

## ALGORITHM FOR THE DETECTION OF ROYA IN CATURRA COFFEE

David Calle González  
Universidad Eafit  
Country  
dcalleg@eafit.edu.co

Julian Ramirez Giraldo  
Universidad EAFIT  
Country  
jramirezg@eafit.edu.co

Mauricio Toro  
Universidad Eafit  
Colombia  
mtorobe@eafit.edu.co

### ABSTRACT

The problem is the necessity of detecting the presence of roya in coffee earlier, in this case, caturra coffee. This can be done by using certain variables like temperature, ph., and other variables obtained by sensors.

the solution of this problem is quite important due to the high loss of coffee regarding to the late detection of this plague in the crops. Also, especially in Colombia that is one of the most important coffee producers in the world, not being able to detect this sickness on time is a very big problem that could affect our coffee production, and hence, our economy.

### KEYWORDS

Data processing, Roya detection

• Information systems~Linked lists • Theory of computation~Oracles and decision trees

### 1. INTRODUCTION

Initially, if we want to understand what conditions lead to the development of roya in coffee we need to understand what the coffee with roya is, known as well as "Roya del cafeto". this is a coffee plants sickness that produces a fungus named "Hemileia vastatrix" which is considered a parasite.

Coffee with roya was not a problem until the first two known epidemics, for example, the roya epidemics in Sri-Lanka could not be controlled, for this reason, the coffee crops had to be removed. The sickness arrival to Colombia took place in the 80, affecting the low areas crops.

Now, the conditions that allow coffee to get infected with roya and ease its propagations are: -The rain: the fungus needs water to exist.

-Temperature: it needs a temperature between 21 and 25 degrees Celsius

-Climatical Alterations: this generates stress in the crops and turns them into a propitious ambient for roya.

### 2. PROBLEM

The objective of this project consists in developing an algorithm that, through decision trees, could let us know whether a batch of coffee is or not infected with roya in an efficient way.

Solving this problem would prevent the loss of big amounts of coffee, something quite important for our country given the fact that we are the third biggest coffee producers in the world, whereby, it could affect our economy as well as it could affect other countries' economies since coffee is a global market .

### 3. RELATED WORK

#### 3.1 C5.0

"This node uses the C5.0 algorithm to generate a decision tree or set of rules. C5.0 models divide the sample according to the field that offers the maximum information gain. The different subsamples defined by the first division are divided again, usually based on another field, and the process is repeated until it is impossible to divide the subsamples again. Finally, the divisions of the lower level are re-examined, and those that do not contribute significantly to the value of the model are eliminated or pruned."

"C5.0 can generate two types of models. A decision tree is a simple description of the divisions that have been found in the algorithm. The different terminal nodes (or "leaf") describe a subset of training data, and each of the cases included in the training data belongs exactly to a terminal node of the tree. In other words, it is possible to make exactly one prediction for each specific data record present in a decision tree."

#### 3.2 QUEST

"QUEST, or rapid and efficient statistical tree, is a binary classification method to generate decision trees. One of the main motivations for its development has been the

reduction of the processing time necessary for large-scale C&RT analyzes with several variables or several cases. A second objective of QUEST is to reduce the tendency of tree classification methods to favor entries that allow more divisions, that is, continuous input fields (numerical range) or those corresponding to various categories.”

“QUEST uses a sequence of rules based on significance checks to evaluate the input fields of a node. For the purpose of selection, you should only perform a single check on the different entries of a node. Unlike what happens with C&RT, not all divisions are examined and, unlike C&RT and CHAID cases, combinations of categories are not checked when evaluating an input field for selection. This increases the speed of the analysis.”

