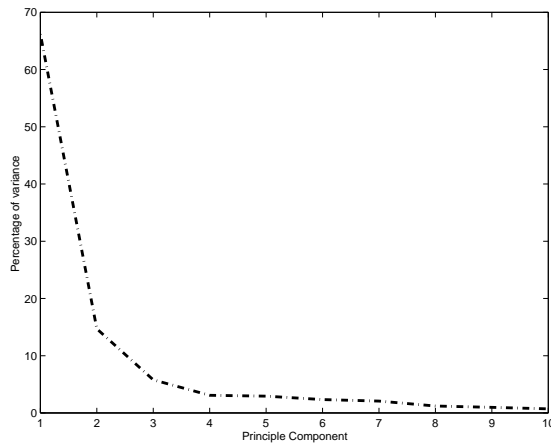


Fig. 10. Percentage of variance captured by principal components for Nasdaq Data.

ure 9(Right) compares the tree complexity of the original C4.5 ensemble with that of the RFT ensemble.

For the orthogonal trees, we project the data along the first 10 most significant principal components. The Figure 10 illustrates the percentage of variance captured by the ten most significant principal components.

Table II(Right) presents the tree-complexity information for this set of experiments. Both the aggregated Fourier tree and the orthogonal trees performed better than the single C4.5 tree or bagging. The tree-complexity result appears to be quite interesting. While a single C4.5 tree had twenty nine nodes in it, the orthogonal tree from the first principal component requires just three

nodes, which is clearly a much more compact representation.

### C. DNA Data Set

The DNA data set<sup>3</sup> is a processed version of the corresponding data set available from UC Irvine repository. The processed StatLog version replaces the symbolic attribute values representing the nucleotides (only A,C,T,G) by 3 binary indicator variables. Thus the original 60 symbolic attributes are changed into 180 binary attributes. The nucleotides A,C,G,T are given indicator values as follows:  $A = 100$ ,  $C = 010$ ,  $G = 001$ ,  $T = 000$ . The data set has three class values 1, 2, and 3 corresponding to exon-intron boundaries (sometimes called **acceptors**), intron-exon boundaries (sometimes called **donors**), and the case when neither is true. We further process the data such that, there are only two class labels i.e. class 1 representing either donors or acceptors, while class 0 representing neither. The training set consists of 2000 instances and 180 attributes of which 47.45% belongs to class 1 while the remaining 52.55% belongs to class 0. The test data set consists of 1186 instances and 180 attributes of which 49.16% belongs to class 0 while the remaining 50.84% belongs to the class 1. Table III(Left) reports the classification error. The root mean squared error for the 10 fold cross validation in the C4.5 experiment is found to be 0.2263 and the standard deviation is 0.6086 .

Method of classification	Error Percentage
C4.5	6.4924 (%)
Bagging	8.9376(%)
Random Forest	4.595275 (%)
Aggregated Fourier Tree(AFT)	8.347(%)
ODT from 1st PC	10.70(%)
ODT-LR	10.70(%)

Method of classification	Tree Complexity
C4.5	131
Bagging (average of 10 trees)	34
Aggregated Fourier Tree(AFT)(10 trees)	3
Orthogonal Decision Tree from 1st PC	25
Orthogonal Decision Trees (average of 5 trees)	7.4

TABLE III

CLASSIFICATION ERROR(LEFT) AND TREE COMPLEXITY(RIGHT) FOR DNA DATA.

It may be interesting to note, that the first five eigenvectors are used in this experiment. Figure 11 shows the variance captured by these components. As before, the redundancy free

<sup>3</sup>Obtained from <http://www.liacc.up.pt/ML/statlog/datasets/dna>

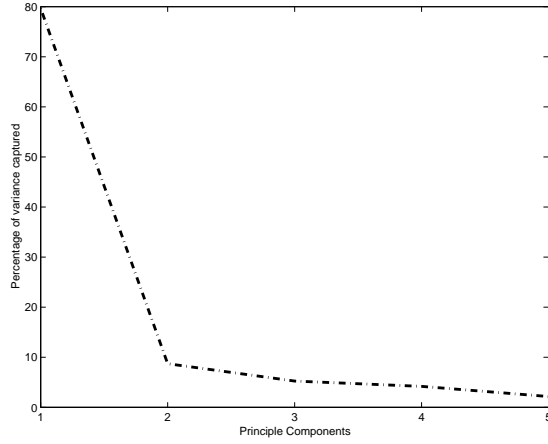


Fig. 11. Percentage of variance captured by principal components for DNA Data.

