

darkred!85!blue!60!black



Felipe de Albuquerque Mello Pereira

Binary Splitting Criteria for Large Categorical Attributes in Decision Trees

Dissertação de Mestrado

Thesis presented to the Informática of the Departamento de
Informática do Centro Técnico Científico da PUC-Rio as partial
fulfillment of the requirements for the degree of Doutor

Advisor: Prof. Eduardo Sany Laber

Rio de Janeiro
February 2018

Felipe de Albuquerque Mello Pereira

**Binary Splitting Criteria for Large Categorical
Attributes in Decision Trees**

Thesis presented to the Postgraduate Program in Informática of
the Departamento de Informática, PUC-Rio as partial fulfillment
of the requirements for the degree of Mestre em Informática.

Prof. Eduardo Sany Laber

Advisor

Departamento de Informática — PUC-Rio

Prof. Hélio Côrtes Vieira Lopes

Departamento de Informática — PUC-Rio

Prof. Marco Serpa Molinaro

Departamento de Informática — PUC-Rio

Prof. Márcio da Silveira Carvalho

Vice Dean of Graduate Studies

Centro Técnico Científico da PUC-Rio

Rio de Janeiro, February 28th, 2018

All rights reserved.

Felipe de Albuquerque Mello Pereira

Bachelor's in Electrical Engineering and Pure Mathematics at the Pontifícia Universidade Católica do Rio de Janeiro (2010 and 2011). Masters' in Mathematics at the Pontifícia Universidade Católica do Rio de Janeiro (2013).

Bibliographic data

Pereira, Felipe de Albuquerque Mello

Binary Splitting Criteria for Large Categorical Attributes in Decision Trees / Felipe de Albuquerque Mello Pereira ; advisor: Eduardo Sany Laber. — 2018.

38 f. : il. ; 30 cm

Dissertação (Mestrado em Informática)-Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2018.

Inclui bibliografia

1. Informática – Teses. 2. Árvores de Decisão; Problema de Corte Máximo; Algoritmos Aproximativos. I. Laber, Eduardo Sany. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

Acknowledgments

TODO: acknowledgment.

Thanks to CNPq for the conceded scholarship during my Masters.

Abstract

Pereira, Felipe de Albuquerque Mello; Laber, Eduardo Sany (advisor).
Binary Splitting Criteria for Large Categorical Attributes in Decision Trees. Rio de Janeiro, 2018. 38p. Dissertação de Mestrado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

TODO: abstract.

Keywords

Decision Trees; Max-cut Problem; Approximated Algorithms

Resumo

Pereira, Felipe de Albuquerque Mello; Laber, Eduardo Sany. **Critérios de Splits Binários para Atributos Categóricos Grandes em Árvores de Decisão..** Rio de Janeiro, 2018. 38p. Dissertação de Mestrado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

TODO: resumobr.

Palavras-chave

Árvores de Decisão; Problema de Corte Máximo; Algoritmos Aproximativos

Contents

1	Introduction	10
1.1	Our Contribution	11
2	Related Work	13
2.1	Organization	13
3	Background	15
3.1	Notation	15
3.2	Impurity Measures	15
3.3	Splitting Criteria	16
3.4	Heuristics for Splitting Decision Tree Nodes	19
4	Framework for Generating Splitting Criteria for Multi-valued Attributes	22
4.1	The Maximum Weighted Cut Problem	22
4.2	A Framework for Generating Splitting Criteria	23
5	Experiments on Splits	27
6	Experiments on Real Datasets	31
6.1	Datasets	31
6.2	Computing the Maximum Cut	33
6.3	Experimental Results	33
7	Conclusions	37

List of Figures

List of Tables

1

Introduction

Decision Trees and Random Forests are among the most popular methods for classification tasks. Decision Trees, specially small ones, are easy to interpret, while Random Forests usually yield more accurate classifications. One of the key issues in these methods is how to select an attribute to associate with a node of the tree/forest. An important related issue is how to split the samples once the attribute is selected.

There is a number of papers discussing aspects related with attribute selection, such as: how to design criteria to evaluate the quality of different types of attributes; whether binary or multi-way splits shall be used and how to remove bias from splitting criteria. For recent surveys on this topic we refer to ???.

Many criteria, with different properties, have been proposed to evaluate the quality of different types of attributes, including continuous and categorical ones. Among the most popular criteria, we have the Gini Gain and the Information Gain .

Despite the large body of work we believe there are still questions to be answered. One of them is to how to properly handle nominal attributes that may assume a large number of values. Before explaining the reason behind our statement we would like to remark that this kind of attribute appears naturally in some applications (e.g.: states of a country or letters from some alphabet). In addition, they may arise as the result of aggregating attributes that have few distinct values with the goal of capturing possible correlation between them, as pointed out by ?. As an example, consider 5 binary attributes (e.g. medical tests) and a target binary variable that has large probability of being positive if at least 3 out of the 5 binary tests are positives. By aggregating the 5 binary variables we obtain a new attribute with $2^5 = 32$ values that captures this relation. If we used the 5 attributes separately we would need 5 levels in the tree to be able to capture the relation between them and the target class, thus incurring a large fragmentation of the set of samples.

To properly face multi-valued nominal attributes we have to deal with the computational time required to compute good splits. Our contribution, explained in the next section, is related with this issue.

A brute force search to compute the best binary split requires $\Omega(2^n)$ time, where n is the number of distinct values the attribute may assume. The computational efficiency can be improved if a n -ary split is used rather than a binary one. However, this may lead to a severe fragmentation of the sample space, which is not desirable: the number of samples available for each of the children of the split node may be small and, as a consequence, the underlying classification tasks may become significantly more difficult. When the target variable is binary, the Gini Gain, proposed in the influential monography by Breiman et al. [1], can be computed efficiently. However, when the number of classes k is larger than 2, most, if not all, of the available exact solutions take exponential time in (n, k) . The Twoing method [2], which is equivalent to Gini Gain when $k = 2$, is an interesting case since its running time is $O(2^{\min\{n, k\}})$ rather than $O(2^n)$.

