

Comparison for DNA Primary Sequence

Bao-Hua Zhang, Hai-Shui Wang, Lu Xu*

Changchun Institute of Applied Chemistry, Chinese Academy of
Sciences, Changchun 130022, Jilin, China

(Received March 19, 2007)

Abstract: In this article, two schemes are suggested based on three exons of β -globin gene belonging to 10 species for comparison of DNA primary sequences. At first, the positions of four nucleic acid bases were extracted, and then based on the information, as the numerical characterization of DNA sequences, the sequence invariants were derived. Sequences comparisons of 10 species selected in this work by using these invariants were performed. The results, especially with scheme 2, are quite satisfactory.

1. Introduction

Deoxyribonucleic acid (DNA) is a long polymer of nucleotides. It is responsible for the genetic propagation of most inherited traits, so comparison of primary sequences of different DNA strands remains one of the important aspects of analysis of DNA data.

There are many species in the world, the differences between their biological traits are enormous, but the differences between their DNA sequences are not enormous as we imagined, so it is difficult to describe effectively the different DNA sequences.

*Corresponding author. Fax:86-431-85685653

E-mail: luxu@ciac.jl.cn

Tel: 86-431-85262239

Postal address: 5625 Renmin Street, Changchun, China.

Postcode: 130022

At present, many schemes such as methods based on matrix ^[1-2], and methods related to graphs ^[3-7] can be used for the comparison of DNA primary sequences.

Formerly, most of the comparisons are based on the first exon of β -globin gene^[8-14]. The first exon of β -globin gene contains 86-94 bases, but there are millions of bases contained in the DNA sequence, so only one exon may be not to give us enough information about a DNA sequence. Thus, many of the results^[1,3] in the literatures did not agree very well with the phylogenetic tree^[15].

In this article, we also follow the rule: important information is contained in the exons. However, our studies are based on three exons of β -globin gene belonging to the 10 species. The suitable invariants were extracted, and the sequence comparisons were made among the 10 species.

1.1 The selected species

We selected 10 species (shown in table1) from phylogenetic tree (figure 1).

Table1. Coding domain sequence (CDS) of ten species

species	Exon1		Exon2		Exon3	
	bases	Region	bases	Region	bases	Region
bovine	86	278-363	223	492-714	129	1613-1741
goat	86	279-364	223	493-715	129	1621-1749
pig	92	871-962	223	1080-1302	129	1944-2072
rabbit	92	480-571	223	698-920	129	1494-1622
rat	92	310-401	223	517-739	129	1377-1505
mouse	92	2718-2809	223	2926-3148	129	3802-3930
gallus	92	465-556	223	649-871	129	1682-1810
geochelone	89	1-89	223	190-412	126	513-638
chimpanzee	105	4189-4293	222	4412-4633	49	5484-5532
gorilla	93	4538-4630	222	4761-4982	49	5833-5881

For facilities to observe, the 10 species are grouped roughly into four classes. The first class includes bovine, goat, rabbit, pig; the second class includes mouse, rat; the third class includes gallus and geochelone(Tortoise), the fourth class includes chimpanzee and gorilla.

index, -0.5. So, we use ${}^m H_t$ to represent our index.

$${}^m H_t = \sum (\delta_1 \cdot \delta_2 \dots \delta_k)^{0.5}$$

In this equation, δ_k represents the members included in a relative position sequence, so ${}^m H_t$ of adenine is:

$${}^1 H_p = (1/33 \times 8/33)^{0.5} + (8/33 \times 13/33)^{0.5} + (13/33 \times 20/33)^{0.5} + \dots$$

The similar indices (${}^1 H_p \sim {}^n H_p$) can be got for the remaining nucleic acids g, c, t. In this article, the indices of ${}^1 H_p \sim {}^5 H_p$ were calculated as invariants for DNA sequence comparison.

(2) Invariants Z

According to step (1), we get 20 invariants belonging to each exon, then, we make quotients between indices of exon 1 and exon 2, as well as exon 2 and exon 3, thus, we got 40 invariants. These are invariants Z, proposed in this paper.

