

# Probabilistic Characterization of Decision Trees

**Amit Dhurandhar**

*Computer and Information Science and Engineering  
University of Florida  
Gainesville, FL 32611, USA*

ASD@CISE.UFL.EDU

**Alin Dobra**

*Computer and Information Science and Engineering  
University of Florida  
Gainesville, FL 32611, USA*

ADOBRA@CISE.UFL.EDU

**Editor:** Leslie Pack Kaelbling

## Abstract

In this paper we use the methodology introduced in Dhurandhar and Dobra (2006) for analyzing the error of classifiers and the model selection measures, to analyze decision tree algorithms. The methodology consists of obtaining parametric expressions for the moments of the Generalization error (GE) for the classification model of interest, followed by plotting these expressions for interpretability. The major challenge in applying the methodology to decision trees, the main theme of this work, is customizing the generic expressions for the moments of GE to this particular classification algorithm. The specific contributions we make in this paper are: (a) we completely characterize a subclass of decision trees namely, Random decision trees, (b) we discuss how the analysis extends to other decision tree algorithms, and (c) in order to extend the analysis to certain model selection measures, we generalize the relationships between the moments of GE and moments of the model selection measures given in Dhurandhar and Dobra (2006) to randomized classification algorithms. An extensive empirical comparison between the proposed method and Monte Carlo, depicts the advantages of the method in terms of running time and accuracy. It also showcases the use of the method as an exploratory tool to study learning algorithms.

## 1. Introduction

Model selection for classification is one of the major challenges in Machine Learning and Data-mining. Given an independent and identically distributed (i.i.d.) sample from the underlying probability distribution, the classification model selection problem consists in building a classifier by selecting among competing models. Ideally the model selected minimizes the Generalization error (GE) – the expected error over the entire input. Since GE cannot be directly computed, part of the sample is used to estimate GE through measures such as Cross-validation, Hold-out-set, Leave-one-out, etc. Though certain rules of thumb are followed by practitioners w.r.t. training size and other parameters specific to the validation measures in evaluating models through empirical studies Kohavi (1995); Blum et al. (1999) and certain asymptotic results exist Vapnik (1998); Shao (1993), the fact remains that most of these models and model selection measures are not well understood in real life (non-asymptotic) scenarios (e.g. what fraction should be test and training, what should be the value  $k$ , in  $k$ -fold cross validation, etc.). This lack of deep understanding limits our ability of utilizing the models most effectively and, maybe more importantly, trusting the models to perform well in a particular application – this is the single most important complaint from users of Machine Learning and Data-mining techniques.

Recently, a novel methodology was proposed in Dhurandhar and Dobra (2006) to study the behavior of models and model selection measures. Since the methodology is at the core of the

current work, we briefly describe it together with the motivation for using this type of analysis for classification in general and decision trees in particular.

### 1.1 *What is the methodology ?*

The methodology for studying classification models consists in studying the behavior of the first two central moments of the GE of the classification algorithm studied. The moments are taken over the space of all possible classifiers produced by the classification algorithm, by training it over all possible datasets sampled i.i.d. from some distribution. The first two moments give enough information about the statistical behavior of the classification algorithm to allow interesting observations about the behavior/trends of the classification algorithm w.r.t. any chosen data distribution.

### 1.2 *Why have such a methodology?*

The answers to the following questions shed light on why the methodology is necessary if tight statistical characterization is to be provided for classification algorithms.

1. *Why study GE ?* The biggest danger of learning is *overfitting* the training data. The main idea in using GE as a measure of success of learning instead on the empirical error on a given dataset is to provide a mechanism to avoid this pitfall. Implicitly, by analyzing GE all the input is considered.
2. *Why study the moments instead of the distribution of GE ?* Ideally, we would study the distribution of GE instead of moments in order to get a complete picture of what is its behavior. Studying the distribution of discrete random variables, except for very simple cases, turns out to be very hard. The difficulty comes from the fact that even computing the pdf in a single point is intractable since all combinations of random choices that result in the same value for GE have to be enumerated. On the other hand, the first two central moments coupled with distribution independent bounds such as Chebychev and Chernoff give guarantees about the worst possible behavior that are not too far from the actual behavior (small constant factor). Interestingly, it is possible to compute the moments of a random variable like GE without ever explicitly writing or making use of the formula for the pdf. What makes such an endeavor possible is extensive use of the linearity of expectation as explained in Dhurandhar and Dobra (2006).
3. *Why characterize a class of classifiers instead of a single classifier ?* While the use of GE as the success measure is standard practice in Machine Learning, characterizing classes of classifiers instead of the particular classifier produced on a given dataset is not. From the point of view of the analysis, without large testing datasets it is not possible to evaluate directly GE for a particular classifier. By considering classes of classifiers to which a classifier belongs, an indirect characterization is obtained for the particular classifier. This is precisely what Statistical Learning Theory (SLT) does; there the class of classifiers consists in all classifiers with the same VC dimension. The main problem with SLT results is that classes based on VC dimension are too large, thus results tend to be pessimistic. In the methodology in Dhurandhar and Dobra (2006), the class of classifiers consists only of the classifiers that are produced by the given classification algorithm from datasets of fixed size from the underlying distribution. This is the probabilistic smallest class in which the particular classifier produced on a given dataset can be placed in.

### 1.3 *How do we implement the methodology ?*

One way of approximately estimating the moments of GE over all possible classifiers for a particular classification algorithm is by directly using Monte Carlo. If we use Monte Carlo directly, we first

need to produce a classifier on a sampled dataset then test on a number of test sets sampled from the same distribution acquiring an estimate of the GE of this classifier. Repeating this entire procedure a couple of times we would acquire estimates of GE for different classifiers. Then by averaging the error of these multiple classifiers we would get an estimate of the first moment of GE. The variance of GE can also be similarly estimated.

Another way of estimating the moments of GE, is by obtaining parametric expressions for them. If this can be accomplished the moments can be computed exactly. Moreover, by dexterously observing the manner in which expressions are derived for a particular classification algorithm, insights can be gained into analyzing other algorithms of interest. Though deriving the expressions may be a tedious task, using them we obtain highly accurate estimates of the moments. In this paper, we propose this second alternative for analyzing a subclass of decision trees. The key to the analysis is focusing on the learning phase of the algorithm. In cases where the parametric expressions are computationally intensive to compute directly, we show that approximating individual terms using Monte Carlo we obtain accurate estimates of the moments when compared to directly using Monte Carlo (first alternative) for the same computational cost.

If the moments are to be studied on synthetic data then the distribution is anyway assumed and the parametric expressions can be directly used. If we have real data an empirical distribution can be built on the dataset and then the parametric expressions can be used.

#### 1.4 Applications of the methodology

It is important to note that the methodology is not aimed towards providing a way of estimating bounds for GE of a classifier on a given dataset. The primary goal is creating an avenue in which learning algorithms can be studied precisely i.e. studying the statistical behavior of a particular algorithm w.r.t. a chosen/built distribution. Below, we discuss the two most important perspectives in which the methodology can be applied.

##### 1.4.1 ALGORITHMIC PERSPECTIVE

If a researcher/practitioner designs a new classification algorithm, he/she needs to validate it. Standard practice is to validate the algorithm on a relatively small (5-20) number of datasets and to report the performance. By observing the behavior of only a few instances of the algorithm the designer infers its quality. Moreover, if the algorithm under performs on some datasets, it can be sometimes difficult to pinpoint the precise reason for its failure. If instead he/she is able to derive parametric expressions for the moments of GE, the test results would be more relevant to the particular classification algorithm, since the moments are over all possible datasets of a particular size drawn i.i.d. from some chosen/built distribution. Testing individually on all these datasets is an impossible task. Thus, by computing the moments using the parametric expressions the algorithm would be tested on a plethora of datasets with the results being highly accurate. Moreover, since the testing is done in a controlled environment i.e. all the parameters are known to the designer while testing, he/she can precisely pinpoint the conditions under which the algorithm performs well and the conditions under which the algorithm under performs.

##### 1.4.2 DATASET PERSPECTIVE

If an algorithm designer validates his/her algorithm by computing moments as mentioned earlier, it can instill greater confidence in the practitioner searching for an appropriate algorithm for his/her dataset. The reason for this being, if the practitioner has a dataset which has a similar structure or is from a similar source as the test dataset on which an empirical distribution was built and favourable results reported by the designer, then this would mean that the results apply not only to that particular test dataset, but to other similar type of datasets and since the practitioner's dataset

belongs to this similar collection, the results would also apply to his. Note that a distribution is just a weighting of different datasets and this perspective is used in the above exposition.

## 1.5 Specific Contributions

In this paper we develop a characterization for a subclass of decision trees. In particular, we characterize Random decision trees which are an interesting variant with respect to three popular stopping criteria namely; fixed height, purity and scarcity (i.e. fewer than some threshold number of points in a portion of the tree). The analysis directly applies to categorical as well as continuous attributes with split points predetermined for each attribute. Moreover, the analysis in Section 2.3 is applicable to even other deterministic attribute selection methods based on information gain, gini gain etc. These and other extensions of the analysis to continuous attributes with dynamically chosen split points is discussed in Section 4. In the experiments that ensue the theory, we compare the accuracy of the derived expressions with direct Monte Carlo on synthetic distributions as well as on distributions built on real data. Notice that using the expressions the moments can be computed without explicitly building the tree. We also extend the relationships between the moments of GE and moments of cross validation error (CE), leave-one-out error (LE) and hold-out-set error (HE) given in Dhurandhar and Dobra (2006) which were applicable only to deterministic classification algorithms, to be applicable to randomized classification algorithms.