When both n and k are large, in the sense that an exhaustive search does not run in a reasonable time, one can rely on heuristics to compute the best binary split. As an example, the GUIDE algorithm [3], the last of a series of algorithms/developments designed by Loh and its contributors, deals with a nominal variable X as follows: if $k = 2$ or $n \leq 11$ the Gini Index is computed; if $k \leq 11$ and $n > 20$ a new variable X' with at most k distinct values is created according to a certain rule and an exhaustive search is performed over it; finally, if $k > 11$ or $n \leq 20$, X is binarized and a Linear Discriminant Analysis (LDA) is employed. These rules reflect the difficulty in dealing with multi-valued nominal attributes. In general, the main drawback of using heuristics is the lack of a theoretical guarantee about their behavior.

1.1

Our Contribution

Given this scenario, in chapter 4 we propose a framework for designing criteria, with nice theoretical properties, for evaluating the quality of multi-valued nominal attributes. Criteria generated according to this framework run in polynomial time in both the number of values and classes and have a theoretical guarantee that they are close to optimal. The key idea consists of formulating the problem of finding the best binary partition for a given attribute A as the problem of finding a cut with maximum weight in a complete graph whose nodes are associated with the values that A may assume and the edges' weights capture the benefit of putting values in different partitions. The motivation behind the use of the max-cut problem is the existence of efficient algorithms with approximation guarantee, in particular the one proposed by [4], with 0.878 approximation, and local search algorithms with 0.5-approximation

?

We discuss two criteria that are derived from this framework: the first one can be seen as a natural variation of the Gini Gain, while the second criterion uses the χ^2 -test to set the edges' weights. For that, each edge e_{ij} , between nodes v_i and v_j , is thought as a binary attribute $A(i, j)$, with values v_i and v_j . After discussing these criteria, we show how to extend them to handle numeric attributes.

We also present a number of experiments that suggest that one of our criteria is competitive with the Twoing method, which is – as far as we know – the only well-established criterion with binary splits that can be optimally computed for large n when $k > 2$. However, in contrast with our methods, Twoing cannot handle datasets that also have a large number of classes. In addition, the experiments also provide evidence of the potential of aggregating attributes for improving the accuracy of decision trees.

2

Related Work

Many splitting criteria have been proposed to deal with continuous and nominal attributes. Arguably, the Gini gain, used by CART, and entropy based measures, such as the Information Gain, adopted by C4.5, are among the most popular ???.

There has been some investigation on methods to compute the best split efficiently ?????. For the 2-class problem, ? proved a theorem which states that an optimal binary partition, for a certain class of splitting criteria, can be determined in linear time on n , the number of distinct values of the attributes, after ordering. The Gini Gain belongs to this class. The other three papers generalize this theorem in different directions and show necessary conditions that are satisfied by optimal partitions for a certain class of splitting criteria. These conditions, though useful to restrict the set of partitions that need to be considered, do not yield a method that is efficient (polynomial time) for large values of n and k . These papers also present heuristics, without approximation guarantee, to obtain good splits.

Other proposals to speed up the attribute selection phase include ??. The first presents a simple heuristic to reduce the number of binary splits considered to choose the best nominal variable among the m available ones. The second extends the method for another class of impurity measures.

In order to properly handle nominal attributes with a large number of values, apart from efficiently computing good splits, it is important to prevent bias in the attribute selection. Indeed, it is widely known that many splitting criteria have bias toward attributes with a large number of values. There are some proposals available to cope with this issue ????. This topic, though relevant, is not the focus of our paper.

2.1

Organization

In chapter 3 we explain how decision trees are used for classification problems and how they are constructed. We also present the main impurity measures and splitting criteria used in the literature, together with their execution-time complexity.

Chapter 4 contains the framework for generating splitting criteria that run in polynomial time. Its relation with the Max-Cut problem and its approximation algorithms are explained and some criteria obtained from this framework are presented.

Later, in chapters ?? and ??, we compare all the existing criteria and see how they perform in practice. In chapter ?? we explore how the many heuristics used to find splits with optimal impurity perform compare with the exact methods. This suggests a couple of criteria that perform better and can be used when the number of values and classes are high. In chapter ?? we analyse these heuristics together with exact methods on real datasets that contain attributes with large number of values and classes. Lastly, in chapter ?? we present our study conclusions.

3 Background

3.1 Notation

We adopt the following notation throughout the dissertation. Let S be a set of N samples and $C = \{c_1, \dots, c_k\}$ be the domain of the class label. In addition, for an attribute A , we use $A(s)$ to denote the value taken by attribute A on sample s ; we use $V = \{v_1, \dots, v_n\}$ to denote the set of values taken by A ; A_{ij} to refer to the number of samples from class c_j for which A takes value v_i ; N_i for the number of samples with value v_i for attribute A and S_j for the number of samples from class c_j . Furthermore, we let $p_j = S_j/N$ and $p_{ij} = \Pr[C = c_j | A = v_i]$. We observe that the estimator of maximum likelihood for p_{ij} is A_{ij}/N_i .

3.2 Impurity Measures

Many of the splitting criteria follow the same algorithm:
end for split S into (L, R) using the attribute A whose best split has the smallest impurity $\text{CreateTree}(L) \text{ CreateTree}(R)$

It is important to note that most impurity calculations on a categorical attribute are done based on what's called a contingency table. It consists of an n matrix where the ij -th entry corresponds to the number of samples that have value i and belong to class k . We'll assume that every node has the contingency table pre-calculated for all nominal attributes. This takes $O(N)$ time for each attribute and cannot be avoided¹. Therefore, whenever we mention the time complexity for a decision-tree constructing algorithm/criterion, this cost will not be mentioned since it does not change between them and don't affect the comparison.

A good place to start is by presenting the two most common impurity measures found in the literature.

¹The only exception is the attribute used to split the parent node, which can be calculated in $O(n \times k)$.

3.2.1

Gini

The Gini Index for a set of samples S is given by

$$Gini(S) = 1 - \sum_{i=1}^k (p_i)^2. \quad (3-1)$$

The Gini Index can be used to generate binary splits and, as a consequence, binary decision trees.

The Gini Gain, Δ_G , induced by a binary partition (L, R) of the set of values V is given by

$$\Delta_G(L, R) = Gini(S) - p_L Gini(S_L) - p_R Gini(S_R), \quad (3-2)$$

where $S_L = \{s \in S | A(s) \in L\}$, $S_R = \{s \in S | A(s) \in R\}$, $p_L = |S_L|/N$ and $p_R = |S_R|/N$. Therefore, the largest the Gini Gain is, the better the partition.

