

Internet Electronic Journal of Molecular Design

December 2002, Volume 1, Number 12, Pages 668–674

Editor: Ovidiu Ivanciuc

Special issue dedicated to Professor Haruo Hosoya on the occasion of the 65th birthday
Part 4

Guest Editor: Jun–ichi Aihara

Numerical Characterization of DNA Primary Sequence

Ping–an He¹ and Jun Wang^{1,2}

¹ Department of Applied Mathematics, Dalian University of Technology, Dalian 116024, P. R. China

² College of Advanced Science and Technology, Dalian University of Technology, Dalian 116024, P. R. China

Received: August 1, 2002; Revised: September 4, 2002; Accepted: November 3, 2002; Published: December 31, 2002

Citation of the article:

P. He and J. Wang, Numerical Characterization of DNA Primary Sequence, *Internet Electron. J. Mol. Des.* 2002, 1, 668–674, <http://www.biochempress.com>.

Numerical Characterization of DNA Primary Sequence[#]

Ping-an He^{1,*} and Jun Wang^{1,2}

¹ Department of Applied Mathematics, Dalian University of Technology, Dalian 116024, P. R. China

² College of Advanced Science and Technology, Dalian University of Technology, Dalian 116024, P. R. China

Received: August 1, 2002; Revised: September 4, 2002; Accepted: November 3, 2002; Published: December 31, 2002

Internet Electron. J. Mol. Des. 2002, 1 (12), 668–674

Abstract

In a previous paper, the authors defined a numerical characterization of DNA primary sequences. A DNA primary sequence was reduced to a few of binary sequences, based on the classifications of the four nucleic acid bases. The reduced sequences are called the characteristic sequences. For each characteristic sequence, we associated two 2×2 matrices, the elements of which are given by the frequency of occurrence of all (0,1) triplets in the characteristic sequence. In this paper, we use eigenvalues of the new matrices to characterize the biological functions of purine–pyrimidine, amino–keto groups and weak–strong H–bonds, respectively.

Keywords. DNA primary sequence; characteristic sequences; leading eigenvalue; similarity; dissimilarity.

1 INTRODUCTION

With the imminent completion of the Human Genome Project and the fast increase of many complete genomes of prokaryote and eukaryote, fundamental questions regarding the characteristics of these sequences arise, the first of which is how to compare genomes. Hence analysis and understanding of the DNA primary sequences are very important tasks in bioinformatics.

Usually, a DNA primary sequence can be taken as a string of letters A, G, C, T, which denote the four nucleic acid bases: adenine, guanine, cytosine and thymine, respectively. Therefore, the analysis and understanding of DNA primary sequences are performed via comparisons of such strings of the four letters. In previous research, the comparisons of DNA primary sequences are mainly to consider the alignment of the DNA primary sequences. The alignment of sequences is performed by the computer to find the smallest number of changes (deletions, insertions, substitutions, shifts) that are necessary to match labels in two DNA primary sequences.

[#] Dedicated to Professor Haruo Hosoya on the occasion of the 65th birthday.

* Correspondence author; E-mail: pinganhe@yahoo.com.cn.

Some researchers consider graphical representations for DNA primary sequences [1–17], in particular, Gates [1], Hamori and Ruskin [2], Leong and Morgenthaler [4], Randić [9–13], Nandy [5–8,13–15], Zhang [16,17] and others, considering a real DNA primary sequence as a curve embedded in 2–D plane or 3–D space. Using research on the graphical representations, we can derive some numerical characterization for DNA primary sequences.

An alternative approach of the comparisons for DNA sequences is suggested by Randić *et al.* [9–13], who considered mathematical invariants of DNA primary sequences rather than the sequences themselves. For chemical structure and chemical graphs we can in fact obtain numerous invariants that are applied for characterization and comparison of structures. There are hundreds of topological indices that have been used in structure–property–activity studies based on molecular graphs. In this way, we can arrive at invariants for DNA primary sequences to associate a matrix with a DNA primary sequence. Once a matrix representation of sequences is given, one can consider suitable matrix invariants as invariants of the comparison of DNA sequences.