## 2. Computing Moments

In this section we first provide the necessary technical groundwork, followed by customization of the expressions for decision trees. We now introduce some notation that is used primarily in this section.  $X$  is a random vector modeling input whose domain is denoted by  $\mathcal{X}$ .  $Y$  is a random variable modeling output whose domain is denoted by  $\mathcal{Y}$  (set of class labels).  $Y(x)$  is a random variable modeling output for input  $x$ .  $\zeta$  represents a particular classifier with its GE denoted by  $GE(\zeta)$ .  $\mathcal{Z}(N)$  denotes a set of classifiers obtained by application of a classification algorithm to different samples of size  $N$ .

### 2.1 Technical Framework

The basic idea in the generic characterization of the moments of GE as given in Dhurandhar and Dobra (2006), is to define a class of classifiers induced by a classification algorithm and an i.i.d. sample of a particular size from an underlying distribution. Each classifier in this class and its GE act as random variables, since the process of obtaining the sample is randomized. Since  $GE(\zeta)$  is a random variable, it has a distribution. Quite often though, characterizing a finite subset of moments turns out to be a more viable option than characterizing the entire distribution. Based on these facts, we revisit the expressions for the first two moments around zero of the GE of a classifier,

$$E_{\mathcal{Z}(N)} [GE(\zeta)] = \sum_{x \in \mathcal{X}} P[X=x] \sum_{y \in \mathcal{Y}} P_{\mathcal{Z}(N)} [\zeta(x)=y] P[Y(x) \neq y] \quad (1)$$

$$E_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [GE(\zeta)GE(\zeta')] = \sum_{x \in \mathcal{X}} \sum_{x' \in \mathcal{X}} P[X=x] P[X=x'] \cdot \sum_{y \in \mathcal{Y}} \sum_{y' \in \mathcal{Y}} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x)=y \wedge \zeta'(x')=y'] \cdot P[Y(x) \neq y] P[Y(x') \neq y'] \quad (2)$$

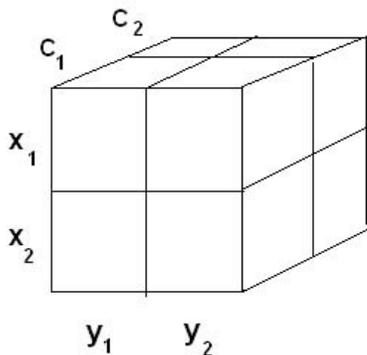


Figure 1: Contingency table with 2 attributes each having 2 values and 2 classes.

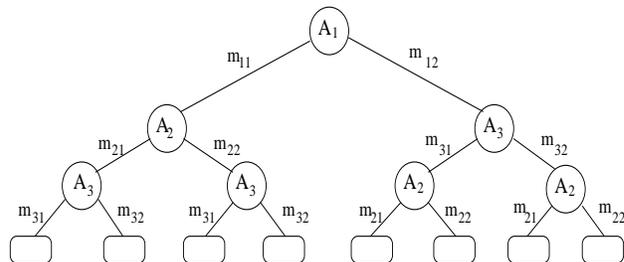


Figure 2: The all attribute tree with 3 attributes  $A_1, A_2, A_3$ , each having 2 values.

From the above equations we observe that for the first moment we have to characterize the behavior of the classifier on each input separately while for the second moment we need to observe its behavior on pairs of inputs. In particular, to derive expressions for the moments of any classification algorithm we need to characterize  $P_{\mathcal{Z}(N)} [\zeta(x) = y]$  for the first moment and  $P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = y \wedge \zeta'(x') = y']$  for the second moment. The values for the other terms denote the error of the classifier for the first moment and errors of two classifiers for the second moment which are obtained directly from the underlying joint distribution. For example, if we have data with a class prior  $p$  for class 1 and  $1-p$  for class 2. Then the error of a classifier classifying data into class 1 is  $1-p$  and the error of a classifier classifying data into class 2 is given by  $p$ . We now focus our attention on relating the above two probabilities, to probabilities that can be computed using the joint distribution and the classification model viz. Decision Trees.

In the subsections that follow we assume the following setup. We consider the dimensionality of the input space to be  $d$ .  $A_1, A_2, \dots, A_d$  are the corresponding discrete attributes or continuous attributes with predetermined split points.  $a_1, a_2, \dots, a_d$  are the number of attribute values/the number of splits of the attributes  $A_1, A_2, \dots, A_d$  respectively.  $m_{ij}$  is the  $i^{\text{th}}$  attribute value/split of the  $j^{\text{th}}$  attribute, where  $i \leq a_j$  and  $j \leq d$ . Let  $C_1, C_2, \dots, C_k$  be the class labels representing  $k$  classes and  $N$  the sample size.

## 2.2 All Attribute Decision Trees (ATT)

Let us consider a decision tree algorithm whose only stopping criteria is that no attributes remain when building any part of the tree. In other words, every path in the tree from root to leaf has all the attributes. An example of such a tree is shown in Figure 2. It can be seen that irrespective of the split attribute selection method (e.g. information gain, gini gain, randomised selection, etc.) the above stopping criteria yields trees with the same leaf nodes. Thus although a particular path in one tree has an ordering of attributes that might be different from a corresponding path in other trees, the leaf nodes will represent the same region in space or the same set of datapoints. This is seen in Figure 3. Moreover, since predictions are made using data in the leaf nodes, any deterministic way of prediction would lead to these trees resulting in the same classifier for a given sample and thus having the same GE. Usually, prediction in the leaves is performed by choosing the most numerous class as the class label for the corresponding datapoint. With this we arrive at the expressions for

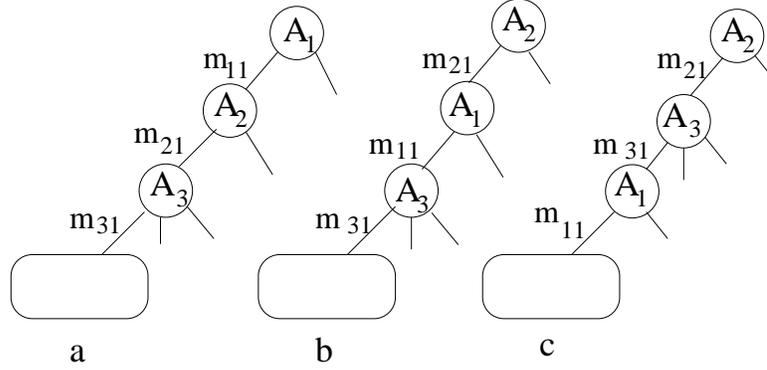


Figure 3: Given 3 attributes  $A_1, A_2, A_3$ , the path  $m_{11}m_{21}m_{31}$  is formed irrespective of the ordering of the attributes. Three such permutations are shown in the above figure.

computing the aforementioned probabilities,

$$P_{\mathcal{Z}(N)} [\zeta(x) = C_i] = \frac{P_{\mathcal{Z}(N)} [ct(m_{p_1}m_{q_2}\dots m_{r_d}C_i) > ct(m_{p_1}m_{q_2}\dots m_{r_d}C_j)]}{\sum_{j \neq i} P_{\mathcal{Z}(N)} [ct(m_{p_1}m_{q_2}\dots m_{r_d}C_i) > ct(m_{p_1}m_{q_2}\dots m_{r_d}C_j)]},$$

$$\forall j \neq i, i, j \in [1, \dots, k]$$

where  $x = m_{p_1}m_{q_2}\dots m_{r_d}$  represents a datapoint which is also a path from root to leaf in the tree.  $ct(m_{p_1}m_{q_2}\dots m_{r_d}C_i)$  is the count of the datapoints specified by the cell  $m_{p_1}m_{q_2}\dots m_{r_d}C_i$ . For example in Figure 1  $x_1y_1C_1$  represents a cell. Henceforth, when using the word "path" we will strictly imply path from root to leaf. By computing the above probability  $\forall i$  and  $\forall x$  we can compute the first moment of the GE for this classification algorithm.

Similarly, for the second moment we compute cumulative joint probabilities of the following form:

$$P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta'(x') = C_v] = \frac{P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(m_{p_1}\dots m_{r_d}C_i) > ct(m_{p_1}\dots m_{r_d}C_j), ct(m_{f_1}\dots m_{h_d}C_v) > ct(m_{f_1}\dots m_{h_d}C_w)]}{\sum_{j \neq i, w \neq v} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(m_{p_1}\dots m_{r_d}C_i) > ct(m_{p_1}\dots m_{r_d}C_j), ct(m_{f_1}\dots m_{h_d}C_v) > ct(m_{f_1}\dots m_{h_d}C_w)]},$$

$$\forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]$$

where the terms have similar conotation as before. These probabilities can be computed exactly or by using fast approximation techniques proposed in Dhurandhar and Dobra (2006).