3.2.2

Entropy

The Entropy for a set of samples S is given by

$$Entropy(S) = - \sum_{i=1}^k p_i \log p_i \quad (3-3)$$

The Entropy can be used to generate binary splits and, as a consequence, binary decision trees.

The Information Gain IG induced by a binary partition (L, R) of the set of values V is given by

$$IG(L, R) = Entropy(S) - p_L Entropy(S_L) - p_R Entropy(S_R), \quad (3-4)$$

where $S_L = \{s \in S | A(s) \in L\}$, $S_R = \{s \in S | A(s) \in R\}$, $p_L = |S_L|/N$ and $p_R = |S_R|/N$. Therefore, the largest the Information Gain is, the better the partition.

3.3

Splitting Criteria

In this section we recall some well-known splitting criteria.

It is important to note that, for numeric attributes, all criteria follow the same algorithm to choose the best split: the values are ordered and the valid splits have the form

$$L = \{v_i : v_i \leq v\}$$

$$R = \{v_i : v_i > v\}$$

for some chosen value v . The criteria evaluate the impurity of these splits and chooses the one with the smallest impurity to split. Since we only have to evaluate at a single value v between each pair of consecutive values v_i, v_{i+1} , it takes $O(n \log n)$ time. Since this is polynomial and quite fast, it is largely ignored throughout this dissertation.

3.3.1

Gini Gain

This criterion generates all 2^n binary values split and the partition with maximum Gini Gain shall be selected. As shown in ?, for the 2-class problem this optimal partition can be computed in $O(n \log n)$ time. First fix one of the two classes and calculate the frequency of samples of each value on this class. If we denote by v_1, \dots, v_n the values in this order, then the best binary split of S will have the form

$$L = \{v_i : v_i \leq v\}$$

$$R = \{v_i : v_i > v\}$$

for some chosen value v . Since we only need to test a single value of v that is between each pair of consecutive values v_i, v_{i+1} , the time complexity follows.

For problems with more than 2 classes, however, there is no efficient procedure with theoretical approximation guarantee to compute the Gini Gain in subexponential time in n .

3.3.2

Twoing

The Twoing criterion for a binary partition (L, R) of the set of values V is given by

$$0.25 \cdot p_L \cdot p_R \cdot \left(\sum_{i=1}^k |p_L^i - p_R^i| \right)^2,$$

where

$$p_L^i = \frac{|\{s \in S_L : s \text{ belongs to class } c_i\}|}{|S_L|} \text{ and } p_R^i = \frac{|\{s \in S_R : s \text{ belongs to class } c_i\}|}{|S_R|}.$$

When the Twoing criterion is used to generate binary decision trees, the binary partition with maximum twoing shall be selected at each node.

As shown in ?, such partition can be calculated in $O(\min\{n(k + \log n)2^k, 2^n\})$ time by considering all possibilities of partitioning the classes into two superclasses and applying the Gini Gain criterion on each of them. We shall remark that, for the 2-class problem, the Twoing criterion and the Gini Gain compute the same binary partitions.

3.3.3

Information Gain

This criterion works exactly the same as the Gini Gain, but replacing the Gini impurity by the Entropy. First it generates all 2^n binary values split and the partition with maximum Information Gain shall be selected. For the 2-class problem, the same result valid for the Gini Gain works here, and the optimal partition can be computed in $O(n \log n)$ time. Once again, when the number of classes is larger than 2 there is no efficient procedure with theoretical approximation guarantee to compute the Information Gain in subexponential time in n .

A related criterion is the Gain Ratio, where the Information Gain of an attribute is normalized by the potential information of that attribute. This is used as a way of decreasing the bias of the k-ary Information Gain criterion towards attributes with larger number of values. Since we are only interested in binary splits in this dissertation, we will not go into its details.

3.3.4

χ^2 criterion

The χ^2 is a popular criterion that was used in ?. It is also the first one shown here not based on impurity measures, and it only generates k-ary (instead of binary) splits. It is mentioned here because of its relation to the framework presented in chapter 4.

For an attribute A the χ^2 criterion is given by

$$\sum_{i=1}^n \sum_{j=1}^k \frac{(A_{ij} - E[A_{ij}])^2}{E[A_{ij}]}, \quad (3-5)$$

where $E[A_{ij}] = N_i p_j$.

3.3.5

Conditional Inference Trees

Conditional Inference Trees are actually a framework of creating criteria that are bias-free when it comes to the number of values in an attribute. It was published by ? and still is the only known method of obtaining criteria that do not have any bias towards attributes with larger number of values.

It first chooses the best attribute to split at the current node and then evaluates all possible binary splits using any given impurity measure, finally using the best one found.

To choose the attribute in which to split, first one has to calculate the conditional expectation $\mu \in \mathbb{R}^{nk}$ and covariance $\Sigma \in \mathbb{R}^{nk \times nk}$ of the permutation

test of every attribute². Then, in order to compare the attributes, we need to calculate the p-value of a univariate test statistics c_{quad} calculated on μ and Σ of every attribute. The only exact form of doing this is by using the quadratic form 3-6, which it follows an asymptotic χ^2 distribution with degrees of freedom given by the rank of Σ . Since this involves the calculation of a pseudoinverse, which is cubic in the dimension of Σ , this criterion is very time consuming.

$$\sum_{i=1}^n \sum_{j=1}^k \frac{(A_{ij} - E[A_{ij}])^2}{E[A_{ij}]}, \quad (3-6)$$

This method, although very complicated and quite slow, is employed by the community when accuracy counts for more than time spent training³. Thus this criterion will be used in our experiments in chapter ?? to compare the different accuracies obtained when changing the splitting criterion used to choose the attribute's values split.

3.4

Heuristics for Splitting Decision Tree Nodes

As seen in the previous section, calculating the optimal split takes exponential time in the number of values or classes. Therefore many heuristics were created to create decision trees in this situation. The most used ones are listed below. All of them work with any impurity measure (e.g.: Gini or Entropy), but some of them work best with one of them. When this is the case, it will be mentioned.

