



TAXONOMIC EVIDENCE OF CLASSIFICATION APPLYING INTELLIGENT DATA MINING. GALACTIC AND GLOBULAR CLUSTERS

Gregorio PERICHINSKY⁽¹⁾ Arturo Carlos SERVETTO⁽¹⁾,
Elizabeth M. JIMENEZ REY ⁽¹⁾, María Delia GROSSI ⁽¹⁾ Félix VALLEJOS⁽¹⁾,
Rosa Beatriz ORELLANA ⁽²⁾, Angel Luis PLASTINO ⁽³⁾

(1) UNIVERSITY OF BUENOS AIRES, FACULTY OF ENGINEERING, DATABASES AND OPERATING SYSTEM LABORATORY COMPUTER SCIENCE DEPARTMENT

(2) UNIVERSITY OF LA PLATA, FACULTY OF ASTRONOMICAL AND GEOPHYSICAL SCIENCES, CELESTIAL MECHANICS DEPARTMENT, LA PLATA - BUENOS AIRES – ARGENTINA

(3) UNIVERSITY OF LA PLATA, FACULTY OF SCIENCES – BUENOS AIRES – ARGENTINA, DEPARTMENT OF PHYSICAL SCIENCES, PROTEM LABORATORY BUENOS AIRES - ARGENTINA

ABSTRACT

Taxonomy aims to group in families, using so-called structure analysis of operational taxonomic units (OTUs or taxons or taxa).

Clusters that constitute families with a new approach, is the purpose of this paper that belong to a series of papers, of this kind, of the authors. In this case of use the objects are stars instead of other celestial bodies, as asteroids or minor planets. But the data of the observed field in galactic and globular clusters, the taxonomic distances are distorted, because of their projection in the celestial sphere, for that reason, using trigonometrical functions the taxonomic distances should be transformed. The original algorithm, that must be modified, is conformed by: Structural analysis, that shows the relationships, in terms of degrees of similarity, through the computation of the Matrix of Similarity, applying the technique of integration dynamic of independent domains, of the semantics of the Dynamic Relational Database Model. The main contribution is to introduce the concept of spectrum of the OTUs, based in the states of their characters. The concept of families' spectra emerges, if the principles of superposition and interference, and the Invariants determined by the maximum of the Bienaymé-Tchebycheff relation, are applied to the spectra of the OTUs.

Through Intelligent Data Mining, we focused our interest on the Quinlan algorithms, applied in classification problems with the Gain of Entropy, we contrast the Computational Taxonomy, obtaining a new criterion and the robustness of the method.

KEYWORDS: classification, cluster (family), spectrum, induction, divide and rule, entropy.

1. INTRODUCTION

Classification is an abstraction technique used to collect objects with common properties.

The following hypothesis: 1) each object belongs to one (and only one) class and 2) for each class at least one object belongs to the classification, allow us to delimit the domain of objects.

The association of concepts in systematic way by recourse to numerical variables has been the source of a great variety of numerical classification techniques that have their origin in Numerical Taxonomy.

The search of classification concepts that facilitate a robust classification structure (not modifiable by the addition of new information and not altered by the incorporation of new entities) constitutes an important endeavor. In such a line, this work develops tools based on Information Theory and, as a result, a new classification technique is found.

A correspondence with ideas pertaining to the field of the Dynamic Databases is also established [2] [3] [4] [5] [8] [9] [10] [11].

Taxonomic objects are here represented by the application of the semantics of the Dynamic Relational Database Model.

Classification of objects to form clusters or families [26] [27] [28] [29] [30].

Families of OTUs are obtained employing as tools i) the Euclidean distance and ii) nearest neighbor techniques. Thus taxonomic evidence is gathered so as to quantify the similarity for each pair of OTUs (pair-group method) obtained from the basic data matrix [7][16]. The main contribution of the series of papers presented until now was to introduce the concept of spectrum of the OTUs, based in the states of their characters. The concept of families' spectra emerges, if the superposition principle is applied to the spectra of the OTUs, and the groups are delimited through the maximum of the Bienaymé-Tchebycheff relation, that determines Invariants [29] [30].

Applying the integrated independent domain technique dynamically to compute the Matrix of Similarity and, by recourse of an iterative algorithm [33], families or clusters are obtained.

A new taxonomic criterion was thereby formulated.

The considerable discrepancies among the incongruities of existing classifications, whose resultant studies in several disciplines, have motivated an interdisciplinary program of research that notices a clustering of objects in stabilized families [35] [36].

In our case, is worked in an interdisciplinary way in Celestial Mechanics [35] [36], Theory of the Information [1][17], Neural Networks[15] and Dynamic Databases [28] and the Algorithmic of the Numerical Taxonomy [7] [34], to achieve the discovery of the depths of the structure formation of the Solar System, an astronomic application is worked out. The result is a new criterion for the classification of Celestial Bodies in the hyperspace of orbital proper elements, Biological Sciences for linguistic and live beings to avoid confusions, uncertainties and ambiguities, in such a way that the classes include the attributes and the relationships [7].

Thus, a new approach to Computational Taxonomy is presented, that has been already employed with reference to Data Mining.