### 2.3 Decision Trees with Non-trivial Stopping Criteria

We just considered decision trees which are grown until all attributes are exhausted. In real life though we seldom build such trees. The main reasons for this could be any of the following: we wish to build small decision trees to save space; certain path counts (i.e. number of datapoints in the leaves) are extremely low and hence we want to avoid splitting further, as the predictions can get arbitrarily bad; we have split on a certain subset of attributes and all the datapoints in that path belong to the same class (purity based criteria); we want to grow trees to a fixed height (or depth). These stopping measures would lead to paths in the tree that contain a subset of the entire set of attributes. Thus from a classification point of view we cannot simply compare the counts in two cells as we did previously. The reason for this being that the corresponding path may not be present in the tree. Hence, we need to check that the path exists and then compare cell counts. Given the classification algorithm, since the  $P_{\mathcal{Z}(N)} [\zeta(x) = C_i]$  is the probability of all possible ways in which an input  $x$  can be classified into class  $C_i$  for a decision tree it equates to finding the following kind

of probability for the first moment,

$$\begin{aligned}
 P_{\mathcal{Z}(N)} [\zeta(x) = C_i] = \\
 \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p exists, \\
 \forall j \neq i, i, j \in [1, \dots, k]]
 \end{aligned} \tag{3}$$

where  $p$  indexes all allowed paths by the tree algorithm in classifying input  $x$ . After the summation, the right hand side term above is the probability that the cell  $path_p C_i$  has the greatest count, with the path " $path_p$ " being present in the tree. This will become clearer when we discuss different stopping criteria. Notice that the characterization for the ATT is just a special case of this more generic characterization.

The probability that we need to find for the second moment is,

$$\begin{aligned}
 P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta'(x') = C_v] = \\
 \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p exists, \\
 ct(path_q C_v) > ct(path_q C_w), path_q exists, \\
 \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]]
 \end{aligned} \tag{4}$$

where  $p$  and  $q$  index all allowed paths by the tree algorithm in classifying input  $x$  and  $x'$  respectively. The above two equations are generic in analyzing any decision tree algorithm which classifies inputs into the most numerous class in the corresponding leaf. It is not difficult to generalize it further when the decision in leaves is some other measure than majority. In that case we would just include that measure in the probability in place of the inequality.

### 2.3.1 CHARACTERIZING *path exists* FOR THREE STOPPING CRITERIA

It follows from above that to compute the moments of the GE for a decision tree algorithm we need to characterize conditions under which particular paths are present. This characterization depends on the stopping criteria and split attribute selection method in a decision tree algorithm. We now look at three popular stopping criteria, namely a) Fixed height based, b) Purity (i.e. entropy 0 or gini index 0 etc.) based and c) Scarcity (i.e. too few datapoints) based. We consider conditions under which certain paths are present for each stopping criteria. Similar conditions can be enumerated for any reasonable stopping criteria. We then choose a split attribute selection method, thereby fully characterizing the above two probabilities and hence the moments.

1. **Fixed Height:** This stopping criteria is basically that every path in the tree should be of length exactly  $h$ , where  $h \in [1, \dots, d]$ . If  $h = 1$  we classify based on just one attribute. If  $h = d$  then we have the all attribute tree.

In general, a path  $m_{i_1} m_{j_2} \dots m_{l_h}$  is present in the tree iff the attributes  $A_1, A_2, \dots, A_h$  are chosen in any order to form the path for a tree construction during the split attribute selection phase. Thus, for any path of length  $h$  to be present we biconditionally imply that the corresponding attributes are chosen.

2. **Purity:** This stopping criteria implies that we stop growing the tree from a particular split of a particular attribute if all datapoints lying in that split belong to the same class. We call such a path pure else we call it impure. In this scenario, we could have paths of length 1 to  $d$  depending on when we encounter purity (assuming all datapoints don't lie in 1 class). Thus, we have the following two separate checks for paths of length  $d$  and less than  $d$  respectively.

a) Path  $m_{i_1} m_{j_2} \dots m_{l_d}$  present iff the path  $m_{i_1} m_{j_2} \dots m_{l_{(d-1)}}$  is impure and attributes  $A_1, A_2, \dots, A_{d-1}$  are chosen above  $A_d$ , or  $m_{i_1} m_{j_2} \dots m_{s_{(d-2)}} m_{l_d}$  is impure and attributes  $A_1, A_2, \dots, A_{d-2}, A_d$

are chosen above  $A_{d-1}$ , or ... or  $m_{j2}...m_{ld}$  is impure and attributes  $A_2, \dots, A_d$  are chosen above  $A_1$ .

This means that if a certain set of  $d - 1$  attributes are present in a path in the tree then we split on the  $d^{th}$  attribute iff the current path is not pure, finally resulting in a path of length  $d$ .

b) Path  $m_{i1}m_{j2}...m_{lh}$  present where  $h < d$  iff the path  $m_{i1}m_{j2}...m_{lh}$  is pure and attributes  $A_1, A_2, \dots, A_{h-1}$  are chosen above  $A_h$  and  $m_{i1}m_{j2}...m_{l(h-1)}$  is impure or the path  $m_{i1}m_{j2}...m_{lh}$  is pure and attributes  $A_1, A_2, \dots, A_{h-2}, A_h$  are chosen above  $A_{h-1}$  and  $m_{i1}m_{j2}...m_{l(h-2)}m_{lh}$  is impure or ... or the path  $m_{j2}...m_{lh}$  is pure and attributes  $A_2, \dots, A_h$  are chosen above  $A_1$  and  $m_{j2}...m_{lh}$  is impure.

This means that if a certain set of  $h - 1$  attributes are present in a path in the tree then we split on some  $h^{th}$  attribute iff the current path is not pure and the resulting path is pure.

The above conditions suffice for "path present" since the purity property is anti-monotone and the impurity property is monotone.

3. **Scarcity:** This stopping criteria implies that we stop growing the tree from a particular split of a certain attribute if its count is less than or equal to some pre-specified pruning bound. Let us denote this number by  $pb$ . As before, we have the following two separate checks for paths of length  $d$  and less than  $d$  respectively.

a) Path  $m_{i1}m_{j2}...m_{ld}$  present iff the attributes  $A_1, \dots, A_{d-1}$  are chosen above  $A_d$  and  $ct(m_{i1}m_{j2}...m_{l(d-1)}) > pb$  or the attributes  $A_1, \dots, A_{d-2}, A_d$  are chosen above  $A_{d-1}$  and  $ct(m_{i1}m_{j2}...m_{l(d-2)}m_{nd}) > pb$  or ... or the attributes  $A_2, \dots, A_d$  are chosen above  $A_1$  and  $ct(m_{i2}m_{j3}...m_{ld}) > pb$ .

b) Path  $m_{i1}m_{j2}...m_{lh}$  present where  $h < d$  iff the attributes  $A_1, \dots, A_{h-1}$  are chosen above  $A_h$  and  $ct(m_{i1}m_{j2}...m_{l(h-1)}) > pb$  and  $ct(m_{i1}m_{j2}...m_{lh}) \leq pb$  or the attributes  $A_1, \dots, A_{h-2}, A_h$  are chosen above  $A_{h-1}$  and  $ct(m_{i1}m_{j2}...m_{l(h-2)}m_{nh}) > pb$  and  $ct(m_{i1}m_{j2}...m_{nh}) \leq pb$  or ... or the attributes  $A_2, \dots, A_h$  are chosen above  $A_1$  and  $ct(m_{i2}m_{j3}...m_{lh}) > pb$  and  $ct(m_{i1}m_{j2}...m_{lh}) \leq pb$ .

This means that we stop growing the tree under a node once we find that the next chosen attribute produces a path with occupancy  $\leq pb$ .

The above conditions suffice for "path present" since the occupancy property is monotone.

We observe from the above checks that we have two types of conditions that need to be evaluated for a path being present namely, i) those that depend on the sample viz.  $m_{i1}m_{j2}...m_{l(d-1)}$  is impure or  $ct(m_{i1}m_{j2}...m_{lh}) > pb$  and ii) those that depend split attribute selection method viz.  $A_1, A_2, \dots, A_h$  are chosen. The former depends on the data distribution which we have specified to be a multinomial. The latter we discuss in the next subsection. Note that checks for a combination of the above stopping criteria can be obtained by appropriately combining the individual checks.

## 2.4 Split Attribute Selection

In decision tree construction algorithms, at each iteration we have to decide the attribute variable on which the data should be split. Numerous measures have been developed Hall and Holmes (2003). Some of the most popular ones aim to increase the purity of a set of datapoints that lie in the region formed by that split. The purer the region, the better the prediction and lower the error of the classifier. Measures such as, i) Information Gain (IG) Quinlan (1986), ii) Gini Gain (GG) Breiman et al. (1984), iii) Gain Ratio (GR) Quinlan (1986), iv) Chi-square test (CS) Shao (2003) etc. aim at realising this intuition. Other measures using Principal Component Analysis Smith (2002), Correlation-based measures Hall (1998) have also been developed. Another interesting yet non-intuitive measure in terms of its utility is the Random attribute selection measure. According to this measure we randomly choose the split attribute from available set. The decision tree that

this algorithm produces is called a Random decision tree (RDT). Surprisingly enough, a collection of RDTs quite often outperform their seemingly more powerful counterparts Liu et al. (2005). In this paper we study this interesting variant. We do this by first presenting a probabilistic characterization in selecting a particular attribute/set of attributes, followed by simulation studies. Characterizations for the other measures can be developed in similar vein by focusing on the working of each measure. As an example, for the deterministic purity based measures mentioned above the split attribute selection is just a function of the sample and thus by appropriately conditioning on the sample we can find the relevant probabilities and hence the moments.

Before presenting the expression for the probability of selecting a split attribute/attributes in constructing a RDT we extend the results in Dhurandhar and Dobra (2006) where relationships were drawn between the moments of HE, CE, LE (just a special case of cross-validation) and GE, to be applicable to randomized classification algorithms. The random process is assumed to be independent of the sampling process. This result is required since the results in Dhurandhar and Dobra (2006) are applicable to deterministic classification algorithms and we would be analyzing RDT's. With this we have the following lemma.

