

Abstract

An important problem faced in archaeology is to establish an efficient computer-based object classification algorithm, however, most of such classifications might be dependent on users' subjectivities, hampering its implementation. Hence, we propose classification criteria based on the object's topological and geometrical properties, such criteria being more robust and objective. This particular project is based on the idea of Euler Characteristic Graph (ECG), which summarizes the topological-geometrical evolution of an object. Finally, based on it, a new possible classification of pre-Columbian masks is suggested.

Introduction

AN IMPORTANT problem faced in archaeology is to establish an efficient computer-based object classification algorithm to further develop an artifacts' database. The assortment of a collection of **pre-Columbian masks** into several groups, each sharing a set of characteristics, arises as a concrete problem. Actual assortments are considerably subjective, making them prone to discrepancies among the archaeology community. Identifying topological-geometrical invariance within masks and employing it as the primary classification criterion is the main goal of this project. To that end, the Instituto Nacional de Antropología e Historia (INAH) provided a 128 different masks test set. After object pre-process, each mask consists of a triangulated 3 dimensional mesh embedded into the $[-1, 1]^3$ cube such that its mass center locates at origin. INAH provided a standard classification of such test set as well, hence establishing a start and comparison point for the suggested classifications. This project was mainly based upon the Euler Characteristic graph (ECG) idea as exposed in [4].

General methodology

CONSIDER a n -dimensional object $X = (V_0, V_1, \dots, V_n)$, with V_k the set of all its k -cells. We define its Euler Characteristic (EC) χ as:

$$\chi = \sum_{k=0}^n (-1)^k |V_k|.$$

Fix a **function** $g : V_0 \rightarrow [a, b]$, with $[a, b]$ a closed interval, and a fixed **number** T of thresholds. These will be the main ingredients to construct the ECG of X . For instance, g could be Gaussian curvature, Shape Index, distance to mass center, etc. From it, auxiliary functions $g_k : V_k \rightarrow [a, b]$ are constructed as follows: for each k -cell $\{v_0, v_1, \dots, v_k\}$ with $v_0, v_1, \dots, v_k \in V_0$, define

$$g_k(\{v_0, v_1, \dots, v_k\}) = \min_{0 \leq i \leq k} \{g(v_i)\}.$$

Once every cell has an associated numerical value, the interval $[a, b]$ is divided into T equally-spaced thresholds $a = t_0 < t_1 < t_2 < \dots < t_T = b$. For each of them, the cardinality of the vertices subset $V_0^{(i)} \subset V_0$ defined as $V_0^{(i)} = \{x \in V_0 : g(x) > t_i\}$ might be computed using a similar method as in *bucket-sort*. Proceed analogously with the subsets $V_k^{(i)}$ of k -cells. Thus, the **EC at i -th threshold** is defined as:

$$\chi_i = \sum_{k=0}^n (-1)^k |V_k^{(i)}|.$$

Notice that if a given vertex x fails to surpass the i -th threshold (in the sense that $g(x) \leq t_i$), then g_k 's definition implies that every k -cell containing x as a vertex will fail to surpass such threshold as well.

Finally, an **object's ECG** is simply defined as the graph given by χ_i versus t_i . Notice that as the threshold value increases, there are less vertices (hence less k -cells) that surpass such threshold. Qualitatively speaking, the ECG is a mathematical tool that summarizes the topological-geometrical evolution of the object as it *disintegrates through time*.

Assuming that the numerical values of function g have already been computed, assigning numerical values to every k -cell is an $O(V_0)$ complexity algorithm; computing and storing every $|V_k^{(i)}|$ value has $O(V_0 + T)$ complexity. Finally, taking advantage of the previously stored values, computing every χ_i has complexity $O(T)$. Thus, the whole algorithm runs with **complexity** $O(V_0 + T)$, demanding $O(V_0 + T)$ memory. Overall, computation of ECGs is a quite efficient algorithm.

It is expected that similar ECGs of two given objects suggest topological and geometrical similarity between these two, thus strongly suggesting that both of them belong into the same category. To determine similarity between graphs, Support Vector Machine (SVM) and k -Nearest Neighbor (KNN) might be used.