On the other hand: (i) the works [35] [36] have clarified subtle points concerning the dynamic evolution in the long-term orbits of the asteroids, whose modeling is an essential prerequisite for the proper elements derived (for the classification in families); (ii) the availability of physical data on sizes, shapes, numerical taxonomy, many hundreds of OTUs has provoked new families analyses [1]; (iii) while the most populous families appear in both criteria in quite homogeneous form, the criterion that take into account the composition of the objects and their precedents, is a criterion with more or less difficulty and the criterion which with less difficulty has identified families is that uses data of scientific attributes; (iv) we do not consider in the transformation isotropic and homogeneous sets, changing the values of the attributes to recompute the values of the zones of inter-gap of the objects (if they exist) in the real space with average values and; (v) elimination of groups with few objects, all of which we consider are outside of a Computational criterion.

2. REQUIREMENTS ENGINEERING

Software engineers have many questions to answer, following Fenton & Pfleeger [12]. We know how to assess our current situation scientifically and to determine the magnitude of change when we manipulate our environment.

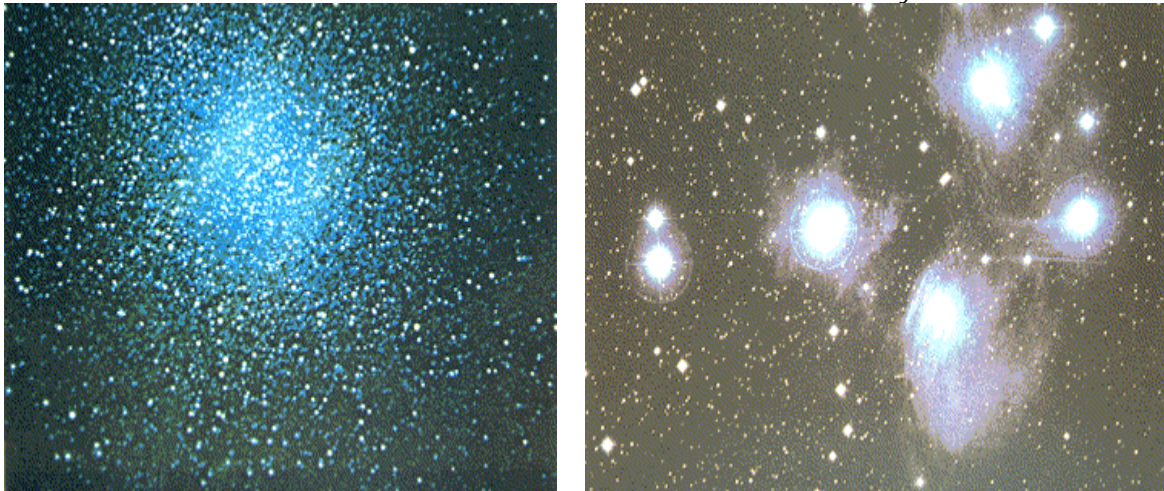
There are key components of empirical investigation in software engineering. Other

methods, including feature analysis, are not addressed in this paper; their descriptions are available in social science research textbooks, and in a series of columns by Kitchenham[23].

2.1. Application in Celestial Bodies.

We will have the classification of celestial bodies in mind, First of all, a fundamental grouping of asteroids was established, the so-called "families of Hirayama" [[19][20][21]]. In Astronomy, taxonomic classification began 50 years ago, used by several authors, for grouping different celestial bodies, among other: Galaxies, Variable Stars, Asteroids, Comets and Cumulous.

This paper is an approach to the problem of Cumulous, because of the considerable discrepancies among the incongruities and existing classifications and in our case we cannot define the invariants and constants with security.



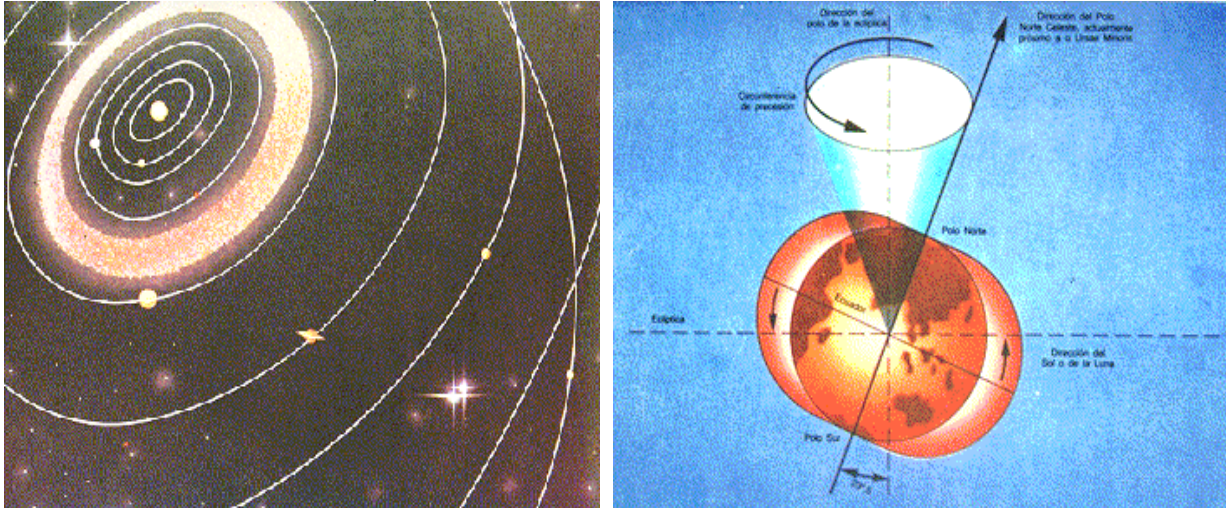
The Cumulous are stars' local condensations united by forces gravitational that they appear in the sky like concentrations of luminous points or, even, as delicate cloudiness. According to their structure they are subdivided in open cumulous and globular clusters. The open clusters, also called Galactic Clusters, they are in the Galactic Disk, that is to say in the central plane of our galaxy, and they are characterized by a stellar density a higher hundred of times that the one that is in the regions that surround in the sun; and however, the stars that compose them are relatively dispersed. The average diameter of the open clusters and the number of stars that contain vary from some dozens to some thousands. They have been observed and classified approximately about 1.000, but it is thought that in our Galaxy it should have at least 15.000. According to the aspect that they present to the telescope, the open clusters are divided in classes.