**Lemma 1** *Let  $D$  and  $T$  be independent discrete random variables, with some distribution defined on each of them. Let  $\mathcal{D}$  and  $\mathcal{T}$  denote the domains of the random variables. Let  $f(d, t)$  and  $g(d, t)$  be two functions such that  $\forall t \in \mathcal{T} E_{\mathcal{D}}[f(d, t)] = E_{\mathcal{D}}[g(d, t)]$  and  $d \in \mathcal{D}$ . Then,  $E_{\mathcal{T} \times \mathcal{D}}[f(d, t)] = E_{\mathcal{T} \times \mathcal{D}}[g(d, t)]$*

**Proof**

$$\begin{aligned} E_{\mathcal{T} \times \mathcal{D}}[f(d, t)] &= \sum_{t \in \mathcal{T}} \sum_{d \in \mathcal{D}} f(d, t) P[T = t, D = d] \\ &= \sum_{t \in \mathcal{T}} \sum_{d \in \mathcal{D}} f(d, t) P[D = d] P[T = t] \\ &= \sum_{t \in \mathcal{T}} E_{\mathcal{D}}[g(d, t)] P[T = t] \\ &= E_{\mathcal{T} \times \mathcal{D}}[g(d, t)] \end{aligned}$$

■

The result is valid even when  $D$  and  $T$  are continuous, but considering the scope of this paper we are mainly interested in the discrete case. This result implies that all the relationships and expressions in Dhurandhar and Dobra (2006) hold with an extra expectation over the  $t$ 's, for randomized classification algorithms where the random process is independent of the sampling process. In equations 1 and 2 the expectations w.r.t.  $Z(N)$  become expectations w.r.t.  $Z(N, t)$ .

## 2.5 Random Decision Trees

In this subsection we explain the randomized process used for split attribute selection and provide the expression for the probability of choosing an attribute/a set of attributes. The attribute selection method we use is as follows. We assume a uniform probability distribution in selecting the attribute variables i.e. attributes which have already not been chosen in a particular branch, have an equal chance of being chosen for the next level. The random process involved in attribute selection is independent of the sample and hence the lemma 1 applies. We now give the expression for the probability of selecting a subset of attributes from the given set for a path. This expression is required in the computation of the above mentioned probabilities used in computing the moments. For the first moment we need to find the following probability. Given  $d$  attributes  $A_1, A_2, \dots, A_d$

the probability of choosing a set of  $h$  attributes where  $h \in \{1, 2, \dots, d\}$  is,

$$P[h \text{ attributes chosen}] = \frac{1}{\binom{d}{h}}$$

since choosing without replacement is equivalent to simultaneously choosing a subset of attributes from the given set.

For the second moment when the trees are different (required in the finding of variance of CE since, the training sets in the various runs in cross validation are different i.e. for finding  $E_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [GE(\zeta)GE(\zeta')]$ ), the probability of choosing  $l_1$  attributes for path in one tree and  $l_2$  attributes for path in another tree where  $l_1, l_2 \leq d$  is given by,

$$P[l_1 \text{ attribute path in tree 1, } l_2 \text{ attribute path in tree 2}] = \frac{1}{\binom{d}{l_1} \binom{d}{l_2}}$$

since the process of choosing one set of attributes for a path in one tree is independent of the process of choosing another set of attributes for a path in a different tree.

For the second moment when the tree is the same (required in the finding of variance of GE and HE i.e. for finding  $E_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [GE(\zeta)^2]$ ), the probability of choosing two sets of attributes such that the two distinct paths resulting from them co-exist in a single tree is given by the following. Assume we have  $d$  attributes  $A_1, A_2, \dots, A_d$ . Let the lengths of the two paths (or cardinality of the two sets) be  $l_1$  and  $l_2$  respectively, where  $l_1, l_2 \leq d$ . Without loss of generality assume  $l_1 \leq l_2$ . Let  $p$  be the number of attributes common to both paths. Notice that  $p \geq 1$  is one of the necessary conditions for the two paths to co-exist. Let  $v \leq p$  be those attributes among the total  $p$  that have same values for both paths. Thus  $p - v$  attributes are common to both paths but have different values. At one of these attributes in a given tree the two paths will bifurcate. The probability that the two paths co-exist given our randomized attribute selection method is computed by finding out all possible ways in which the two paths can co-exist in a tree and then multiplying the number of each kind of way by the probability of having that way. A detailed proof is given in the appendix. The expression for the probability based on the attribute selection method is,

$$P[l_1 \text{ and } l_2 \text{ length paths co-exist}] = \sum_{i=0}^v vPr_i (l_1 - i - 1)! (l_2 - i - 1)! (p - v) prob_i$$

where  $vPr_i = \frac{v!}{(v-i)!}$  denotes permutation and  $prob_i = \frac{1}{d(d-1)\dots(d-i)(d-i-1)^2\dots(d-l_1+1)(d-l_1)\dots(d-l_2+1)}$  is the probability of the  $i^{th}$  possible way. For fixed height trees of height  $h$ ,  $(l_1 - i - 1)!(l_2 - i - 1)!$  becomes  $(h - i - 1)!^2$  and  $prob_i = \frac{1}{d(d-1)\dots(d-i)(d-i-1)^2\dots(d-h+1)^2}$ .

## 2.6 Putting things together

We now have all the ingredients that are required for the computation of the moments of GE. In this subsection we combine the results derived in the previous subsections to obtain expressions for  $P_{\mathcal{Z}(N)} [\zeta(x) = C_i]$  and  $P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta'(x') = C_v]$  which are vital in the computation of the moments.

Let *s.c.c.s.* be an abbreviation for stopping criteria conditions that are sample dependent. Conversely, *s.c.c.i.* be an abbreviation for stopping criteria conditions that are sample independent or conditions that are dependent on the attribute selection method. We now provide expressions for the above probabilities categorized by the 3 stopping criteria.

## 2.6.1 FIXED HEIGHT

The conditions for "path exists" for fixed height trees depend only on the attribute selection method as seen in subsection 2.3.1. Hence the probability used in finding the first moment is given by,

$$\begin{aligned}
 P_{\mathcal{Z}(N)} [\zeta(x) = C_i] &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{ exists}, \forall j \neq i, i, j \in [1, \dots, k]] \\
 &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), s.c.c.i., \forall j \neq i, i, j \in [1, \dots, k]] \\
 &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), \forall j \neq i, i, j \in [1, \dots, k]] P_{\mathcal{Z}(N)} [s.c.c.i.] \quad (5) \\
 &= \sum_p \frac{P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), \forall j \neq i, i, j \in [1, \dots, k]]}{dC_h}
 \end{aligned}$$

where  $dC_h = \frac{d!}{h!(d-h)!}$  and  $h$  is the length of the paths or the height of the tree. The probability in the last step of the above derivation can be computed from the underlying joint distribution. The probability for the second moment when the trees are different is given by,

$$\begin{aligned}
 P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta'(x') = C_v] &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{ exists}, ct(path_q C_v) > ct(path_q C_w), path_q \text{ exists}, \\
 &\quad \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \\
 &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \cdot \\
 &\quad P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [s.c.c.i.] \\
 &= \sum_{p,q} \frac{P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]]}{dC_h^2} \quad (6)
 \end{aligned}$$

where  $h$  is the length of the paths. The probability for the second moment when the trees are identical is given by,

$$\begin{aligned}
 P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta(x') = C_v] &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{ exists}, ct(path_q C_v) > ct(path_q C_w), path_q \text{ exists}, \\
 &\quad \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \\
 &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \cdot \\
 &\quad P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [s.c.c.i.] \\
 &= \sum_{p,q} \sum_{t=0}^b bPr_t (h-t-1)!^2 (r-v) prob_t P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), \\
 &\quad \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \quad (7)
 \end{aligned}$$

where  $r$  is the number of attributes that are common in the 2 paths,  $b$  is the number of attributes that have the same value in the 2 paths,  $h$  is the length of the paths and  $prob_t = \frac{1}{d(d-1)\dots(d-t)(d-t-1)^2\dots(d-h+1)^2}$ . As before, the probability comparing counts can be computed from the underlying joint distribution.

### 2.6.2 PURITY AND SCARCITY

The conditions for "path exists" in the case of purity and scarcity depend on both the sample and the attribute selection method as can be seen in 2.3.1. The probability used in finding the first moment is given by,

$$\begin{aligned}
 P_{\mathcal{Z}(N)} [\zeta(x) = C_i] &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{ exists}, \forall j \neq i, i, j \in [1, \dots, k]] \\
 &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), s.c.c.i, s.c.c.s., \forall j \neq i, i, j \in [1, \dots, k]] \\
 &= \sum_p P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), s.c.c.s., \forall j \neq i, i, j \in [1, \dots, k]] P_{\mathcal{Z}(N)} [s.c.c.i.] \\
 &= \sum_p \frac{P_{\mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), s.c.c.s., \forall j \neq i, i, j \in [1, \dots, k]]}{dC_{h_p-1}(d-h_p+1)}
 \end{aligned} \tag{8}$$

where  $h_p$  is the length of the path indexed by  $p$ . The joint probability of comparing counts and *s.c.c.s.* can be computed from the underlying joint distribution. The probability for the second moment when the trees are different is given by,

$$\begin{aligned}
 P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta'(x') = C_v] &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{ exists}, ct(path_q C_v) > ct(path_q C_w), path_q \text{ exists}, \\
 &\quad \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \\
 &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), s.c.c.s., \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \cdot \\
 &\quad P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [s.c.c.i.] \\
 &= \sum_{p,q} \frac{P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), s.c.c.s., \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]]}{dC_{h_p-1} dC_{h_q-1} (d-h_p+1)(d-h_q+1)}
 \end{aligned} \tag{9}$$