Particular methodology

A C/C++ based algorithm was implemented to compute the ECG of every object. *Point Cloud Library (PCL)* was used to compute the principal curvature values at each vertex [5] and the magnitude of projections from each vertex to each cube's face. Once every ECG was computed, they were assorted into

9 different families using the standard SVM with linear kernel and KNN algorithms as readily available in *Python* [2, 3]. It is worth noting that both SVM and KNN supervised classification methods, thus, are dependent on a training set previously defined by a user.

For each mask, several ECGs were computed experimenting with different functions $g : V \rightarrow [a, b]$ such as Shape Index, (normalized) mean curvature, and length of projections from each vertex to right, top and frontal cube's faces. New ECGs were obtained by concatenating two or more previous ECGs. Several T values were employed as well, experimenting with 32, 64, and 256 different thresholds.

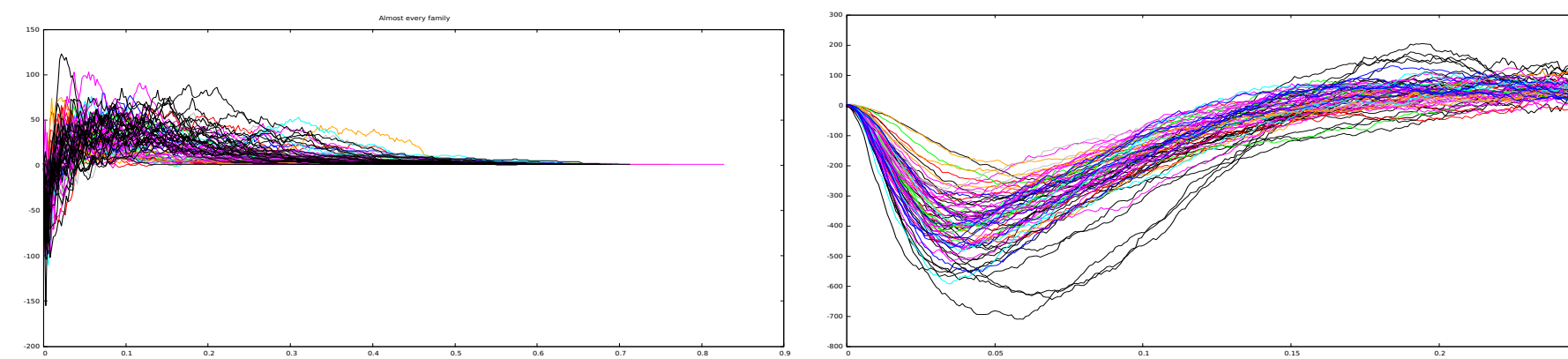
Establishing a new classification of a 128 pre-Columbian masks test set was the main goal. According to INAH, these are sorted into 9 different families; thus, the goal was to specifically assort 128 masks into 9 different groups. The main hindrance faced was the heterogeneity of number of representatives of each original group; most of the groups (6 out of 9) had less than 10 items each, while a particular group had more than 50 members. The lack of members hinders a possible characterization of its family, which in turn hinders the development of a reasonable classification.

Another problem aroused by such lack of characterization was the practical impossibility of employing unsupervised classification methods such as k -Means Clustering (KMC) [1]. In this case, the largest family will tend to absorb the rest, making the classification a unique large cluster along some scattered individuals.

6 out of 9 original families were included completely into the training set due to lack of items in each of them, while 10 representatives were included from the remaining three families. In total, the training set was made of 58 masks divided into 9 families; to classify the remaining 70 a linear kernel SVM method was used.

Results

COMPUTING each mask's ECG based on normalized mean curvature and 256 thresholds was one of the first experiments. Every graph was plotted on the same plane, and colored according to their original family: two ECGs will be the same color if their corresponding masks belong to the same original family. The same procedure was repeated using Shape Index.



(a) Mean curvature (b) Shape Index

Figure 1: ECGs with $T = 256$ thresholds. Each of the 9 original families correspond to a different color.