The Globular clusters they are distributed in a region with form spheroid and called galactic haul, they are characterized by a high stellar density and for a high concentration of stars in the central part of the cluster, until the point that is impossible in many cases, even with a potent telescope, to distinguish each star of those that appear like an only light source. These are less numerous than the open, but bigger and richer cluster in stars. The stellar cluster in general, they have been revealed as a crucible that contains stars of all the types and ages and, therefore, they are object of fundamental study for the investigations on the stellar evolution. The globular clusters on the other hand, are of old formation: about ten thousand million years. Some open clusters are also called in movement, because the stars that compose them are encouraged by an evident movement toward a common apex.

With this term the apparent flight of the distant galaxies is indicated, determined thanks to the Doppler effects. The astronomer E. Hubble realized that the velocity of move away or recession, like one says with the most appropriate word, of the galaxies they increased with the growth of its distances. This discovery gives origin to the

cosmological theory of the Big Bang. As consequence of the explosion of this, it began to expand. The expansion would continue at the moment and it is the one that the astronomers measure under the displacement from toward the red of the spectrum bands of distant galactic sources.

Then it seems that the stars are all on a spherical surface of radio infinite that, with the passing of the hours, it rotates and in fact the celestial bodies occupy different distances with regard to the observer, taking place deformations in the algorithm that makes that the cut is an ellipse instead of a circumference.

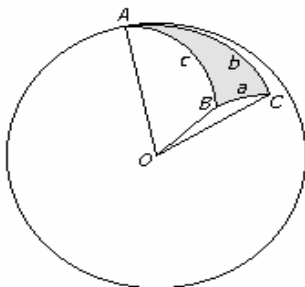


Given a Trihedron or trihedral Angle, is an angle polyhedron of three faces. It also has three dihedral.

The faces and the dihedral of a trihedron complete the following estates:

- ✚ Each face is smaller than the sum of the other ones two.
- ✚ The sum of the three faces is smaller than 360° .
- ✚ The dihedral sum of the three is bigger than 180° and smaller than 540° .

The intersection of a trihedron with a spherical surface with the center in its vertex is a spherical triangle:



The sides of the triangle, a , b , c , are arcs of maximum circumference whose measures coincide with those of the respective faces of the trihedron. The angles of the triangle are the corresponding ones dihedral of the trihedron.

In the algorithm when having stars of the field of the celestial sphere the stars that belong to a cluster to be able to establish similar movements, for the stars of the same cluster, it is necessary to apply trigonometrical functions to homogenize their taxonomic distances.

In many cluster they give a sphere correctly but there are deformations that produce elliptic courts instead of circular, that it is the current situation and it is for that reason that this work is presented that although it is similar to that carried out up to now you should continue with the approaches of Stars in CLUSTERING of Invariants that we denominate of 4th Dispersion and applying the maximum of Tchebycheff (see eq. 3.2.5).

2.2. Spectral analysis classification criterion

With these motivations, we have decided to accomplish with our spectral analysis criterion, the classifications extended to the proper elements database of stars as we did with asteroids in families [12] [13] [14] [17] [18] [19] [20] [21]. We recognize that the works of Zappala [35][36] are very important (automatic classification and hierarchic method), and a point of inflection in the early 90's but is different the approach because we work in computational taxonomy, in a taxonomic hyperspace and not in a transformed space not clearly univocal.

We do not consider in the transformation of isotropic and homogeneous sets, changing the values of the attributes to recompute the values of the regions that emerges from the clustering eliminating groups of few objects, all of which we consider are outside a Computational criterion.

Thus, a new approach to Computational Taxonomy is presented, that has been already employed with reference to Data Mining.

2.3. Intelligent Data Mining Introduction

Machine Learning is the field dedicated to the development of computational methods underlying learning processes and to applying computer-based learning systems to practical problems. Data Mining tries to solve those problems related to the search of interesting patterns and important regularities in large databases [24] [25]. Data Mining uses methods and strategies from other areas, including Machine Learning. When we apply Machine Learning techniques to solve a Data Mining problem, we refer to it as an Intelligent Data Mining.