where  $h_p$  and  $h_q$  are the lengths of the paths indexed by  $p$  and  $q$ . The probability for the second moment when the trees are identical is given by,

$$\begin{aligned}
 & P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [\zeta(x) = C_i \wedge \zeta(x') = C_v] \\
 &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), path_p \text{exists}, ct(path_q C_v) > ct(path_q C_w), path_q \text{exists}, \\
 &\quad \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]] \\
 &= \sum_{p,q} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), ct(path_q C_v) > ct(path_q C_w), s.c.c.s., \forall j \neq i, \forall w \neq v, \\
 &\quad i, j, v, w \in [1, \dots, k]] P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [s.c.c.i.] \\
 &= \sum_{p,q} \sum_{t=0}^b \frac{bPr_t(h_p - t - 2)!(h_q - t - 2)!(r - v)prob_t}{(d - h_p + 1)(d - h_q + 1)} P_{\mathcal{Z}(N) \times \mathcal{Z}(N)} [ct(path_p C_i) > ct(path_p C_j), \\
 &\quad ct(path_q C_v) > ct(path_q C_w), s.c.c.s., \forall j \neq i, \forall w \neq v, i, j, v, w \in [1, \dots, k]]
 \end{aligned} \tag{10}$$

where  $r$  is the number of attributes that are common in the 2 paths sparing the attributes chosen as leaves,  $b$  is the number of attributes that have the same value,  $h_p$  and  $h_q$  are the lengths of the 2 paths and without loss of generality assuming  $h_p \leq h_q$   $prob_t = \frac{1}{d(d-1)\dots(d-t)(d-t-1)^2\dots(d-h_p)^2(d-h_p-1)\dots(d-h_q)}$ . As before, the probability of comparing counts and s.c.c.s. can be computed from the underlying joint distribution.

Using the expressions for the above probabilities the moments of GE can be computed. In next section we perform experiments on synthetic as well as distributions built on real data to portray the efficacy of the derived expressions.

### 3. Experiments

To exactly compute the probabilities for each path the time complexity for fixed height trees is  $O(N^2)$  and for purity and scarcity based trees is  $O(N^3)$ . Hence, computing exactly the probabilities and consequently the moments is practical for small values of  $N$ . For larger values of  $N$ , we propose computing the individual probabilities using Monte Carlo (MC). In the empirical studies we report, we show that the accuracy in estimating the error (i.e. the moments of GE) by using our expressions with MC is always greater than by directly using MC for the same computational cost. In fact, the accuracy of using the expressions is never worse than MC even when MC is executed for 10 times the number of iterations as those of the expressions. The true error or the golden standard against which we compare the accuracy of these estimators is obtained by running MC for a week, which is around 200 times the number of iterations as those of the expressions.

**Notation:** In the experiments, AF refers to the estimates obtained by using the expressions in conjunction with Monte Carlo. MC- $i$  refers to simple Monte Carlo being executed for  $i$  times the number of iterations as those of the expressions. The term True Error or TE refers to the golden standard against which we compare AF and MC- $i$ .

**General Setup:** We perform empirical studies on synthetic as well as real data. The experimental setup for synthetic data is as follows: We fix  $N$  to 10000. The number of classes is fixed to two. We observe the behavior of the error for the three kinds of trees with the number of attributes fixed to  $d = 5$  and each attribute having 2 attribute values. We then increase the number of attribute values to 3, to observe the effect that increasing the number of split points has on the performance of the estimators. We also increase the number of attributes to  $d = 8$  to study the effect that increasing the number of attributes has on the performance. With this we have a  $d + 1$  dimensional contingency table whose  $d$  dimensions are the attributes and the  $(d + 1)^{th}$  dimension represents the class labels. When each attribute has two values the total number of cells in the table is  $c = 2^{d+1}$  and with three values the total number of cells is  $c = 3^d \times 2$ . If we fix the probability of observing a datapoint in

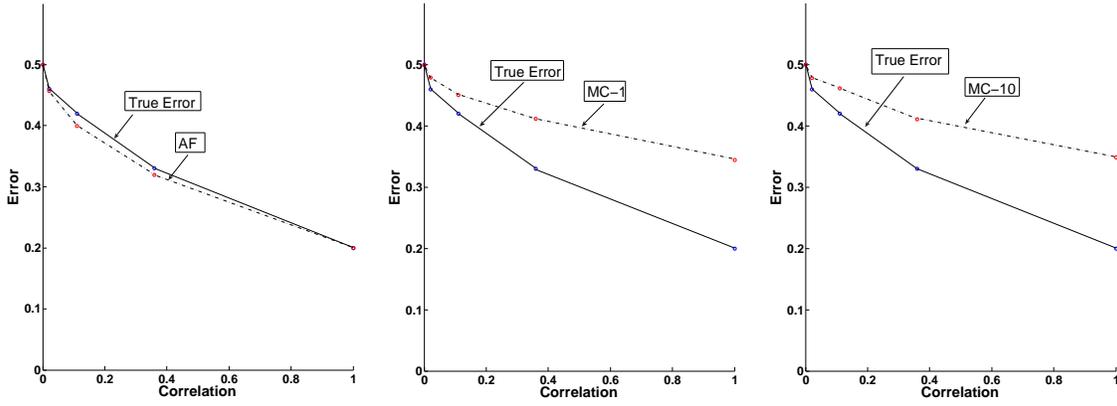


Figure 4: Fixed Height trees with  $d = 5$ ,  $h = 3$  and attributes with binary splits.

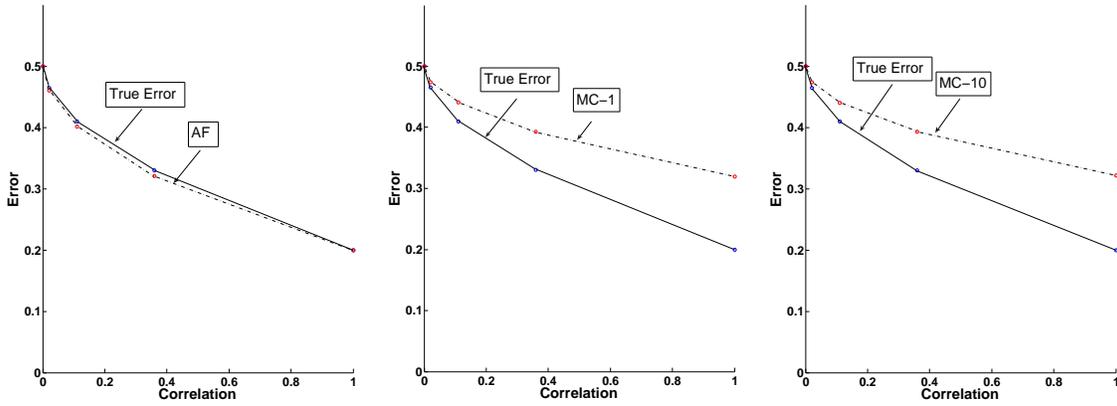


Figure 5: Fixed Height trees with  $d = 5$ ,  $h = 3$  and attributes with ternary splits.

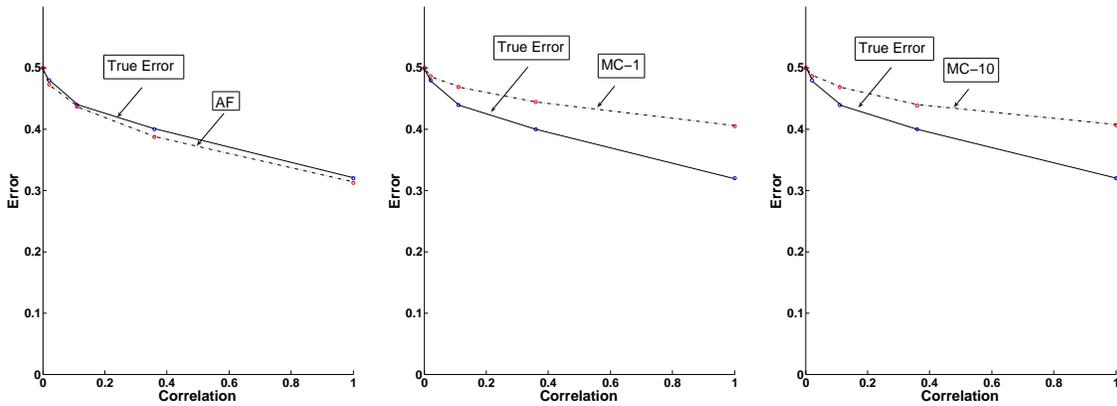


Figure 6: Fixed Height trees with  $d = 8$ ,  $h = 3$  and attributes with binary splits.

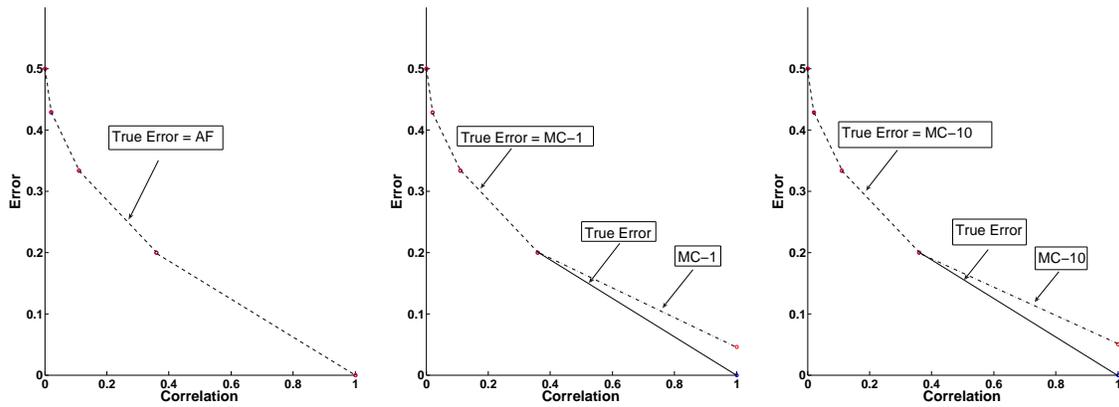


Figure 7: Purity based trees with  $d = 5$  and attributes with binary splits.