3.4.1

SLIQ and SLIQ-ext

SLIQ was presented in ? and it's a very simple greedy heuristic. Given an attribute, one starts with all the values going to the left split, and none on the right split. We then choose a value to go from the left to the right split. This value is the one that, when changing from the left to the right sides, decreases the impurity (increases the impurity gain) the most. This is repeated until there is no way of moving a value from the left to the right and decreasing the impurity.

SLIQ-ext is a simple extension, where we keep changing values from the left to the right until the left side is empty (that is, we move from the left to the right even if that increases the impurity). Once again the value to move

²Since the formulae are very complicated and won't be used throughout this dissertation, they are going to be omitted. The interested reader can obtain them in the section 3 of ?.

³Since this method does not have any bias when choosing which attribute to split, the accuracy of these trees tend to be higher than the trees obtained by biased criteria.

is chosen in a greedy fashion. SLIQ-ext returns the values' split seen that had the lowest impurity.

3.4.2

PC and PC-ext

These heuristics are based on the Principal Component of the contingency table and were presented in ?. One first calculates the class probability distribution of every value, which is done by normalizing the contingency table rows to have 1 in the sum norm. We then group values into “supervalues” where each value in the same supervalue has the same class probability distribution. Now we get the contingency table of these supervalues and calculate the first principal component of this matrix. One then calculates the inner product of each class probability vector of the supervalues with the principal component and sort the supervalues by it. Denote by v_1, \dots, v_{n^*} the n^* supervalues sorted in this manner and denote by pc the first principal component. We then calculate the impurity gain of the supervalues' splits of the form

$$L = \{v_i : v_i \cdot pc \leq t\}$$

$$R = \{v_i : v_i \cdot pc > t\}$$

where t is a chosen threshold. Once we find the supervalues split with the largest impurity gain, we translate the supervalues into original values to obtain a valid partition.

PC-ext is a simple extension of this algorithm, where instead of only testing the supervalues splits given by (TODO: ex ref), we also test the split given by it and exchanging the last supervalue on the left with the first supervalue on the right (where first and last are given by the order after calculating the inner product).

3.4.3

Largest Alone

First one calculates the most frequent class and group the other classes in a single superclass. We then apply the Gini Gain criterion on this two-class problem. Since calculating the class frequencies can be done using the contingency table, this heuristic takes $O(N + nk + n \log n)$ time in total.

Another advantage of the Largest Alone heuristic is that it has an approximation guarantee of 2 (TODO: conferir) to the best gini impurity of the original problem. It is also proved that there is no other way of grouping classes into superclasses that has a smaller approximation guarantee for the

Gini impurity. It can also be used with the Information Gain, instead of Gini Gain, but its approximation guarantee increases to 3 (TODO: conferir). These bounds are proved in (TODO: ver como colocar referencia).

3.4.4

List Scheduling

First one calculates the frequency of every class. Then, we use a List Scheduling algorithm to group the classes into 2 superclasses as balanced as possible (in terms of number of samples). Lastly we apply the Information Gain criterion on this two-class problem. Again, since calculating the class frequencies can be done using the contingency table, this heuristic and the List Scheduling algorithm is linear in the number of classes, this heuristic also takes $O(N + nk + n \log n)$ time in total.

Similarly to the Largest Alone heuristic, the List Scheduling heuristic has an approximation guarantee of 2 (TODO: conferir) to the best entropy impurity of the original problem. This is proved in (TODO: ver como colocar referencia). It is also proved that, for the entropy impurity, the best form of grouping classes into superclasses is by balancing them the best way possible.

4

Framework for Generating Splitting Criteria for Multi-valued Attributes

First we recall some definitions and results for the Max-Cut problem. These definitions will be used in the following section, when we define our framework.

4.1

The Maximum Weighted Cut Problem

We recall some definitions from graph theory. A cut X in a weighted graph $G = (V, E)$ is a subset of vertexes of V . The weight of a cut X , denoted here by $w(X)$, is the sum of the weights of the edges that have one endpoint in X and the other one in $V - X$.

The problem of computing the cut X^* with maximum weight in a graph with non-negative weights is NP-Hard. However, there are good approximation algorithms available. A remarkable one is the randomized algorithm proposed in ? that relies on a formulation of the max-cut problem via semidefinite programming (SDP). This algorithm, denoted here by GW, returns a cut X that satisfies $E[w(X)] \geq 0.878w(X^*)$. It involves solving an SDP on the graph weights matrix, calculating the Cholesky decomposition of it and then generating a random partition of the values based on the inner product of the decompositions column vectors with a randomly generated vector on the sphere of dimension n . As solving an SDP takes $O(n^4)$ arithmetic operations (see ?) and calculating the Cholesky decomposition takes $O(n^3)$ operations, the time complexity is high but polynomial.

Another possibility to solve the Max-Cut problem is by using the GreedyCut algorithm, presented in Algorithm ???. It obtains a cut X such that $w(X) \geq 0.5w(X^*)$?. The algorithm starts with two empty sets X and X' . Then, it scans the nodes and assigns each of them to the set that provides the maximum improvement on the weight of the current cut (ties are broken arbitrarily). Is it easy to see that the time complexity of this greedy algorithm is $O(n^2)$.

The solutions obtained by both GW and GreedyCut can be improved via a local search. In its simplest version, it moves a node from one group to

the other while some improvement on the cut weight is possible. Although this algorithm is not polynomial in the worst case, it has polynomial behavior in the smoothed analysis framework ?. In addition, it is always possible to set a limit on the number of moves. A more refined version allows exchanging a pair of nodes as long as the weight of the cut is improved. In our experiments we use the version presented at Algorithm ??.

4.2

A Framework for Generating Splitting Criteria

In this section we explain our approach to building binary splitting criteria for multi-valued nominal attributes.

Let A be a nominal attribute that takes values in the domain $V = \{v_1, \dots, v_n\}$. Our framework to produce a splitting criterion I consists of three steps:

1. Create a complete graph $G = (V, E)$ with n vertexes.
2. Assign a non-negative weight w_{ij} to the edge that connects v_i to v_j . This value shall reflect the benefit of putting v_i and v_j in different partitions. Different definitions of w_{ij} yield to different criteria, as we explain further.
3. Ideally, the value of the criterion I for attribute A is the weight of the cut with maximum weight in G . However, this is not a reasonable possibility for large n since the problem of computing the cut X^* with maximum weight in a graph with non-negative weights is NP-Hard. Thus, the value of criterion I is given by the weight of the cut obtained by some algorithm, with approximation guarantee, for the maximum cut problem in G .