The figures 1(a) and 1(b) discourage the search of a coherent classification based on ECGs and either mean curvature or Shape Index. Notice that each family are thoroughly mixed with the rest of families.

On the other hand, the graphic based on the three principal projections (from each vertex to rightmost, top and frontal cube's faces) encourages the development of a classification based on it.

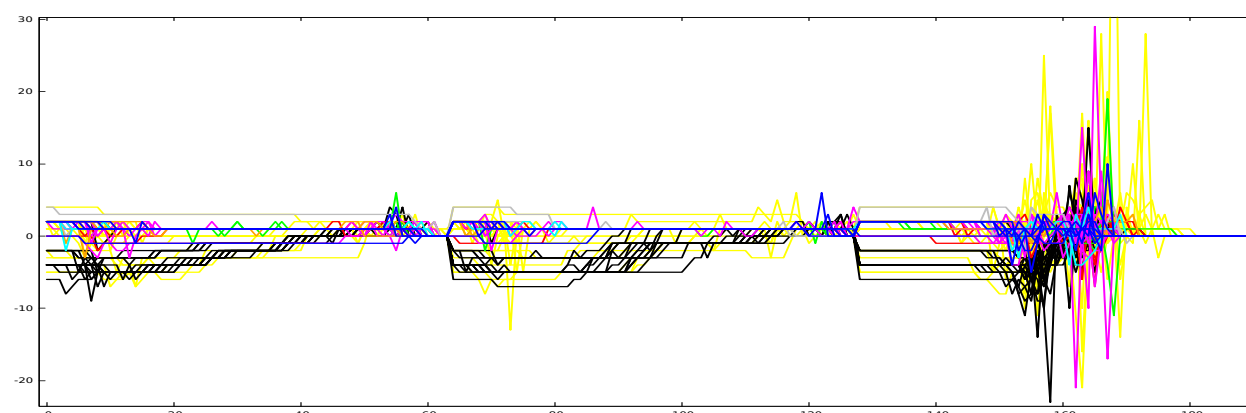


Figure 2: ECGs obtained by concatenating three ECGs, one for each principal projection. Different colors refer to different families.

After running the training phase, a new classification was obtained. The first striking difference between this assortment and the original one is the homogenization of number of specimens per family, since only two out of nine families possess less than 10 items. Out of the remaining seven families, 10 items were chosen randomly and each family was plotted.