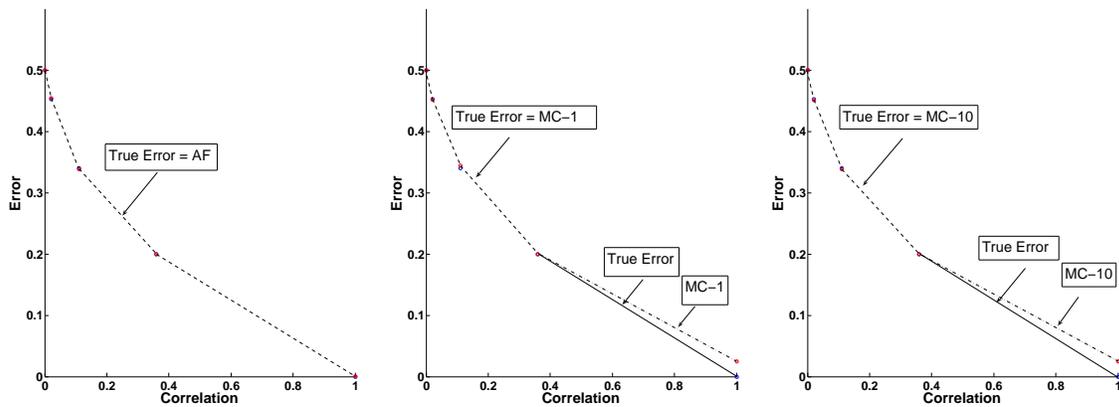


Figure 8: Purity based trees with  $d = 5$  and attributes with ternary splits.

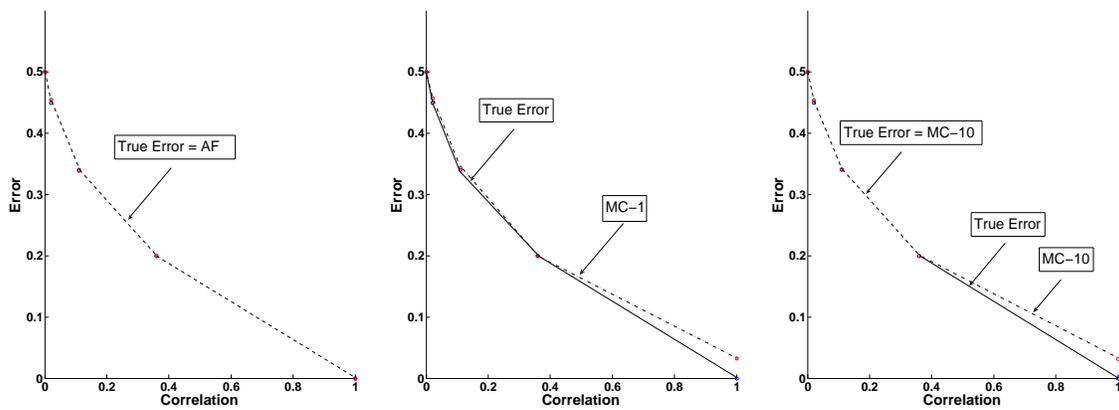


Figure 9: Purity based trees with  $d = 8$  and attributes with binary splits.

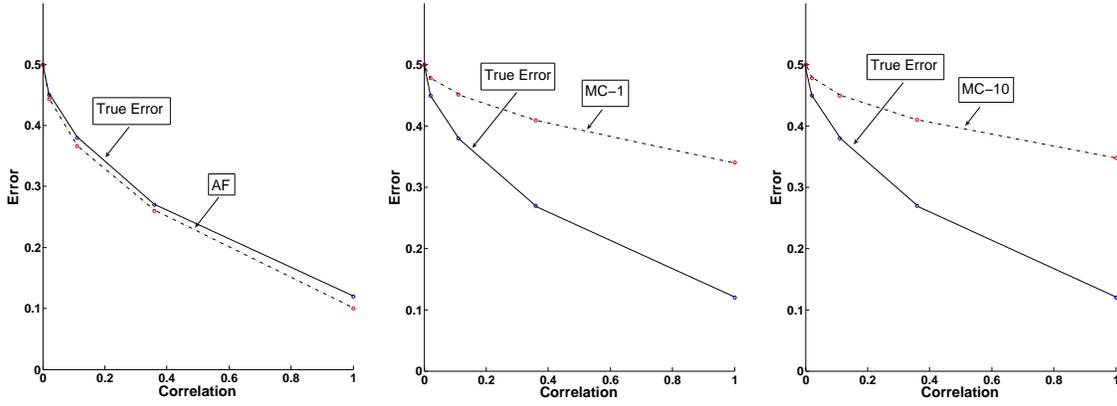


Figure 10: Scarcity based trees with  $d = 5$ ,  $pb = \frac{N}{10}$  and attributes with binary splits.

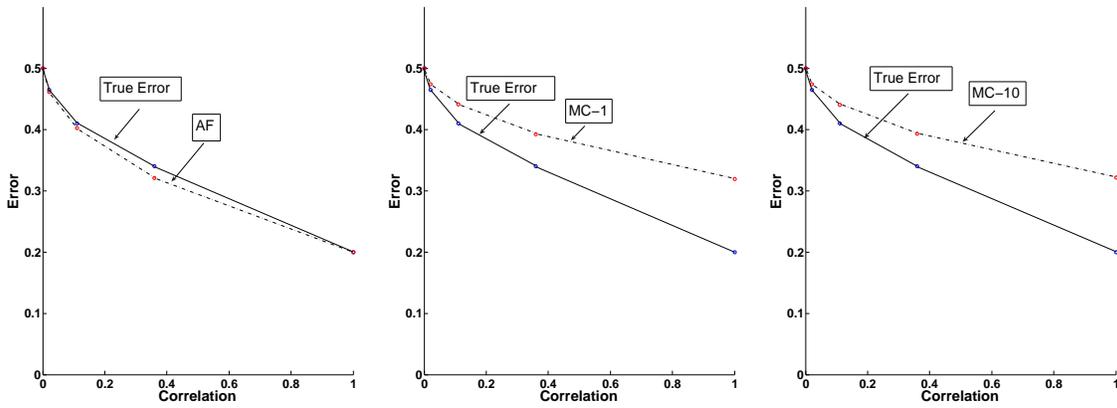


Figure 11: Scarcity based trees with  $d = 5$ ,  $pb = \frac{N}{10}$  and attributes with ternary splits.

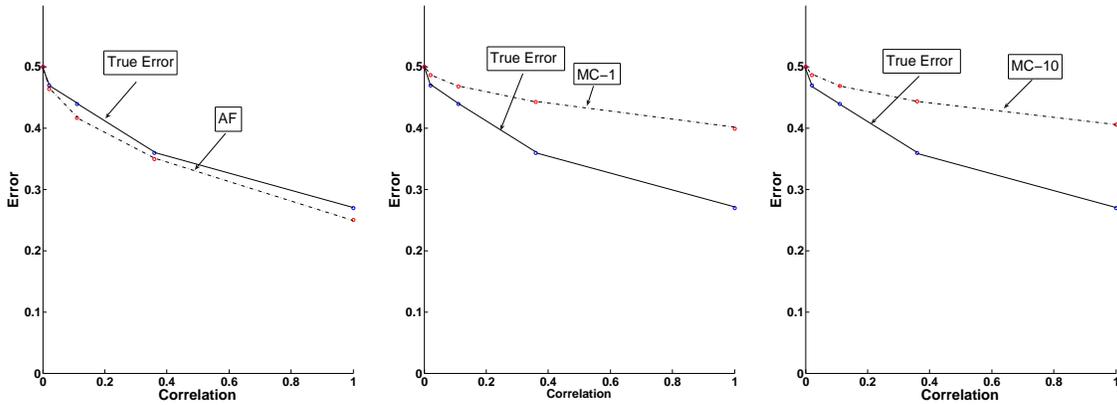


Figure 12: Scarcity based trees with  $d = 8$ ,  $pb = \frac{N}{10}$  and attributes with binary splits.