### 3.3 C&R

“The Classification and Regression Tree (C&R) node is a tree-based prediction and classification method. Similar to C5.0, this method uses repeated partition to divide training records into segments with similar output field values. The C&RT node begins by examining the input fields to find the best division, which has been measured by reducing the impurity index resulting from the division. The division defines two subgroups, which continue to be divided into two other subgroups successively until a stop criterion is activated. All divisions are binary (only two subgroups are created)”

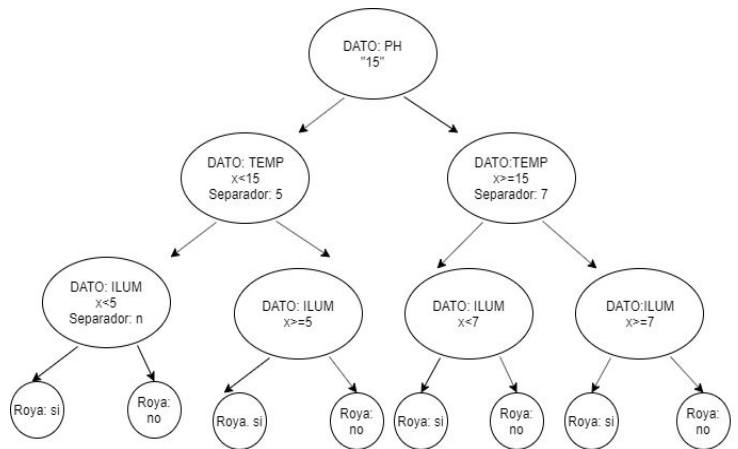
### 3.4 CHAID

“CHAID, or automatic detection of chi-square interactions (Chi-squared Automatic Interaction Detection), is a classification method to generate decision trees using chisquare statistics to identify optimal divisions.”

“CHAID first examines the cross-tabulation tables between the input fields and the results, and then controls the significance through a chi-square independence check. If several of these relationships are statistically important, CHAID will select the most relevant input field (the smallest P value). If an input has more than two categories, compare these categories and counteract those that do not show differences in results. To do this, the pair of categories that present the least difference will be put together, and so on. This process of merging categories stops when all remaining categories differ from each other at the specified check level. In the case of nominal input fields, all categories can be merged. However, in ordinal sets, they can access the adjacent categories.”

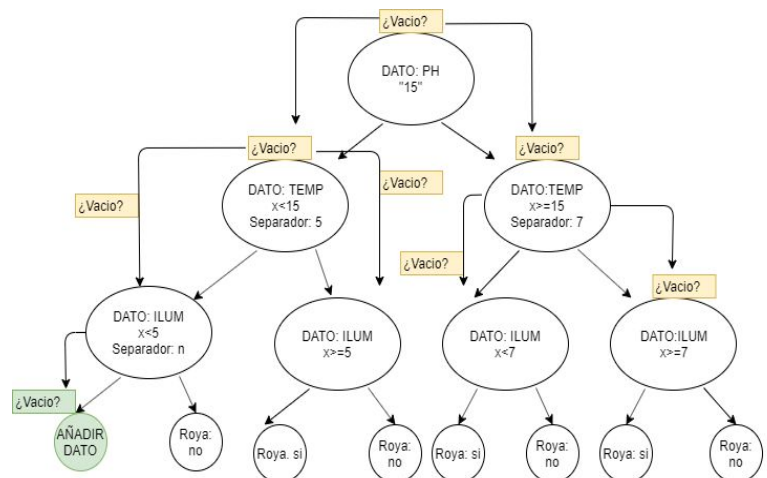
## 4. DECISION TREES

The data structure chosen is a decision tree created by us that follows the following structure:



It should be clarified that not only will work with the 5 generations of data present, they will be more but it is the initial model of the structure. It is about each node containing one of the dataset data and defining for each generation a value called "value that best separates" that is responsible for sending the decision to the left or right node.

### 4.1 DATA STRUCTURE OPERATIONS

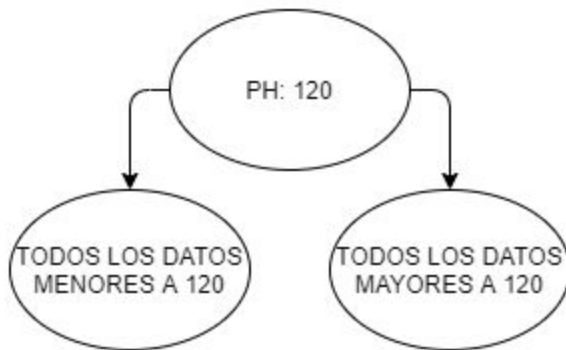


#### 4.1.1 ADD

The algorithm is responsible for going node by node in the tree and asking “Is it empty or not?” Until you find an empty one (represented in green in the graph) and adds the data in the position.