In a previous paper [3], we introduced another representation for DNA sequences, which is based on the idea of the coarse–grained description of the DNA primary sequence: we classify the four nucleic acid bases into two groups, purine–pyrimidine, amino–keto groups and weak–strong H–bonds, respectively, and then label the bases of purine, amino and weak H–bonds by 1, and the bases of pyrimidine, keto and strong H–bonds by 0, respectively. Thus, from a DNA primary sequence we obtained three (0,1)–sequences, which are called the characteristic sequences of the DNA primary sequence. For each characteristic sequence we constructed a set of 2×2 matrices, which are based on counting of the frequency of occurrence of all (0,1) triplets of the characteristic sequence. The leading eigenvalues of these matrices are computed and considered as invariants for the comparison of DNA primary sequences.

Table 1. Exon–1 of the β–Globin genes for Eight Species.

Species	Sequence	Length
Human	ATGGTGCACCTGACTCCTGAGGAGAAGTCTGCCGTTACTGCCCTGTG GGGCAAGGTGAACGTGGAGTAAGTTGGTGGTGAGGCCCTGGGCAG	92
Goat	ATGCTGACTGCTGAGGAGAAGGCTGCCGTCACCGCTTCTGGGG CAAGGTGAAAGTGGATGAAGTTGGTGCTGAGGCCCTGGGCAG	86
Gallus	ATGGTGCACCTGGACTGCTGAGGAGAAGCAGCTCATCACCGCCTCTG GGGCAAGGTCAATGTGGCCGAATGTGGGGCCGAAGCCCTGGCCAG	92
Opossum	ATGGTGCACCTGACTTCTGAGGAGAAGAAGTGCATCACTACCATCTG GTCTAAGGTGCAGGTTGACCAGACTGGTGGTGAGGCCCTTGGCAG	92
Lemur	ATGACTTTGCTGAGTGCTGAGGAGAATGCTCATGTCCACTCTCTGTG GGGCAAGGTGGATGTAGAGAAAGTTGGTGGCGAGGCCTTGGGCAG	92
Mouse	ATGGTGCACCTGACTGATGCTGAGAAGGCTGCCGTTACTGCCCTGTG GGGCAAGGTGAACGTGGATGAAGTTGGTGGTGAGGCCCTGGGCAGG	93
Rabbit	ATGGTGCATCTGTCCAGTGAGGAGAAGTCTGCGGTCACTGCCCTGT GGGGCAAGGTGAATGTGGAAGAAGTTGGTGGTGAGGCCCTGGGC	90
Rat	ATGGTGCACCTAACTGATGCTGAGAAGGCTACTGTTAGTGGCCTGTG GGGAAAGGTGAACCTGATAATGTTGGCGCTGAGGCCCTGGGCAG	92

In this paper, through comparison of characteristic sequences we try to find the biological functions of purine–pyrimidine, amino–keto groups and weak–strong H–bonds, respectively. In Table 1, the exon–1 of the β –globin gene for eight species are listed, which were reported by Randić [9].

2 CONSTRUCTION OF THE CHARACTERISTIC SEQUENCES

Nucleic acids and proteins are all linear macromolecules. However, comparison of DNA primary sequences should be considered not only the string structures, but also their chemical structures. In DNA primary sequences, the four bases A, C, G, T can be divided into two classes according to their chemical structures: purine $R = \{A, G\}$ and pyrimidine $Y = \{C, T\}$, or amino group $M = \{A, C\}$ and keto group $K = \{G, T\}$. Besides these, the division can be also made according to the strength of the hydrogen bond, *i.e.*, weak H–bonds $W = \{A, T\}$ and strong H–bonds $S = \{G, C\}$.

Let $S = a_1a_2a_3\cdots$ be a DNA primary sequence. Using above classifications, we can transform a DNA primary sequence into three (0,1) sequences by three homomorphic maps $\phi_i, i = 1, 2, 3$, $\phi_i(S) = \phi_i(a_1)\phi_i(a_2)\cdots$, as follows:

$$\phi_1 = \begin{cases} 1 & \text{if } a_i \in R \\ 0 & \text{if } a_i \in Y \end{cases}$$

$$\phi_2 = \begin{cases} 1 & \text{if } a_i \in M \\ 0 & \text{if } a_i \in K \end{cases}$$

$$\phi_3 = \begin{cases} 1 & \text{if } a_i \in W \\ 0 & \text{if } a_i \in S \end{cases}$$

Thus, we obtain three (0,1) sequences corresponding to the same DNA primary sequence, and we call them as (R,Y) –, (M,K) – and (W,S) –characteristic sequences of the DNA primary sequence, respectively.