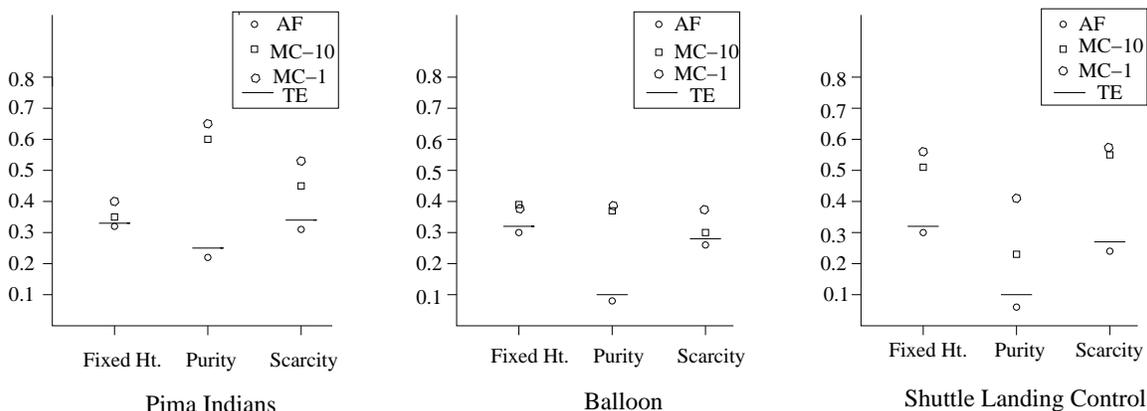


Figure 13: Comparison between AF and MC on three UCI datasets for trees pruned based on fixed height ( $h = 3$ ), purity and scarcity ( $pb = \frac{N}{10}$ ).

cell  $i$  to be  $p_i$  such that  $\sum_{i=1}^c p_i = 1$  and the sample size to  $N$  the distribution that perfectly models this scenario is a multinomial distribution with parameters  $N$  and the set  $\{p_1, p_2, \dots, p_c\}$ . In fact, irrespective of the value of  $d$  and the number of attribute values for each attribute the scenario can be modelled by a multinomial distribution. In the studies that follow the  $p_i$ 's are varied and the amount of dependence between the attributes and the class labels is computed for each set of  $p_i$ 's using the Chi-square test Connor-Linton (2003). More precisely, we sum over all  $i$  the squares of the difference of each  $p_i$  with the product of its corresponding marginals, with each squared difference being divided by this product i.e. correlation =  $\sum_i \frac{(p_i - p_{im})^2}{p_{im}}$ , where  $p_{im}$  is the product of the marginals for the  $i^{th}$  cell. The behavior of the error for trees with the three aforementioned stopping criteria is seen for different correlation values and for a class prior of 0.5.

In case of real data, we perform experiments on distributions built on three UCI datasets. We split the continuous attributes at the mean of the given data. We thus can form a contingency table representing each of the datasets. The counts in the individual cells divided by the dataset size provide us with empirical estimates for the individual cell probabilities ( $p_i$ 's). Thus, with the knowledge of  $N$  (dataset size) and the individual  $p_i$ 's we have a multinomial distribution. Using this distribution we observe the behavior of the error for the three kinds of trees with results being applicable to other datasets that are similar to the original.

**Observations:** Figures 4, 5 and 6 depict the error of fixed height trees with the number of attributes being 5 for the first two figures and 8 for the third figure. The number of attribute values increases from 2 to 3 in figures 4 and 5 respectively. We observe in these figures that AF is significantly more accurate than both MC-1 and MC-10. In fact the performance of the 3 estimators namely, AF, MC-1 and MC-10 remains more or less unaltered even with changes in the number of attributes and in the number of splits per attribute. A similar trend is seen for both purity based trees i.e. figures 7, 8 and 9 as well as scarcity based trees 10, 11 and 12. Though in the case of purity based trees the performance of both MC-1 and MC-10 is much superior as compared with their performance on the other two kinds of trees, especially at low correlations. The reason for this being that, at low correlations the probability in each cell of the multinomial is non-negligible and with  $N = 10000$  the event that every cell contains atleast a single datapoint is highly likely. Hence, the trees we obtain with high probability using the purity based stopping criteria are all ATT's. Since in an ATT all the leaves are identical irrespective of the ordering of the attributes in any path, the randomness in the classifiers produced, is only due to the randomness in the data generation process and not because of the random attribute selection method. Thus, the space of classifiers

over which the error is computed reduces and MC performs well even for a relatively fewer number of iterations. At higher correlations and for the other two kinds of trees the probability of smaller trees is reasonable and hence MC has to account for a larger space of classifiers induced by not only the randomness in the data but also by the randomness in the attribute selection method.

In case of real data too figure 13, the performance of the expressions is significantly superior as compared with MC-1 and MC-10. The performance of MC-1 and MC-10 for the purity based trees is not as impressive here since the dataset sizes are much smaller (in the tens or hundreds) compared to 10000 and hence the probability of having an empty cell are not particularly low. Moreover, the correlations are reasonably high (above 0.6).

**Reasons for superior performance of expressions:** With simple MC, trees have to be built while performing the experiments. Since, the expectations are over all possible classifiers i.e. over all possible datasets and all possible randomizations in the attribute selection phase, the exhaustive space over which direct MC has to run is huge. No tree has to be explicitly built when using the expressions. Moreover, the probabilities for each path can be computed parrallely. Another reason as to why calculating the moments using expressions works better is that the portion of the probabilities for each path that depend on the attribute selection method are computed *exactly* (i.e. with no error) by the given expressions and the inaccuracies in the estimates only occur due to the sample dependent portion in the probabilities.

## 4. Discussion

In the previous sections we derived the analytical expressions for the moments of the GE of decision trees and depicted interesting behavior of RDT's built under the 3 stopping criteria. It is clear that using the expressions we obtain highly accurate estimates of the moments of errors for situations of interest. In this section we discuss issues related to extension of the analysis to other attribute selection methods and issues related to computational complexity of algorithm.

### 4.1 Extension

The conditions presented for the 3 stopping criteria namely, fixed height, purity and scarcity are applicable irrespective of the attribute selection method. Commonly used deterministic attribute selection methods include those based on Information Gain (IG), Gini Gain (GG), Gain ratio (GR) etc. Given a sample the above metrics can be computed for each attribute. Hence, the above metrics can be implemented as corresponding functions of the sample. For e.g. in the case of IG we compute the loss in entropy ( $q \log q$  where the  $q$ 's are computed from the sample) by the addition of an attribute as we build the tree. We then compare the loss in entropy of all attributes not already chosen in the path and choose the attribute for which the loss in entropy is maximum. Following this procedure we build the path and hence the tree. To compute the probability of *path exists*, we add these sample dependent conditions in the corresponding probabilities. These conditions account for a particular set of attributes being chosen, in the 3 stopping criteria. In other words, these conditions quantify the conditions in the 3 stopping criteria that are attribute selection method dependent. Similar conditions can be derived for the other attribute selection methods (attribute with maximum gini gain for GG, attribute with maximum gain ratio for GR) from which the relevant probabilities and hence the moments can be computed. Thus, while computing the probabilities given in equations 3 and 4 the conditions for *path exists* for these attribute selection methods depend totally on the sample. This is unlike what we observed for the randomized attribute selection criterion where the conditions for *path exists* depending on this randomized criterion, were sample independent while the other conditions in purity and scarcity were sample dependent. Characterizing these probabilities enables us in computing the moments of GE for these other attribute selection methods.

In the analysis that we presented, we assumed that the split points for continuous attributes were determined apriori to tree construction. If the split point selection algorithm is dynamic i.e. the

split points are selected while building the tree, then in the *path exists* conditions of the 3 stopping criteria we would have to append an extra condition namely, the split occurs at "this" particular attribute value. In reality, the value of "this" is determined by the values that the samples attain for the specific attribute in the particular dataset, which is finite.<sup>1</sup> Hence, while analyzing we can choose a set of allowed values for "this" for each continuous attribute. Using these updated set of conditions for the 3 stopping criteria the moments of GE can be computed.

## 4.2 Scalability

The time complexity of implementing the analysis is proportional to the product of the size of the input/output space <sup>2</sup> and the number of paths that are possible in the tree while classifying a particular input. To this end, it should be noted that if a stopping criterion is not carefully chosen and applied, then the number of possible trees and hence the number of allowed paths can become exponential in the dimensionality. In such scenarios, studying small or at best medium size trees is feasible. For studying larger trees the practitioner should combine stopping criteria (e.g. pruning bound and fixed height or scarcity and fixed height) i.e. combine the conditions given for each individual stopping criteria or choose a stopping criterion that limits the number of paths (e.g. fixed height). Keeping these simple facts in mind and on appropriate usage, the expressions can assist in delving into the statistical behavior of the errors for decision tree classifiers.

## 5. Conclusion

In this paper we have developed a general characterization for computing the moments of the GE for decision trees. In particular we have specifically characterized RDT's for three stopping criteria namely, fixed height, purity and scarcity. Being able to compute moments of GE, allows us to compute the moments of the various validation measures and observe their relative behavior. Using the general characterization, characterizations for specific attribute selection measures (e.g. IG, GG etc.) other than randomized can be developed as described before. As a technical result, we have extended the theory in Dhurandhar and Dobra (2006) to be applicable to randomized classification algorithms; this is necessary if the theory is to be applied to random decisions trees as we did in this paper. The experiments reported in section 3 had two purposes: (a) portray the manner in which the expressions can be utilized as an exploratory tool to gain a better understanding of decision tree classifiers, and (b) show conclusively that the methodology in Dhurandhar and Dobra (2006) together with the developments in this paper provide a superior analysis tool when compared with simple Monte Carlo.

More work needs to be done to explore the possibilities and test the limits of the kind of analysis that we have performed. However, if learning algorithms are analyzed in the manner that we have shown, it would aid us in studying them more precisely, leading to better understanding and improved decision-making in the practice of model selection.