What distinguishes the criteria generated by our framework is how the weights of the edges are set and the method employed to compute the cut on graph G . Here, we discuss two ways to set the weights: the first one yields to criteria that are related with the Gini Gain, while the second is built upon some given splitting criterion that works well for binary attributes.

4.2.1

The Squared Gini Criterion

Here, we discuss how to set the weights so that we obtain a criterion that can be seen as a variation of the Gini Gain discussed in Section ??.

In fact, Lemma 1 below show that it is possible to define the weights of the edges so that

$$w(S_L) = Gini(S) - p_L^2 \cdot Gini(S_L) - p_R^2 \cdot Gini(S_R) \quad (4-1)$$

for every partition (L, R) of V .

Note that the weight of the cut S_L in the above identity is similar to the expression for the Gini Gain given by equation (3-2). The difference is that p_L and p_R are replaced with p_L^2 and p_R^2 , respectively. Because of the squares, this new criterion tends to favor more balanced partitions.

For that the proof of Lemma 1, recall that A_{ix} is the number of samples of class x that have value v_i , and C is used to denote the set of classes.

Lemma 1 *For every i, j , with $i \neq j$ and $i, j \in \{1, \dots, n\}$, let*

$$w_{ij} = \frac{2 \sum_{\substack{x, y \in C \\ x \neq y}} A_{ix} A_{jy}}{N^2}.$$

Then, for every partition (L, R) of V we have $w(S_L) = Gini(S) - p_L^2 \cdot Gini(S_L) - p_R^2 \cdot Gini(S_R)$.

Proof: Let $S_{x,L}$ and $S_{y,R}$ be the number of samples of classes x and y in groups L and R , respectively. Moreover, let N_L and N_R be the number of samples in L and R , respectively. It follows from equation (3-1) that

$$N^2 Gini(S) = N^2 - \sum_{x=1}^k (S_{x,L} + S_{x,R})^2$$

$$N_L^2 Gini(S_L) = N_L^2 - \sum_{x=1}^k S_{x,L}^2$$

and

$$N_R^2 Gini(S_R) = N_R^2 - \sum_{x=1}^k S_{x,R}^2.$$

Since $N = (N_L + N_R)$ it follows that

$$N^2 Gini(S) - N_L^2 Gini(S_L) - N_R^2 Gini(S_R) =$$

$$2N_L N_R - 2 \sum_{x \in C} S_{x,L} S_{x,R} =$$

$$2 \sum_{x \in C} S_{x,L} \sum_{x \in C} S_{x,R} - 2 \sum_{x \in C} S_{x,L} S_{x,R} =$$

$$2 \sum_{\substack{x \neq y \\ x, y \in C}} S_{x,L} S_{y,R} = 2 \sum_{\substack{x \neq y \\ x, y \in C}} \left(\sum_{i \in L} \sum_{j \in R} A_{ix} A_{jy} \right) =$$

$$N^2 \sum_{i \in L} \sum_{j \in R} w_{ij} = N^2 w(S_L)$$

Dividing the first term and the last term by N^2 in the above expression and, using $N_L = p_L \cdot N$ and $N_R = p_R \cdot N$, we establish the lemma. ■

It is worth mentioning that symmetric misclassification costs can be easily introduced in this case. In fact, let $mix(x, y)$ be the cost of mixing samples from classes x and y . We can define

$$w_{ij} = \sum_{\substack{x, y \in C \\ x \neq y}} mix(x, y) p_{ix} p_{jy}.$$

This measure favors the separation of the classes that incur a large cost in the case they are mixed.

4.2.2

Setting weights according to other splitting criteria

Our second way of defining the weights makes use of some given splitting criterion for binary nominal attributes. Such criterion is used to measure the quality of separating samples with value v_i from those with value v_j , for each i and j , and, thus, defining the edges' weights. Here, we investigate the criterion obtained by defining w_{ij} as the value of the χ^2 -test for the attribute A when evaluated over the restricted dataset that contains only the samples of S with values v_i and v_j . In formulae,

$$w_{ij} = \sum_{\ell=1}^k \frac{(A_{i\ell} - E[A_{i\ell}])^2}{E[A_{i\ell}]} + \sum_{\ell=1}^k \frac{(A_{j\ell} - E[A_{j\ell}])^2}{E[A_{j\ell}]}$$

where $E[A_{i\ell}] = N_i p_\ell$ and $E[A_{j\ell}] = N_j p_\ell$.

In order to reduce the bias towards attributes with many values, we divide w_{ij} by $n - 1$, for every pair i, j . We make this adjustment because each value contributes to the weight of $n - 1$ edges.

We shall remark that, although not explored in this work, other criteria such as Information Gain or Gini Index could be used, instead of χ^2 test, to set the weights.

4.2.3

Handling Numeric Attributes

We observe that criteria from our framework can handle a numeric attribute A with t distinct values v_1, \dots, v_t by considering it as collection of $t - 1$ binary attributes, where the j -th attribute, A^j , splits the samples into

the groups $\{s|A(s) \leq v_j\}$ and $\{s|A(s) > v_j\}$. The split obtained by criterion I on a numeric attribute A matches the split of the best attribute among A^1, \dots, A^{t-1} , according to I .

5 Experiments on Splits

In this chapter we compare the ability of different heuristics in finding the values' split with lowest impurity.

We report a number of experiments with the heuristics proposed/analyzed in the previous sections. Our experiments are very similar to those proposed in ? except for a few details. All experiments are Monte Carlo simulations with 10,000 runs, each using a randomly-generated contingency table for the given number of values n and classes k . By a contingency table we mean a matrix where each row corresponds to a distinct vector of the input V . Each table was created by uniformly picking a number in $\{0, \dots, 7\}$ for each entry. This guarantees a substantial probability of a row/column having some zero frequencies, which is common in practice. Differing from ?, if all the entries corresponding to a value or a class are zero, we re-generate the contingency table, otherwise the number of actual values and classes would not match n and k .