#### 4.1.2 SEPARATE DATA

**DATO BASE PARA REPARTIR: 120**



This method allows to distribute the data depending on a data called "base data" which is what defines the separation of the data to each node (left and right), better called the data of the variable that best separates the data then, yes a data is less than that base data goes to the left side, if the data is larger than the base data goes to the right side.

It should be clarified that this function works in conjunction with the function mentioned in **4.1.1 ADD** to define where there is an empty space to add the data and how to separate them

#### 4.2 CRITERIA FOR WHICH THE DATA STRUCTURE WAS CHOSEN

The decision to take the data structure of “decision trees” was given because it was concluded that the ease of data search is better in a binary structure than in a linked list because to search for an element in a linked list we must go through one to one the elements to find what is sought instead in a binary decision tree that by the same structure the data is already ordered in the case but the

complexity is  $O(1)$  or rather, it is reduced as minimum search at 50% compared to the linked list

#### 4.2 COMPLEXITY ANALYSIS

<b>Print</b>	<b><math>O(n)</math></b>
<b>Tree</b>	<b><math>O(2n)</math></b>
<b>Scroll (Recorrer)</b>	<b><math>n\log(n)</math></b>
<b>Add</b>	<b><math>O(1)</math></b>

#### 4.3 TIMES OF EXECUTION AND USE OF

	Operación	Mmax	Mmin	Tmax	Tmin	Mav	Tav
Datos	imprimir	522240	490776	3	0	1013000	0.46
	generar arbol	522240	0	1	0	10444,8	0.02
	repartir datos	1048576	522240	3	1	632871	0,31
	Operación	Mmax	Mmin	Tmax	Tmin	Mav	Tav
Data_train	imprimir	2048	2048	3	0	2048	0.5
	generar arbol	522240	0	1	0	10444,8	0.02
	repartir datos	526336	293664	3	1	396742	0,2
	Operación	Mmax	Mmin	Tmax	Tmin	Mav	Tav
Test	imprimir	524288	524288	3		524288	0.42
	generar arbol	522240	0	1	0	10444,8	0.02
	repartir datos	526336	293688	2	1	323486	0,2

#### MEMORY

The graph shows different columns defined as:

Operation: The evaluated method

Mmax: Maximum memory used

Mmin: Minimum memory used

Tmax: Maximum time

Tmin; Minimum time

Mav: Average Memory

Tav: Average time

#### REFERENCES

IBM, Nodo CHAID:

[https://www.ibm.com/support/knowledgecenter/es/SS3RA7\\_sub/modeler\\_mainhelp\\_client\\_ddita/clementine/chaidnode\\_general.html](https://www.ibm.com/support/knowledgecenter/es/SS3RA7_sub/modeler_mainhelp_client_ddita/clementine/chaidnode_general.html)

IBM, Nodo Árbol C&R:

[https://www.ibm.com/support/knowledgecenter/es/SS3RA7\\_sub/modeler\\_mainhelp\\_client\\_ddita/clementine/cartnode\\_general.html](https://www.ibm.com/support/knowledgecenter/es/SS3RA7_sub/modeler_mainhelp_client_ddita/clementine/cartnode_general.html)

IBM, Nodo QUEST:

[https://www.ibm.com/support/knowledgecenter/es/SS3RA7\\_sub/modeler\\_mainhelp\\_client\\_ddita/clementine/questnode\\_general.html](https://www.ibm.com/support/knowledgecenter/es/SS3RA7_sub/modeler_mainhelp_client_ddita/clementine/questnode_general.html)

IBM, Nodo C5.0

[https://www.ibm.com/support/knowledgecenter/es/SS3RA7\\_sub/modeler\\_mainhelp\\_client\\_ddita/clementine/c50node\\_general.html](https://www.ibm.com/support/knowledgecenter/es/SS3RA7_sub/modeler_mainhelp_client_ddita/clementine/c50node_general.html)