## 6. Appendix

The probability that two paths of lengths  $l_1$  and  $l_2$  ( $l_2 \geq l_1$ ) co-exist in a tree based on the randomized attribute selection method is given by,

$$P[l_1 \text{ and } l_2 \text{ length paths co-exist}] = \sum_{i=0}^v vPr_i(l_1 - i - 1)!(l_2 - i - 1)!(r - v)prob_i$$

---

1. Since dataset is finite.

2. In case of continuous attributes the size of the input/output space is the size after discretization

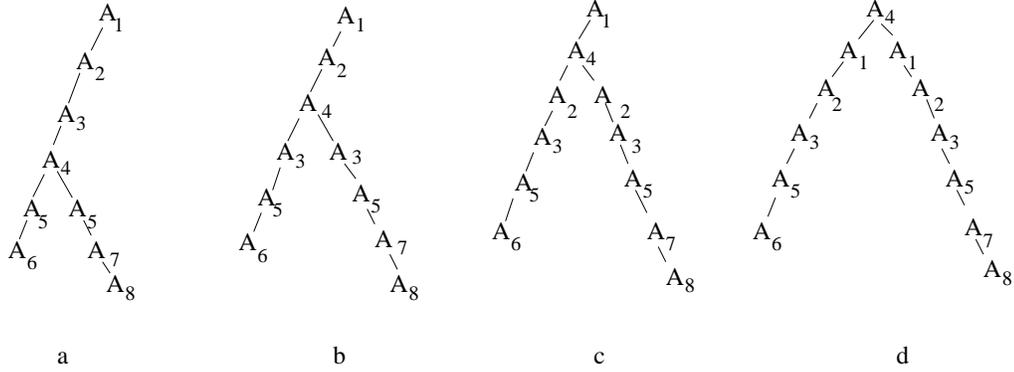


Figure 14: Instances of possible arrangements.

where  $r$  is the number of attributes common in the two paths,  $v$  is the number attributes with the same values in the two paths,  $vPr_i = \frac{v!}{(v-i)!}$  denotes permutation and

$$prob_i = \frac{1}{d(d-1)\dots(d-i)(d-i-1)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}.$$

We prove the above result with the help of an example. The derivation of the above result will become clearer through the following example. Consider the total number of attributes to be  $d$  as usual. Let  $A_1$ ,  $A_2$  and  $A_3$  be three attributes that are common to both paths and also having the same attribute values. Let  $A_4$  and  $A_5$  be common to both paths but have different attribute values for each of them. Let  $A_6$  belong to only the first path and  $A_7$ ,  $A_8$  to only the second path. Thus, in our example  $l_1 = 6$ ,  $l_2 = 7$ ,  $r = 5$  and  $v = 3$ . For the two paths to co-exist notice that atleast one of  $A_4$  or  $A_5$  has to be at a lower depth than the non-common attributes  $A_6$ ,  $A_7$ ,  $A_8$ . This has to be true since, if a non-common attribute say  $A_6$  is higher than  $A_4$  and  $A_5$  in a path of the tree then the other path cannot exist. Hence, in all the possible ways that the two paths can co-exist, one of the attributes  $A_4$  or  $A_5$  has to occur at a maximum depth of  $v + 1$  i.e. 4 in this example. Figure 14a depicts this case. In the successive tree structures i.e. Figure 14b, Figure 14c the common attribute with distinct attribute values ( $A_4$ ) rises higher up in the tree (to lower depths) until in Figure 14d it becomes the root. To find the probability that the two paths co-exist we sum up the probabilities of such arrangements/tree structures. The probability of the subtree shown in Figure 14a is  $\frac{1}{d(d-1)(d-2)(d-3)(d-4)^2(d-5)^2(d-6)}$  considering that we choose attributes w/o replacement for a particular path. Thus the probability of choosing the root is  $\frac{1}{d}$ , the next attribute is  $\frac{1}{d-1}$  and so on till the subtree splits into two paths at depth 5. After the split at depth 5 the probability of choosing the respective attributes for the two paths is  $\frac{1}{(d-4)^2}$ , since repetitions are allowed in two separate paths. Finally, the first path ends at depth 6 and only one attribute has to be chosen at depth 7 for the second path which is chosen with a probability of  $\frac{1}{d-6}$ . We now find the total number of subtrees with such an arrangement where the highest common attribute with different values is at depth of 4. We observe that  $A_1$ ,  $A_2$  and  $A_3$  can be permuted in whichever way w/o altering the tree structure. The total number of ways of doing this is  $3!$  i.e.  $3Pr_3$ . The attributes below  $A_4$  can also be permuted in  $2!3!$  w/o changing the tree structure. Moreover,  $A_4$  can be replaced by  $A_5$ . Thus, the total number of ways the two paths can co-exist with this arrangement is  $3Pr_3 2!3!2$ . The probability of the arrangement is hence given by,  $\frac{3Pr_3 2!3!2}{d(d-1)(d-2)(d-3)(d-4)^2(d-5)^2(d-6)}$ . Similarly, we find the probability of the arrangement in Figure 14b where the common attribute with different values is at depth 3 then at depth 2 and finally at the root. The probabilities for the successive arrangements are  $\frac{3Pr_3 3!4!2}{d(d-1)(d-2)(d-3)^2(d-4)^2(d-5)^2(d-6)}$ ,  $\frac{3Pr_1 4!5!2}{d(d-1)(d-2)^2(d-3)^2(d-4)^2(d-5)^2(d-6)}$  and  $\frac{3Pr_0 5!6!2}{d(d-1)^2(d-2)^2(d-3)^2(d-4)^2(d-5)^2(d-6)}$  respectively. The total probability for the paths to co-exist is given by the sum of the probabilities of these individual arrangements.

In the general case, where we have  $v$  attributes with the same values the number of arrangements possible is  $v+1$ . This is because the depth at which the two paths separate out lowers from  $v+1$  to 1. When the bifurcation occurs at  $v+1$  the total number of subtrees is  $vPr_v(l_1-v-1)!(l_2-v-1)!(r-v)$  with this arrangement.  $vPr_v$  is the permutations of the common attributes with same values.  $(l_1-v-1)!$  and  $(l_2-v-1)!$  are the total permutations of the attributes in path 1 and 2 respectively after the split.  $r-v$  is the number of choices for the split attribute. The probability of any one of the subtrees is  $\frac{1}{d(d-1)\dots(d-v)(d-v-1)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}$  since until a depth of  $v+1$  the two paths are the same and then from  $v+2$  the two paths separate out. The probability of the first arrangement is thus,  $\frac{vPr_v(l_1-v-1)!(l_2-v-1)!(r-v)}{d(d-1)\dots(d-v)(d-v-1)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}$ . For the second arrangement with the bifurcation occurring at a depth of  $v$ , the number of subtrees is  $vPr_{v-1}(l_1-v)!(l_2-v)!(r-v)$  and the probability of any one of them is  $\frac{1}{d(d-1)\dots(d-v+1)(d-v)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}$ . The probability of the arrangement is thus  $\frac{vPr_{v-1}(l_1-v)!(l_2-v)!(r-v)}{d(d-1)\dots(d-v+1)(d-v)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}$ . Similarly, the probabilities of the other arrangements can be derived. Hence the total probability for the two paths to co-exist which is the sum of the probabilities of the individual arrangements is given by,

$$P[l_1 \text{ and } l_2 \text{ length paths co-exist}] = \sum_{i=0}^v \frac{vPr_i(l_1-i-1)!(l_2-i-1)!(r-v)}{d(d-1)\dots(d-i)(d-i-1)^2\dots(d-l_1+1)^2(d-l_1)\dots(d-l_2+1)}$$

## References

- Avrim Blum, Adam Kalai, and John Langford. Beating the hold-out: Bounds for k-fold and progressive cross-validation. In *Computational Learning Theory*, 1999. URL [citeseer.ist.psu.edu/blum99beating.html](http://citeseer.ist.psu.edu/blum99beating.html).
- L. Breiman, J. Friedman, R. Olshen, and C. Stone. *Classification and Regression Trees*. Wadsworth and Brooks, 1984.
- Jeff Connor-Linton. Chi square tutorial. [http://www.georgetown.edu/faculty/ballc/webtools/web\\_chi\\_tut.html](http://www.georgetown.edu/faculty/ballc/webtools/web_chi_tut.html), 2003.
- Amit Dhurandhar and Alin Dobra. Semi-analytical method for analyzing models and model selection measures based on moment analysis. [www.cise.ufl.edu/submit/ext\\_ops.php?op=list&type=report&by\\_tag=REP-2007-296&display\\_level=full](http://www.cise.ufl.edu/submit/ext_ops.php?op=list&type=report&by_tag=REP-2007-296&display_level=full), 2006.
- M. Hall. Correlation-based feature selection for machine learning, 1998. URL [citeseer.ist.psu.edu/hall99correlationbased.html](http://citeseer.ist.psu.edu/hall99correlationbased.html).
- Mark A. Hall and Georey Holmes. Benchmarking attribute selection techniques for discrete class data mining. *IEEE TRANSACTIONS ON KDE*, 2003.
- R. Kohavi. A study of cross-validation and bootstrap for accuracy estimation and model selection. In *In Proceedings of the Fourteenth IJCAI*, 1995. URL [overcite.lcs.mit.edu/kohavi95study.html](http://overcite.lcs.mit.edu/kohavi95study.html).
- Fei Tony Liu, Kai Ming Ting, and Wei Fan. Maximizing tree diversity by building complete-random decision trees. In *PAKDD*, pages 605–610, 2005.
- J.R. Quinlan. Induction of decision trees. *Machine Learning*, 1(1):81–106, 1986.
- Jun Shao. Linear model selection by cross validation. *JASA*, 88, 1993.
- Jun Shao. *Mathematical statistics*. Springer-Verlag, 2003.

Lindsay I. Smith. A tutorial on principal components analysis.  
[www.csnet.otago.ac.nz/cosc453/student\\_tutorials/principal\\_components.pdf](http://www.csnet.otago.ac.nz/cosc453/student_tutorials/principal_components.pdf), 2002.

Vapnik. *Statistical Learning Theory*. Wiley & Sons, 1998.