This paper analyses the TDIDT (Top Down Induction Trees) induction family, and in particular to the C4.5 algorithm. We tried to determine the degree of efficiency achieved by the C4.5 algorithm when applied in data mining to generate valid models of the data in classification problems with the Gain of Entropy.

The C4.5 algorithm generates decision trees and decision rules from pre-classified data. The “divide and rule” method is used to build the decision trees. This method divides the input data in subsets according to some pre-established criteria. Then it works on each of these subsets dividing them again, until all the cases present in one subset belong to the same class.

2.3.1. Constructing the decision trees

2.3.1.1. ID3

The Induction Decision Trees algorithm was developed as a supervised learning method, for build decision trees from a set of examples. The examples must have a group of attributes and a class. The attributes and classes must be discrete, and the classes must be disjoint. The first versions of this algorithm allowed just two classes: positive and negative. This restriction was eliminated in later releases, but the disjoint class restriction was preserved. The descriptions generated by ID3 cover each one of the examples in the training set.

2.3.1.2. C4.5

The C4.5 algorithm is a descendant of the ID3 algorithm, and solves many of its predecessor's limitations. For example, the C4.5 works with continuous attributes, by dividing the possible results in two branches: one for those values $A_i \leq N$ and another one for $A_i > N$. Moreover, the trees are less bushy because each leaf covers a distribution of classes and not one class in particular as the ID3 trees, this makes trees less profound and more understandable[24][25]. C4.5 generates a decision tree partitioning the data recursively, according to the depth-first strategy. Before making each partition, the system analyses all the possible tests that can divide the data set and selects the test with the higher information gain or the higher gain ratio. For discrete attributes, it considers a test with n possible outcomes, n being the amount of possible values that the attribute can take. For continuous attribute, a binary test is performed on each of the values that the attribute can take.

2.3.1.3. Decision trees

The trees TDIDT, to those which belong generated them by the ID3 and post C4.5, are built from method of Hunt. The ID3 and C4.5 algorithms use the “divide and rule” strategy to build the initial decision tree from the training data [23].

The form of this method to build a decision tree as of a set T of training data, divides the data in each step according to the values of the “best” attribute. Any test that divides T in a non trivial manner, as long as two different $\{T_i\}$ are not empty, is very simple.

They will be the classes $\{C_1, C_2, \dots, C_k\}$. T contains cases belonging to several classes, in this case, the idea is to refine T in subsets of cases that tend, or seem to tend toward a collection of cases belonging to an only class. It is chosen a test based on an only attribute, that has one or more resulted, mutually excluding $\{O_1, O_2, \dots, O_n\}$. T is partition of the subsets T_1, T_2, \dots, T_n where T_i contains all the cases of T that have the result O_i for the elected test. The decision tree for T consists in a node of decision identifying the test, with a branch for each possible result. The construction mechanism of the tree is applied recursively to each subset of training data, so that the i -th branch carries to the decision tree built by the subset T_i of training data.

Still, the ultimate objective behind the process of constructing the decision tree isn't just to find any decision tree, but to find a decision tree that reveals a certain structure of the domain, that is to say, a tree with predictive power. That is the reason why each leave must cover a large number of cases, and why each partition must have the smallest possible number of classes. In an ideal case, we would like to choose in each step the test that generates the smallest decision tree.

Basically, what we are looking for is a small decision tree consistent with the training data. We could explore and analyze all the possible decision trees and choose the simplest one. However, the searching and hypothesis space has an exponential number of trees that would have to be explored. The problem of finding the smallest decision tree consistent with the training data has NP-complexity.

To calculate which is the "best" attribute to divide the data in each step, both the information gain and the gain ratio were used. Moreover, the trees generated with the C4.5 algorithm were pruned according to the method, this post-pruning was made in order to avoid the overfitting of the data.

2.3.1.4. Transforming decision trees to decision rules

Decision trees that are too big or too bushy are somewhat difficult to read and understand because each node must be interpreted in the context defined by the previous branches. In any decision tree, the conditions that must be satisfied when classifying a case can be found following a trail from the root to the leaf to which that case belongs. If that trail was transformed directly into a production rule, the antecedent of the rule would be the conjunction of all the tests in the nodes that must be traversed to reach the leaf. All the antecedents of the rules built this way are mutually exclusive and exhaustive.

To transform a tree to decision rules, the C4.5 algorithm traverses the decision tree in preorder (from the root to the leaves, from left to right) and constructs a rule for each path from the root to the leaves. The rule's antecedent is the conjunction of the value tests belonging to each of the visited nodes, and the class is the one corresponding to the leaf reached.

For the Evaluation of the TDIDT family is used a crossed-validation approach to evaluate the decision trees and the production rules obtained. Each dataset was divided into two sets with proportions **2:3** and **1:3**. We used two thirds of the original data as a training set and one third to evaluate the results [31].

3. NUMERICAL TAXONOMY

We infer an **analogy** of the **taxonomic representation** [26] [27] [28] [29] [30] **in dynamic relational database**.