We evaluated Twoing, and PCext. The latter is a method proposed in ? that defines the partition of V by using a hyperplane in k whose normal direction is the principal direction of the contingency table associated with the instance. According to the experiments reported in ? PCext consistently outperformed SliqExt and the Flip Flop method in terms of speed and the impurity of the partitions found.

Table ?? and ?? show, for different values of n and k , the percentage of times that Twoing outperformed/ was outperformed by PCext for Gini and Entropy, respectively. Note that the percentages do not necessarily sum up 100% because there were ties. We only show results for $k \leq 9$ because for larger values of k Twoing becomes non practical due to its running time. In addition, we do not present results for small values of n because in this case the optimal partition can be found reasonably quick by testing all possible partitions so that there is no motivation for heuristics.

In general, we observe an advantage of Twoing for both criteria, being more clear for Entropy impurity. The maximum relative excess between the impurities of the partitions found by Twoing and PCext was 0.9% for Gini and 1.4% for Entropy. On the other hand, the maximum relative excess between

the impurities of the partitions found by PCext and Twoing was 3.48% for Gini and 6.23% for Entropy. In terms of speed, as expected, Twoing was faster up to $k = 7$ and then PCext becomes faster.

We do not report the results for because it was not competitive with the other two heuristics. However, due to its running time it might be used when both n and k are very large and speed is an issue. When $n = 200$ and $k = 100$, using one core, is almost 50 times faster than PCext, with the latter using 8 cores. In addition, could be used together with PCext, incurring a negligible overhead, to guarantee that the ratio between the impurity of the partition found and the optimal one is bounded.

Taking into account these experiments, those reported in ? and the theoretical properties of the available algorithms, Table ?? suggests guidelines on how to solve the problem of finding the binary partition of minimum impurity in practical situation. Of course small, medium and large depend on the available hardware and the time one is up to wait to train/test classification models.

All experiments in this section are Monte Carlo simulations with 10,000 runs, each using a randomly-generated contingency table for the given number of values n and classes k . The contingency tables were created by uniformly picking a number in $\{0, \dots, 7\}$ for each entry, as done in Copersmith et al (TODO: ref). This guarantees a substantial probability of a row/column having some zero frequencies, which is common in practice. If all the entries corresponding to a value or a class are zero, we re-generate the contingency table, otherwise the number of actual values and classes would not match n and k .

We ran the experiments for two different impurity measures: the Gini Gain and the Information Gain. We measured the performance of the same criteria for both of them. All of the non-exact criteria studied are based on the superclass trick, using the theorem proved by Breiman et al (TODO: ref). To separate the classes into a pair of superclasses we use two heuristics: (where the class with the largest frequency is separated from the others) and ListScheduling (balancing the superclasses using a polynomial algorithm with a $4/3$ -approximation, as explained in TODO: ref). Since (ListScheduling) should be the best heuristic for the Gini (InformationGain) impurity, we compare against ListScheduling () as a baseline. Given a partition of classes into superclasses, we use the theorem (TODO:ref) to find the optimal partition of values in linear time after sorting.

We compare those methods with the exact and the Twoing criteria, which

run in exponential time in n and k , respectively. For all the non-exact criteria, we measured the impurity with respect to all the k classes—instead of the impurity with respect to the superclasses—since that’s what we are optimizing. We also studied a criterion that picks a random partition of values into two non-empty groups, but it was significantly worse than all the other methods, thus we chose not to report its results.

We are interested in choosing what criterion to use when the exact ones don’t run in reasonable time. In some experiments we choose values for n that are not very large because we want to be able to compare the other criteria with the exact ones. In others, we choose n to be very large, and don’t analyze the exact criteria because they don’t finish executing.

The first set of experiments uses $k = 3$ classes and $n = 6, 12$ values. Since the Twoing criterion behaves almost perfectly, both in terms of execution time and impurity found, we don’t compare it against the and ListScheduling criteria.

From tables ??, ??, ?? and ??, it is striking how well Twoing performs. In almost all the simulations the impurity found is exactly the same as the exact methods. This suggests that, in cases when the exact methods don’t run in reasonable time but Twoing does (that is, large n and small k), we lose almost nothing by using it.

The second set of experiments uses $k = 9$ classes and $n = 6, 12$ values. The results are shown in tables ??, ??, ?? and ?. The comparison between and ListScheduling is shown in tables ?? and ?.

Once again Twoing behaves remarkably well when it comes to getting the lowest impurity, but its execution time is already much larger than the criteria that use a single superclass split.

Comparing the heuristics between themselves we can see that, in terms of the Gini impurity, the heuristic is the best among them, while for the Entropy impurity, the best one is the ListScheduling method. This is consistent with the theory developed earlier in this paper.

The last set of experiments is focused on studying what happens when the number of values is very large ($n \geq 30$) and the number of classes increases.

First we study what happens when the number of classes is not very large ($k = 3, 9$). In this comparison we use Twoing as a baseline, since its behavior is so close to the exact methods in the other experiments and since it is still feasible to run it, while the exact methods aren’t. Note that, since Twoing evaluates all the possible superclass partitions, the impurity it finds is never worse than that of and ListSchedule. The results are shown in tables ??, ??, ??, ??, ?? and ?.

Note that the results for the and the ListScheduling methods are the same when $k = 3$ because they give the same superclass partitions. We can also notice that the number of times that Twoing finds a partition with impurity smaller than the other methods increases with k . This is explained by the exponential growth of the number of possible superclass partitions with k . What's surprising is that the method gives slightly better partitions than ListSchedule for the Information Gain impurity. Since both methods are within the approximation bounds set in this paper, this does not contradict our results.

In the next experiments, we use large number of values and classes. This makes both Twoing and the exact methods infeasible to be executed, so we only compare with ListSchedule. We use the minimum impurity between the two as a baseline.

Comparing the results in tables ??, ??, ??, ??, ?? and ?? we can see that, for both the Gini and Information Gain impurity measures, the is usually the best method. This is specially interesting for the Information Gain impurity, since the approximation guarantees for the method that uses a balanced superclass partition is better. Furthermore, the results for the criterion improve as the number of values and classes increase. This suggests that, in practice, one should prefer the method when n and k are very large, even though the approximation guarantee for the ListSchedule method is tighter.