2.2 Coding scheme 2

(1) Graphical representation of DNA sequences

In this part, we consider graphical representation of DNA primary sequences. At first we transform a sequence into a graph, then, extract invariants from this graph. The approach is illustrated on coding domain sequence of goat β -globin gene. The coding domain sequence (CDS) of goat β -globin gene is:

```
atgctgactgctgaggagaaggctgccctcaccggcttctggggcaaggtgaaagtgatgaagttg
gtgctgaggccctgggcaggctgctggtgtctaccctggactcagaggttctttgacactttgggg
actgtcctctgctgatgctgttatgaacaatgctaaggtgaaggcccatggcaagaaggtgctagactc
ctftagtaacggcatgaagcatttgacgacctcaagggcaccttctgctcagctgagtgagctgactgt
gataagctcacgtgatcctgagaactcaagctcctgggcaacgtgctggtggtgctgctgctcgc
caccatggcagtgaaattcaccctgctgctgaggtgagttcagaaggtggtgctggtgttccaat
gcctggccccacagatatcactaa
```

We assign to each nucleic acid base two number sequences [17]: one numerical sequence is the position of a base in the DNA sequence, we use m to represent it; the

other is the position of a base in the subsequence of nucleic acid bases of the same kind, we use n to represent it. Then, the CDS sequence mentioned above corresponding to adenine(a) leads to the following two numerical sequences: $m=[1, 7, 14, 17, 19, 20, 31, 59, \dots, 434, 437, 438]$; $n=[1, 2, 3, 4, 5, 6, 7, \dots, 88, 89]$.

Based on m, n , i.e., to take n as the x-coordinate, m as y-coordinate, we can drive a curve in 2-D plane (see Figure 2).

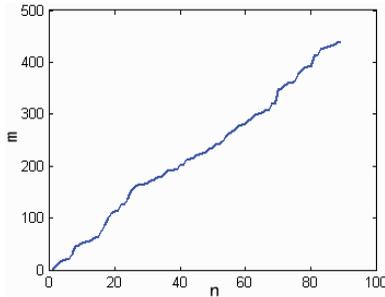


Figure 2 Graphical representation of DNA sequence

Similarly, graphic representation can be constructed for the remaining nucleic acids, guanine (g), cytosine(c), thymine (t).

(2) Extracting invariant from the graph

We also take base adenine in the domain sequence mentioned above as an example to explain the extraction of the sequence invariants. From 0 to 89, we equally take k points from the x-coordinate, that the values of y-coordinate corresponding to the points of x-coordinate are the invariants of base adenine. The magnitude of k is arbitrary, but it should be large enough. Because of only when the k value is large enough, the DNA sequence can be described sufficiently.

(3) Proportional characterization of exons

According to step (2) above, we get $(4 \times k)$ invariants for each exon, then, we make quotients between invariants of exon1 and exon2, as well as exon2 and exon3. thus, we got $(8 \times k)$ invariants. Using such invariants, comparisons were made for coding domain sequence of the 10 species.

3. Similarity comparison

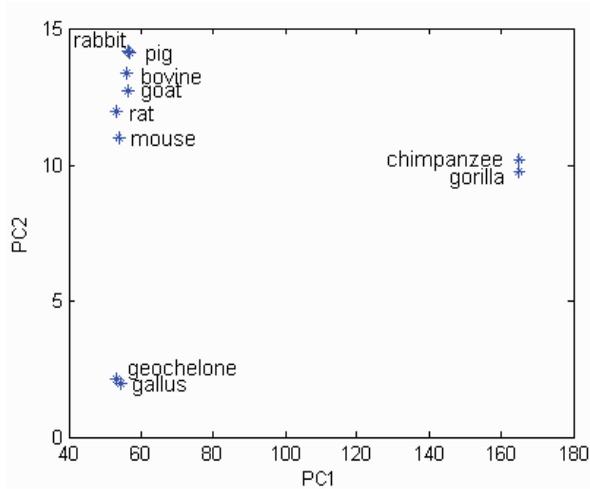


Figure 3 Projected graph of the 10 species with principal component analysis

4. Concluding remarks

All the invariants proposed in this article contain proportional relationships of three exons in the coding domains of the 10 species. In which, based on the scheme1, the results are basically agree with phylogenetic tree, but there are exceptions. Whereas, based on the scheme2, the results agree with phylogenetic tree very well. Therefore, the following conclusions could be given out:

- (1) The primary sequence of a DNA looks to be simple, but it is difficult to describe it effectively. We would say that among many methods those that indicate great dissimilarity among two species (while other approaches may not show such big difference) are to be more trusted even if other methods do not show such difference. Converse, of course is not true: two species showing apparent similarity (having small difference in their selected invariants used for similarity analysis) need not be similar at all - but if they show great similarity within a number of different representations - they are likely to be similar.
- (2) In appearance, schem2 also is a method of graph transformation, but it has no the disadvantages of the type approaches, such as those methods can not be used to the longer sequences; loss of information associated with repeating moves that overlap; the choice of axes for various bases is arbitrary and so on. Theoretically, the scheme in this research can be used for any lengthy sequences.
- (3) It is expected that those invariants can be applied to similarity analysis of RNA

sequences.

REFERENCES

- (1) Randić, M. Condensed Representation of DNA Primary Sequences. *J.Chem.Inf. Comput. Sci.* **2000**, 40, 50-56.
- (2) Randić, M. On Characterization of DNA Primary Sequences by Condensed Matrix. *Chem. Phys. Lett.* **2000**, 317, 29-34.
- (3) Nandy, A. A New Graphical Representation and Analysis of DNA Sequence Structure. I. Methodology and Application to Globin Genes. *Curr. Sci.* **1994**, 66, 309-314.
- (4) Randić, M.; Balaban, A. T. On a Four-Dimensional Representation of DNA Primary Sequences. *J. Chem. Inf. Comput. Sci.* **2003**, 43, 532-539.
- (5) Randić, M.; Vračko, M.; Lerš, N.; Plavšić, D. Novel 2-D graphical representation of DNA sequences and their numerical characterization. *Chem. Phys. Lett.* **2003**, 368, 1-6.
- (6) Randić, M.; Vračko, M.; Lerš, N.; Plavšić, D. Analysis of Similarity/Dissimilarity of DNA Sequences Based on Novel 2-D Graphical Representation. *Chem. Phys. Lett.* **2003**, 371, 202-207.
- (7) Zupan, J.; Randić, M. Algorithm for DNA Sequences into "Spectrum-like" and "Zigzag" Representations. *J. Chem. Inf. Comput. Sci.* **2005**, 45, 309-313.
- (8) Randić, M.; Vračko, M. On Similarity of DNA Primary Sequences. *J.Chem.Inf. Comput. Sci.* **2000**, 40, 599-606.
- (9) Randić, M.; Balaban, A. T.; Basak, S. C. On Structural Interpretation of Several Distance Related Topological Indices, *J. Chem. Inf. Comput. Sci.* **2001**, 41, 593-601.
- (10) Randić, M.; Guo, X. F.; Basak, S. C. On the Characterization of DNA Primary Sequences by Triplet of Nucleic Acid Bases *J. Chem. Inf. Comput. Sci.* **2001**, 41, 619-626.
- (11) Randić, M.; Vračko, M.; Nandy, A.; Basak, S. C. On 3-D Graphical Representation of DNA Primary Sequences and Their Numerical Characterization. *J. Chem. Inf. Comput. Sci.* **2001**, 41, 1235-1244.

- (12) Guo, X. F.; Randić, M.; Basak, S. C. A Novel 2-D Graphical Representation of DNA Sequences of Low Degeneracy. *Chem. Phys. Lett.* **2002**, 350, 106-112.
- (13) Randić, M.; Zupan, J. Highly Compact 2-D Graphical Representation of DNA Sequences. *SAR QSAR Environ Res.* **2004**, 15, 147-157.
- (14) Randić, M. 2-D Graphical Representation of Proteins Based on Virtual Genetic Code. *SAR QSAR Environ. Res.* **2004**, 15, 191-205.
- (15) Wang, J. Y.; Zhu, S. G.; Xu, C. F. (Eds.), *Biochemistry*, Chapter 4, High Education, Beijing, **2003**, p. 183.
- (16) Randić, M.; Characterization of Molecular Branching. *J Am Chem Soc.* **1975**, 97, 6609-6615.
- (17) Randić, M.; Basak, S. C. Characterization of DNA Primary Sequences Based on the Average Distances between Bases. *J. Chem. Inf. Comput. Sci.* **2001**, 41, 561-568.