We explain the theoretical development of a domain's structured Database and how they can be represented in a Dynamic Database.

Immediately we apply our model to the structural aspects of the taxonomy, applying Scaling Methods for domains [7] [34].

We define numerical methods used for establishing and defining clusters by their taxonomic distances.

We shall let C_{jk} stand for a general dissimilarity coefficient of which taxonomic distance, d_{jk} , is a special example. Euclidean distances will be used in the explanation of clustering techniques.

In discussing clustering procedures we make a useful distinction between three types of measure.

We use clustering strategy of space-conserving or the space-distorting strategies that appears as though the space in the immediate vicinity of a cluster has been contracted or dilated and if we return to the criterion of admission for a candidate joining an extant cluster, this is constant in all **pair-group** method.

Thus we can represent the **data matrix** and to compute the **resemblance of normalized domains**.

The steps of clustering are the **recomputation** of the coefficient of similarity for future admission followed by the **admission criterion** for new members to an established cluster.

The strategies of both **space-conserving** and **space-distorting** that appear in the immediate vicinity of a cluster either contract or dilate the space, and this is constant in all **pair-group** methods [7] [34].

3.1. Calculation of the Average and of the Standard Deviation for the Normalization

In normalizing characters we compute the average value and the standard deviation of each string (the states of each character) and express each state as a deviation of the average in standard deviation units. The normalization of the states of the character makes the average of all characters to vanish. Likewise, variances adopt the value unity. We have

$$\bar{X}_j = (\sum_i^n X_{ij}) / n \quad (3.1.1)$$

$$\sigma_j = ((\sum_i^n (X_{ij} - \bar{X}_j)^2) / (n - 1))^{1/2} \quad (3.1.2)$$

$$\bar{X}'_{ij} = (X_{ij} - \bar{X}_j) / \sigma_j \quad (3.1.3)$$

For the normalized domains we calculate both the average difference among characters (its absolute value) and the concomitant taxonomic distances. For the latter we consider two metrics: that of Minkowski and the so-called Manhattan one [7]. The quantity

$$\bar{D} = (\sum_i^n | X_{ij} - X_{ik} |) / n \quad (3.1.4)$$

is the mean difference among characters,

$$\Delta_{jk} = [\sum_i^n (X_{ij} - X_{ik})^2]^{1/2} \quad (3.1.5)$$

is the distance Δ_{jk} among OTUs, and we consider further the average value

$$d_{jk} = ((\sum_i^n (X_{ij} - X_{ik})^2 / n))^{1/2} \quad (3.1.6)$$

due to the fact that Δ_{jk} grows with the number of characters.

The expectation value (d) of the d_{jk} for a normal distribution of zero and variance unity is:

$$E(d) = ((n - 1)! (\pi / n)^{1/2}) / (2^{n-2} [((n/2) - 1)!]^2) \quad (3.1.7)$$

After using the Stirling approximation we get

$$E(d) \approx \sqrt{2} (1 - 1/n)^{1/2} ((1 + (1 / (n - 2)))^{1/2} (1 / e)) \quad (3.1.8)$$

and the expectation value of the variance for (d) turns out to be

$$E(\sigma^2_d) = 2 - [E(d)]^2 \approx 1/n. \quad (3.1.9)$$

3.2. Dispersion

The variance is a moment of second order and represents to the moment of inertia of the distribution of objects (masses) with respect to their center of gravity (the so-called centroid) [6].

$\overline{X'_{ij}} = (X_{ij} - \overline{X_j}) / \sigma_j$ (3.2.1) is a normalized variable that represents the deviation of the X_{ij} with respect to their mean (in units of σ_j).

As usual, we take the dispersion to be given by the variance σ^2_d . The mean-squares method is now to be applied.

Let $g(X_{ij})$ be a not negative function of the variable X_{ij} . For all $k > 0$ will have the probability function:

$$P [g (X_{ij}) \geq K] \leq (E (g (X_{ij})) / K) \quad (3.2.2).$$

Theorem of Tchebycheff

Let S be the set of all the X_{ij} that satisfy the inequality $g(X_{ij}) \geq K$. The truth of the theorem stems from the relationship (valid in any number of dimensions):

$$Eg(X_{ij}) = \int_{-\infty}^{\infty} g(X_{ij})dF \geq K \int_S dF = KP(S) \quad (3.2.3)$$

If $g(X_{ij}) = (X_{ij} - \overline{X_j})^2$, $K = k^2 \sigma_j^2$, which leads, for all $k > 0$, to the inequality of Bienaymé-Tchebycheff:

$$P (| X_{ij} - \overline{X_j} | \geq k \cdot \sigma_j) \leq 1 / k^2 \quad (3.2.4)$$

This inequality shows that

$$\overline{X_j} - k \cdot \sigma_j < X_{ij} < \overline{X_j} + k \cdot \sigma_j \quad (3.2.5)$$

(maximal value equal to $1 / k^2$).

In particular, for an average value $\overline{X_j}$ and deviation σ_j with a mass $1 / 2 \cdot k^2$ located at each the points $X_{ij} = \overline{X_j} \pm k \cdot \sigma_j$ one has

$$P (| X_{ij} - \overline{X_j} | \geq k \cdot \sigma_j) = 1 / k^2 \quad (3.2.6)$$

a maximal limit value that can not be improved upon.