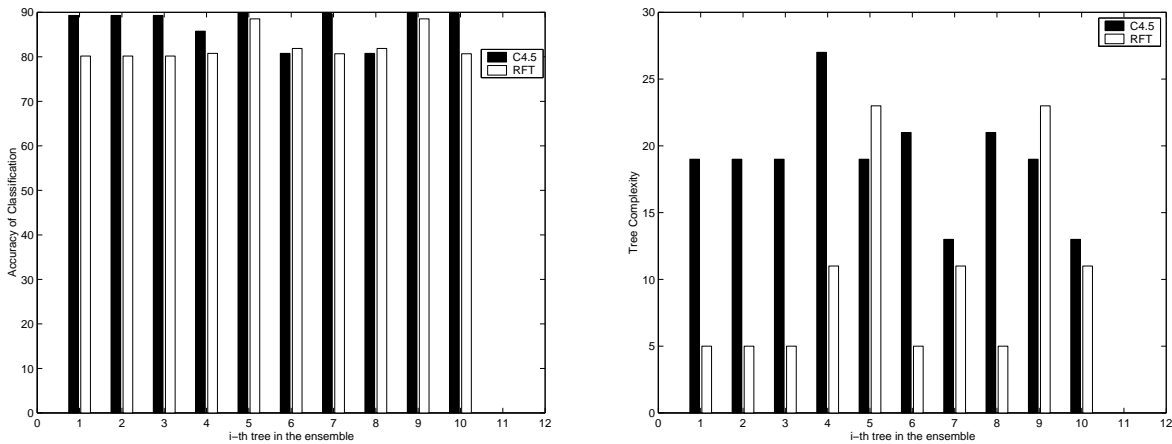


Fig. 12. The accuracy and tree complexity of C4.5 and RFT for DNA data

trees are combined by the weights obtained from Least Square Regression. Table III(Right) reports the tree complexity for this data set.

Figure 12(Left) compares the accuracy of the original C4.5 ensemble with that of the Reconstructed Fourier Tree(RFT) ensemble preserving 90% of the energy of the spectrum. Figure 12(Right) compares the tree complexity of the original C4.5 ensemble with that of the RFT ensemble.

#### D. House of Votes Data

The 1984 United States Congressional Voting Records Database is obtained from the University of California, Machine Learning Repository. This data set includes votes for each of the U.S. House of Representatives Congressmen on the 16 key votes identified by the CQA including water project cost sharing, adoption of budget resolution, mx-missile, immigration etc. It has 435 instances, 16 boolean valued attributes and a binary class label(democrat or republican).Our experiments use the first 335 instances for training and the remaining 100 instances for testing. In our experiments, missing values in the data are replaced by one.

The results of classification are shown in the Table IV(Left) while the tree complexity is shown in Table IV(Right). The root mean squared error for the 10 fold cross validation in the C4.5 experiment is found to be 0.2634 and the standard deviation is 0.3862. For Bagging, fifteen trees are constructed using the dataset, since this produced the best classification results. The size of each bag was 20% of the training data set. Random Forest of fifteen trees, each constructed by considering 8 random features produces an average out of bag error of 0.05502. The accuracy of classification and the tree complexity of the original C4.5 and RFT ensemble are illustrated in the left and right hand side of Figure 13 respectively.

For orthogonal trees, the coefficient matrix is projected onto the first five most significant principal components. Figure 14(Left) illustrates the amount of variance captured by each of the principal components.

Method of classification	Error Percentage
C4.5	8.0 (%)
Bagging	11.0(%)
Random Forest	5.6(%)
Aggregated Fourier Tree(AFT)	11(%)
ODT from 1st PC	11(%)
ODT-LR	11(%)

Method of classification	Tree Complexity
C4.5	9
Bagging (average of 15 trees)	5.266
Aggregated Fourier Tree (AFT)(15 trees)	5
Orthogonal Decision Tree from 1st PC	5
Orthogonal Decision Trees (average of 5 trees)	3

TABLE IV

CLASSIFICATION ERROR(LEFT) AND TREE COMPLEXITY(RIGHT) FOR HOUSE OF VOTES DATA.

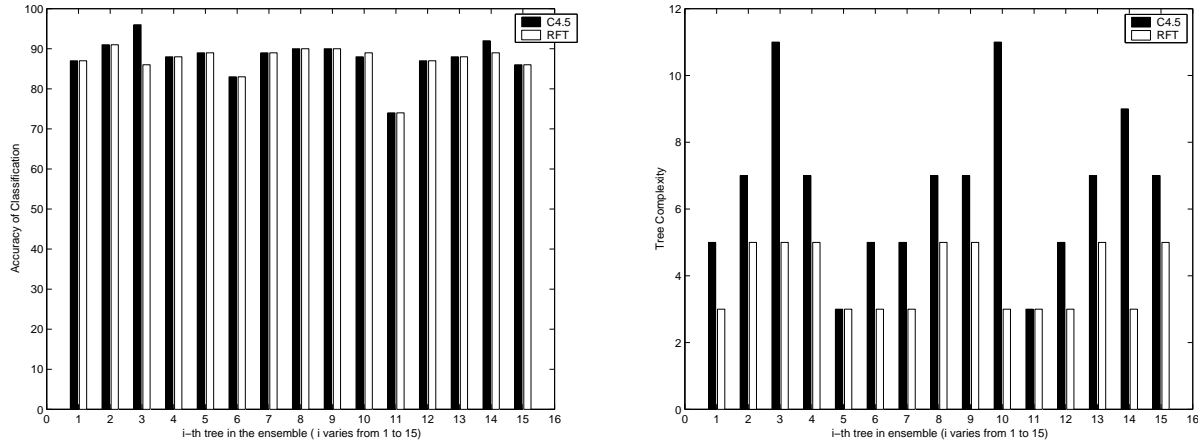


Fig. 13. The accuracy and tree complexity of C4.5 and RFT for House of Votes data