The above experiments strongly suggest that, when we can't run the exact criteria, the best choice we can make is to choose the Twoing criterion, if its execution time is feasible. If that's also not possible, for the Gini impurity, the best choice is always to use the criterion. For the Information Gain, we should choose between ListSchedule, if n and k are not significantly large, or , for larger number of values and classes.

In all the simulations the criteria satisfied the approximation guarantees given in this paper. Since they all performed much better than their approximation bounds, the criterion with the best guarantee is not necessarily the best in terms of the expected impurity of the chosen partition.

6

Experiments on Real Datasets

In this chapter we describe our experimental study on real datasets. First, we describe the chosen datasets. Next, we discuss the max-cut algorithms employed and, then, we present our results.

All experiments described in the following sections were executed on a machine with the following settings: Intel(R) Core(TM) i7-4790 CPU @ 3.60GHz with 32 GB of RAM. The code was developed using Python 3.6.1 with the libraries numpy, scipy, scikit-learn and cvxpy. The project can be accessed in github.com/felipeamp/max_cut_paper. It includes the code, the datasets and the results of our experiments.

6.1

Datasets

We employed 11 datasets in total. Eight of them are from the UCI repository: Mushroom, KDD98, Adult, Nursery, Covertype, Cars, Contraceptive and Poker ?. Two others are available in Kaggle: San Francisco Crime and Shelter Animal Outcome ?Austin Animal Center (2016). The last dataset was created by translating texts from the Reuters database ? into phonemes, using the CMU pronouncing dictionary ?.

We chose these datasets because they have at least 1000 samples and they either contain multi-valued attributes or attributes that can be naturally aggregated to produce multi-valued attributes. From the KDD98 dataset we derived the datasets KDD98- k , for $k = 2$ and 9. These datasets contain only the positive samples (people that donate money) of KDD98 and the target attribute, Target_D, is split into k classes, where the i -th class correspond to the i -th quantile in terms of amount of money donated. For the Reuters Phonemes dataset, we extracted 10000 samples containing the 15 most common phonemes as class and try to predict when they are about to happen given the 3 preceding phonemes. This dataset is motivated by Spoken Language Recognition problems, where phonotactic models are used as an important part of the classification system ?. For the San Francisco Crime dataset, we give the month, day of the week, police department district and latitude/longitude and try to predict the crime category. Lastly, for the Shelter Animal Outcomes dataset,

we converted the age into a numeric field containing the number of days old and separated the breed into two categorical fields, repeating the breed in both in case there was only one originally. We also removed the AnimalID, Name and the DateTime. For this dataset we try to predict the outcome type and subtype (concatenated into a single categorical field). For both San Francisco Crime and Shelter Animal Outcomes datasets we created a version of them (**S.F. Crime-15** and **Shelter-15**), containing only 15 classes, instead of the 39 and 22 original ones, respectively. This was done by grouping the rarest classes into a single one.

We also created extended versions of some of the above datasets by adding nominal attributes obtained by aggregating some of the original ones, as we detail below. Our goals are examining the impact of multi-valued attributes in the classification performance and also understanding how the different splitting criteria handle them.

Table ?? illustrates this construction.

- **Nursery-Ext.** This dataset is obtained by adding three new attributes to dataset Nursery. The first attribute has 15 distinct values and it is constructed through the aggregation of 2 attributes from group EMPLOY, one with 5 values and the other with 3 values. The second attribute has 72 distinct values corresponding to the aggregation of attributes from the attributes in group STRUCT.FINAN. The third attribute, with 9 distinct values, is the combination of the attributes in group SOC.HEALTH.
- **Coverttype-Ext.** We combined 40 binary attributes related with the soil type into a new attribute with 40 distinct values. The same approach was employed to combine the 4 binary attributes related with the wilderness area into a new attribute with 4 distinct values. This is an interesting case because, apparently, the 40 (soil type) binary attributes as well as the 4 (wilderness area) binary attributes were derived from a binarization of two attributes, one with 40 distinct value and the other with 4 distinct values.
- **Cars-Ext.** To obtain this dataset, the 2 attributes related with the concept PRICE, `buying` and `maint`, were combined into an attribute with 16 distinct values. Moreover, the 3 attributes related with concept CONFORT were combined into an an attribute with 36 distinct values.
- **Contraceptive-Ext.** The 2 attributes related with the couple’s education were combined into an attribute with 16 distinct values. Moreover, the 3 attributes related with the couple’s occupations and standard of living were aggregated into a new attribute with 32 distinct values.

Samples with missing values were removed from the datasets. Table ?? provides some statistics.

6.2

Computing the Maximum Cut

The GW algorithm requires the solution of a semidefinite program (SDP), which may be computationally expensive despite its polynomial time worst case behavior. As an example, for an attribute with 100 distinct values, the solution of the corresponding SDP takes in average about 2 second in our machine. On the one hand, this is a tiny amount of time compared with that required to perform an exhaustive search on the 2^{100} possible binary partitions. On the other hand, faster alternatives are desirable, even at the cost of losing part of the theoretical approximation guarantee.

To avoid solving a SDP, we also evaluated a procedure that first executes the GreedyCut algorithm presented in Section 4.1 and then runs a local search as described in the Algorithm ?? in the same section. The use of this approach combined with the two ways of setting the edges' weights lead to Greedy LocalSearch SquaredGini (GLSG) and Greedy LocalSearch χ^2 ($GL\chi^2$) criteria, respectively. For attributes with 100 distinct values this approach is 60-70 times faster than the one based on the GW algorithm.

In fact, experiments similar to those described in the next section show that GL approach consistently obtains better results than the GW algorithm in terms of both speed and accuracy. This advantage is likely related with the fact that the weights of the cuts computed by the GL approach are, in general, larger than those obtained by the GW algorithm. This pattern was observed in a set of experiments with the max-cut instances induced by the attributes of the datasets used in the experiments section. Due to these results, we proceed our investigation using only the GL approach.