This inequality shows that the quantity of mass of the distribution is to be found in the interval

$$\overline{X_j} - k \cdot \sigma_j < X_{ij} < \overline{X_j} + k \cdot \sigma_j \quad (3.2.7).$$

The inequality permits one to fix both distribution levels and the radius of a cluster.

If we take $k=2$ then we obtain $k \cdot \sigma_j = \sqrt{2} \cdot \sigma_j$ as the maximal value (3.2.8).

As are taken the stars of the field, there are some of them that belong to different clusters and, as are not parallel the sides of the triangles, the value of the corresponding attribute should be adjusted, after applying the normalization and the algorithm that it calculates matrix of similarity. Applying the theorem of the sine for spherical triangles, it is modified the value of the taxonomic distance.

3.3. Clusters and Spectra.

In discussing Sequential, Agglomerative, Hierarchic and Non-overlapping (SAHN) [34] clustering procedures we make a useful distinction between the three types of measure.

We shall be concerned with clusters J, K and L containing t_j , t_k and t_l OTUs, respectively, where t_j , t_k and t_l all ≥ 1 . OTUs j and k are contained in clusters J and K , and $l \in L$, respectively. Given two clusters J and K that are to be joined, the problem is to evaluate the dissimilarity between the resulting joint cluster and additional candidates L for further fusion. The fused cluster is denoted (J, K) , with $t_{j,k} = t_j + t_k$ OTUs.

The cluster center or centroid represents an average object, which is simply a mathematical construct that permits the characterization of the Density, the Variance, the taxon radius and the range as **INVARIANT** quantities.

The states of the taxonomic characters in a class, defined ordinarily with reference to the set of their properties, allow one to calculate the distances between the members of the class. The distances can be established by the similarity relationship among individuals (obtaining a matrix of similarity that has been computed).

Considering characteristic spectra [13][14][18][32], in addition to the states of the characters or attributes of the OTUs, we introduce here the new **SPECTRAL** concepts of i) **OBJECTS** and ii) **FAMILY SPECTRA**.

Within the taxonomic space this method of clustering delimits taxonomic groups in such a manner that they can be visualized as characteristic spectra of an OTU and characteristic spectra of the families.

We define an individual spectral metric for the set of distances between an OTU and the other OTUs of the set. Each one provides the states of the characters and, therefore, is constant for each OTU, if the taxonomic conditions do not change (in analogy with the fasors) having an individual taxonomic spectrum (ITS).

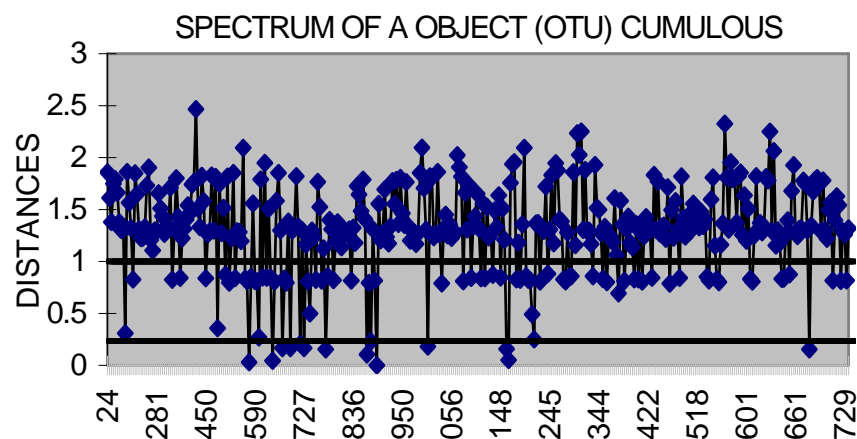
The spectrum of taxonomic similarity is the set of distances between the OTUs of the set that determine the constant characteristics of a cluster or family, for a given type of taxonomic conditions. Multiplying the values of the domain related with the attributes of the sides of the spherical triangle (see figure in 2.1. paragraph) by the cosine of the opposed side.

Invariants are found that characterize each cluster. Among them we mention the variance, the radius, the density and the centroid.

These invariants are associated with the spectra of taxonomic similarity that identify each family.

3.3.1. Variation Range Normalization.

There exist sound reasons for considering that the weight of a character should be inversely proportional to its variability. For normally distributed quantitative characters their information content (in the information theory sense) is proportional to the variance. If the variances are made equal, then each character contributes an equal informational amount. Such an uniform probability yields, of course, the maximum possible entropy.



One observes in the graph, for the line of equal Invariant (ordinate unity), a region that clearly shows the objects that constitute it. Objects belonging to other regions are to be found above such a line. Below the 0.2343 (*RED*) line at ordinate one sees objects of a family. Above these two lines we encounter other objects. A more detailed analysis is required in order to ascertain to which family these objects belong.

3.4. Tests of Intelligent Data Mining

A software system was constructed to evaluate the C4.5 algorithm. This system takes the training data as an input and allows the user to choose whether he wants to

construct a decision tree according to the C4.5. If the user chooses the C4.5, the decision tree is generated, then it is pruned and the decision rules are built.