In [3], we constructed a set of 2×2 matrix and computed their leading eigenvalues for three characteristic sequences. Using the leading eigenvalues, we compared the similarities and dissimilarities for eight species in Table 1. The results in [3] coincide with the result of Randić's papers. It demonstrates that the comparison of three characteristic sequences is the same as the comparison of DNA primary sequences.

As we have seen in [3], the three characteristic sequences contain all information of the primary sequence. On the other hand, each characteristic sequence is a coarse–grained description for the DNA primary sequence, *i.e.*, some information for DNA primary sequence may be lost in a characteristic sequence so that different DNA primary sequences may have certain similar characteristic sequences. This just reflects the functions of the classifications. Therefore, comparing

each characteristic sequences has special significance. In Table 2, we list the characteristic sequences of the eight DNA sequences of Table 1. In the next section we will compare each characteristic sequence and get some conclusions that cannot be obtained from direct comparison of DNA primary sequences.

Table 2. Characteristic Sequences of the Eight DNA Sequence from Table 1

human
1011010100011000001111111100010010010010000101111011110111010111101110011011011110000111011 1000001111001101100100101100100110001100111000000011100001110000100110000000001001110000110 11001001001010100101001011010100001110100001010000011001011001001011101100100101000001000010
goat
10100110010011111111110010010010011000001111011110111110111011110011010011110000111011 10010011001001001011001001100111100100100000111000011100001001100000001001001110000110 11001010100101001011000100001010000011010000011001011101001101101100100101000001000010
gallus
10110101001110010011111111101100010010011000001111011110011010110011101011110011110000110011 10000011100011001001001011011010110111100110100000111000111000001101100000001101101110001110 11001001010010100101001011001001011010000001010000011001011101000001110100000001100001000010
opossum
1011010100011000001111111111001010010010010001100011110101110011001110011011011110000011011 10000011100011001001001011011100110111011110100001011000011000001111011000000001001110000110 11001001011010110101001011011010011010110011010010111001001001101001010100100101000001100010
lemur
1011000010011101001111111101000101001000000010111101111011 1010111111110011011011110000111011 10011000010010001001001011001011000111101010000000111000001000101011100000000101001100000110 11010111001010100101001011100101101010010101010000011001001101101011101100100001000011000010
mouse
101101010001100110100111111100100100100100001011110111101110101110111100110110111100001110111 10000011110011001001001011001001100011001110000000111000011100001001100000000010011100001100 1100100100101010110010110001000011101000010100000110010110010011011011001001010000010000100
rabbit
1011010100010001101111111100010110010010000101111011110111010111111100110110111100001110 100000110100011100010010110010010001110011100000001110000110000011011000000000010011100001 11001001101010010101001011010100001010100001010000011001 0111010011011011001001010000010000
rat
1011010100011001101001111111001001001101100010111111111011 1000011011010011010011110000111011 10000011110111001001001011001011000010000110000000111000011111001011000000101001001110000110 11001001001110101100101011000110101110100001010000111001011000101111101100000101000001000010

3 COMPARISON OF CHARACTERISTIC SEQUENCES

In [3], we also introduced a $2 \times 2 \times 2$ cubic matrix with 8 entries $f_{ijk}^X = (100m_{ijk}^X)/(N-2)$, where m_{ijk}^X is the enumeration of the (0,1) triplet ijk in characteristic sequence X and N is the length of X . Clearly, it represents the 100 times of the frequency of occurrence of the (0,1) triplet ijk in X . That we take the 100 times is for convenience of tabulation and computation. By F^R , F^M and F^W we denote the cubic matrices for the (R,Y)-, (M,K)- and (W,S)-characteristic sequences, respectively. We partition each of the cubic matrices into a pair of 2×2 condensed matrices F_0^X and F_1^X , where

$F_0^X = (f_{0,jk}^X)$ and $F_1^X = (f_{1,jk}^X)$ with X being R , M or W . In [3], we computed the leading eigenvalues condensed matrices as above. In Table 3, all leading eigenvalues of characteristic sequences are listed, as reported in [3].