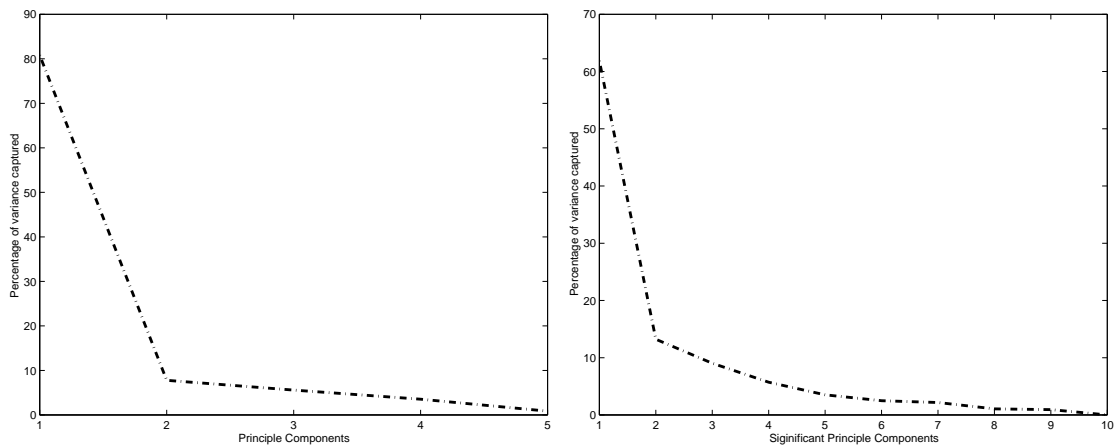


Fig. 14. Percentage of variance captured by principal components for (Left) House of Votes Data and (Right) Contraceptive Method Usage data.

### E. Contraceptive Method Usage Data

This dataset, is obtained from the University of California Irvine, Machine Learning Repository and is a subset of the 1987 National Indonesia Contraceptive Prevalence Survey. The samples are married women who are either not pregnant or do not know if they are at the time of interview. The problem is to predict the current contraceptive method choice of a woman based on her demographic and socio-economic characteristics. There are 1473 instances, and 10 attributes including a binary class label. All attributes are processed so that they are binary. Our experiments use 1320 instances for the training set while the rest form the test data set.

The results of classification are tabulated in the Table V(Left) while Table V(Right) shows the tree complexity. The root mean squared error for the 10 fold cross validation in the C4.5 experiment is found to be 0.5111 and the standard deviation is 1.8943. Random Forest built with 10 trees, considering 5 random features produces an average error in classification of about 45.88% and an average out of bag error of 0.42556. Figure 15(Left) compares the accuracy of the original C4.5 ensemble with that of the Reconstructed Fourier Tree(RFT) ensemble preserving 90% of the energy of the spectrum. Figure 15(Right) compares the tree complexity of the original C4.5 ensemble with that of the RFT ensemble.

For ODTs, the data is projected along the first ten principal components. Figure 14(Right) shows the amount of variance captured by each principal component. It is interesting to note that the first principal component captures only about 61.85% of the variance and thus the corresponding ODT generated from the first principal component has a relatively high tree complexity.

Method of classification	Error Percentage
C4.5	49.6732(%)
Bagging	52.2876(%)
Random Forest	45.88234 (%)
Aggregated Fourier Tree(AFT)	33.98(%)
ODT from 1st PC	46.40(%)
ODT-LR	46.40(%)

Method of classification	Tree Complexity
C4.5	27
Bagging(average of 10 trees)	24.8
Aggregated Fourier Tree(AFT)(10 trees)	55
Orthogonal Decision Tree from 1st PC	15
Orthogonal Decision Trees (average of 10 trees)	6.6

TABLE V

CLASSIFICATION ERROR(LEFT) AND TREE COMPLEXITY(RIGHT) FOR CONTRACEPTIVE METHOD USAGE DATA.

## VIII. CONCLUSIONS

This paper introduced the notion of orthogonal decision trees and offered a methodology to construct them. Orthogonal decision trees are functionally orthogonal to each other and they provide an efficient redundancy-free representation of large ensembles that are frequently produced by techniques like Boosting [2], [3], Bagging[4], Stacking [5], and random forests [6]. The proposed technique is also likely to be very useful in ensemble-based mining of distributed [10] and stream data [7], [8].

