



## Genome analysis with distance to the nearest dissimilar nucleotide

Vera Afreixo<sup>a,\*</sup>, Carlos A.C. Bastos<sup>b,c</sup>, Armando J. Pinho<sup>b,c</sup>, Sara P. Garcia<sup>b</sup>, Paulo J.S.G. Ferreira<sup>b,c</sup>

<sup>a</sup> Department of Mathematics, University of Aveiro, 3810-193 Aveiro, Portugal

<sup>b</sup> Signal Processing Lab, IEETA, University of Aveiro, 3810-193 Aveiro, Portugal

<sup>c</sup> Department of Electronics, Telecommunications and Informatics, University of Aveiro, 3810-193 Aveiro, Portugal

### ARTICLE INFO

#### Article history:

Received 6 September 2010

Received in revised form

24 January 2011

Accepted 24 January 2011

Available online 2 February 2011

#### Keywords:

Alignment-free genome comparison

Inter-nucleotide distances

Nearest dissimilar distances

DNA sequences

### ABSTRACT

DNA may be represented by sequences of four symbols, but it is often useful to convert those symbols into real or complex numbers for further analysis. Several mapping schemes have been used in the past, but most of them seem to be unrelated to any intrinsic characteristic of DNA. The objective of this work was to study a mapping scheme that is directly related to DNA characteristics, and that could be useful in discriminating between different species.

Recently, we have proposed a methodology based on the inter-nucleotide distance, which proved to contribute to the discrimination among species. In this paper, we introduce a new distance, the distance to the nearest dissimilar nucleotide, which is the distance of a nucleotide to first occurrence of a different nucleotide. This distance is related to the repetition structure of single nucleotides. Using the information resulting from the concatenation of the distance to the nearest dissimilar and the inter-nucleotide distance, we found that this new distance brings additional discriminative capabilities. This suggests that the distance to the nearest dissimilar nucleotide might contribute with useful information about the evolution of the species.

© 2011 Elsevier Ltd. All rights reserved.

### 1. Introduction

DNA sequences have been converted to numerical signals using different mappings. A commonly used mapping is to consider binary sequences that describe the position of each symbol (Voss, 1992). The binary representation is certainly one of the earliest and one of the most popular mappings of DNA. However, several other different mappings have been proposed (see for example Silverman and Linsker, 1986; Jeffrey, 1990; Zhang and Zhang, 1994; Buldyrev et al., 1995; Anastassiou, 2001; Cristea, 2003; Ning et al., 2003; Brodzik and Peters, 2005; Liao et al., 2005; Akhtar et al., 2007; Randic, 2008; Nair and Mahalakshmi, 2005; Afreixo et al., 2009).

Some of the mappings used in DNA processing do not have a simple numerical interpretation and others do not have biological motivation. Also, some of the representations are not reversible and do not take into account the sequence structure. Currently, there is no ideal mapping to analyze every type of correlation in DNA sequences.

In a previous work, we explored the inter-nucleotide (IN) distance, the distance to the first occurrence of the same symbol, to perform a comparative analysis between species (Afreixo et al.,

2009). In this work, we present a new DNA numerical profile and a new mapping to explore the correlation structure of DNA: the distance to the nearest dissimilar (ND) nucleotide. This representation converts any DNA sequence into a unique numerical sequence with lower length, where each number represents the distance of a symbol to the next occurrence of a different symbol. We introduced also four sequences, one for each nucleotide, to represent the ND distances. This allows to perform comparative analysis between the behavior of the four nucleotides distance sequences and the global sequence.

From the perspective of molecular evolution, DNA sequences may reflect both the results of random mutation and selective evolution. One should subtract the random background from the simple counting result in order to highlight the contribution of selective evolution (Qi et al., 2004; Ding et al., 2010). Therefore, we present an analysis of the relative error to highlight the contribution of selective evolution of the DNA of each species. This residual analysis may be used, for example, to perform multiple organism comparisons.

Phylogenetic trees reproduce the evolutionary tree that represents the historical relationships between the species. Recent phylogenetic tree algorithms use nucleotide sequences. Typically, these trees are constructed with multiple sequence alignment (Hodge and Cope, 2000), which is a computationally demanding task. Recently, alignment-free methods have been proposed and present some advantages over multiple sequences alignment

\* Corresponding author.

E-mail address: vera@ua.pt (V. Afreixo).

methods (see for example Sims et al., 2009; Vinga and Almeida, 2003). The distance that we address in this paper seems to possess discriminating properties that might be helpful in inferring phylogenies. We do believe that this claim is supported by the examples of trees that are provided. However, we also believe that, by itself, this distance measure does not convey all necessary information for building phylogenies. Instead, it should be regarded as potentially useful for working in cooperation and complementing other measures.

## 2. Materials and methods

### 2.1. DNA sequences

In this study, we used the complete DNA sequences of 29 species: 27 were obtained from the National Center for Biotechnology Information (NCBI) (<ftp://ftp.ncbi.nih.gov/genomes/>); *Populus trichocarpa* (California poplar) obtained from the Joint Genome Institute (<http://genome.jgi-psf.org/>) and *Xenopus tropicalis* (Western clawed frog) from Xenbase (<http://www.xenbase.org/>). The species used in this work are listed in Table 1.

