

000 Clustering DNA sequences by relative compression

001

002 Morteza Hosseini
003 seyedmorteza@ua.pt

004 Diogo Pratas
005 pratas@ua.pt

006 Armando J. Pinho
007 ap@ua.pt

IEETA/DETI,
University of Aveiro

008

009 Abstract

010

011 This document demonstrates the format requirements for papers submitted to the Portuguese Conference on Pattern Recognition. The format is designed for easy on-screen reading, and to print well at one or two pages per sheet. Additional features include: pop-up annotations for citations [? ?]; a margin ruler for reviewing; and a greatly simplified way of entering multiple authors and institutions.

016

017 1 Introduction

018

019 2 Results

020

021 The proposed method is implemented and publicly available at `github.com/smortezah/Clusico`, under GPLv3 license. The machine used for the tests had an 8-core 3.40 GHz Intel® Core™ i7-6700 CPU with 32 GB RAM.

025 For the experiments, we have used 30 mitochondrial DNA (mtDNA) sequences from three groups of Actinopterygii (Ray-finned fishes), Chondrichthyes (Cartilaginous fishes) and Mammalia, that can be downloaded from `www.ncbi.nlm.nih.gov/nuccore`. Each groups contains 10 sequences and their sizes varies from 16,189 to 18,431 bases.

029 In order to classify the sequences, we first ran GeCo [2] on all sequences, considering them as references as well as targets, to find similarity of sequences to each other. For measuring the similarity, normalized relative compression (NRC) were used, that can be calculated as [3]

$$033 \text{NRC}(x|y) = \frac{C(x|y)}{|x| \log_2 |\Phi|}, \quad (1)$$

035 in which $C(x|y)$ is the information in the sequence x and is obtained by compressing x relatively to the sequence y , $|x|$ is the size of sequence x and $|\Phi|$ is the cardinality of input DNA sequences, i.e. $\text{size}(\{A, C, G, T\}) = 4$. Values of NRC falls within the range $(0, 1]$ and the more similar two sequences are, the less is this value.

040 In the next step, we used weighted pair group method with arithmetic mean (WPGMA), which is a bottom-up hierarchical clustering method, to classify the sequences based on NRC values. The WPGMA algorithm employs a similarity matrix to construct a rooted tree (dendrogram) [1, 4].

044 Figure 1a shows similarity of different sequences (NRC values), obtained by GeCo. As is show, when a sequence is compressed relatively to itself, the NRC value will be approximately 0. These cases are shown with red squares.

047 Figure 1b demonstrates the result of classification of the sequences, which is obtained by WPGMA algorithm. The sequences in Mammalia, Chondrichthyes and Actinopterygii groups are shown with red, green and blue colors, respectively. On top of this figure, the dendrogram is plotted, which shows similarity of different sequences within each group and also, similarity of different groups. As it is show, the two groups of Chondrichthyes and Actinopterygii, that are fishes, are more similar to each other, in comparison with Mammalia.

054

055 References

056

- 057 [1] Harry Clifford, Frank Wessely, Satish Pendurthi, and Richard D
058 Emes. Comparison of clustering methods for investigation of
059 genome-wide methylation array data. *Frontiers in genetics*, 2:88,
060 2011.
061 [2] Diogo Pratas, Armando J Pinho, and Paulo JSG Ferreira. Efficient
062 compression of genomic sequences. In *Data Compression Conference (DCC)*, pages 231–240. IEEE, 2016.

- [3] Diogo Pratas, Raquel Silva, and Armando Pinho. Comparison of compression-based measures with application to the evolution of primate genomes. *Entropy*, 20(6):393, 2018.
[4] R. R. Sokal and C. D. Michener. A statistical method for evaluating systematic relationships. *University of Kansas Science Bulletin*, 38: 1409–1438, 1958.

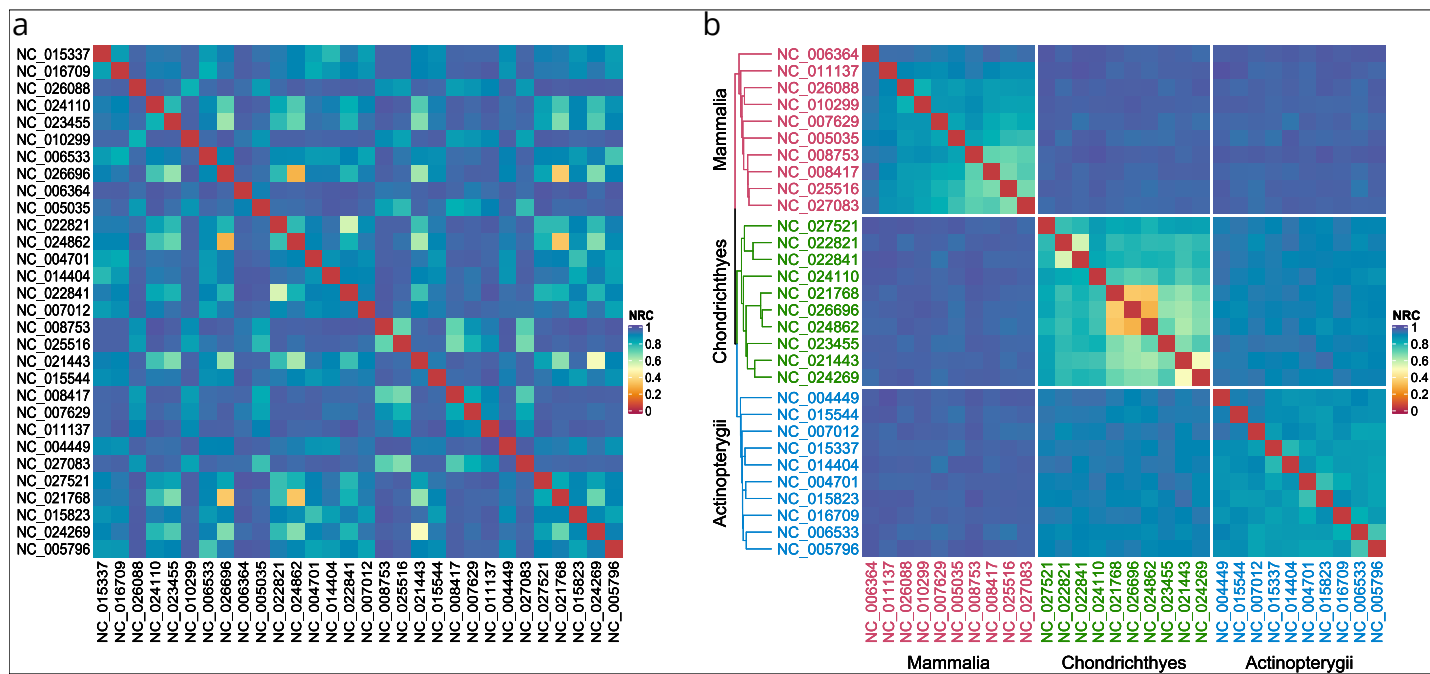


Figure 1: NRC.