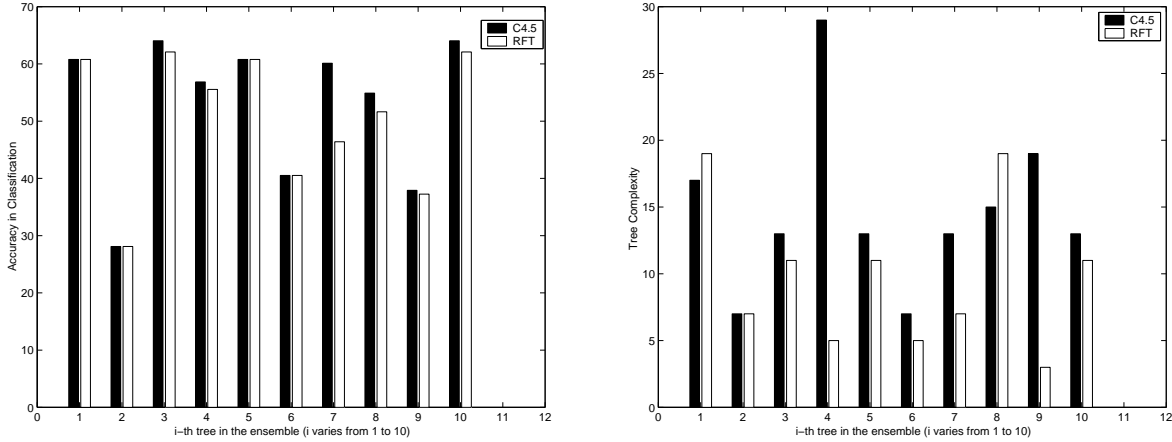


Fig. 15. The accuracy and tree complexity of C4.5 and RFT for Contraceptive Method Usage data

The proposed approach exploits the earlier work done by the first author and his colleagues [20], [9] which showed that the Fourier transform of decision trees can be efficiently computed. This work also shows that we can compute the tree back from its Fourier spectrum. The paper also offered a collection of new results regarding the properties of the multi-variate Fourier spectrum of decision trees. Although, the paper considers the Fourier representation, this is clearly not the only available linear representation around. However, our work shows that it is particularly suitable for representing decision trees.

This work also opens up several new possibilities. Linear systems theory offers many tools for analyzing properties like stability and convergence. For example, eigenvalues of a linear system are directly associated with the stability of the system. Similar concepts may be useful in understanding the behavior of large ensembles. We plan to explore these issues in the future.

#### ACKNOWLEDGMENTS

The authors acknowledge supports from NSF CAREER award IIS-0093353, NSF grant IIS-0203958, and NASA grant NAS2-37143. The work of B.-H. Park was partially funded by the Scientific Data Management Center (<http://sdmcenter.lbl.gov>) under the Department of Energy's Scientific Discovery through Advanced Computing (DOE SciDAC) program (<http://www.scidac.org>).

#### REFERENCES

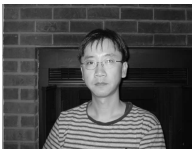
- [1] J. R. Quinlan, "Induction of decision trees," *Machine Learning*, vol. 1, no. 1, pp. 81–106, 1986.



- [2] Y. Freund, "Boosting a weak learning algorithm by majority," *Information and Computation*, vol. 121, no. 2, pp. 256–285, 1995.
- [3] H. Drucker and C. Cortes, "Boosting decision trees," *Advances in Neural Information Processing Systems*, vol. 8, pp. 479–485, 1996.
- [4] L. Breiman, "Bagging predictors," *Machine Learning*, vol. 24, no. 2, pp. 123–140, 1996.
- [5] D. Wolpert, "Stacked generalization," *Neural Networks*, vol. 5, pp. 241–259, 1992.
- [6] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [7] W. Fan, S. Stolfo, and J. Zhang, "The application of adaboost for distributed, scalable and on-line learning," in *Fifth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Diego, California, 1999.
- [8] W. N. Street and Y. Kim, "A streaming ensemble algorithm (sea) for large-scale classification," in *Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, San Francisco, CA, 2001.
- [9] H. Kargupta and B. Park, "A fourier spectrum-based approach to represent decision trees for mining data streams in mobile environments," *IEEE Transactions on Knowledge and Data Engineering*, vol. 16, no. 2, pp. 216–229, 2002.
- [10] B. Park, A. R., and H. Kargupta, "A fourier analysis-based approach to learn classifier from distributed heterogeneous data," in *Proceedings of the First SIAM International Conference on Data Mining*, Chicago, US, 2001.
- [11] B. H. Park and H. Kargupta, "Constructing simpler decision trees from ensemble models using fourier analysis," in *Proceedings of the 7th Workshop on Research Issues in Data Mining and Knowledge Discovery, ACM SIGMOD*, 2002, pp. 18–23.
- [12] F. Chung, *Spectral Graph Theory*. Providence, Rhode Island, USA: American Mathematical Society, 1994.
- [13] H. Kargupta and B. Park, "Mining time-critical data stream using the Fourier spectrum of decision trees," in *Proceedings of the IEEE International Conference on Data Mining*. IEEE Press, 2001, pp. 281–288.
- [14] H. Kargupta, B. Park, S. Pittie, L. Liu, D. Kushraj, and K. Sarkar, "Mobimine: Monitoring the stock market from a PDA," *ACM SIGKDD Explorations*, vol. 3, no. 2, pp. 37–46, January 2002.
- [15] N. Linial, Y. Mansour, and N. Nisan, "Constant depth circuits, fourier transform, and learnability," *Journal of the ACM*, vol. 40, pp. 607–620, 1993.
- [16] E. Kushilevitz and Y. Mansour, "Learning decision trees using the Fourier spectrum," *SIAM Journal on Computing*, vol. 22, no. 6, pp. 1331–1348, 1993.
- [17] D. Goldberg, "Genetic algorithms and Walsh functions: Part I, a gentle introduction," *Complex Systems*, vol. 3, no. 2, pp. 129–152, 1989.
- [18] C. J. Merz and M. J. Pazzani, "A principal components approach to combining regression estimates," *Machine Learning*, vol. 36, no. 1–2, pp. 9–32, 1999.
- [19] C. Merz and M. Pazzani, "A principal components approach to combining regression estimates," *Machine Learning*, vol. 36, pp. 9–32, 1999.
- [20] H. Kargupta, B. Park, D. Hershberger, and E. Johnson, "Collective data mining: A new perspective towards distributed data mining," in *Advances in Distributed and Parallel Knowledge Discovery, Eds: Kargupta, Hillol and Chan, Philip*. AAAI/MIT Press, 2000.