We continue having contour problems in some galactic clusters, should find invariants and constant of widespread Tchebycheff for all the cases.

3.4.1. Compute of the Information Gain

In the cases, in those which the set contains examples belonging to different classes, is accomplished a test on the different attributes and is accomplished a partition according to the "better" attribute. To find the "better" attribute, is used the theory of the information, that supports that the information is maximized when the entropy is minimized. The entropy determines the randomness or disorder of a set.

We suppose that we have negative and positive examples. In this context the entropy of the subset S_i , $H(S_i)$, it can be calculated as:

$$H(S_i) = -p_i^+ \log p_i^+ - p_i^- \log p_i^- \quad (3.4.1)$$

Where p_i^+ is the probability of a example is taken in random mode of S_i will be positive. This probability may be calculated as

$$p_i^+ = \frac{n_i^+}{n_i^+ + n_i^-} \quad (3.4.2)$$

Being n_i^+ the quantity of positives examples of S_i , and n_i^- the quantity of negatives examples.

The probability p_i^- is calculated in analogous form to p_i^+ , replacing the quantity of positives examples by the quantity of negatives examples, and conversely.

Generalizing the expression (3.4.1) for any type of examples, we obtain the general formulation of the entropy:

$$H(S_i) = \sum_{i=1}^n -p_i \log p_i \quad (3.4.3)$$

In all the calculations related to the entropy, we define $0 \log 0$ equal to 0.

If the attribute at divide the set S in the subsets S_i , $i = 1, 2, \dots, n$, then, the total entropy of the system of subsets will be:

$$H(S, at) = \sum_{i=1}^n P(S_i) \cdot H(S_i) \quad (3.4.4)$$

Where $H(S_i)$ is the entropy of the subset S_i and $P(S_i)$ is the probability of the fact that an example belong to S_i . It can be calculate, used the relative sizes of the subsets, as:

$$P(S_i) = \frac{|S_i|}{|S|} \quad (3.4.5)$$

The gain of information may be calculate as the decrease in entropy. Thus:

$$I(S, at) = H(S) - H(S, at) \quad (3.4.6)$$

Where $H(S)$ is the value of the entropy a priori, before accomplishing the subdivision, and $H(S, at)$ is the value of the entropy of the subsets system generated by the partition according to at .

The use of the entropy to evaluate the best attribute is not the only one existing method or used in Automatic Learning. However, it is used by Quinlan upon developing the ID3 (paragraph 2.3.1.1.) and his succeeding the C4.5 (paragraph 2.3.1.2.).

3.4.2. Numerical Data

The decision trees can be generated so much as discrete attributes as continuous attributes. When it is worked with discrete attributes, the partition of the set according to the value of an attribute is simple.

To solve this problem, it can be appealed to the binary method. This method consists in forming two ranges of agreement values to the value of an attribute, that it can be taken as symbolic.

4. RESULTS AND CONCLUSIONS

The hypothesis space for this algorithm is complete according to the available attributes. Because any value test can be represented with a decision tree, this algorithm avoids one of the principal risks of inductive method that works reducing the spaces of the hypothesis.

The C4.5 with post-pruning results in trees smaller and less bushy. If we analyze the trees obtained in the domain, we'll see that the percentages of error obtained with the C4.5 are between a 3% and a 3.7%, since that the C4.5 generate smaller trees and smaller rulesets. Derivative of the fact that each leaf in a tree generated covers a distribution of classes.

Nevertheless, we can state that the results show that the proportion of error depends on the data domain like is established in the development of the paper.

For future study, we suggest an analysis the input datasets and the Invariants and Constants as is mentioned in the previous paragraphs (2.1.) – (3.2.).

REFERENCES

- [1.] Abramson, N., "Information Theory and Coding". McGraw Hill. Paraninfo. Madrid. 1966.
- [2.] Batini, C., Ceri, S., Navathe, S.B. "Conceptual Databases Design" Addison Wesley. 1998.
- [3.] Codd E. F. "Relational Completeness of Data Base Sublanguages". Database Systems, Courant Computer Science Symposia Series 6, Englewood Cliffs, New Jersey, Prentice-Hall. 1972.
- [4.] Codd E. F. "How Relational is your Database Management System?". Computer World. 1985.
- [5.] Codd E. F. "The Relational Model for Database Management: Version 2". Addison Wesley. 1990.
- [6.] Cramer, Harald. "Mathematics Methods in Statistics". Aguilar Edition. Madrid. Spanish. 1958.
- [7.] Crisci, J.V. , Lopez Armengol, M.F. "Introduction to Theory and Practice of the Numerical Taxonomy", A.S.O. Regional Program of Science and Technology for Development. Washington D.C.Spanish. 1983.
- [8.] Date, C.J. "An Introduction to Dabase Systems Vol. I". 6th Ed. Addison Wesley. 1995.
- [9.] Date, C.J. "Relational Database: Selected Writings". Addison Wesley. 1986.
- [10.] de Miguel, A., Piattoni, M. "Concepts and Design of Databases." Addison Wesley. 1994. Spanish.
- [11.] Elmasri, R., Navathe, S. "Fundamentals of Database Systems". The Benjamin/Cummings Publishing Company and Addison Wesley. 1997.
- [12.] Fenton, N.E., Pfleeger, Sh.L. "Software Metrics". PWS Publishing Company. 1997.
- [13.] Feynman, R.P., Leighton, R.B. & Sands, M. "Lectures on physics, Mainly Mechanics, Radiation and Heat". pp. 25-2 ff, 28-6 ff, 29-1 ff, 37-4. 1971.
- [14.] Frank, N.H. "Introduction to Mechanics and Heat". Science Service. Washington. Editorial Atlante. Spanish. 1949.
- [15.] Freeman, J.A., Skapura, D.M. "Neural Networks. Algorithms, applications and techniques of programming". Addison Wesley. Iberoamericana. Spanish. 1991.
- [16.] Gennari, J.H. "A Survey of Clustering Methods" (b). Technical Report 89-38. Department of Computer Science and Informatics. University of California., Irvine, CA 92717. 1989.
- [17.] Hamming, R.W. "Coding and information theory". Englewood Cliffs, NJ: Prentice Hall. 1980.
- [18.] Hetcht, E. and Zajac, A., "Optic". Fondo Educativo Interamericano. pp. 5-11-206-207-293-297-459-534. Spanish 1977.
- [19.] Hirayama, K. "Groups of Asteroids Probably Common Origin". Proceeding of Physics-Mathematics Society. Japan II:9. pp 354-351. 1918.