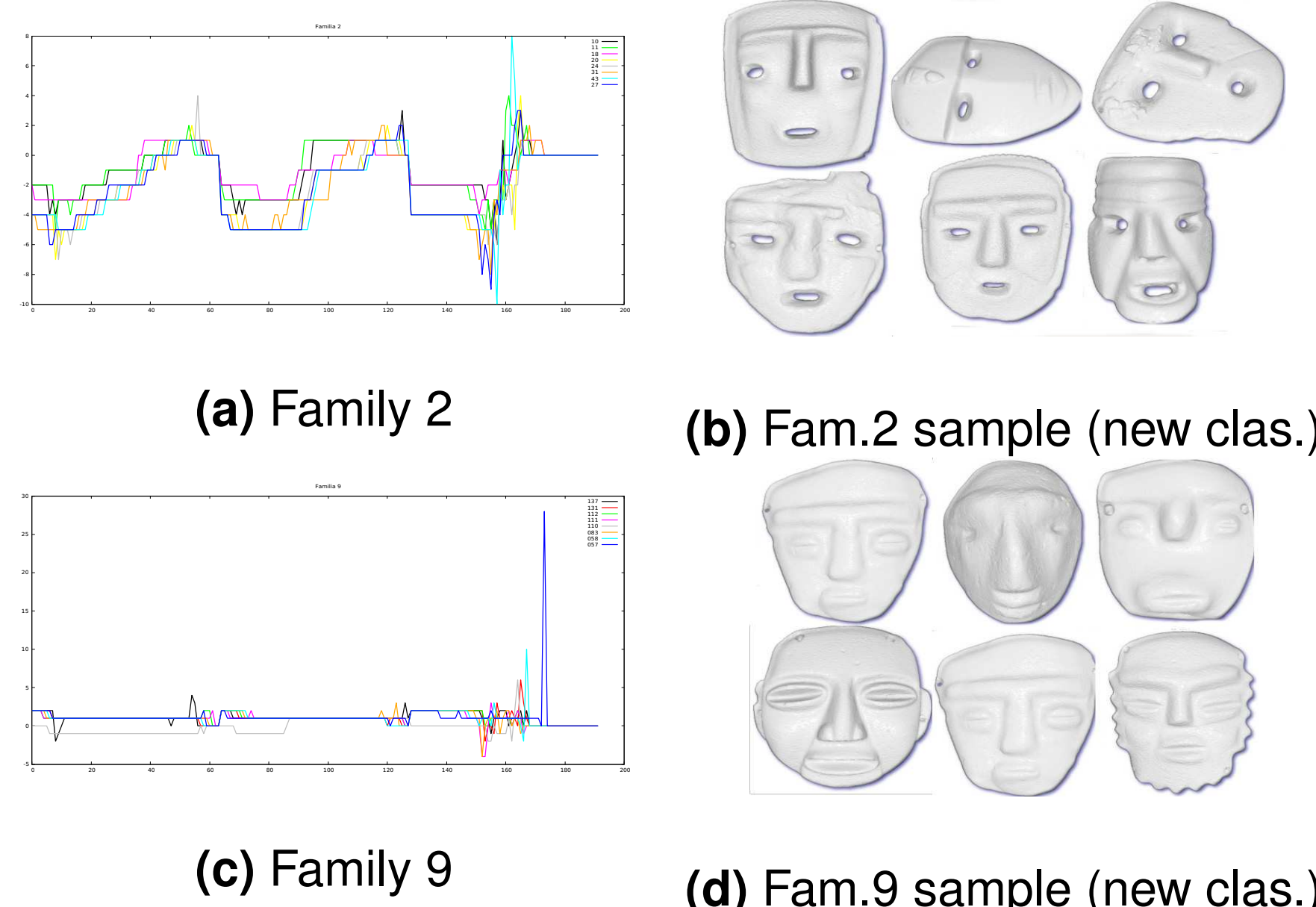
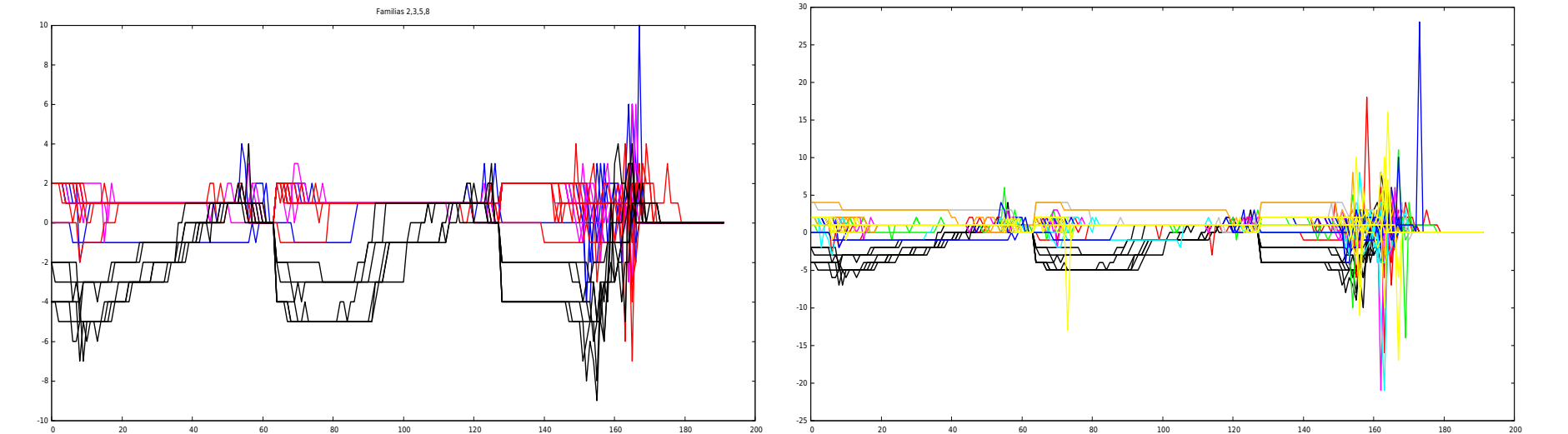


Figure 3: ECGs of different masks within the same family according to the new classification. Different colors refer to different items.