**Hillol Kargupta** is an Associate Professor in the Department of Computer Science and Electrical Engineering, University of Maryland, Baltimore County. He received the PhD degree in computer science from the University of Illinois at Urbana-Champaign in 1996. He is also a cofounder of Agnik LLC, an ubiquitous intelligence company. His research interests include mobile and distributed data mining and computation in biological process of gene expression. Dr. Kargupta won a US National Science Foundation CAREER award in 2001 for his research on ubiquitous and distributed data mining. He along with his coauthors received the best paper award at the 2003 IEEE International Conference on Data Mining for a paper on privacy-preserving data mining. He won the 2000 TRW Foundation Award and the 1997 Los Alamos Award for Outstanding Technical Achievement. His research has been funded by the US National Science Foundation, US Air Force, Department of Homeland Security, NASA, and various other organizations. He has published more than 90 peer-reviewed articles in journals, conferences, and books. He has coedited two books: *Advances in Distributed and Parallel Knowledge Discovery*, AAAI/MIT Press, and *Data Mining: Next Generation Challenges and Future Directions*, AAAI/MIT Press. He is an associate editor of the *IEEE Transactions on Knowledge and Data Engineering* and *IEEE Transactions on Systems, Man, and Cybernetics, Part B*. He regularly serves in the organizing and program committee of many data mining conferences. More information about him can be found at <http://www.cs.umbc.edu/hillol>.



**Byung-Hoon Park** received the MS and PhD degrees in computer science both from the Washington State University in 1996 and 2001, respectively. He is currently a research scientist at the Computer Science and Mathematics Division of Oak Ridge National Laboratory (ORNL). His research areas include distributed data mining, computational biology, genetic computing, data stream analysis, and text mining. His research activities have been supported by the Genomes-to-Life program of Department of Energy (DOE), Scientific Data Management (SDM) of DOE SciDAC program, and the Biodefense Knowledge Center projects of Department of Homeland Security. Before joining the ORNL, Dr. Park was with University of Maryland Baltimore County as a postdoctoral research associate, where he was involved in NASA EOS distributed data mining project. He served on the program committees of several data mining conferences and workshops. He also serves as a reviewer of numerous journals and conferences.



**Haimonti Dutta** received her BS degree in Computer Science from Jadavpur University, Kolkata, India in 1999 and the MS degree in Computer and Information Science from Temple University, Philadelphia in 2002. She worked for an year as a Software Consultant at iGate Global Solutions, Bangalore. She is currently a Phd student in the Department of Computer Science and Electrical Engineering at the University of Maryland, Baltimore County. Her research interests include distributed data mining, data stream monitoring, grid mining and medical informatics.

## APPENDIX

### APPENDIX: THE EXPONENTIAL DECAY PROPERTY OF THE NON-BOOLEAN FOURIER SPECTRUM

Let us consider  $l$ -dimensional discrete domain  $\Lambda = \prod_{j=1}^l \{0, \dots, \lambda_j - 1\}$ , where each element  $\mathbf{x} \in \Lambda$  is denoted as  $(x_1, x_2, \dots, x_l)$ . Note that  $\lambda_1, \lambda_2, \dots, \lambda_l$  denote the cardinalities of  $x_1, x_2, \dots, x_l$  respectively. That is, the component  $x_j$  can take only the values  $0, 1, \dots, \lambda_j - 1$ .