- [20.] Hirayama,K. "Groups of Asteroids Probably Common Origin". The Astronomical Journal: 31, pp 185-188. 1918.
- [21.] Hirayama,K. "Present State of the Families of Asteroids". Proceeding of Physics-Mathematics Society. Japan II:9. pp 482-485. 1933.
- [22.] Hunt, E.B., Marin, J., Stone, P.J. 1966 (1995-AI). Experiments in Induction. New York: Academic Press, USA.
- [23.] Kitchenham,B., Pickard,L., Pfleeger, S.L. "Case studies for method and tool evaluation". IEEE Software, 12(4) pp 52-62. 1995.
- [24.] Quinlan, J.R. 1986. Induction of Decision Trees. In Machine Learning, Ch. 1, p.81-106. Morgan Kaufmann.
- [25.] Quinlan, J.R. 1987. Generating Production Rules from Decision trees. Proceeding of the Tenth International Joint Conference on Artificial Intelligence, p. 304-307. San Mateo, CA., Morgan Kaufmann, USA.
- [26.] Perichinsky, G., Servetto, A., Crocco, E. "Relational Data Bases Structured on Dynamic Domains of Attributes".18th Sessions. Operative Research and Informatic's Argentine Society. 1989.b.
- [27.] Perichinsky, G., Servetto, A. "Dynamically Integrated Independent Domains on Data Bases".19th Sessions Operations Research and Informatic's Argentine Society. 1990.
- [28.] Perichinsky, G., Feldgen, M., Clúa, O. "Conceptual Contrast of Dynamic Data Bases with the Relational Model" in Proceedings International Association of Science and Technology for Development. 14th Applied Informatics Conference. Innsbruck, Austria. 1996.
- [29.] Perichinsky, G., Orellana, R. and Plastino, A.L. "Spectra of Taxonomic Evidence in Databases." Proceedings of International Conference on Computer Science, Software Engineering, Information Technology, e-Business and Applications (CESITeA '02). Foz do Iguazu. Brazil. 2002.
- [30.] Perichinsky, G., Orellana, R., Plastino, A.L., Garcia Martinez, R., Servente, M. and Servetto, A. C. "Taxonomic Evidence Applying Algorithms of Intelligent Data Mining" Proceedings of International Conference on Computer Science, Software Engineering, Information Technology, e-Business and Applications (CESITeA '03). Rio de Janeiro. Brazil. 2003.
- [31.] Perichinsky G., Jiménez Rey E. M., Grossi M. D., Vallejos F. A., Servetto A. C., Orellana R. B., Plastino A. L. Taxonomic evidence applying intelligent information algorithm and the principle of maximum entropy: the case study of asteroids families. Electronic magazine of Systems of Information, RESI. ISSN 1677-3071. Edición 6–Año IV–Volumen IV–Número 2. Departamento de Informática y Estadística. Universidad Federal de Santa Catarina. Brasil. 2005.
- [32.] Sawyer, R.A. "Experimental Spectroscopy". Dover Publication. New York. 1963.
- [33.] SEI "Rationale for SQL Ada module Description language (SAMeDL)" Ver. 2.0 CMU/SEI-92-TR-16, oct 1992.
- [34.] Sokal, R.R., Sneath, P.H.A. "Numerical Taxonomy". W.H. Freeman and Company. 1973.
- [35.] Zappala, V., Cellino,A., Farinella,P., Knêzevic,Z., "Asteroid Families. I. Identification by Hierarchical Clustering and Reliability Assessment". The Astronomical Journal, 100, 2030. 1990.
- [36.] Zappala, V., Cellino,A., Farinella,P., Milani,A., "Asteroid Families. II. Extension to Unnumbered Multiopposition Asteroids" The Astronomical Journal, 107, 772. 1994.