Table 3. Leading eigenvalues of the 6 matrices F_0^X and F_1^X for the eight DNA sequences of Table 1.

Species	F_0^R	F_1^R	F_0^M	F_1^M	F_0^W	F_1^W
Human	21.7	28.9	31.3	18.3	30.0	22.2
Goat	20.3	30.8	30.2	21.7	30.4	21.4
Gallus	22.6	28.2	26.5	23.2	33.2	20.7
Opossum	23.0	26.6	27.6	21.8	26.6	25.0
Lemur	21.1	30.3	33.2	20.4	26.5	22.9
Mouse	21.6	29.6	31.5	19.9	29.6	23.4
Rabbit	20.5	31.7	34.7	18.6	29.2	22.2
Rat	22.8	28.3	30.3	21.3	28.2	22.3

For each characteristic sequence, we take the leading eigenvalue as a two-dimensional vector (F_0^X, F_1^X) , by which we compare the (R,Y) -, (M,K) - and (W,S) -characteristic sequences of DNA primary sequences based on the Euclidean distance between the end points of the two-dimensional vectors, respectively. The results of comparisons are listed on the three tables, where Table 4 reveals the information of purine-pyrimidine group, Table 5 the information of amino-keto group, and Table 6 the information of weak-strong H-bonds, respectively.

Table 4. Similarity/dissimilarity table for the eight DNA sequences of Table 1 based on their (R,Y) characteristic sequences.

Species	Goat	Gallus	Opossum	Lemmur	Mouse	Rabbit	Rat
Human	2.36008	1.14018	2.64197	1.52315	2.86007	3.04631	1.25300
Goat		3.47131	4.99300	0.94339	5.22015	0.92195	3.53553
Gallus			1.64924	2.58070	1.78885	4.08167	0.22361
Opossum				4.15933	0.40000	5.67979	1.71172
Lemmur					4.3566	1.52315	2.62488
Mouse						5.86686	1.80278
Rabbit							4.10488

Table 5. Similarity/dissimilarity table for the eight DNA sequences of Table 1 based on their (M,K) characteristic sequences.

Species	Goat	Gallus	Opossum	Lemmur	Mouse	Rabbit	Rat
Human	3.57351	6.8593	5.02096	2.83196	1.71172	3.41321	3.16228
Goat		3.99249	2.50200	3.26956	2.14009	5.46443	0.41231
Gallus			1.84391	7.26154	5.93633	9.40213	4.24853
Opossum				5.67539	4.20476	7.69675	2.64764
Lemmur					1.74642	2.34307	3.03645
Mouse						3.49285	1.76918
Rabbit							5.16236

Observing Tables 4, 5 and 6, we can obtain some information for each characteristic sequence. For example, the species gallus is the most dissimilarly with others in Table 6. However, we do not see the same result from Table 4 and 5, even the value of gallus-rat pair is the least in Table 4. The

species gallus is the only non-mammalian species among these considered species. Whether or not this means that the essential nature of the mammalian species may be revealed mainly in the characteristic sequence of weak-strong H-bonds group. On the other hand, the results in Tables 5 and 6 are very similar to that of the comparison of DNA primary sequences. This means that the information of the similarities for eight sequences may contain mainly in the reduce sequences of amino-keto groups and weak-strong H-bonds groups.

Table 6. Similarity/dissimilarity table for the eight DNA sequences of Table 1 based on their (W,S) characteristic sequences.

Species	Goat	Gallus	Opossum	Lemmur	Mouse	Rabbit	Rat
Human	0.894427	3.53412	4.40454	3.56931	1.96977	0.80000	1.80278
Goat		2.88617	5.23450	4.17852	2.72029	1.44222	2.37697
Gallus			7.87718	7.05195	5.50364	4.2720	5.24976
Opossum				2.10238	2.56125	3.82099	3.13847
Lemmur					1.70294	2.78927	1.80278
Mouse						1.28062	0.70000
Rabbit							1.00499

Furthermore, we can observe the least value in each table: the gallus-rat pair in Table 4, the goat-rat pair in Table 5, and the mouse-rat pair in Table 6, respectively. Whether these results imply that the three characteristic sequences reflect some intrinsic essence of species rat from different aspect.