The Fourier basis function that corresponds to the partition  $\mathbf{j}$  is,

$$\psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x}) = \prod_{m=1}^l \exp\left(\frac{2\pi i}{\lambda_m} x_m j_m\right)$$

The Fourier transform that corresponds to the partition  $\mathbf{j}$  is then,

$$w_{\mathbf{j}} = \prod_{i=1}^l \frac{1}{\lambda_i} \sum_{\mathbf{x}} \psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x}) \phi(\mathbf{x})$$

The inverse Fourier transform for an instance vector  $\mathbf{x}$  is,

$$f(\mathbf{x}) = \sum_{\mathbf{j}} w_{\mathbf{j}} \bar{\psi}_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$$

where  $\bar{\psi}_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$  is the complex conjugate of  $\psi_{\mathbf{j}}^{\bar{\lambda}}(\mathbf{x})$ .

Now we prove the exponential decay property of Fourier spectrum in a non-Boolean domain. For the sake of simplicity, let us further assume that each  $\lambda_j = 2^{q_j}$ , where  $q_j$  is any non-negative integer, so that the  $j$ -th non-binary variable can be represented using  $q_j$  bits. The results stated here are general and they extend easily to the case where the cardinalities are not powers of two. In the following proof, we treat a schema<sup>4</sup> as both a subset of  $\Lambda$  and a vector in  $\bar{\Lambda} = \prod_{j=1}^l \{*, 0, \dots, \lambda_j - 1\}$ .  $\bar{\Lambda}$  is the set of all schemata on  $\Lambda$ . The proof is based on transforming each non-Boolean vector in the  $l$ -dimensional space  $\Lambda$  to a (longer) Boolean vector in  $\Lambda' = \{0, 1\}^Q$ , where  $Q = \sum_{j=1}^l q_j$ . This is done by replacing each feature of a vector in  $\Lambda$  with its binary expansion. Following arguments establish a correspondence between the energy contained in corresponding sets of Fourier coefficients in the spectrum of the function defined on  $\Lambda$  and the function induced on  $\Lambda'$  by the transformation. Thus, since the exponential decay property holds for the Boolean case, it must hold for the general discrete case as well.

<sup>4</sup>A schema is a hyperplane that denotes a subset of domain members. It is essentially a similarity based equivalence class. Schemata is its plural form.

We next define the transformation

$$\kappa : \bar{\Lambda} \mapsto \bar{\Lambda}',$$

formally.  $\bar{\Lambda}'$  refers to the set of all schemata on  $\Lambda'$ , that is,  $\bar{\Lambda}' = \{*, 0, 1\}^Q$ . We first define feature-wise transformations

$$b_j : \{*, 0, 1, \dots, \lambda_j - 1\} \mapsto \{*, 0, 1\}^{q_j}$$

by

$$b_j(x_j) = \begin{cases} \overbrace{** \dots *}^{q_j} & \text{if } x_j = * \\ x'_j & \text{if } x_j \in \{0, 1, \dots, \lambda_j - 1\} \end{cases}$$

where  $x'_j$  is the  $q_j$ -bit binary representation of  $x_j$ . Now for any schema  $\mathbf{s} = (s_1, s_2, \dots, s_l) \in \bar{\Lambda}$ ,  $\kappa$  is defined as

$$\kappa : (s_1, s_2, \dots, s_l) \mapsto (b_1(s_1), b_2(s_2), \dots, b_l(s_l))$$

$\kappa$  is essentially a map from an  $l$ -feature schema in an arbitrary discrete domain to a  $Q$ -feature schema in a binary domain. We note here that we treat  $\Lambda$  as a subset of  $\bar{\Lambda}$  and thus can apply  $\kappa$  to elements of  $\Lambda$ . For a subset  $A \subseteq \Lambda$  we use the notation  $\kappa(A)$  to denote the set  $\{\kappa(\mathbf{x}) \mid \mathbf{x} \in A\}$ .

Let us further assume that  $f(\mathbf{x})$  is a functional representation of a decision tree  $T$  whose domain of definition is  $\Lambda$ . We establish some further before we proceed. We use  $\Gamma$  to denote the set of *schemata*  $\{*, 0\}^l$ . That is, schemata in  $\Gamma$  have only zeroes at their fixed (non-wildcard) features. We also define a set of schemata associated with any fixed schema  $\mathbf{s} \in \Gamma$ :

$$S(\mathbf{s}) = \{\mathbf{h} \in \bar{\Lambda} \mid h_k = *, \text{ if } s_k = 0 \text{ and } h_k \in \{0, 1, \dots, \lambda_k - 1\}, \text{ otherwise}\}$$

where  $h_k$  denotes the  $k$ -th feature in the schema  $\mathbf{h}$ , and similarly for  $s_k$ . Thus elements of  $S(\mathbf{s})$  can have wildcards at those positions where  $\mathbf{s}$  has zeroes.