### 2.2. Distance to the nearest dissimilar

Consider the alphabet  $\mathcal{A} = \{A, C, G, T\}$  and let  $s = (s_k)_{k \in \{1, \dots, N\}}$  be a symbolic sequence defined in  $\mathcal{A}$ . Consider a numerical sequence,  $w^x$ , that represents the distance to the nearest dissimilar of symbol  $x \in \mathcal{A}$ . As an example, the four ND distance sequences for the short DNA fragment CAAACCGTTAAGTAACAGGGA-TATTGGCCC are

$$w^A = (3, 2, 2, 1, 1, 1), \quad w^C = (1, 2, 1, 3),$$

$$w^G = (1, 1, 3, 2), \quad w^T = (3, 1, 1, 2).$$

**Table 1**

List of DNA builds used for each species.

Species	Reference
<i>Homo sapiens</i> (human)	Build 36.3
<i>Pan troglodytes</i> (chimpanzee)	Build 2.1
<i>Macaca mulatta</i> (Rhesus macaque)	Build 1.1
<i>Mus musculus</i> (mouse)	Build 37.1
<i>Rattus norvegicus</i> (brown rat)	Build 4.1
<i>Equus caballus</i> (horse)	Build 2.1
<i>Cannis familiaris</i> (dog)	Build 2.1
<i>Bos taurus</i> (cow)	Build 4.1
<i>Ornithorhynchus anatinus</i> (platypus)	Build 1.1
<i>Monodelphis domestica</i> (opossum)	Build 2
<i>Gallus gallus</i> (chicken)	Build 2.1
<i>Xenopus tropicalis</i> (Western clawed frog)	Build 4.1
<i>Danio rerio</i> (zebrafish)	Build 3.1
<i>Apis mellifera</i> (honey bee)	Build 4.1
<i>Caenorhabditis elegans</i> (nematode)	NC003279
<i>Vitis vinifera</i> (grape vine)	Build 1.1
<i>Populus trichocarpa</i> (California poplar)	Build 1.0
<i>Arabidopsis thaliana</i> (thale cress)	AGI 7.2
<i>Saccharomyces cerevisiae</i> str.S228C (budding yeast)	SGD 1
<i>Schizosaccharomyces pombe</i> (fission yeast)	Build 1.1
<i>Dictyostelium discoideum</i> str.AX4 (amoeba)	Build 2.1
<i>Plasmodium falciparum</i> 3D7 (protozoon)	Build 2.1
<i>Escherichia coli</i> str.K12 substr.MG1655 (bacterium)	NC000913
<i>Bacillus subtilis</i> str.168 (bacterium)	NC000964
<i>Chlamydia trachomatis</i> str.D/UW-3/CX (bacterium)	NC000117
<i>Mycoplasma genitalium</i> str.G37 (bacterium)	NC000908
<i>Streptococcus mutans</i> str.UA159 (bacterium)	NC004350
<i>Streptococcus pneumoniae</i> str.ATCC 700669 (bacterium)	NC011900
<i>Aeropyrum pernix</i> str.K1 (archaeota)	NC000854

The global sequence of ND distances for this example is

$$w = (1, 3, 2, 1, 3, 2, 1, 1, 2, 1, 1, 3, 1, 1, 1, 2, 2, 3).$$

Note that the ND distance of each nucleotide corresponds to the repeat length of that nucleotide.

**Algorithm 1.** Computation of  $w$  for sequence  $s$ .

```

p := 1
p' := 1
while p' ≤ N do
  i := 0
  while sp' = sp'+i do
    i := i + 1
  end while
  wp := i
  p := p + 1
  p' := p' + i
end while
    
```

**Algorithm 2.** Computation of  $w^x$  with  $x \in \mathcal{A}$ .

```

p := 1
p' := 1
while p' ≤ N do
  i := 0
  while sp'+i = 'x' do
    i := i + 1
  end while
  if i = 0 then
    p' := p' + 1
  else
    wp := i
    p := p + 1
    p' := p' + i
  end if
end while
    
```

Note that

$$\sum_{i=1}^L w_i = N, \tag{1}$$

where  $L$  is the length of  $w$ . Let  $n_i$  be the number of occurrences of ND distance  $i$ , then

$$\sum_{i=1}^K n_i i = N, \tag{2}$$

where  $K$  is the largest ND distance present in the data sequence. The mean distance is

$$\frac{\sum_{i=1}^L w_i}{L} = \frac{N}{L}. \tag{3}$$

### 2.2.1. Relationship between ND and IN distances

In general, the ND distance distribution complements the information that is accumulated in the first IN distance. We recall that the IN distance is the distance to the first occurrence of the same nucleotide. Let  $n'_i$  be the absolute frequency of the  $i$ th IN distance,

$$n'_1 = \sum_{i=2}^K n_i(i-1) + \delta = N - L + \delta \tag{4}$$

Download English Version:

<https://daneshyari.com/en/article/4497161>

Download Persian Version:

<https://daneshyari.com/article/4497161>

[Daneshyari.com](https://daneshyari.com)