It is worth mentioning that the greedy algorithm that precedes the local search guarantees that we still have a theoretical guarantee of being at least 0.5 of the optimal solution. The local search may be not

In addition, it was recently proved that a local search runs in polynomial time TODO: melhorar

6.3

Experimental Results

We performed a number of experiments to evaluate how the proposed methods behave with real datasets. All experiments consist of building decision trees with a predefined maximum depth. In addition, to prevent the selection

of non-informative nominal attributes, we used a χ^2 -test for each attribute at every node of the tree: if the χ^2 -test on the contingency table of attribute A has p -value larger than 10% at a node ν , then A is not used in ν . Furthermore, attributes with less than 15 samples associated with its second most frequent value are also not considered for splitting. This helps avoid data overfitting.

Table ??.(a) presents the results of an experiment to compare the accuracy of Decision Trees built by our methods with those built by Twoing. In this experiment, the maximum depth was set to 16 and we considered just the nominal attributes of the datasets. Each accuracy is the average of 20 stratified 3-fold cross-validations, each generated with a different seed. The entry associated with (\mathcal{D}, I) has two pieces of information: the average accuracy of criterion I on dataset \mathcal{D} and the number of criteria with accuracy statistically lower than that of I on dataset \mathcal{D} . The statistical test used for criteria comparison is a one-tailed paired t -student test with a 95% confidence level. In general there was a balance between Twoing and $\text{GL}\chi^2$, with GLSG being slightly worse. Running the same experiment using maximum depth 5 instead of 16 showed more balanced results. This suggests that our criteria should also be useful in boosting tree methods.

The results of Table ??.(a) also provide evidence of the potential of considering aggregated attributes. The accuracy obtained for the extended versions of datasets **Nursery**, **Cars** and **CoverType** are considerably higher than those obtained for the original versions. For **Contracep**, the effect is not clear.

Table ??.(b) presents the comparison between our methods and Twoing in another scenario, where we use c_{quad} , one of the bias-free criterion proposed in ?, to select the attribute at each node of the tree. Then, both Twoing and our methods are used only for splitting the chosen attribute, which allows for a more direct comparison of their splitting ability. Again, we observed a balance between Twoing and $\text{GL}\chi^2$. This experiment also showed that it is not possible to run c_{quad} in reasonable time for the **Shelter-15** dataset. This happens because it calculates a pseudoinverse of a matrix whose dimension grows with the number of values and classes, which is infeasible for large n and k .

Another key aspect to discuss is the computational cost of the proposed criteria. Table ?? shows the running time of Twoing, $\text{GL}\chi^2$ and GLSG when they are used for both selecting and splitting purposes (the experiment of Table ??.(a)). When the number of classes is small all the criteria have very similar execution time, with Twoing being faster only on the **KDD98-2** dataset. As the number of classes increases, the GL-based methods become much faster

than Twoing, with the turning point being around $k = 7$. For datasets with 15 classes our criteria are 30-300 times faster. We also ran experiments using all the classes available in both the **S.F. Crime** and **Shelter** datasets (39 and 22, respectively). Twoing can not be executed in a reasonable time with that many classes, while GLSG and $GL\chi^2$ ran in approximately 100 seconds on the **S.F. Crime** dataset and 300 seconds on the **Shelter** dataset. This behavior for the Twoing criterion is not surprising, since its running time has an exponential dependence of the number of classes k . Nonetheless, since the execution time for our criteria in this experiment grew in an approximately linear fashion with k , it suggests that they can also be used with datasets that have a much larger number of classes. It is also interesting to note that the aggregated attributes usually appeared at or near the root of the decision trees.

Table ?? shows experiments similar to those presented at Table ??(a), but now using also the numeric attributes. We observed a significant gain in terms of accuracy for almost all datasets. The performance of $GL\chi^2$ was competitive with that of Twoing for all datasets but **KDD98-9** and **CoverType**, where it was considerably better and worse, respectively.

6.3.1 Depth 5

In Table ?? we can see a balance among the different methods. In this experiment we also note that GLSG is more competitive with Twoing and $GL\chi^2$ than in the experiment with depth 16 described in the paper. This suggests that GLSG loses competitiveness as the tree depth increases. Another interesting observation is that the GW-based criteria were about equal or slightly inferior to their GL-based counterparts. This advantage is likely related with the fact that the weights of the cuts computed by the GL approach in this experiment are, in general, larger than those obtained by the GW algorithm.

For the experiment with the bias-free criteria, shown in table ?? we once again observed a balance among the different criteria, with a slight advantage towards our methods. Perhaps surprisingly, the bias free approach had significantly worse results for the datasets with extended attributes.

Table ?? shows the running time of each criterion when used for both selecting and splitting purposes. Twoing is the fastest method when the number of classes is small and the GL-based methods become competitive and even the fastest when the number of classes gets larger, as illustrated by the results on **KDD98** dataset. We also observe that the running time for the GW-based criteria were usually one or two orders of magnitude larger than the others. The only clear exception was in the **CoverType** dataset, where the number of

samples is very large while the attributes' number of values is much smaller.

Finally, Table ?? shows experiments similar to those presented at Table ??, but now using also the numeric attributes. We observed a significant gain in terms of accuracy for all datasets except for KDD98-2. Also note that GLSG was inferior to both Twoing and $GL\chi^2$.

7

Conclusions

In this paper we proposed a framework for designing splitting criteria for handling multi-valued nominal attributes. Criteria derived from our framework can be implemented to run in polynomial time in n and k , with theoretical guarantee of producing a split that is close to the optimal one.

Experiments over 11 datasets suggest that the $GL\chi^2$ criterion, obtained from our framework, is competitive with the well-established Twoing criterion in terms of both accuracy and speed for datasets with a small number of classes ($k \leq 7$). It is also much faster than Twoing when the number of classes is greater than 10, while keeping a comparable accuracy.

Therefore, our methods are an interesting alternative to deal with datasets with a large number of classes that contain nominal attributes with a large number of different values, since those cannot be properly handled by Twoing due to its exponential running time dependence on the number of classes.

Furthermore, our experiments also reinforce the potential of aggregating attributes as a tool for improving the accuracy of decision trees. An interesting topic for future research is evaluating the behavior of our criteria in boosted tree methods. Another direction for future work is developing new methods for automatic aggregating attributes, or improving the available ones.

Bibliography

Austin Animal Center (2016), 'Shelter animal outcomes dataset', *URL:*
kaggle.com/c/shelter-animal-outcomes.