It's worth to notice that, with exception of at most one mask per family, every item follows a defined pattern within its family. Removing the outlying masks and plotting the same four families simultaneously, as in figure 4(a), it is reasonable to claim a different pattern for each different family. Regarding this fact, 10 items of each family are randomly chosen (if a family has less than 10 items, all its members are chosen) and their respective ECGs are plotted on the same plane. The resulting figure, 4(b), colored according to the new classification, shows a reasonable pattern for each color.



(a) Families 2,3,5,9 (b) New classification

Figure 4: ECGs based on the three main projections, now colored according to the new classification

An SVM was also run using ECGs based upon mean curvature, Shape Index, and concatenations of curvatures and projections. As foretold in figure 1, none of such classifications yielded coherent results, since visually, is difficult to recognize some resemblance among the masks in each family.

Also, a KMC algorithm was run using ECGs based on curvatures, projections, or concatenations. All of these resulted in a unique large cluster and several small ones as forewarned, hence, there was not a serious consideration of it as an assortment.

Conclusions

The computation of the ECG associated to a given object, especially if it is based on projections, is a simple algorithm of linear complexity and memory. It is possible to process a large number of objects quickly, even when these are composed of a large number of vertex. This efficiency suggests stronger applications in a wide range of real-time algorithms, especially pattern and object recognition in general.

The ECG algorithm is quite general and is open for further experimentation with other functions $g : V \rightarrow [a, b]$, such as distance from each vertex to mass center, or other curvatures obtained from the principal curvature values.

A larger database, with more specimens per family might solve the lack of characterization problem faced throughout the project. Dealing with such database might allow an even more objective assortment, since it would allow to run unsupervised classification algorithms, diminishing the subjectivity of picking a training set.

Figure 4(b) shows a visually recognizable pattern followed by each of the nine determined families, strongly suggesting that it is possible to obtain a reasonably objective final classification employing the aforementioned methods.

As it was pointed out, the use of a supervised classification algorithm such as SVM implies a dependency upon a training set. One possible way to hone such classification consists on honing the training set. To that end, it might be useful to repeat the idea behind figure 3: plot every ECG associated to every mask of a fixed original family on the same plane; such family will be represented in the training set by the masks whose ECGs follow all a same pattern. Such classification ought to consider archaeological guidance into determining the most representative masks of every family. Finally, such obtained classification might established a more robust and objective criteria to assort pre-Columbian masks, especially those which pose discussion and discrepancies among the archaeological community.

Acknowledgements

The authors would like to thank Instituto Nacional de Antropología e Historia through its project *Desarrollo de Aplicaciones de Computación en Arqueología* and Red Temática CONACYT de Tecnologías Digitales para la Difusión del Patrimonio Cultural for providing the 3D scanned pre-Columbian masks which made possible this project.

References

- [1] D. Arthur, S. Vassilvitskii, "k-means++: The advantages of careful seeding", *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms Society for Industrial and Applied Mathematics* pp.1027-1035, 2007.
- [2] C. Chang, C. Lin, "LIBSVM: A library for support vector machines", *ACM Transactions on Intelligent Systems and Technology* Vol. 2, No. 3, Article No. 27, 2011.
- [3] M. Kuhn, K. Johnson, *Applied Predictive Modeling*. Springer, p. 159, 2013.
- [4] E. Richardson, M. Weirman, "Efficient classification using the Euler Characteristic". *Pattern Recognition Letters* Vol.49, pp.99-106, 2014.
- [5] R.B. Rusu, *Semantic 3D Object Maps for Everyday Manipulation in Human Living Environments*. Dissertation presented at Computer Science Department, Technische Universitaet Muenchen, Germany, 2009.