*Lemma 6:* For any schema  $\mathbf{s} \in \Gamma$ ,

$$\sum_{\mathbf{i} \in \mathbf{s}} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}^{\bar{\Lambda}}(\mathbf{x}) = \sum_{\mathbf{j} \in \kappa(\mathbf{s})} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{x}))$$

**Proof:** For any  $\mathbf{i} \in \mathbf{s}$ ,

$$\begin{aligned}\eta_{\mathbf{i}} &= \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in \Lambda} f(\mathbf{x}) \psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x}) \\ &= \frac{1}{|S(\mathbf{s})|} \sum_{\mathbf{h} \in S(\mathbf{s})} \frac{1}{|\mathbf{h}|} \sum_{\mathbf{x} \in \mathbf{h}} f(\mathbf{x}) \psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x})\end{aligned}$$

where  $|\mathbf{h}|$  and  $|S(\mathbf{s})|$  denote the sizes of  $\mathbf{h}$  and  $S(\mathbf{s})$  respectively. Now for any  $\mathbf{i} \in \mathbf{s}$  and  $\mathbf{h} \in S(\mathbf{s})$ ,  $\psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x})$  is invariant over  $\mathbf{x} \in \mathbf{h}$ . Let us denote this value by  $\bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h})$ . Then we get

$$\begin{aligned}\eta_{\mathbf{i}} &= \frac{1}{|S(\mathbf{s})|} \sum_{\mathbf{h} \in S(\mathbf{s})} \frac{1}{|\mathbf{h}|} \sum_{\mathbf{x} \in \mathbf{h}} f(\mathbf{x}) \psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x}) \\ &= \frac{1}{|S(\mathbf{s})|} \sum_{\mathbf{h} \in S(\mathbf{s})} \frac{\psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h})}{|\mathbf{h}|} \sum_{\mathbf{x} \in \mathbf{h}} f(\mathbf{x}) \\ &= \frac{1}{|S(\mathbf{s})|} \sum_{\mathbf{h} \in S(\mathbf{s})} \phi(\mathbf{h}) \psi_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h})\end{aligned}$$

where  $\phi(\mathbf{h})$  is the average of  $f(\mathbf{x})$ , for  $\mathbf{x} \in \mathbf{h}$ . Now  $\phi(\mathbf{h})$  (by inverse Fourier transform) for any  $\mathbf{h} \in S(\mathbf{s})$  is

$$\phi(\mathbf{h}) = \frac{1}{|\mathbf{h}|} \sum_{\mathbf{x} \in \mathbf{h}} \sum_{\mathbf{i} \in \Lambda} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x}) = \frac{1}{|\mathbf{h}|} \sum_{\mathbf{i} \in \Lambda} \eta_{\mathbf{i}} \sum_{\mathbf{x} \in \mathbf{h}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x}) = \sum_{\mathbf{i} \in \mathbf{s}} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h})$$

since  $\sum_{\mathbf{x} \in \mathbf{h}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h})$  equals 0 if  $\mathbf{i} \notin \mathbf{s}$  and  $|\mathbf{h}|$  otherwise. Similarly,  $\phi(\kappa(\mathbf{h}))$  for any  $\mathbf{h} \in S(\mathbf{s})$  is

$$\phi(\kappa(\mathbf{h})) = \sum_{\mathbf{j} \in \mathbf{s}} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{h}))$$

Since for any  $\mathbf{h}$ ,  $\phi(\mathbf{h}) = \phi(\kappa(\mathbf{h}))$ ,

$$\sum_{\mathbf{i} \in \mathbf{s}} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{h}) = \sum_{\mathbf{j} \in \kappa(\mathbf{s})} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{h}))$$

Therefore,

$$\sum_{\mathbf{i} \in \mathbf{s}} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}^{\bar{\lambda}}(\mathbf{x}) = \sum_{\mathbf{j} \in \kappa(\mathbf{s})} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{x})), \text{ for all } \mathbf{x} \in \mathbf{h}.$$

This completes the proof of the lemma.

■

Now let us define  $\theta(\mathbf{s})$  as,

$$\theta(\mathbf{s}) = \{\mathbf{i} \in \mathbf{s} \mid o(\mathbf{i}) = o(\mathbf{s})\}$$

where  $o(\mathbf{i})$  and  $o(\mathbf{s})$  denote the orders of  $\mathbf{i}$  and  $\mathbf{s}$  respectively.  $\theta(\mathbf{s})$  is a subset of  $\mathbf{s}$  which only includes partitions whose orders are the same as that of  $\mathbf{s}$ . Now consider the following corollary.