Generally, we can also observe the least value of all species in each table, so that we can obtain information of the (R,Y)-, (M,K)- and (W,S)-characteristic sequences, respectively. For example, the mouse species, the least value in Tables 4–6 are the opossum, human and rat, respectively. These results illuminate that the three characteristic sequences reflect some essence of mouse species from different aspect.

4 CONCLUSIONS

Comparing characteristic sequences, we can get some information that cannot be obtained from the direct comparison of DNA primary sequences and observe some special nature in species from different aspect. Although some information may be lost in characteristic sequences, we can focus our attention on the information of our interest. This is the advantage of our approach.

Acknowledgment

This work is supported in part by the National Natural Science Foundation of China.

5 REFERENCES

- [1] M. A. Gates, A Simple way to look a DNA, *J. Theor. Biol.* **1986**, *119*, 319–328.
- [2] E. Hamori and J. Ruskin, H curves, a novel method of representation of nucleotide series especially suited for long DNA sequences, *J. Biol. Chem.* **1983**, *258*, 1318–1327.

- [3] P. He and J. Wang, Characteristic sequences for DNA primary sequence, *J. Chem. Inf. Comput. Sci.* **2002**, *42*, 1080–1085.
- [4] P. M. Leong and S. Morgenthaler, Random walk and gap plots of DNA sequences, *Comput. Applic. Biosc.* **1995**, *11*, 503–507.
- [5] A. Nandy, A new graphical representation and analysis of DNA sequence structure: I. Methodology and application to globin genes, *Curr. Sci.* **1994**, *66*, 309–314.
- [6] A. Nandy and P. Nandy, Graphical analysis of DNA sequences structure: II. Relative abundance of nucleotides in DNAs, gene evolution and duplication, *Curr. Sci.* **1995**, *68*, 75–85.
- [7] A. Nandy, Graphical analysis of DNA sequence structure: III. Indications of evolutionary distinctions and characteristics of introns and exons, *Curr. Sci.* **1996**, *70*, 661–668.
- [8] A. Nandy, P. Nandy, and S. C. Basak, Quantitative Descriptor for SNP Related Gene Sequences, *Internet Electron. J. Mol. Des.* **2002**, *1*, 367–373, <http://www.biochempress.com>.
- [9] M. Randić, Condensed Representation of DNA Primary Sequences, *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 50–56.
- [10] M. Randić and M. Vračko, On the Similarity of DNA Primary Sequences, *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 599–606.
- [11] M. Randić, On characterisation of DNA primary sequences by a condensed matrix, *Chem. Phys. Lett.* **2000**, *317*, 29–34.
- [12] M. Randić, X. Guo, and S. C. Basak, On the characterization of DNA primary sequences by triplet of nucleic acid bases, *J. Chem. Inf. Comput. Sci.* **2001**, *41*, 619–626.
- [13] M. Randić, M. Vračko, A. Nandy, and S. C. Basak, On 3–D representation of DNA primary sequences, *J. Chem. Inf. Comput. Sci.* **2000**, *40*, 1235–1244.
- [14] C. Raychaudhury and A. Nandy, Indexing scheme and similarity measures for macromolecular sequences, *J. Chem. Inf. Comput. Sci.* **1999**, *39*, 243–247.
- [15] A. Nandy, Investigations on Evolutionary Changes in Base Distributions in Gene Sequences, *Internet Electron. J. Mol. Des.* **2002**, *1*, 545–558, <http://www.biochempress.com>.
- [16] R. Zhang and C. T. Zhang, Z–curve, An Intuitive Tool for Visualizing and Analyzing the DNA sequences, *J. Biomol. Str. Dyn.* **1994**, *11*, 767–782.
- [17] C. T. Zhang, A Symmetrical Theory of DNA sequences and Its Application, *J. Theor. Biol.* **1997**, *187*, 297–306.

Biographies

Ping-an He is a PhD student of Applied Mathematics at the Dalian University of Technology. His main research interests include combinatorics, graph theory and bioinformatics.

Jun Wang is a Professor of Applied Mathematics at the Dalian University of Technology, the advisor of the first author.