*Corollary 1:* For any  $\theta(\mathbf{s})$  and  $\kappa(\theta(\mathbf{s}))$ ,

$$\sum_{\mathbf{i} \in \theta(\mathbf{s})} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}(\mathbf{x}) = \sum_{\mathbf{j} \in \kappa(\theta(\mathbf{s}))} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{x}))$$

**Proof:** Let  $P_{\mathbf{s},n}$  be the set of all schemata that are obtained by replacing  $n$  of the  $*$ -s in  $\mathbf{s}$  with zero. Then,

$$\begin{aligned} \sum_{\mathbf{i} \in \theta(\mathbf{s})} \eta_{\mathbf{i}} \psi_{\mathbf{i}}(\mathbf{x}) \bar{\psi}_{\mathbf{i}}(\mathbf{x}) &= \sum_{\mathbf{k} \in \mathbf{s}} \eta_{\mathbf{k}} \bar{\psi}_{\mathbf{k}}(\mathbf{x}) + \sum_{n=1}^{o(\mathbf{s})} (-1)^n \sum_{\mathbf{r} \in P_{\mathbf{s},n}} \sum_{\mathbf{k} \in \mathbf{r}} \eta_{\mathbf{k}} \bar{\psi}_{\mathbf{k}}(\mathbf{x}) \\ &= \sum_{\mathbf{m} \in \kappa(\mathbf{s})} w_{\mathbf{m}} \psi_{\mathbf{m}}(\kappa(\mathbf{x})) + \sum_{n=1}^{o(\mathbf{s})} (-1)^n \sum_{\mathbf{r} \in P_{\mathbf{s},n}} \sum_{\mathbf{m} \in \kappa(\mathbf{r})} w_{\mathbf{m}} \psi_{\mathbf{m}}(\kappa(\mathbf{x})) \\ &= \sum_{\mathbf{j} \in \kappa(\theta(\mathbf{s}))} w_{\mathbf{j}} \psi_{\mathbf{j}}(\kappa(\mathbf{x})) \end{aligned}$$

The second equality follows from Lemma 6.

■

Now let us rephrase Linial's original lemma as follows. That is, in the Boolean Fourier basis notation,

$$\sum_{o(\mathbf{j}) \geq k} w_{\mathbf{j}}^2 \leq \varphi(k)$$

where  $\varphi(k)$  decreases exponentially in  $k$ . Now consider the main theorem.

*Theorem 1:* In a non-Boolean  $l$ -dimensional discrete domain  $\mathbf{D}$ , for any non-negative integer  $k \leq l$  and Fourier spectrum  $\eta_{\mathbf{i}}$ -s of a decision tree defined over  $\mathbf{D}$ ,

$$\sum_{o(\mathbf{i}) \geq k} \|\eta_{\mathbf{i}}\|^2 \leq \varphi(k)$$

where  $\|\eta_{\mathbf{i}}\|$  denotes the magnitude of  $\eta_{\mathbf{i}}$ .

**Proof:** For the sake of convenience, let us define,

$$\begin{aligned} f_{\mathbf{s}}(\mathbf{x}) &= \sum_{\mathbf{i} \in \theta(\mathbf{s})} \eta_{\mathbf{i}} \bar{\psi}_{\mathbf{i}}(\mathbf{x}) \\ f_{\kappa(\mathbf{s})}(\kappa(\mathbf{x})) &= \sum_{\mathbf{j} \in \kappa(\theta(\mathbf{s}))} w_{\mathbf{j}} \psi_{\mathbf{j}}(\mathbf{x}) \end{aligned}$$

Then, by Corollary 1,  $f_{\mathbf{s}}(\mathbf{x}) = f_{\kappa(\mathbf{s})}(\kappa(\mathbf{x}))$ . Following *Parseval's Identity*,

$$\begin{aligned} \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in \Lambda} f_{\mathbf{s}}^2(\mathbf{x}) &= \frac{1}{|\Lambda|} \sum_{\mathbf{x} \in \Lambda} f_{\kappa(\mathbf{s})}^2(\kappa(\mathbf{x})) \\ &= \sum_{\mathbf{i} \in \theta(\mathbf{s})} \|\eta_{\mathbf{i}}\|^2 \\ &= \sum_{\mathbf{j} \in \kappa(\theta(\mathbf{s}))} w_{\mathbf{j}}^2 \end{aligned}$$

Since, for any  $\mathbf{i} \in \mathbf{s}$  and  $\mathbf{j} \in \kappa(\mathbf{s})$ ,  $o(\mathbf{i}) \leq o(\mathbf{j})$ ,

$$\begin{aligned} \sum_{o(\mathbf{s}) \geq k} \sum_{\mathbf{i} \in \theta(\mathbf{s})} \|\eta_{\mathbf{i}}\|^2 &= \sum_{\kappa(\mathbf{s})} \sum_{\mathbf{j} \in \kappa(\theta(\mathbf{s}))} w_{\mathbf{j}}^2 \\ &\leq \sum_{o(\mathbf{j}) \geq k} w_{\mathbf{j}}^2 \\ &\leq \varphi(k) \end{aligned}$$

Thus, the non-Boolean Fourier spectrum of a decision tree also has the exponential decay property.

■